

# 601 Appendices

602	<b>A Broader impact</b>	<b>17</b>
603	<b>B Limitations</b>	<b>17</b>
604	<b>C The scaling toolbox: practical methods for enabling systematic scaling</b>	<b>18</b>
605	C.1 The continuous time/SDE perspective . . . . .	18
606	C.2 Scaling rules for optimization . . . . .	19
607	C.3 Commonly used values of hyperparameters at different batch sizes . . . . .	20
608	C.4 Progressive scaling . . . . .	20
609	<b>D EMA approximation theorems with SDEs</b>	<b>21</b>
610	D.1 SGD with model EMA . . . . .	21
611	D.2 Adaptive gradient methods with model EMA . . . . .	23
612	<b>E Additional proofs</b>	<b>24</b>
613	E.1 Iterations of SGD + EMA . . . . .	24
614	E.2 Limiting behavior of Polyak-Ruppert averaging . . . . .	25
615	<b>F Additional details and results for Polyak-Ruppert averaging</b>	<b>26</b>
616	F.1 Noisy parabola . . . . .	27
617	F.2 Image Classification . . . . .	29
618	F.3 Applying the EMA Scaling Rule to Batch Normalization . . . . .	30
619	<b>G Additional details and results for Automatic Speech Recognition (ASR)</b>	<b>31</b>
620	G.1 Detailed results . . . . .	34
621	G.2 Scaling to $\kappa = 16$ with Progressive Scaling . . . . .	36
622	<b>H Additional details and results for self-supervised image representation learning</b>	<b>37</b>
623	H.1 Components of self-supervised learning . . . . .	37
624	H.2 A Vision Transformer recipe for BYOL . . . . .	38
625	H.3 The role of Batch Normalization and Layer Normalization in BYOL with ViTs . . . . .	39
626	H.4 Longer training duration with incremental Progressive Scaling . . . . .	40
627	H.5 Building intuition around Progressive Scaling and momentum sensitivity . . . . .	40
628	H.6 Scaling a ResNet-50 BYOL using LARS and Progressive Scaling . . . . .	41
629	H.7 Preventing collapse phenomena in DINO at scale . . . . .	42

630 **A Broader impact**

631 This work shows how to adapt Machine Learning (ML) optimization in the presence of a model  
632 Exponential Moving Average (EMA). There are a number of benefits to this:

- 633 1. Scaling rules democratize the training of ML models: they give ML researchers the ability to  
634 replicate the optimization of large scale systems, even if those researchers *do not* have access  
635 to i) significant parallel computational resources *or* ii) the technical tooling to do so.
- 636 2. Our EMA Scaling Rule lowers compute usage as it removes the necessity for a hyperparameter  
637 search over momenta; in the case where our scaling assumptions hold, if we know the value  
638 of the optimal momentum  $\rho_B$  at some batch size  $B$  (for example, the momentum that gives the  
639 best transfer performance), then the optimal value at another batch size  $\hat{B}$  is exactly the one  
640 given by the EMA Scaling Rule  $\hat{\rho} = \rho_B^{\kappa}$ , for scaling  $\kappa = \hat{B}/B$ .
- 641 3. Our EMA Scaling Rule enables researchers to more quickly iterate through experimental ideas,  
642 and opens up access to large-scale training (for example, larger models and larger datasets) for  
643 Pseudo-Labeling and Self-Supervised Learning (SSL) techniques.

644 These points have potential negative consequences:

- 645 1. As our EMA Scaling Rule enables researchers to iterate the same experiments more quickly,  
646 and perform large-scale training with EMA-based methods, this may encourage a greater number  
647 of experiments, or the training of larger models. Either of these possibilities leads to greater  
648 energy consumption.
- 649 2. As the need to determine momentum hyperparameters has now been removed, researchers who  
650 were previously discouraged from attempting to scale these methods due to an *extra* hyperpa-  
651 rameter to tune may begin to perform such experiments, leading, once more, to greater energy  
652 consumption.

653 The environmental impact of each of these two points may be significant.

654 **B Limitations**

655 The EMA Scaling Rule provides a recipe for producing training dynamics independent of the batch  
656 size used in stochastic optimization. The technology underpinning it will not *always* give the desired  
657 behavior, however.

658 The first issue occurs with the wording present in the EMA Scaling Rule: [...] *and scale other*  
659 *optimizers according to their own scaling rules* (Definition 1.1):

- 660 1. This statement requires that the given Stochastic Differential Equation (SDE) approximation  
661 we are using for the model optimizer is itself providing well-behaved scaling, that is, that in the  
662 *absence* of a model EMA, the model optimization trajectories at the batch sizes  $B$  and  $\kappa B$ , with  
663 optimizer hyperparameters appropriately scaled, are close. In general we know this is not true.  
664 First, we know that the SDE approximation for Stochastic Gradient Descent (SGD) breaks at a  
665 given  $\kappa$  due to discretization error (Li et al., 2021). Second, we know that if the gradient noise  
666 is not sufficiently large, the SDE approximation for Adam does not exist (Malladi et al., 2022),  
667 i.e. an SDE motivated scaling rule has no meaning.
- 668 2. This statement requires knowledge of how to scale the corresponding model optimizer. We  
669 have principled ways to achieve this for SGD (Li et al., 2021), and for the adaptive optimization  
670 methods RMSProp and Adam (Malladi et al., 2022). Empirically, a square-root scaling  
671 law for LAMB (You et al., 2020) has been observed, however, it has not been derived formally.  
672 Problematically, there is no known hyperparameter scaling law or SDE approximation known  
673 for LARS (You et al., 2017), which has been used in Bootstrap Your Own Latent (BYOL)  
674 (Grill et al., 2020) and many other large-scale training procedures for convolution-based archi-  
675 tectures. Despite this, we are able to demonstrate in Appendix H.6 that a combination of the  
676 EMA Scaling Rule and progressive scaling can match, or surpass baseline BYOL performance  
677 at a batch size of 32,768 using LARS, indicating that although the theoretical guarantees may  
678 not hold, there is still practical utility in the tools we provide in this work.

679 3. It may be the case that the optimal performance attainable by a given model setup exists at a  
680 level of discretization/gradient noise where no SDE exists. In this case, SDE-derived scaling  
681 rules can never be valid, and no scaling of this dynamics can be achieved with known tools.

682 The second issue is related to the case when the optimizer scaling rule is valid. In this case, the error  
683 for the EMA Scaling Rule at finite learning rate  $\eta$  at large  $\kappa$  can be considerable. In cases where the  
684 model EMA plays a role in the overall optimization, the error introduced by the EMA Scaling Rule  
685 can break the preservation of model dynamics.

686 Put another way, an optimizer scaling rule and the EMA Scaling Rule each introduce their own dis-  
687 cretization errors. In the case where EMA plays a role in optimization, as soon as the discretization  
688 error of *either* the optimizer scaling rule *or* the EMA Scaling Rule is large, the error for the joint  
689 optimization procedure is large. This is *at least* as bad as cases that *do not* use a model EMA during  
690 the optimization process.

## 691 C The scaling toolbox: practical methods for enabling systematic scaling

692 There are many different components involved in preserving optimization dynamics at different  
693 batch sizes. In this appendix we collect into a single place the different concepts and values that we  
694 found useful in practice, in an attempt to make the practice of scaling as accessible as possible.

### 695 C.1 The continuous time/SDE perspective

696 Here we discuss the mindset difference required when trying to preserve training dynamics. In  
697 ML we typically use stochastic optimization, leading us to think of the optimization in terms of  
698 *performing updates*, or *stepping the optimizer*. This notion has become more common in the era of  
699 large datasets, where it may be the case that we only see a fraction of the dataset during optimization.

700 For dynamics preservation under scaling, we suggest that it is simpler to consider the *amount of data*  
701 seen by the training process, or alternatively, the amount of *continuous time* in the discretization of  
702 SDEs view. The reason is the following. The SDE scaling rule results (Definition 1.1, Li et al.  
703 (2019, 2021); Malladi et al. (2022)) follow from showing that different discretizations of the SDE  
704 are close to that SDE, providing we appropriately scale hyperparameters (see Section 2.2). Each of  
705 these discretizations shares the *total continuous time*  $T = \hat{\eta} \times \hat{N}_{\text{iter}}$ <sup>7</sup> of the underlying SDE, but each  
706 discretization has a *different* number of iterations  $\hat{N}_{\text{iter}} = N_{\text{iter}}/\kappa$ .

707 This perspective is already adopted, perhaps by accident in some domains. For example, in Com-  
708 puter Vision (CV), it is typical to compare model performance after optimization on ImageNet1k  
709 after a *number of epochs*, whilst also specifying a learning rate warmup after a *number of epochs*.  
710 This transforms the schedule into the form *wait until the process meets [condition]*, where here  
711 *[condition]* is *when the process has seen sufficiently many samples*.

712 More generally, we can specify any *condition* that is not a property of the discretization procedure  
713 itself. Instead, the discretization procedure should be viewed as a numerical approximation method  
714 for the SDE we are evolving, and the properties of that discretization process (like *number of steps*)  
715 are not *of specific interest* in the world view where we do decouple optimization from the batch size.  
716 A specific example of this more general case is present in Section 3.3, where for scaling  $\kappa > 2$  we  
717 wait until the pre-training Word Error Rate (WER) is sufficiently low.

718 There may be cases where one is working with a setup that is explicitly defined in terms of quantities  
719 related to the discretization process. Indeed, the optimizer hyperparameters are examples of these,  
720 and need to be scaled accordingly with  $\kappa$ . The other typical example of this is conditions based on  
721 the *number of optimizer steps*, rather than the number of epochs. In this case, these quantities should  
722 be scaled to achieve the desired condition in the same amount of time, i.e. as above  $\hat{N}_{\text{iter}} = N_{\text{iter}}/\kappa$ ,  
723 where  $N_{\text{iter}}$  is the number of iterations specified at the base batch size  $B$ . Concretely, if training is  
724 specified in a number of steps, then doubling the batch size implies you should train for half the  
725 number of steps.

---

<sup>7</sup>This is in the case of SGD, for RMSProp and Adam one should use  $T = \hat{\eta}^2 \times \hat{N}_{\text{iter}}$  (Malladi et al., 2022).

726 **C.2 Scaling rules for optimization**

727 For ease of reference, we collect all the scaling rules related to batch size modification we are aware  
728 of. We begin with the most well-known, the SGD Scaling Rule (Definitions 2.2 and C.1).

729 **Definition C.1** (SGD Scaling Rule). *When running SGD (Definition 2.1) with batch size  $\hat{B} = \kappa B$ ,  
730 use a learning rate  $\hat{\eta} = \kappa \eta$  (Krizhevsky, 2014; Goyal et al., 2017).*

731 The SGD Scaling Rule is also known as the Linear Scaling Rule (LSR), although for clarity, this  
732 work adopts the naming convention [*Algorithm Name*] *Scaling Rule*, which means all parameters of  
733 those algorithms are appropriately scaled from batch size  $B$  to  $\kappa B$ .

734 Next, we give the two scaling rules known for the adaptive optimizers RMSProp (Tieleman et al.,  
735 2012) and Adam (Kingma & Ba, 2015) in Definition C.2 and Definition C.3 respectively.

736 **Definition C.2** (RMSProp Scaling Rule). *When running RMSProp (Tieleman et al., 2012) with  
737 batch size  $\hat{B} = \kappa B$ , use a learning rate  $\hat{\eta} = \sqrt{\kappa} \eta$ , beta coefficient  $\hat{\beta} = 1 - \kappa \times (1 - \beta)$ , and adaptivity  
738 parameter  $\hat{\epsilon} = \frac{\epsilon}{\sqrt{\kappa}}$  (Malladi et al., 2022).*

739 **Definition C.3** (Adam Scaling Rule). *When running Adam (Kingma & Ba, 2015) with batch size  
740  $\hat{B} = \kappa B$ , use a learning rate  $\hat{\eta} = \sqrt{\kappa} \eta$ , beta coefficients  $\hat{\beta}_1 = 1 - \kappa \times (1 - \beta_1)$ ,  $\hat{\beta}_2 = 1 - \kappa \times (1 - \beta_2)$ ,  
741 and adaptivity parameter  $\hat{\epsilon} = \frac{\epsilon}{\sqrt{\kappa}}$  (Malladi et al., 2022).*

742 Next, we present a contribution of this work, the EMA Scaling Rule (Definitions 1.1 and C.4), which  
743 extends the above scaling rules to allow the presence of a model EMA which is able to contribute to  
744 the overall optimization (see Appendices D and E.1 for derivations).

745 **Definition C.4** (EMA Scaling Rule). *When computing the EMA update (Definition 2.3) of a model  
746 undergoing stochastic optimization with batch size  $\hat{B} = \kappa B$ , use a momentum  $\hat{\rho} = \rho^\kappa$  and scale other  
747 optimizers according to their own scaling rules.*

748 Concretely, if we are using SGD in the presence of a model EMA, Definitions C.1 and C.4 state that  
749 we should take  $\hat{\eta} = \kappa \eta$  and  $\hat{\rho} = \rho^\kappa$  when scaling by  $\kappa = \hat{B}/B$ .

750 The final scaling rule is for weight decay, and follows from the scaling logic discussed in Ap-  
751 pendix C.1 and Krizhevsky (2014). If we take the weight decay regularization penalty  $\lambda$  defined at  
752 batch size  $B$ , what should the weight decay  $\hat{\lambda}$  be for batch size  $\hat{B} = \kappa B$ ? For simplicity, consider  $\kappa$   
753 updates of optimization of parameters  $\theta_t$  in the presence of weight decay only

$$\theta_{t+\kappa} = \theta_{t+\kappa-1} - \eta \lambda \theta_{t+\kappa-1} = (1 - \eta \lambda) \theta_{t+\kappa-1} = (1 - \eta \lambda)^\kappa \theta_t. \quad (11)$$

754 Therefore, to match the effect of weight decay with a single iteration step, we need to match

$$1 - \hat{\eta} \hat{\lambda} = (1 - \eta \lambda)^\kappa. \quad (12)$$

755 Solving for  $\hat{\lambda}$  and expanding around  $\eta \approx 0$  gives

$$\hat{\lambda} = \frac{1 - (1 - \eta \lambda)^\kappa}{\hat{\eta}} \approx \frac{\eta}{\hat{\eta}} \times \kappa \lambda + O(\eta). \quad (13)$$

756 This leads to the Weight Decay Scaling Rule (Definition C.5).

757 **Definition C.5** (Weight Decay Scaling Rule). *When using weight decay with batch size  $\hat{B} = \kappa B$ , use  
758 a penalty term  $\hat{\lambda} = (\kappa \hat{\eta} / \eta) \lambda$ , where  $\hat{\eta}$  and  $\eta$  represent the scaled and unscaled learning rates of the  
759 corresponding optimizer (Krizhevsky, 2014; Li et al., 2018; Loshchilov & Hutter, 2019).*

760 The Weight Decay Scaling Rule implies that using *linear* scaling for the learning rate  $\eta$  then the  
761 weight decay penalty is automatically scaled, and when using *square-root* scaling for the learning  
762 rate  $\eta$  (e.g. in the case of the Adam Scaling Rule (Definition C.3)) then the weight decay penalty  
763 should also be scaled with a *square-root* as is proposed in Loshchilov & Hutter (2019).

764 Finally, we see that if the implementation of weight decay does not have an update scaled by the  
765 learning rate, i.e. the update is  $\theta_{t+1} = (1 - \lambda) \theta_t$ , then the scaling rule is optimizer-independent, and  
766 becomes linear for small weight decay, i.e.  $\hat{\lambda} = \kappa \lambda$ , and for arbitrary  $\lambda$  takes the form  $\hat{\lambda} = 1 - (1 - \lambda)^\kappa$ .

Table 2: Scaled learning rates  $\hat{\eta}$  at different batch sizes  $\hat{B} = \kappa B$  given reference learning rates  $\eta$  defined at batch size  $B$ . The reference values of each column are boldened. Note that this is only valid when there is a notion of *single sample*. In the sequence learning setup (for example, in Section 3.3), the notion of batch size should be appropriately replaced with the *dynamic batch size*, i.e. total sequence length.

Batch size $\hat{B}$	$\hat{\eta} = \kappa\eta$ [SGD]			$\hat{\eta} = \sqrt{\kappa}\eta$ [RMSProp, Adam]		
	$B = 256$		$B = 512$	$B = 256$	$B = 4096$	
	$\eta = 0.1$	$\eta = 0.3$	$\eta = 0.1$	$\eta = 10^{-3}$	$\eta = 4.8$	$\eta = 10^{-3}$
32	0.0125	0.0375	0.00625	0.00035	0.42426	0.00009
64	0.025	0.075	0.0125	0.0005	0.6	0.00013
128	0.05	0.15	0.025	0.00071	0.84853	0.00018
256	<b>0.1</b>	<b>0.3</b>	<b>0.05</b>	<b>0.001</b>	1.2	0.00025
512	0.2	0.6	<b>0.1</b>	0.00141	1.69706	0.00035
1024	0.4	1.2	0.2	0.002	2.4	0.0005
2048	0.8	2.4	0.4	0.00283	3.39411	0.00071
4096	1.6	4.8	0.8	0.004	<b>4.8</b>	<b>0.001</b>
8192	3.2	9.6	1.6	0.00566	6.78823	0.00141
16384	6.4	19.2	3.2	0.008	9.6	0.002
32768	12.8	38.4	6.4	0.01131	13.57645	0.00283
65536	25.6	76.8	12.8	0.016	19.2	0.004

Table 3: Scaled EMA momenta  $\hat{\rho} = \rho^\kappa$  at different batch sizes  $\hat{B} = \kappa B$  given reference momenta  $\rho$  defined at batch size  $B$ . The reference values of each column are boldened. Again in the sequence learning setup, batch size should be appropriately replaced with a notion of sequence length.

Batch size $\hat{B}$	$B = 256$			$B = 4096$			
	$\rho = 0.9999$	$\rho = 0.999$	$\rho = 0.99$	$\rho = 0.996$	$\rho = 0.992$	$\rho = 0.99$	$\rho = 0.97$
32	0.99999	0.99987	0.99874	0.99997	0.99994	0.99992	0.99976
64	0.99997	0.99975	0.99749	0.99994	0.99987	0.99984	0.99952
128	0.99995	0.9995	0.99499	0.99987	0.99975	0.99969	0.99905
256	<b>0.9999</b>	<b>0.999</b>	<b>0.99</b>	0.99975	0.9995	0.99937	0.9981
512	0.9998	0.998	0.9801	0.9995	0.999	0.99874	0.9962
1024	0.9996	0.99601	0.9606	0.999	0.99799	0.99749	0.99241
2048	0.9992	0.99203	0.92274	0.998	0.99599	0.99499	0.98489
4096	0.9984	0.98412	0.85146	<b>0.996</b>	<b>0.992</b>	<b>0.99</b>	<b>0.97</b>
8192	0.9968	0.96849	0.72498	0.99202	0.98406	0.9801	0.9409
16384	0.99362	0.93798	0.5256	0.9841	0.96838	0.9606	0.88529
32768	0.98728	0.8798	0.27625	0.96844	0.93776	0.92274	0.78374
65536	0.97472	0.77405	0.07632	0.93788	0.8794	0.85146	0.61425

### 767 C.3 Commonly used values of hyperparameters at different batch sizes

768 In the literature it is common to give a base learning rate  $\eta$  defined at batch size 256, implicitly  
 769 using the SGD Scaling Rule, even when using the Adam optimizer. Because the scaling of other  
 770 optimization hyperparameters was not understood until recently, it is also common to just present  
 771 these *for the experiment*, e.g. the Adam betas and epsilon, and the EMA momentum, implicitly  
 772 defined at the scale of the experiment, for example at batch size 4096. One way to deal with this  
 773 in practice is to define a single reference batch size  $B$  at which *all* hyperparameters are defined, and  
 774 then scale from there. In this case, it is easiest to compute *using linear scaling* the learning rate at  
 775 the redefined base batch size  $\eta = \tilde{\kappa} \eta_{\text{orig}}$ , where  $\tilde{\kappa} = B/B_{\text{orig}}$ , and then scale this new reference  $\eta$  as  
 776  $\hat{\eta} = \kappa\eta$ ,  $\kappa = \hat{B}/B$ , along with e.g. the momentum defined at  $B$ .

777 As this process can be slightly frustrating, we have provided tables of typical learning rates in Table 2  
 778 and momenta in Table 3.

### 779 C.4 Progressive scaling

780 In Section 3.4 we introduced Progressive Scaling (Definition 3.2) to test our hypothesis that early  
 781 in the BYOL training procedure, there are dynamics that are challenging to replicate at larger batch

---

**Algorithm 1** Stochastic Gradient Descent with Progressive Scaling

---

**Require:** Base learning rate  $\eta$ , base momentum  $\rho$  for base batch size  $B$   
**Require:** Initial target model parameters  $\theta$  and model EMA parameters  $\zeta$   
**Require:** Epochs  $E$  and schedule of batch sizes  $\mathcal{B} = B_1, B_2, \dots, B_E$   
**Require:** Loss function  $\mathcal{L}$

```
for  $e$  in  $1, 2, \dots, E$  do
   $\hat{B} \leftarrow \mathcal{B}[e]$                                 ▶ Get current batch size
   $\kappa \leftarrow \hat{B}/B$                                 ▶ Compute scaling factor
   $\hat{\eta} \leftarrow \kappa\eta$                             ▶ Get scaled learning rate
   $\hat{\rho} \leftarrow \rho^\kappa$                             ▶ Get scaled momentum
  for  $b$  in  $1, 2, \dots, \text{floor}(E/\hat{B})$  do
    Sample a minibatch of  $\hat{B}$  samples  $\mathcal{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\hat{B})}\}$ 
     $\theta \leftarrow \theta - (\hat{\eta}/\hat{B}) \sum_{\mathbf{x} \in \mathcal{X}} \nabla_{\theta} \mathcal{L}(\mathbf{x}; \theta, \zeta)$     ▶ SGD Update
     $\zeta \leftarrow \hat{\rho} \zeta + (1 - \hat{\rho}) \theta$                 ▶ EMA Update
  end for
end for
```

---

782 sizes. To remove ambiguity, in Algorithm 1 we provide pseudo-code for how to use Progressive  
783 Scaling.

784 In Algorithm 1, the prefactor of the SGD update could also have been written  $\eta/B$ , although an  
785 equivalent use of the base momentum is not possible.

786 Finally, we outline how to extend Algorithm 1 to more complex setups, like those presented in  
787 Section 3.4:

- 788 1. Optimizer scaling rules are used appropriately, for example the Adam scaling rule in case of  
789 using the Adam optimizer to update parameters  $\theta$ .
- 790 2. Schedules for hyperparameters are computed using the base hyperparameters, and are then  
791 modified by application of the scaling law in epoch (outer) loop.
- 792 3. Schedules for hyperparameters at the *step* rather than epoch level can be achieved in practice  
793 through recomputing the schedule and updating the notion of minibatch index appropriately  
794 throughout training.

795 All of the above techniques are used in Section 3.4. In addition, scheduling batch sizes within epoch  
796 is possible, providing one maintains a notion of computation within some fixed continuous time  
797  $T_{\text{fixed}}$ . We did not investigate this scenario.

## 798 D EMA approximation theorems with SDEs

### 799 D.1 SGD with model EMA

800 We will now derive the EMA scaling rule when tracking model parameters and the model is trained  
801 using SGD. We employ a strategy similar to Malladi et al. (2022), where we associate to each  
802 iterative process a Stochastic Differential Equation (SDE). In order to control the distance between  
803 the SDE and the discrete process, we use the tools from Li et al. (2019).

804 **Definition D.1** (Polynomial growth, Definition 1 in (Li et al., 2019)). *The set  $G$  is the set of con-*  
805 *tinuous functions  $\mathbb{R}^d \rightarrow \mathbb{R}$  with at most polynomial growth, i.e., for  $g \in G$  there exists two scalars*  
806  *$\kappa_1, \kappa_2 > 0$  such that for all  $\mathbf{x} \in \mathbb{R}^d$ , we have  $|g(\mathbf{x})| \leq \kappa_1(1 + \|\mathbf{x}\|^{\kappa_2})$ .*

807 *For an integer  $\alpha > 0$ ,  $G^\alpha$  is the set of functions  $\mathbb{R}^d \rightarrow \mathbb{R}$  that are  $\alpha$ -times continuously differentiable*  
808 *and such that all their derivatives up to order  $\alpha$  are in  $G$ .*

809 Similarly to Malladi et al. (2022), we use Noisy Gradient Oracle with Scale Parameter (NGOS) to  
810 define the update rules on the parameters.

811 **Definition D.2** (Noisy Gradient Oracle with Scale Parameter (NGOS), adaptation of (Malladi et al.,  
812 2022)). *A NGOs is a tuple  $\mathcal{G}_\sigma = (f, \Sigma, \mathcal{Z}_\sigma)$ . Given a noise scale parameter  $\sigma > 0$ , the NGOs  $\mathcal{G}_\sigma$*

813 takes as input the parameters  $\theta$  and outputs a random vector  $\mathbf{g} = \nabla f(\theta, \zeta) + \sigma \epsilon$  where  $\nabla f(\theta, \zeta)$  is  
814 the gradient of  $f$  with respect to  $\theta$  at  $(\theta, \zeta)$ , and  $\epsilon$  is a random vector drawn from the distribution  
815  $\mathcal{Z}_\sigma(\theta, \zeta)$  with zero mean and covariance  $\Sigma(\theta, \zeta)$ .

816 Note that in the above definition, the probability distribution  $\mathcal{Z}_\sigma(\theta, \zeta)$  is allowed to change with the  
817 scale  $\sigma$ , but its first two moments — its mean and its covariance — are fixed with  $\sigma$ . We have the  
818 following theorem for model EMA under optimization with SGD:

819 **Theorem D.1** (SDE for SGD + EMA). *Consider the couple  $\mathbf{x}_k = (\theta_k, \zeta_k)$  where  $\theta_k$  are the iterates  
820 of SGD with a NGOS (Definition D.2) and  $\zeta_k$  is an EMA of  $\theta_k$ , defined, starting from  $\mathbf{x}_0 = \mathbf{x}_0$ , by*

$$\theta_{k+1} = \theta_k - \eta \mathbf{g}_k, \text{ with } \mathbf{g}_k = \nabla f(\theta_k, \zeta_k) + \sigma \epsilon_k, \text{ and } \epsilon_k \sim \mathcal{Z}_\sigma(\theta_k, \zeta_k), \quad (14)$$

$$\zeta_{k+1} = \rho \zeta_k + (1 - \rho) \theta_k. \quad (15)$$

821 Define  $\beta_0 = (1 - \rho)/\eta$ ,  $\sigma_0 = \sigma\sqrt{\eta}$ , and define the SDE for  $X_t = (\Theta_t, Z_t)$ , starting from  $X_0 = \mathbf{x}_0$ , by

$$d\Theta_t = -\nabla f(\Theta_t, Z_t)dt + \sigma_0 \Sigma(\Theta_t, Z_t)^{\frac{1}{2}} dW_t, \text{ with } W_t \text{ a Wiener process} \quad (16)$$

$$dZ_t = \beta_0(\Theta_t - Z_t)dt. \quad (17)$$

822 Assume that  $f$  is continuously differentiable, with  $f \in G^3$  and  $\Sigma^{\frac{1}{2}} \in G^2$  (Definition D.1). Then, for  
823 any time horizon  $T > 0$  and test function  $g \in G^2$ , there exists a constant  $c > 0$  such that

$$\max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}[g(X_{\eta k})] - \mathbb{E}[g(\mathbf{x}_k)]| \leq c \times \eta. \quad (18)$$

824 *Proof.* The proof uses the same tools as in Li et al. (2019). Define  $\Delta(\theta, \zeta) = \eta(-\nabla f(\theta, \zeta) +$   
825  $\sigma \epsilon, \beta_0(\theta - \zeta))$  with  $\epsilon \sim \mathcal{Z}_\sigma(\theta, \zeta)$  the one-step update for the SGD + EMA update, such that  
826  $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta(\mathbf{x}_k)$ . We have the first two moments:

$$\mathbb{E}[\Delta(\theta, \zeta)] = \eta(-\nabla f(\theta, \zeta), \beta_0(\theta - \zeta)) \quad (19)$$

$$\mathbb{V}[\Delta(\theta, \zeta)] = \eta \sigma_0^2 \begin{bmatrix} \Sigma(\theta, \zeta) & 0 \\ 0 & 0 \end{bmatrix} \quad (20)$$

827 and the higher-order moments are  $O(\eta^2)$ . Similarly, let  $\tilde{\Delta}(\theta, \zeta)$  be the solution at time  $\eta$  of the SDE  
828 defined by Equation (6) starting from  $X_0 = (\theta, \zeta)$ . From Ito's formula, we also obtain

$$\mathbb{E}[\tilde{\Delta}(\theta, \zeta)] = \eta(-\nabla f(\theta), \beta_0(\theta - \zeta)) \quad (21)$$

$$\mathbb{V}[\tilde{\Delta}(\theta, \zeta)] = \eta \sigma_0^2 \begin{bmatrix} \Sigma(\theta, \zeta) & 0 \\ 0 & 0 \end{bmatrix} \quad (22)$$

829 and the higher-order moments are  $O(\eta^2)$ . Hence, the moments of the discrete iteration and of the  
830 SDE match up to second order. Following the same proof technique as in Li et al. (2019) then leads  
831 to the advertised theorem.  $\square$

832 This theorem is a simple adaptation of the results of Li et al. (2019). Intuitively, it is expected that  
833  $X_t$  and  $\mathbf{x}_k$  are close since  $\mathbf{x}_k$  is the Euler-Maruyama discretization of  $X_t$  with learning rate  $\eta$ . We  
834 then have the corollary.

835 **Corollary D.1.1** (Validity of the EMA Scaling Rule). *Assume that  $f$  is continuously differentiable,  
836 with  $f \in G^3$  and  $\Sigma^{\frac{1}{2}} \in G^2$ . Let  $\theta_k^B, \zeta_k^B$  the iterates of the Equation (5) with batch size  $B$  and  
837 hyperparameters  $\eta, \rho$ . Let  $\theta_k^{\kappa B}, \zeta_k^{\kappa B}$  be iterates with batch size  $\kappa B$ , learning rate  $\eta$  determined by  
838 the SGD Scaling Rule (Definition 2.2) and momentum determined by the EMA Scaling Rule, linear  
839 version (Definition 1.1). Then, for any time horizon  $T > 0$  and function  $g \in G^2$ , there exists a  
840 constant  $d > 0$  such that*

$$\max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}[g(\theta_{\lfloor k/\kappa \rfloor}^{\kappa B}, \zeta_{\lfloor k/\kappa \rfloor}^{\kappa B})] - \mathbb{E}[g(\theta_k, \zeta_k)]| \leq d \times \eta. \quad (23)$$

841 *Proof.* The proof is similar to Malladi et al. (2022). Under the scaling rule, both  $\mathbf{x}_k = (\theta_k, \zeta_k)$  and  
842  $\hat{\mathbf{x}}_{\lfloor k/\kappa \rfloor} = (\theta_{\lfloor k/\kappa \rfloor}^{\kappa B}, \zeta_{\lfloor k/\kappa \rfloor}^{\kappa B})$  have the same limiting SDE. Hence we have from the previous theorem  
843 that for all test function  $g$ , we can find  $c, c'$  such that

$$\max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}[g(X_{\eta k})] - \mathbb{E}[g(\mathbf{x}_k)]| \leq c \times \eta \text{ and } \max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}[g(X_{\eta k})] - \mathbb{E}[g(\hat{\mathbf{x}}_{\lfloor k/\kappa \rfloor})]| \leq c' \times \eta. \quad (24)$$

844 The triangle inequality then gives

$$\max_{k=0, \dots, \lfloor T/\eta \rfloor} |\mathbb{E}[g(\hat{\mathbf{x}}_{\lfloor k/\kappa \rfloor})] - \mathbb{E}[g(\mathbf{x}_k)]| \leq (c + c') \times \eta. \quad (25)$$

845 Hence, taking  $d = c + c'$  gives the expected result.  $\square$

## 846 D.2 Adaptive gradient methods with model EMA

847 We now turn to the case where one uses an adaptive gradient method rather than SGD to train the  
848 model. We follow derivations similar to those of Malladi et al. (2022), with an added EMA. Like  
849 above, we consider that the loss function  $f$  also depends on the EMA tracking parameter  $\zeta_k$ . We  
850 begin with RMSProp with EMA, which iterates:

$$\mathbf{v}_{k+1} = \gamma \mathbf{v}_k + (1 - \gamma) \mathbf{g}_k^2, \quad \text{with } \mathbf{g}_k = \nabla f(\boldsymbol{\theta}_k, \zeta_k) + \sigma \boldsymbol{\epsilon}_k, \quad \text{and } \boldsymbol{\epsilon}_k \sim \mathcal{Z}_\sigma(\boldsymbol{\theta}_k, \zeta_k), \quad (26)$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta(\sqrt{\mathbf{v}_k} + \varepsilon)^{-1} \times \mathbf{g}_k \quad (27)$$

$$\zeta_{k+1} = \rho \zeta_k + (1 - \rho) \boldsymbol{\theta}_k. \quad (28)$$

851 Like in Malladi et al. (2022), we place ourselves in the high noise regime, in which the term  $\mathbf{g}_k^2$  in  
852 Equation (26) is approximated by  $\mathbf{g}_k^2 \approx \sigma^2 \text{diag}(\Sigma(\boldsymbol{\theta}_k, \zeta_k))$ . We use the same scaling rules, with an  
853 additional one for  $\rho$ :

$$\gamma_0 = (1 - \gamma)/\eta^2, \quad \sigma_0 = \sigma\eta, \quad \varepsilon_0 = \varepsilon\eta, \quad \text{and } \beta_0 = (1 - \rho)/\eta^2, \quad (29)$$

854 and we let  $\mathbf{u}_k = \mathbf{v}_k/\sigma^2$ . The equations for RMSProp with EMA then become, using only these new  
855 variables and  $\eta$ :

$$\mathbf{u}_{k+1} - \mathbf{u}_k = \eta^2 \gamma_0 (\text{diag}(\Sigma(\boldsymbol{\theta}_k, \zeta_k)) - \mathbf{u}_k), \quad (30)$$

$$\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k = -(\sqrt{\mathbf{u}_k} + \varepsilon_0)^{-1} (\eta^2 \nabla f(\boldsymbol{\theta}_k, \zeta_k) + \eta \boldsymbol{\epsilon}_k) \quad (31)$$

$$\zeta_{k+1} - \zeta_k = \eta^2 \beta_0 (\boldsymbol{\theta}_k - \zeta_k). \quad (32)$$

856 This formulation makes it clear that these iterations can be seen as the discretization of the SDE

$$dU_t = \gamma_0 (\text{diag}(\Sigma(\boldsymbol{\Theta}_t, Z_t)) - U_t) dt, \quad (33)$$

$$d\boldsymbol{\Theta}_t = -(\sigma_0 \sqrt{U_t} + \varepsilon_0)^{-1} (\nabla f(\boldsymbol{\Theta}_t, Z_t) dt + \sigma_0 \Sigma(\boldsymbol{\Theta}_t, Z_t)^{1/2} dW_t) \quad (34)$$

$$dZ_t = \beta_0 (\boldsymbol{\Theta}_t - Z_t) dt, \quad (35)$$

with step size  $\eta^2$ . Of course, we recover the SDE of Malladi et al. (2022) in the case where  $\beta_0 = 0$ . A formal proof of closeness between the iterates and the SDE trajectory is out of the scope of the present paper since it would imply redoing much of the theoretical work developed in Malladi et al. (2022). Still, the previous informal analysis hints that for RMSProp, the scaling rule in Equation (29) should be used. In other words, given a certain set of hyperparameters  $\gamma, \eta$  and  $\rho$ , if the batch size goes from  $B$  to  $\hat{B} = \kappa \times B$ , the noise level becomes  $\hat{\sigma} = \sigma/\sqrt{\kappa}$ , and keeping the quantities in Equation (29) constant means that we should use as new hyperparameters

$$\hat{\gamma} = 1 - (1 - \gamma) \times \kappa, \quad \hat{\eta} = \eta \times \sqrt{\kappa}, \quad \text{and } \hat{\rho} = 1 - (1 - \rho) \times \kappa.$$

857 The linear rule  $\hat{\rho} = 1 - (1 - \rho) \times \kappa$  is at the first order equivalent to the exponential scaling rule  
858  $\hat{\rho} = \rho^\kappa$ . Hence, even though the limiting SDE differs greatly from that of SGD, and even though the  
859 scaling rule regarding the learning rate differs, we recover for the momentum term  $\rho$  the exact same  
860 scaling rule as for SGD.

861 We finish the discussion with the case of Adam, which leads once again to the same rule as for SGD.  
862 Adam with EMA tracking of the network parameters iterates

$$\mathbf{m}_{k+1} = \beta_1 \mathbf{m}_k + (1 - \beta_1) \mathbf{g}_k, \quad \text{with } \mathbf{g}_k = \nabla f(\boldsymbol{\theta}_k, \zeta_k) + \sigma \boldsymbol{\epsilon}_k, \quad \text{and } \boldsymbol{\epsilon}_k \sim \mathcal{Z}_\sigma(\boldsymbol{\theta}_k, \zeta_k), \quad (36)$$

$$\mathbf{v}_{k+1} = \beta_2 \mathbf{v}_k + (1 - \beta_2) \mathbf{g}_k^2 \quad (37)$$

$$\tilde{\mathbf{m}}_{k+1} = \mathbf{m}_{k+1}/(1 - \beta_1^{k+1}) \quad (38)$$

$$\tilde{\mathbf{v}}_{k+1} = \mathbf{v}_{k+1}/(1 - \beta_2^{k+1}) \quad (39)$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta(\sqrt{\tilde{\mathbf{v}}_k} + \varepsilon)^{-1} \times \tilde{\mathbf{m}}_{k+1} \quad (40)$$

$$\zeta_{k+1} = \rho \zeta_k + (1 - \rho) \boldsymbol{\theta}_k. \quad (41)$$



863 Here, we use the same minor modification of the iterations as in Malladi et al. (2022), where we use  
 864  $\mathbf{v}_k$  instead of  $\mathbf{v}_{k+1}$  in the denominator of the  $\Theta_k$  update.

865 We consider the following scaling for the hyperparameters

$$c_1 = (1 - \beta_1)/\eta^2, \quad c_2 = (1 - \beta_2)/\eta^2, \quad \sigma_0 = \sigma\eta, \quad \varepsilon_0 = \varepsilon\eta, \quad \text{and } \beta_0 = (1 - \rho)/\eta^2, \quad (42)$$

866 and  $\gamma_1(t) = 1 - \exp(-c_1 t)$ ,  $\gamma_2(t) = 1 - \exp(-c_2 t)$ , and  $\mathbf{u}_k = \mathbf{v}_k/\sigma^2$ . The SDE for Adam + EMA is  
 867 given by

$$dM_t = c_1 \left( (\nabla f(\Theta_t, Z_t) - M_t) dt + \sigma_0 \Sigma(\Theta_t, Z_t)^{1/2} dW_t \right) \quad (43)$$

$$dU_t = c_2 (\text{diag}(\Sigma(\Theta_t, Z_t)) - U_t) dt \quad (44)$$

$$d\Theta_t = -\frac{\sqrt{\gamma_2(t)}}{\gamma_1(t)} (\sigma_0 \sqrt{U_t} + \varepsilon_0 \sqrt{\gamma_2(t)})^{-1} \times M_t dt \quad (45)$$

$$dZ_t = \beta_0 (\Theta_t - Z_t) dt. \quad (46)$$

This is once again the same SDE as in Malladi et al. (2022) with the added EMA term. Like previously, this SDE hints at the fact that the scaling rule in eq. (42) should be used. In other words, given a set of hyperparameters  $\beta_1, \beta_2, \eta$ , and  $\rho$ , if the batch size goes from  $B$  to  $\kappa \times B$ , then the noise level becomes  $\hat{\sigma} = \sigma/\sqrt{\kappa}$  and keeping quantities in eq. (42) constant means that we should use as new hyperparameters

$$\hat{\beta}_1 = 1 - (1 - \beta_1) \times \kappa, \quad \hat{\beta}_2 = 1 - (1 - \beta_2) \times \kappa, \quad \hat{\eta} = \eta \times \sqrt{\kappa}, \quad \text{and } \hat{\rho} = 1 - (1 - \rho) \times \kappa.$$

868 We once again recover a linear rule for  $1 - \rho$  which is equivalent to the exponential scaling rule  
 869  $\hat{\rho} = \rho^\kappa$  in the limit  $\rho \rightarrow 0$ .

## 870 E Additional proofs

### 871 E.1 Iterations of SGD + EMA

872 Here we derive a critical component of the EMA Scaling Rule, the matrix equation of Equation (4)  
 873 from which the EMA Scaling Rule (Definition 1.1) follows.

874 **Theorem E.1** (Iterations of SGD + EMA). *Assuming that gradients change slowly over iterations*  
 875 *of SGD (Definition 2.1) and EMA (Definition 2.3):  $\nabla_{\theta} \mathcal{L}(x; \theta_{t+j}, \zeta_{t+j}) \approx \nabla_{\theta} \mathcal{L}(x; \theta_t, \zeta_t) \approx \mathbf{g}$ , for*  
 876  *$j = 1, 2, \dots, \kappa$  and representative gradient  $\mathbf{g}$ , iterating over  $\kappa$  independent minibatches produces*  
 877 *model states*

$$\begin{bmatrix} \theta_{t+\kappa} \\ \zeta_{t+\kappa} \\ \mathbf{g} \end{bmatrix} = \begin{bmatrix} 1 & 0 & -\eta \\ 1 - \rho & \rho & 0 \\ 0 & 0 & 1 \end{bmatrix}^{\kappa} \cdot \begin{bmatrix} \theta_t \\ \zeta_t \\ \mathbf{g} \end{bmatrix} = \begin{bmatrix} \theta_t - \eta \kappa \mathbf{g} \\ \rho^{\kappa} \zeta_t + (1 - \rho^{\kappa}) \theta_t + O(\eta \times \beta_{\rho}) \\ \mathbf{g} \end{bmatrix}. \quad (47)$$

878 *Proof.* First note that for matrices of the form

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & a_{0,2} \\ 1 - a_{1,1} & a_{1,1} & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (48)$$

879 their multiplication follows

$$\begin{aligned} \mathbf{A} \mathbf{B} &= \begin{bmatrix} 1 & 0 & a_{0,2} \\ 1 - a_{1,1} & a_{1,1} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & b_{0,2} \\ 1 - b_{1,1} & b_{1,1} & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & a_{0,2} + b_{0,2} \\ 1 - a_{1,1} b_{1,1} & a_{1,1} b_{1,1} & (1 - a_{1,1}) b_{0,2} \\ 0 & 0 & 1 \end{bmatrix}, \end{aligned} \quad (49)$$

880 and

$$\begin{aligned} \mathbf{A} \mathbf{B} \mathbf{C} &= \begin{bmatrix} 1 & 0 & a_{0,2} + b_{0,2} \\ 1 - a_{1,1} b_{1,1} & a_{1,1} b_{1,1} & (1 - a_{1,1}) b_{0,2} \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & c_{0,2} \\ 1 - c_{1,1} & c_{1,1} & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & a_{0,2} + b_{0,2} + c_{0,2} \\ 1 - a_{1,1} b_{1,1} c_{1,1} & a_{1,1} b_{1,1} c_{1,1} & (1 - a_{1,1}) b_{0,2} + (1 - a_{1,1} b_{1,1}) c_{0,2} \\ 0 & 0 & 1 \end{bmatrix}. \end{aligned} \quad (50)$$

881 By induction

$$\mathbf{A}^\kappa = \begin{bmatrix} 1 & 0 & \kappa \times a_{0,2} \\ 1 - a_{1,1}^\kappa & a_{1,1}^\kappa & \delta(a_{0,2}, a_{1,1}, \kappa) \\ 0 & 0 & 1 \end{bmatrix}, \quad (51)$$

882 where

$$\delta(a_{0,2}, a_{1,1}, \kappa) = a_{0,2} \sum_{i=1}^{\kappa-1} (1 - a_{1,1}^i) = a_{0,2} \left( \kappa - \frac{1 - a_{1,1}^\kappa}{1 - a_{1,1}} \right), \quad \text{for } a_{1,1} \neq 1. \quad (52)$$

883 It follows that

$$\begin{bmatrix} 1 & 0 & -\eta \\ 1 - \rho & \rho & 0 \\ 0 & 0 & 1 \end{bmatrix}^\kappa = \begin{bmatrix} 1 & 0 & -\kappa \eta \\ 1 - \rho^\kappa & \rho^\kappa & \delta(\eta, \rho, \kappa) \\ 0 & 0 & 1 \end{bmatrix} \quad (53)$$

884 where the EMA Scaling Rule error

$$\delta(\eta, \rho, \kappa) = (-\eta) \left( \kappa - \frac{1 - \rho^\kappa}{1 - \rho} \right) \approx (-\eta) (\kappa - \kappa + \mathcal{O}(\beta_\rho)) = 0 + \mathcal{O}(\eta \times \beta_\rho), \quad (54)$$

885 where  $\beta_\rho \equiv 1 - \rho$  and the approximation is around  $\rho = 1$ .  $\square$

## 886 E.2 Limiting behavior of Polyak-Ruppert averaging

887 Here we sketch the asymptotic behavior of a target model  $\theta$  and its EMA  $\zeta$ . Let us assume that  $\theta$   
 888 converges to the stationary distribution  $\lim_{t \rightarrow \infty} \theta_t = \theta^*$ ,  $\theta^* \sim p_\infty(\theta)$ . We are interested in statistical  
 889 properties of  $\zeta^* = \lim_{t \rightarrow \infty} \zeta_t$ , as this will formalize the notion of how the EMA depends on the a  
 890 time-horizon defined by its momentum  $\rho$  as discussed in Table 1.

891 As a warm-up, for  $n$  independent random variables  $x_1, \dots, x_n$ , we know that the sample mean  $\bar{x} =$   
 892  $\frac{1}{n}(x_1, x_2, \dots, x_n)$  has the statistical properties

$$\mathbb{E}[\bar{x}] = \mu, \quad \text{Var}[\bar{x}] = \frac{\sigma^2}{n}, \quad (55)$$

893 where  $\mu$  and  $\sigma$  are the population mean and variance. This gives us an idea of what to expect. As we  
 894 will now show, the expectation of  $\zeta^*$  should have no time-horizon dependence, whereas the variance  
 895 of  $\zeta^*$  will depend on its time horizon (i.e. the number of samples it integrates over) which is defined  
 896 by  $\rho$ .

897 In the case of a weighted sum

$$\bar{x}^{(w)} = \sum_{i=1}^n w_i x_i, \quad (56)$$

898 then if the  $x_i$  are Independent and Identically Distributed (i.i.d.), then

$$\mathbb{E}[\bar{x}^{(w)}] = \sum_{i=1}^n w_i \mathbb{E}[x_i] = n \bar{w} \mu, \quad \bar{w} = \frac{1}{n} \sum_{i=1}^n w_i, \quad (57)$$

899 and for the variance (Kish, 1965)

$$\text{Var}[\bar{x}^{(w)}] = n \cdot \bar{w}^2 \cdot \sigma^2 \quad \bar{w}^2 = \frac{1}{n} \sum_{i=1}^n w_i^2, \quad \sigma^2 = \text{Var}[x_i]. \quad (58)$$

900 We can verify that we reproduce the well-known result in Equation (55) in the case where all weights  
 901 are equal to  $\frac{1}{n}$  as follows

$$\forall i : w_i = \frac{1}{n} \implies \bar{w}^2 = \frac{1}{n} \cdot \sum_{i=1}^n \left( \frac{1}{n} \right)^2 = \frac{1}{n^2} \implies \text{Var}[\bar{x}^{(w)}] = n \cdot \frac{1}{n^2} \cdot \sigma^2 = \frac{\sigma^2}{n}. \quad (59)$$

902 In the case of an exponential moving average we have

$$\zeta_{t+1} = \rho \zeta_t + (1 - \rho) \theta_t = \rho^t \zeta_1 + (1 - \rho) \sum_{i=0}^{t-1} \rho^i \theta_{t-i}. \quad (60)$$

903 Let's consider the specific case where we are at iteration  $k$  which is sufficiently large that  $\zeta$  and  $\theta$   
 904 have converged to their stationary distributions. From  $k$ , the iterations unfold as

$$\zeta_{t+1} = \rho^{t+1-k} \zeta_k + (1 - \rho) \sum_{i=0}^{t-k} \rho^i \theta_{t-i}. \quad (61)$$

905 We rearrange for terms in  $\zeta$

$$\zeta_{t+1} - \rho^{t+1-k} \zeta_k = (1 - \rho) \sum_{i=0}^{t-k} \rho^i \theta_{t-i}, \quad (62)$$

906 and before proceeding to the final result, using  $n = t + 1 - k$ , we compute the convenient quantities

$$\bar{\rho} = \frac{1}{n} \sum_{i=0}^{n-1} \rho^i = \frac{1}{n} \times \frac{1 - \rho^n}{1 - \rho} \quad (63)$$

$$\overline{\rho^2} = \frac{1}{n} \sum_{i=0}^{n-1} \rho^{2i} = \frac{1}{n} \times \frac{1 - \rho^{2n}}{1 - \rho^2}. \quad (64)$$

907 Taking expectation of Equation (62) and setting statistics to their stationary values, we have

$$(1 - \rho^n) \mathbb{E}[\zeta^*] = (1 - \rho) n \bar{\rho} \mathbb{E}[\theta^*] = (1 - \rho^n) \mathbb{E}[\theta^*], \quad (65)$$

908 where we have used the result in Equation (57). It follows that for  $\rho \neq 1$  we have

$$\mathbb{E}[\zeta^*] = \mathbb{E}[\theta^*], \quad (66)$$

909 independent of  $\rho$ . Finally, we can take the variance of Equation (62). First the left hand side

$$\text{Var} [\zeta_{t+1} - \rho^n \zeta_k] = \text{Var} [\zeta_{t+1}] + \rho^{2n} \text{Var} [\zeta_k] = (1 + \rho^{2n}) \text{Var} [\zeta^*]. \quad (67)$$

910 Next the right hand side

$$\text{Var} \left[ (1 - \rho) \sum_{i=0}^{n-1} \rho^i \theta_{t-i} \right] = (1 - \rho)^2 \text{Var} \left[ \sum_{i=0}^{n-1} \rho^i \theta_{t-i} \right] = (1 - \rho)^2 \cdot \left( \frac{1 - \rho^{2n}}{1 - \rho^2} \right) \cdot \text{Var}[\theta^*]. \quad (68)$$

911 Finally, equating left and right hand sides and rearranging for  $\text{Var}[\zeta^*]$  gives

$$\text{Var} [\zeta^*] = \frac{1 - \rho^{2n}}{1 + \rho^{2n}} \cdot \frac{1 - \rho}{1 + \rho} \cdot \text{Var} [\theta^*] \quad (69)$$

912 In the limit  $t \rightarrow \infty$ , the momentum-dependent prefactor becomes

$$\lim_{t \rightarrow \infty} \left( \frac{1 - \rho^{2n}}{1 + \rho^{2n}} \cdot \frac{1 - \rho}{1 + \rho} \right) = \frac{1 - \rho}{1 + \rho} \implies \lim_{t \rightarrow \infty} \text{Var} [\zeta^*] = \frac{1 - \rho}{1 + \rho} \cdot \text{Var} [\theta^*]. \quad (70)$$

913 Equations (69) and (70) validate our intuition. When  $\rho \rightarrow 0$ , then  $\zeta$  behaves like  $\theta$  independent of  
 914  $T$ , with their variance and expectation matching. When  $\rho > 0$ , the momentum-dependent prefactor  
 915 serves as an aggregator over the history when  $t$  is sufficiently large compared to  $k$ , reducing the  
 916 variance  $\text{Var}[\zeta^*]$  but preserving its expectation. This formalizes the notion of time horizon discussed  
 917 in Table 1.

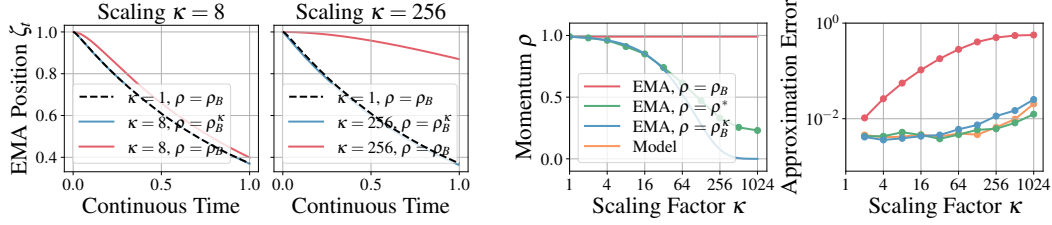
## 918 F Additional details and results for Polyak-Ruppert averaging

919 **Additional background** Polyak-Ruppert averaging (Definition 3.1) is a simplification of Stochastic  
 920 Weight Averaging (SWA) (Izmailov et al., 2018) which uses a more complex multi-cycle sched-  
 921 ule based weighting of the model parameters. Both Definition 3.1 and SWA present similar favor-  
 922 able properties like wider minima and better generalization (Izmailov et al., 2018). For example,  
 923 He et al. (2022) observed that a supervised ViT-H/14 overfits on ImageNet1k (Russakovsky et al.,  
 924 2014) without a model EMA, achieving an accuracy of 80.9%. Equipping a Polyak-Ruppert average  
 925 ( $\rho = 0.9999$ ) alleviated overfitting and gave a 83.1% accuracy.

926 **Organization** In this appendix, we look at additional momenta for one-dimensional noisy  
 927 parabola, as well as extensions to  $D$ -dimensions (Appendix F.1), provide a more detailed view  
 928 of the results of Section 3.2 (Appendix F.2), and investigate the scenario where the EMA Scal-  
 929 ing Rule (Definition 1.1) is applied to batch normalization (Ioffe & Szegedy, 2015) coefficients  
 930 (Appendix F.3).

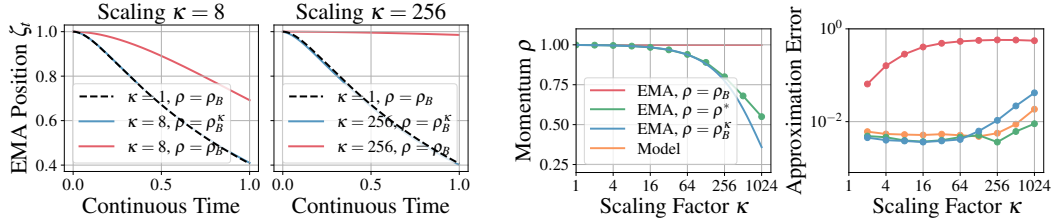
## 931 F.1 Noisy parabola

932 **Additional one-dimensional examples** First we consider additional one-dimensional examples,  
 933 investigating the effect of modifying the base momentum  $\rho_B$ . We present  $\rho_B = 0.99$  in Figure 7, and  
 934  $\rho_B = 0.999$  in Figure 8. The results for  $\rho_B = 0.9999$  are presented in main text in Figure 1.



(a) Trajectory of the model EMA  $\zeta$  under different scalings  $\kappa$ , with  $\rho_B = 0.99$ ,  $\eta_B = 10^{-4}$ . (b) Choices for momentum (left) with corresponding approximation errors (Equation (10)) (right).

Figure 7: (a) We show the effect of scaling by comparing model EMA trajectories of the baseline ( $\kappa = 1$ , black dashed) to  $\kappa = 8$  (left) and  $\kappa = 256$  (right), with ( $\rho = \rho_B^k$ , blue) and without ( $\rho = \rho_B$ , red) the EMA Scaling Rule. (b, left) The momentum according for different scaling rules and the empirically optimal  $\rho^*$  (Equation (10)). (b, right) The approximation error (Equation (10)) of trajectories in (b, left) and the target model (orange). Error for  $\rho^*$  is computed using a hold-out to mitigate overfitting.



(a) Trajectory of the model EMA  $\zeta$  under different scalings  $\kappa$ , with  $\rho_B = 0.999$ ,  $\eta_B = 10^{-4}$ . (b) Choices for momentum (left) with corresponding approximation errors (Equation (10)) (right).

Figure 8: (a) We show the effect of scaling by comparing model EMA trajectories of the baseline ( $\kappa = 1$ , black dashed) to  $\kappa = 8$  (left) and  $\kappa = 256$  (right), with ( $\rho = \rho_B^k$ , blue) and without ( $\rho = \rho_B$ , red) the EMA Scaling Rule. (b, left) The momentum according for different scaling rules and the empirically optimal  $\rho^*$  (Equation (10)). (b, right) The approximation error (Equation (10)) of trajectories in (b, left) and the target model (orange). Error for  $\rho^*$  is computed using a hold-out to mitigate overfitting.

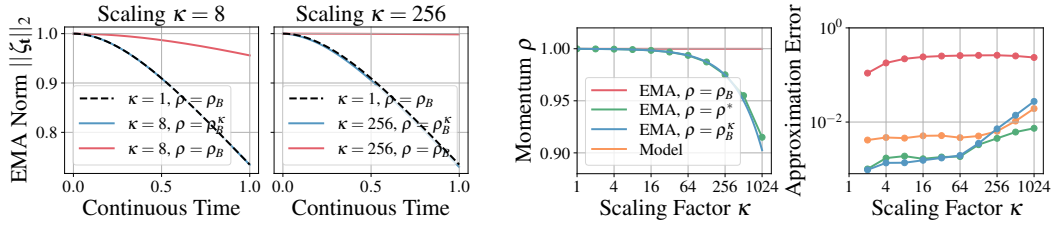
935 As described by the scaling error term in Equation (54), the approximation error at a given  $\kappa$  is  
 936 higher for lower momenta  $\rho$ . For a large range of scalings  $\kappa$ , the EMA Scaling Rule and the optimal  
 937 momenta  $\rho^*$  are consistent. In summary, we see the synthetic experiments validate the results of  
 938 Section 3.1 for a range of momenta  $\rho$ .

939 **Examples in higher dimensions** Our final use of the synthetic *noisy* parabola will consider an  
 940 extension to  $D$  dimensions. Consider the optimization of  $\theta \in \mathbb{R}^D$  in a *noisy parabola* at the origin:

$$\mathcal{L}(\theta) = \frac{a}{2} \theta^\top \theta, \quad \theta_{k+1} = \theta_k - \eta \mathbf{g}_k, \quad \mathbf{g}_k = a \theta_k + \epsilon_k, \quad \epsilon_k \sim \mathcal{N}\left(\mathbf{0}, \frac{b \mathbf{g}_k^2 + c}{\kappa}\right), \quad (71)$$

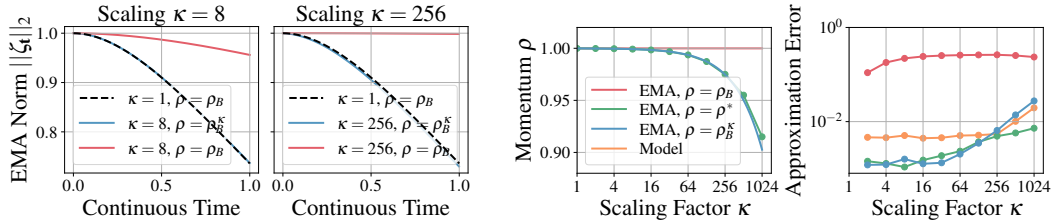
941 for curvature  $a > 0$ , scaled additive  $b > 0$ , and additive  $c > 0$  noise coefficients. The scaling factor  
 942  $\kappa$  in the covariance denominator implements the reduction in gradient noise as the scaling (i.e., the  
 943 batch size) increases (Jastrzebski et al., 2017). Let  $\theta \in \mathbb{R}^D$  be optimized with SGD (Definition 2.1)

944 and let there be a Polyak-Ruppert average (Definition 3.1)  $\zeta \in \mathbb{R}^D$  with momentum  $\rho = 1 - \beta$  for  $\theta$ .  
 945 We consider dimensionalities  $D = 2$  (Figure 9),  $D = 16$  (Figure 10), and  $D = 100$  (Figure 11). We  
 946 observe no significant differences in the EMA scaling behavior as we vary dimensions.



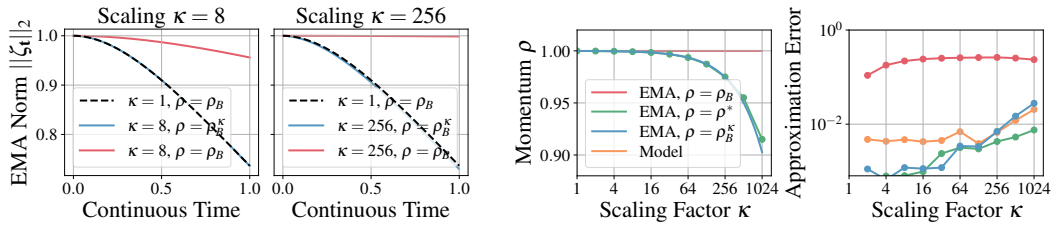
(a) Norm of the model EMA  $\zeta$  under different scalings  $\kappa$ , with  $\rho_B = 0.9999$ ,  $\eta_B = 10^{-4}$ ,  $D = 2$ . (b) Choices for momentum (left) with corresponding approximation errors (Equation (10)) (right).

Figure 9: (a) We show the effect of scaling by comparing model EMA trajectories of the baseline ( $\kappa = 1$ , black dashed) to  $\kappa = 8$  (left) and  $\kappa = 256$  (right), with ( $\rho = \rho_B^{\kappa}$ , blue) and without ( $\rho = \rho_B$ , red) the EMA Scaling Rule. (b, left) The momentum according for different scaling rules and the empirically optimal  $\rho^*$  (Equation (10)). (b, right) The approximation error (Equation (10)) of trajectories in (b, left) and the target model (orange). Error for  $\rho^*$  is computed using a hold-out to mitigate overfitting.



(a) Norm of the model EMA  $\zeta$  under different scalings  $\kappa$ , with  $\rho_B = 0.9999$ ,  $\eta_B = 10^{-4}$ ,  $D = 16$ . (b) Choices for momentum (left) with corresponding approximation errors (Equation (10)) (right).

Figure 10: (a) We show the effect of scaling by comparing model EMA trajectories of the baseline ( $\kappa = 1$ , black dashed) to  $\kappa = 8$  (left) and  $\kappa = 256$  (right), with ( $\rho = \rho_B^{\kappa}$ , blue) and without ( $\rho = \rho_B$ , red) the EMA Scaling Rule. (b, left) The momentum according for different scaling rules and the empirically optimal  $\rho^*$  (Equation (10)). (b, right) The approximation error (Equation (10)) of trajectories in (b, left) and the target model (orange). Error for  $\rho^*$  is computed using a hold-out to mitigate overfitting.



(a) Norm of the model EMA  $\zeta$  under different scalings  $\kappa$ , with  $\rho_B = 0.9999$ ,  $\eta_B = 10^{-4}$ ,  $D = 100$ . (b) Choices for momentum (left) with corresponding approximation errors (Equation (10)) (right).

Figure 11: (a) We show the effect of scaling by comparing model EMA trajectories of the baseline ( $\kappa = 1$ , black dashed) to  $\kappa = 8$  (left) and  $\kappa = 256$  (right), with ( $\rho = \rho_B^{\kappa}$ , blue) and without ( $\rho = \rho_B$ , red) the EMA Scaling Rule. (b, left) The momentum according for different scaling rules and the empirically optimal  $\rho^*$  (Equation (10)). (b, right) The approximation error (Equation (10)) of trajectories in (b, left) and the target model (orange). Error for  $\rho^*$  is computed using a hold-out to mitigate overfitting.

Table 4: Supervised ResNet-v2 50 hyperparameters used in Polyak-Ruppert Averaging experiments.

Supervised ResNet-v2 50	
ImageNet1k Test Top-1	76.27 ± 0.10%
ImageNet1k EMA Test Top-1	76.55 ± 0.07%
Weight initialization	kaiming_normal (relu)
Backbone normalization	BatchNorm
Synchronized BatchNorm over replicas	No
Learning rate schedule	Multi step: ×0.1 at [30, 60, 80] epochs
Learning rate warmup (epochs)	5
Learning rate minimum value	1 × 10 <sup>-6</sup>
Training duration (epochs)	90
Optimizer	SGD + Momentum
SGD momentum	0.9
Optimizer scaling rule	Linear
Base learning rate	0.4
Base batch size	1024
Base Polyak momentum	0.9999
Weight decay	1 × 10 <sup>-4</sup>
Weight decay scaling rule	None
Weight decay skip bias	Yes
Numerical precision	bf16
Augmentation stack	ImageNet
Label smoothing rate	0.1

947 **F.2 Image Classification**

948 **Hyperparameters** We present the base hyperparameters for our image experiments in Table 4.

949 **Data** For large scale vision evaluation, we use the ImageNet1k dataset (Russakovsky et al., 2014),  
 950 a widely used dataset containing approximately 1.2 million labeled images, distributed almost uni-  
 951 formly across 1000 different object classes, like animals, plants, and vehicles.

952 The images in ImageNet1k are not consistent in resolution. To handle this, they are resized and  
 953 cropped to a standard size (in our case, 224 × 224), before further processing. This is part of the  
 954 standard ImageNet augmentation stack for convolutional networks mentioned in Table 4.

955 **Compute** *[This section has been redacted to preserve anonymity during the peer-review process.*  
 956 *If this work is accepted, the full details compute used for these experiments, including: the experi-*  
 957 *ments presented, hyperparameter optimization, and the development process, will be provided.]*

958 **Additional results** In Figure 12 we present a more detailed view of the results in Section 3.2. First,  
 959 we see that for all train metrics, model trajectories match, and that a learning rate step schedule after  
 960 warmup is present. As discussed in Figure 12, a gap in EMA Test Top-1 trajectories begins at scaling  
 961  $\kappa = 4$ , with a more pronounced effect visible at  $\kappa = 8$ . From Figure 12 it is clear that the (non-EMA)  
 962 Test Top-1 performance trajectory is no longer matching at these scalings, demonstrating that the  
 963 problem is not due to a breakdown of the EMA Scaling Rule, but rather, that the model is overfitting  
 964 at larger batch sizes due to batch normalization (Ioffe & Szegedy, 2015).

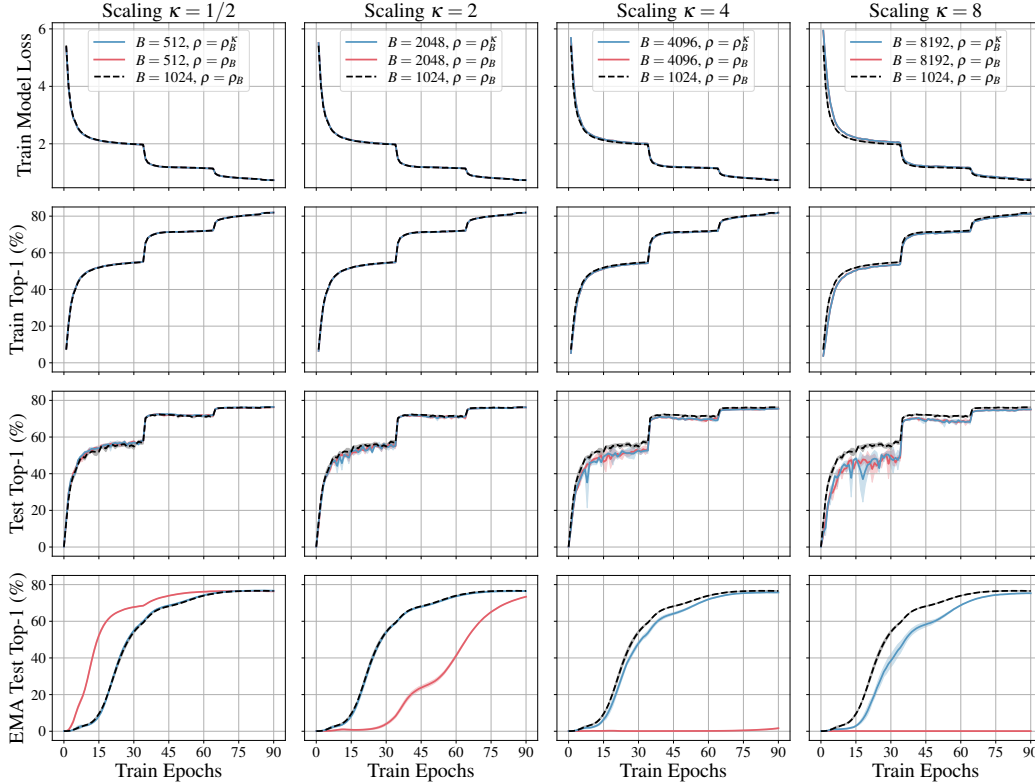


Figure 12: *ResNetv2-50 Polyak-Ruppert averaging on ImageNet1k* for different scalings  $\kappa$ . The baseline model ( $\kappa = 1$ , black dashed) uses batch size 1024 and momentum  $\rho_B = 0.9999$ , is scaled down to a batch size of 512 (left), and up to a batch size of 4096 (right) with (blue,  $\rho = \rho_B^\kappa$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule (Definition 1.1). Bands indicate the mean and standard deviation across three runs.

### 965 F.3 Applying the EMA Scaling Rule to Batch Normalization

966 In Section 3.2 and Appendix F.2, we investigated a range of scalings  $\kappa$ , *with* and *without* applying  
 967 the EMA Scaling Rule to the Polyak momentum. In those experiments, we maintained batch normal-  
 968 ization (Ioffe & Szegedy, 2015) coefficients of  $\rho_{\text{BN}} = 0.9$  throughout<sup>8</sup>, i.e. the EMA Scaling  
 969 Rule is not applied. Yet, the running statistics of Batch Normalization *are* an EMA with values  
 970 determined by  $\rho_{\text{BN}}$  and so it is reasonable to suspect we should apply the EMA Scaling Rule to  $\rho_{\text{BN}}$   
 971 also.

972 In Figure 13 we investigate the effect of applying the EMA Scaling Rule to Batch Normalization  
 973 coefficients, using  $\hat{\rho}_{\text{BN}} = \rho_{\text{BN}}^\kappa$ . We observe that the Test Top-1 trajectories *with* the EMA Scaling  
 974 Rule applied are slightly closer to the reference trajectories for scalings  $\kappa \geq 2$  than those trajectories  
 975 *without* the EMA Scaling Rule. As the effect is not particularly large, at least in this setup, we do  
 976 pursue further ablating applications of the EMA Scaling Rule to batch normalization coefficients,  
 977 and always use  $\rho_{\text{BN}} = 0.1$  for Batch Normalization, independent of  $\kappa$ .

<sup>8</sup>In many ML frameworks, this value is defined using  $\beta_\rho = 1 - \rho$ , i.e. the default is 0.1 and corresponds to  $\beta_{\text{BN}}$  rather than 0.9 corresponding to  $\rho_{\text{BN}}$ . We use  $\rho_{\text{BN}}$  to maintain consistency across this work.

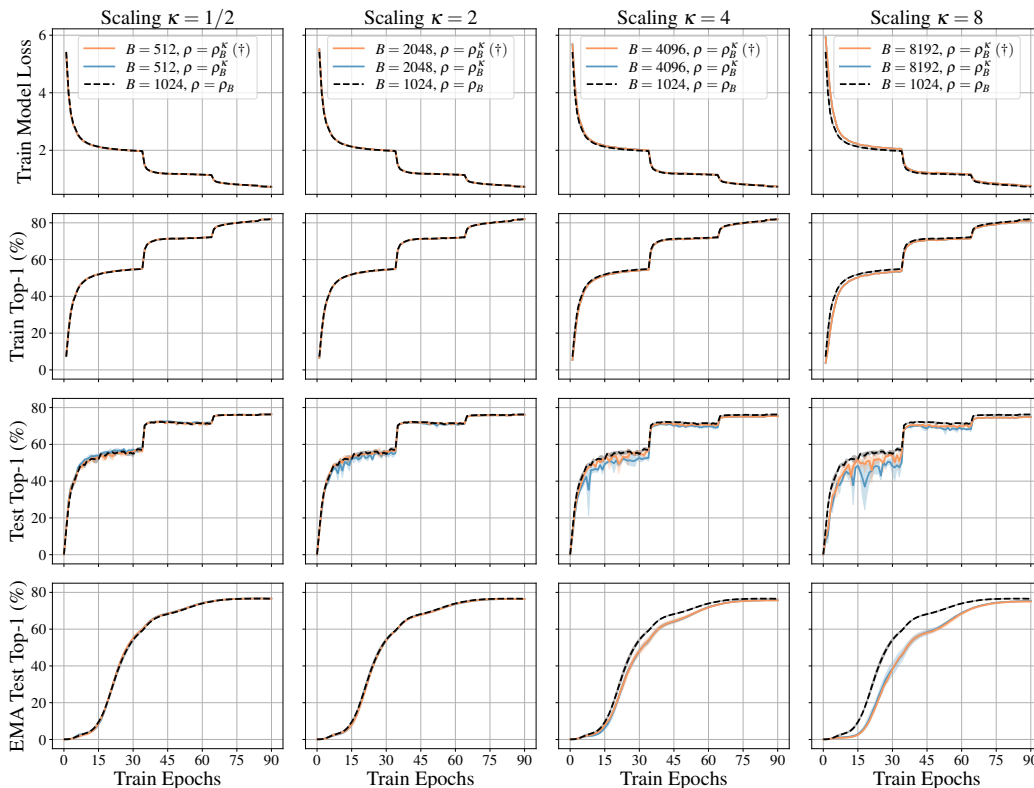


Figure 13: *ResNet2-50 Polyak-Ruppert averaging on ImageNet1k* for different scalings  $\kappa$ . The baseline model ( $\kappa = 1$ , black dashed) uses batch size 1024 and momentum  $\rho_B = 0.9999$ , is scaled down to a batch size of 512 (left), and up to a batch size of 4096 (right) with the EMA Scaling Rule applied to *only* model parameters (blue,  $\rho = \rho_B^\kappa$ ), and model parameters *and* buffers (orange,  $\rho = \rho_B^\kappa (\dagger)$ ). Bands indicate the mean and standard deviation across three runs.

## 978 G Additional details and results for Automatic Speech Recognition (ASR)

979 In this section we provide additional details for the speech recognition experiments in both the  
 980 supervised and semi-supervised case.

981 **Data** We use the LibriSpeech dataset (Panayotov et al., 2015) which is a dataset of audio-  
 982 transcription pairs. For supervised Polyak-Ruppert averaging experiments, we use *train-clean-100*  
 983 as training data, and for semi-supervised pseudo-labeling experiments, we use *train-clean-100* as  
 984 the labeled and *train-clean-360* and *train-other-500* as the unlabeled data. The standard LibriSpeech  
 985 validation sets (*dev-clean* and *dev-other*) are used to tune all hyperparameters, as well as to select  
 986 the best models. Test sets (*test-clean* and *test-other*) are only used for reporting final model per-  
 987 formance, measured in WER without an external language model. We maintain the original 16kHz  
 988 sampling rate, and compute log-mel filterbanks with 80 coefficients for a 25ms sliding window,  
 989 strided by 10ms, later normalized to zero mean and unit variance for each input sequence.

990 **Acoustic model** We employ a vanilla encoder-based only transformer model trained with the Con-  
 991 nectionist Temporal Classification (CTC) loss (Graves et al., 2006). We use the training configura-  
 992 tion from Likhomanenko et al. (2021a), which has three stages: i) 1D convolutions to perform strid-  
 993 ing (kernel of 7 with stride of 3); ii) a Transformer encoder with 36 layers, post-LayerNorm, four  
 994 attention heads, an embedding dimension of 768, an MLP dimension of 3072, a dropout frequency  
 995 of 0.3, and a layer drop frequency of 0.3; and iii) a linear layer to map to the target vocabulary<sup>9</sup>. To  
 996 reduce model training time by a factor of approximately 2 – 3 $\times$ , and to reduce memory footprint,

<sup>9</sup>The token set of this vocabulary consists of the 26 English alphabet letters augmented with the apostrophe and a word boundary token.



Table 5: Hyperparameters summary for speech recognition task for supervised (left) and semi-supervised pseudo-labeling (right) training with a vanilla transformer. The  $0.3 \rightarrow 0.1$  in the dropout and layer drop rates indicates that a rate of 0.3 is used during pre-training, and a rate of 0.1 is used during pseudo-labeling.

	Supervised	Pseudo-Labeling
Librispeech test-clean / test-other WER	7.8/19.1	4.8/11.5
Optimizer	Adam	Adam
Optimizer scaling rule	Adam	Adam
Base $(\beta_1, \beta_2)$	(0.995, 0.999)	(0.995, 0.999)
Base learning rate	0.0001	0.0001
Base learning rate warmup (steps)	64k	64k
Learning rate schedule	Fixed (no decay)	Fixed (no decay)
Learning rate minimum value	0	0
Base training duration (steps)	400k	500k
Base batch size (dynamic)	$8 \times 290s$	$8 \times 290s$
Base teacher momentum	0.99995	0.9999
Weight decay	None	None
Numerical precision	bf16	bf16
Augmentation stack	SpecAug	SpecAug
Dropout	0.3	$0.3 \rightarrow 0.1$
Layer drop	0.3	$0.3 \rightarrow 0.1$
Gradient clipping	1	1
Labeled:unlabeled data ratio	N/A	1:3
Base pre-training steps	N/A	20k
Base start of EMA accumulation (steps)	N/A	19k

997 we use CAPE positional embeddings (Likhomanenko et al., 2021b) instead of relative positional  
 998 embeddings (Shaw et al., 2018): both models perform similarly.

999 **Training** Here we discuss our training procedure for base batch size  $B = 8 \times 290s$ , which is adapted  
 1000 from Likhomanenko et al. (2021a), and is summarized in Table 5. We use SpecAugment (Park et al.,  
 1001 2019) activated after 5k steps of training: two frequency masks with frequency mask parameter  
 1002  $F = 30$ , ten time masks with maximum time-mask ratio  $p = 0.1$  and time mask parameter  $T = 50$  are  
 1003 used; time warping is not used.

1004 One difference in setup is we use the Adam optimizer, whereas Likhomanenko et al. (2021a) used  
 1005 Adagrad (Duchi et al., 2010). Even though Adagrad can be viewed as a particular limit ( $\beta_1 = 0$  and  
 1006  $\beta_2 \rightarrow 1$ ) of Adam (Kingma & Ba, 2015), we were unable to produce reasonable optimization in  
 1007 practice when applying the Adam Scaling Rule of Malladi et al. (2022) in this limit. As a conse-  
 1008 quence, we chose to work with the Adam optimizer, where its scaling rule has been shown to work  
 1009 (Malladi et al., 2022), and we take  $\beta_1 = 0.995$ ,  $\beta_2 = 0.999$ , and  $\epsilon = 10^{-8}$ . We obtained similar results  
 1010 for  $\beta_1 = 0.99$ . Finally, we use a linear learning rate warmup (64k steps) after which the learning rate  
 1011 is kept constant until convergence. This performance can be improved further by using a step decay  
 1012 schedule as shown in prior work. We also apply gradient clipping of 1, and do not use weight decay.

1013 **Pseudo-Labeling** The pseudo-labeling process comprises of two stages: i) The pre-training phase,  
 1014 where we train model on labeled data for 20k steps with model EMA accumulation starting after  
 1015 19k steps; and ii) the pseudo-labeling phase, where we involve unlabeled data by generating pseudo-  
 1016 labels from the model EMA (teacher) and provide them to the model (student) as if they were  
 1017 ground-truth labels. Pseudo-labels are generated without any dropout applied to the teacher, and  
 1018 no data augmentation is applied for the corresponding inputs. To produce the pseudo-label, we use  
 1019 *hard transcription* (Definition G.1)

1020 **Definition G.1** (Hard Transcription). *For a sequence of frames, select the most probable token*  
 1021 *per frame, removing repetitions and the CTC blank token. For example, “h##eelll##ll###oo” is*  
 1022 *transformed into “hello”, where “#” is the CTC blank token.*

1023 These hard transcriptions are then used as transcription for student optimization. We use a 1:3  
 1024 proportion of labeled to unlabeled data as this was found to be optimal in Likhomanenko et al.  
 1025 (2021a), and we decrease model dropout and layer drop rates to 0.1 after pre-training phase. As  
 1026 we have access to the ground-truth labels on the data being treated as unlabeled, we can track

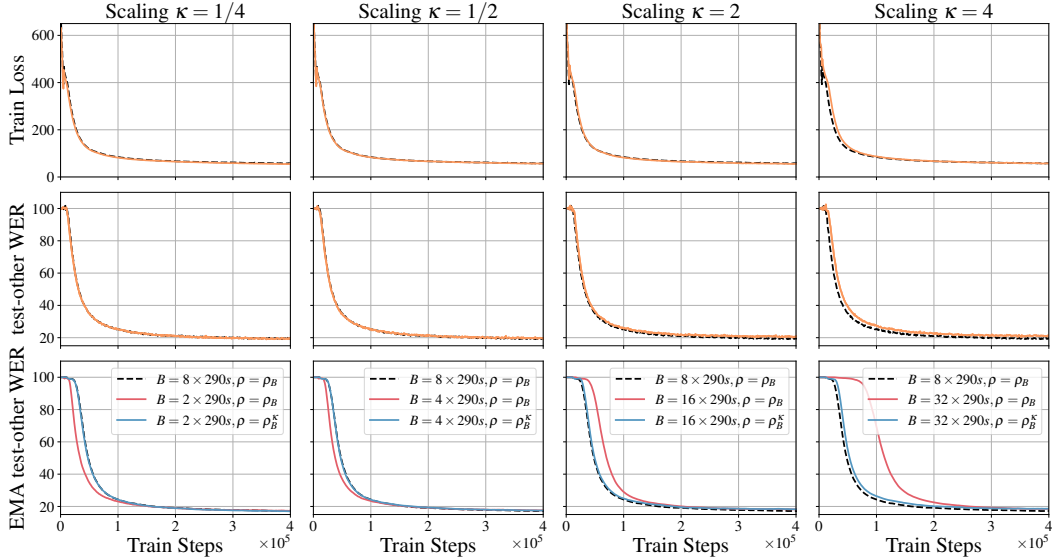


Figure 14: *Transformer Polyak-Ruppert averaging on LibriSpeech (trained on train-clean-100)* with different scalings  $\kappa$ . The baseline ( $\kappa = 1$ , black dashed) is trained with Adam and momentum  $\rho_B = 0.99995$  at a *dynamic batch size*  $B = 8 \times 290s$ , which corresponds to a single train step on the  $x$ -axis. We investigate dynamic batch sizes down to  $B = 2 \times 290s$  (left) and up to  $B = 32 \times 290s$  (right), with (blue,  $\rho = \rho_B^\kappa$ ), and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule (model non-EMA is marked by orange). The Adam Scaling Rule (Malladi et al. (2022), Definition C.3) is used throughout. For momentum  $\rho_B = 0.9999$  we observe similar trajectories behaviour for all models.

1027 pseudo-label quality by computing pseudo-labels on this data, and compute the WER against their  
 1028 ground-truth. Pseudo-label quality is the primary metric to evaluate progress on unlabeled data, as  
 1029 loss on pseudo-labeled data is unreliable when a teacher model and pseudo-labels are evolving with  
 1030 each time step.

1031 **Scaling of batch size** Sequential data is typically processed using dynamic batching as it is more  
 1032 computationally efficient than using a fixed number of sequences (Ott et al., 2019). In our work, we  
 1033 use dynamic batching of  $\sim 290s$  audio per GPU. Moreover, for CTC we do not apply any additional  
 1034 sequence normalization. We experimented with fixed batching, but did not observe any significant  
 1035 differences in conclusions compared with the dynamic batching.

1036 We note that dynamic batching is a more challenging setting for achieving systematic scaling, as the  
 1037 number of independent sequences in any given batch may change, and the i.i.d. assumption does  
 1038 not hold at the frame level. Despite these violations of the assumptions of Section 2.2, our results  
 1039 demonstrate that the Adam Scaling Rule (Definition C.3, Malladi et al. (2022)) holds in the case of  
 1040 dynamic batches, as does our EMA Scaling Rule (Definition 1.1).

1041 The base batch size is set to  $B = 8 \times 290s$ , and in our experiments we scale down to batch size of  
 1042  $B = 2 \times 290s$  and up to batch size of  $B = 128 \times 290s$ . The number of warmup and pre-training  
 1043 steps, steps before SpecAugment is turn on and model EMA is accumulated are scaled according to  
 1044 Appendix C.1.

1045 **Compute** [This section has been redacted to preserve anonymity during the peer-review process.  
 1046 If this work is accepted, the full details compute used for these experiments, including: the experi-  
 1047 ments presented, hyperparameter optimization, and the development process, will be provided.]

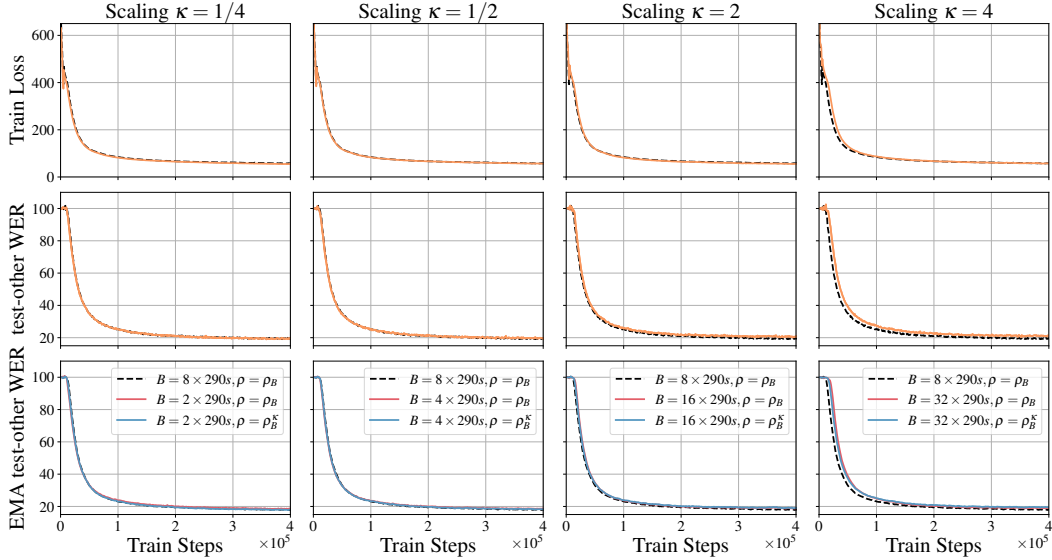


Figure 15: *Transformer Polyak-Ruppert averaging on LibriSpeech (trained on train-clean-100) with different scalings  $\kappa$ . The baseline ( $\kappa = 1$ , black dashed) is trained with Adam and momentum  $\rho_B = 0.999$  at a dynamic batch size  $B = 8 \times 290s$ , which corresponds to a single train step on the  $x$ -axis. We investigate dynamic batch sizes down to  $B = 2 \times 290s$  (left) and up to  $B = 32 \times 290s$  (right), with (blue,  $\rho = \rho_B^\kappa$ ), and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule (model non-EMA is marked by orange). The Adam Scaling Rule (Malladi et al. (2022), Definition C.3) is used throughout. If momentum  $\rho_B$  is small and accumulation history is short we observe no any significant difference between models which all are matching the reference trajectory despite scaling  $\kappa$ .*

## 1048 G.1 Detailed results

1049 We present detailed comparison between models trained with and without EMA Scaling Rule in Fig-  
 1050 ures 14 and 15 for supervised training and in Figures 16 and 17 for semi-supervised training.

1051 First, we observe that if the Adam Scaling Rule does not hold perfectly<sup>10</sup> (there is a mismatch  
 1052 between trajectories for the model before pseudo-labels are involved) the EMA Scaling Rule also  
 1053 gives discrepancies with the reference trajectory, however they are negligible compared to models  
 1054 trained without EMA Scaling Rule. For the semi-supervised training, to alleviate the difficulties with  
 1055 a breakdown of the Adam Scaling Rule for large  $\kappa$  we postpone the pseudo-labeling process until  
 1056 the model reaches similar WER as the baseline. This allows us to align the initial model conditions  
 1057 for pseudo-labeling. In this scenario we are able to match the reference trajectory up to  $\kappa = 8$ .

1058 We note that this result reveals that errors for the Adam Scaling Rule *and* the EMA Scaling Rule  
 1059 are contributing, although the way in which they contribute is different, and one can dominate the  
 1060 other. We observe in Figure 16 that if the initial conditions of the models are similar (attained by  
 1061 using the same WER as a condition to begin pseudo-labeling) then the error from the EMA Scaling  
 1062 Rule dominates over that of the Adam Scaling Rule, causing a divergence in training dynamics.

1063 Second, we observe in practice that the EMA Scaling Rule holds for both fixed batching (a sequence  
 1064 length in the batch can vary significantly) and for dynamic batching (when total number of frames  
 1065 in the batch is fixed, though padding still is accounted to the this amount). This shows that EMA  
 1066 Scaling Rule is applicable to sequential data too.

1067 Third, we observe in Figures 15 and 17 that for smaller values of  $\rho_B$ , scaling with or without  
 1068 EMA Scaling Rule behave similarly, and reference trajectories match in the supervised and semi-  
 1069 supervised cases. However, if the momentum is too large, the *teacher* moves slowly and is uninfor-  
 1070 mative, whereas if the momentum is too low, the *teacher* and the *student* are effectively be the same  
 1071 model, implying: i) the student will self-predict with high confidence, removing any benefits of dis-

<sup>10</sup>See Malladi et al. (2022) for a discussion on scenarios that lead to a breakdown of the Adam Scaling Rule.

1072 tillation<sup>11</sup>; and ii) training instability or model divergence will happen in the low-resource settings  
 1073 (Likhomanenko et al., 2021a; Higuchi et al., 2022).

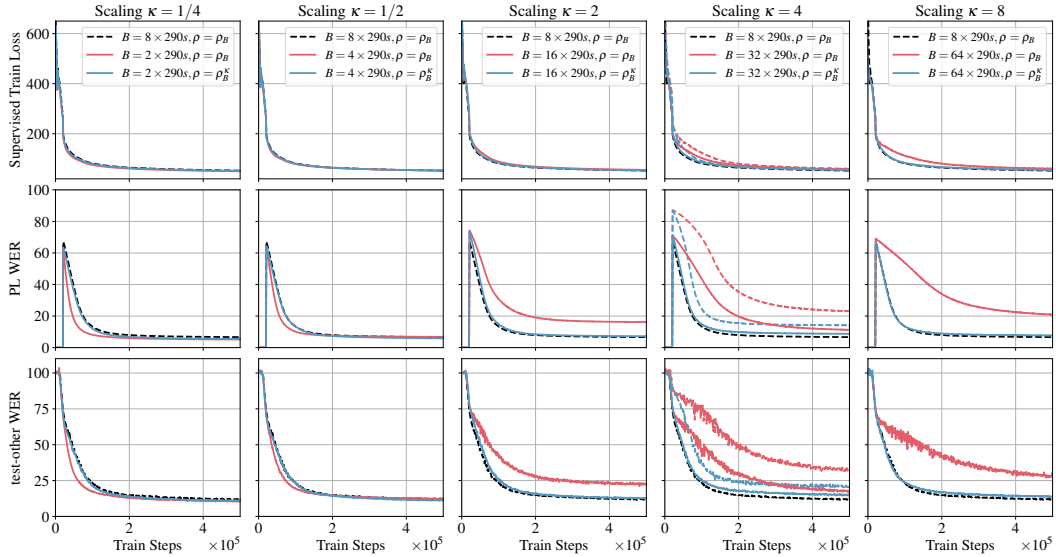


Figure 16: Transformer pseudo-labeling on LibriSpeech (trained on train-clean-100 as labeled and the rest of LibriSpeech as unlabeled) with different scalings  $\kappa$ . The baseline ( $\kappa = 1$ , black dashed) is trained with Adam at a dynamic batch size of  $8 \times 290$  seconds, which corresponds to a single train step on the x-axis. The model EMA (teacher) is updated with momentum  $\rho_B = 0.9999$ . We investigate dynamic batch sizes down to  $B = 2 \times 290s$  (left) and up to  $B = 64 \times 290s$  (right), with (blue,  $\rho = \rho_B^\kappa$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule. The Adam Scaling Rule (Malladi et al. (2022), Definition C.3) is used throughout. For  $\kappa \leq 2$ , we start pseudo-labeling after  $20k/\kappa$  training steps; while for  $\kappa > 2$ , we start when pre-training WER matches the baseline WER ( $24k/\kappa$  for  $\kappa = 4$  and  $29k/\kappa$  for  $\kappa = 8$ ). For  $\kappa = 4$  we experimented with both variants: we start pseudo-labeling after  $20k/\kappa$  (dashed) and when pre-training WER matches the baseline WER (solid,  $24k/\kappa$ ).

<sup>11</sup>He et al. (2020) alleviated the problem with the proper amount of noise during student model training, whilst Xu et al. (2020) used beam-search decoding with a language model.

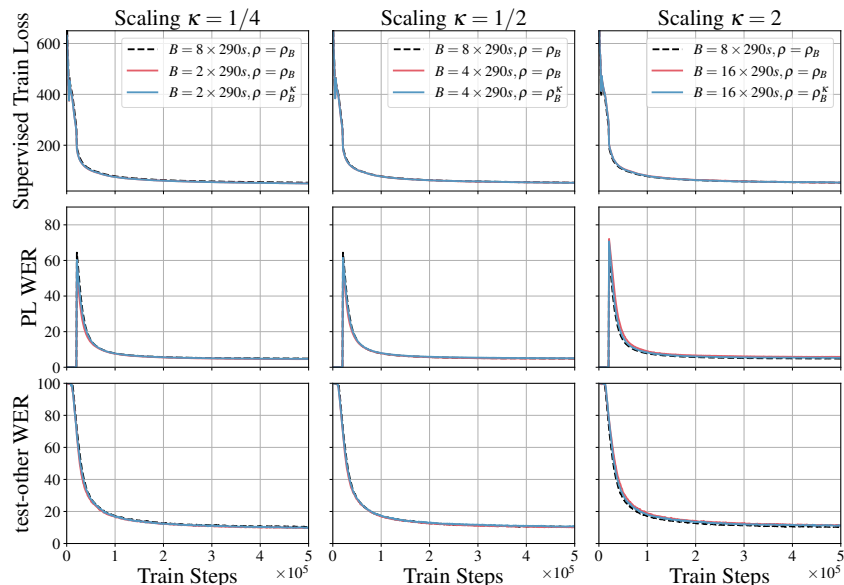


Figure 17: Transformer pseudo-labeling on LibriSpeech (trained on train-clean-100 as labeled and the rest of LibriSpeech as unlabeled) with different scalings  $\kappa$ . The baseline ( $\kappa = 1$ , black dashed) is trained with Adam at a dynamic batch size of  $8 \times 290$  seconds, which corresponds to a single train step on the x-axis. The model EMA (teacher) is updated with momentum  $\rho_B = 0.999$ . We investigate dynamic batch sizes down to  $B = 2 \times 290s$  (left) and up to  $B = 16 \times 290s$  (right), with (blue,  $\rho = \rho_B^k$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule. The Adam Scaling Rule is used throughout. In case of short history accumulation for the momentum (compared to Figure 14) we observe similar to supervised training (Figure 15) no significant different between all models trajectories throughout the training while matching the reference one.

## 1074 G.2 Scaling to $\kappa = 16$ with Progressive Scaling

1075 Finally, we aim to scale semi-supervised pseudo-labeling further to  $\kappa = 16$ . In this case we observe  
 1076 that Adam Scaling Rule does not hold in the pre-training phase and there is no model convergence.  
 1077 To overcome this, we apply Progressive Scaling (Definition 3.2). We pre-train models on supervised  
 1078 data with  $\kappa = 8$  for 29k of reference steps (model EMA accumulation starts at 28k steps). We then  
 1079 scale to  $\kappa = 16$  and begin pseudo-labeling. We see in Figure 18 that Progressive Scaling enables us  
 1080 to scale pseudo-labeling to  $\kappa = 16$  with (middle) and without (left) the EMA Scaling Rule. Second,  
 1081 models *with* the EMA Scaling Rule track the baseline much closer than models without the EMA  
 1082 Scaling Rule, although a small gap is present. We further experimented with Progressive Scaling,  
 1083 postponed the transition condition to the  $\kappa = 16$  until 75k reference steps. In Figure 18 (right), we  
 1084 see this scaled model tracks the reference trajectory, and so using a combination of the EMA Scaling  
 1085 Rule and Progressive Scaling, we are able to scale pseudo-labeling to  $\kappa = 16$ , corresponding to a  
 1086 dynamic batch size of  $128 \times 290s$ .

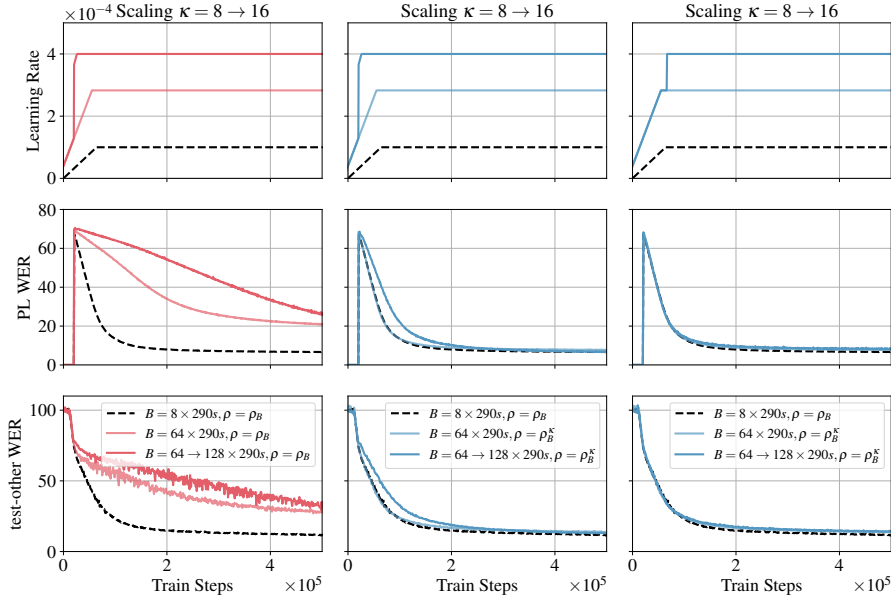


Figure 18: Transformer pseudo-labeling on LibriSpeech (trained on train-clean-100 as labeled and the rest of LibriSpeech as unlabeled) with different Progressive Scaling from  $\kappa = 8$  to  $\kappa = 16$  ( $\kappa = 8 \rightarrow 16$ ). The baseline ( $\kappa = 1$ , black dashed) is trained with Adam at a *dynamic batch size* of  $8 \times 290$  seconds, which corresponds to a single train step on the  $x$ -axis. The model EMA (*teacher*) is updated with momentum  $\rho_B = 0.9999$ . The scaling with  $\kappa = 8$  is shown with lighter color for reference from Figure 16. We investigate dynamic batch sizes progressively from  $B = 64 \times 290s$  to  $B = 128 \times 290s$ , with (blue,  $\rho = \rho_B^\kappa$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule. For reference (top) we show the learning rate schedule with Progressive Scaling. The Adam Scaling Rule (Malladi et al. (2022), Definition C.3) is used throughout. Left and middle correspond to Progressive Scaling with scale from  $\kappa = 8$  to  $\kappa = 16$  at 29k steps, while right corresponds to 75k steps.

## 1087 H Additional details and results for self-supervised image representation 1088 learning

1089 **Organization** This appendix is structured into three sections. We first give an overview of our  
1090 chosen SSL method BYOL (Appendix H.1), our recipe for training BYOL using Vision Transfor-  
1091 mers (ViTs) (Appendix H.2), ablations of normalization approaches that lead to the development of  
1092 this recipe (Appendix H.3), and additional results corresponding to longer training duration (Ap-  
1093 pendix H.4) and further understanding the impact of Progressive Scaling (Appendix H.5).

1094 Second, we demonstrate that the EMA Scaling Rule combined with Progressive Scaling can scale  
1095 a ResNet-50 BYOL model trained with LARS to batch size 32,768 without performance drop,  
1096 demonstrating the empirical utility of the tools we provide outside of their theoretical validity (Ap-  
1097 pendix H.6).

1098 Finally, we show that it is possible to systematically scale DINO (Caron et al., 2021) using a com-  
1099 bination of Progressive Scaling and the EMA Scaling Rule, providing a solution for researchers and  
1100 practitioners wanting to train DINO at scale.

### 1101 H.1 Components of self-supervised learning

1102 First, a key component of many SSL methods is the *stop-gradient* or StopGrad (Definition H.1).

1103 **Definition H.1** (Stop Gradient/StopGrad( $\cdot$ )). *The stop-gradient operator StopGrad( $\cdot$ ) prevents the*  
1104 *flow of gradient information*

$$\frac{df(\text{StopGrad}(h(x; \omega)); \theta)}{d\omega} \equiv 0 \quad (72)$$

1105 *for all parameteric functions  $h$  and  $f$  and for all parameters  $\theta$  and  $\omega$ .*

1106 Applying a *stop-gradient* is sometimes called *detaching* in the literature. Now, we introduce the  
1107 update rule of our representative SSL method BYOL in Definition H.2.

1108 **Definition H.2** (BYOL Update). *BYOL learns unsupervised features by minimizing the cosine dis-*  
 1109 *tance between the predictions of a student backbone  $f(\cdot; \theta)$  (typically a ResNet or Vision Trans-*  
 1110 *former), projected through  $h(\cdot; \omega)$  (typically a Multi-Layer Perceptron (MLP)), and the predictions*  
 1111 *of an EMA teacher  $f(\cdot; \zeta)$  (Grill et al., 2020). The update for the parameters of BYOL is then*

$$(\theta_{t+1}, \omega_{t+1}) = (\theta_t, \omega_t) - \eta \times \frac{1}{B} \sum_{x \in \mathbb{B}} \nabla_{(\theta, \omega)} \mathcal{L}(x; \theta_t, \omega_t, \zeta_t) \quad (73)$$

$$\zeta_{t+1} = \rho \zeta_t + (1 - \rho) \theta_{t+1} \quad (74)$$

$$\text{with } \mathcal{L}(x; \theta_t, \omega_t, \zeta_t) = \frac{1}{2} \cos [h(f(x_1; \theta_t); \omega_t), \text{StopGrad}(f(x_2; \zeta_t))] + (x_1 \leftrightarrow x_2), \quad (75)$$

1112 where  $\cos(\mathbf{a}, \mathbf{b}) \equiv 1 - \mathbf{a} \cdot \mathbf{b} / (|\mathbf{a}| |\mathbf{b}|)$  is the cosine distance, and  $x_1$  and  $x_2$  are two views of a single  
 1113 variate  $x$ , often produced by augmentations, and  $x_1 \leftrightarrow x_2$  denotes symmetrization over  $x_1$  and  $x_2$ .

1114 As noted in Section 3.4, the BYOL EMA update (Equation (74)) uses  $\theta_{t+1}$  instead of our analyzed  
 1115  $\theta_t$  (Equation (4)). The effect upon the overall EMA update is  $\mathcal{O}(\eta \times \beta_\rho)$  and so is captured by the  
 1116 EMA Scaling Rule (Definition 1.1).

1117 One more piece of technology typically employed in SSL is a *tracking probe* (Definition H.3) which  
 1118 we will use to evaluate the performance of BYOL on downstream tasks of interest, for example,  
 1119 image classification.

1120 **Definition H.3** (Tracking Probe/Linear Probe). *When optimizing model parameters  $\omega_t$  of an SSL*  
 1121 *method, simultaneously optimize the parameters  $\xi$  of a probe model  $r(\cdot; \xi)$  under a downstream*  
 1122 *objective  $\mathcal{L}^{(d)}$ . For example, in classification, with data  $x$  and samples  $y$*

$$\mathcal{L}^{(d)}(x, y, \theta_t, \xi_t) = -\log P(y|r(\text{StopGrad}(h(x; \omega_t)); \xi)) \quad (76)$$

$$\mathcal{L}^{(total)}(x, y; \theta_t, \omega_t, \zeta_t, \xi_t) = \mathcal{L}(x; \theta_t, \omega_t, \zeta_t) + \mathcal{L}^{(d)}(x, y, \omega_t, \xi_t), \quad (77)$$

1123 *This is a probe for the teacher, which is typically the better choice due to Polyak-Ruppert averaging*  
 1124 *effects (see Section 3.2). When the  $r$  is a linear model, the tracking probe is called a linear probe.*

1125 It is also typical to use a Batch Normalization layer *without* trainable affine terms before this linear  
 1126 layer as in He et al. (2022) to stabilize probe training. In this case, the running statistics can be  
 1127 absorbed into a definition of the linear layer weights and biases, and so this is still a *linear probe*,  
 1128 although we will call this a *pre-bn linear probe* to remove ambiguity.

## 1129 H.2 A Vision Transformer recipe for BYOL

1130 **Hyperparameters** We present the base hyperparameters for training BYOL with a ViT-B/16 back-  
 1131 bone in Table 6. This recipe was developed by starting from a well-known supervised ViT-B/16  
 1132 recipe (He et al., 2022) and performing a search over weight decay and learning rate hyperparame-  
 1133 ter choices. We find that BYOL performs well with heavy weight decay ( $\lambda = 0.3$ ) and a low learning  
 1134 rate ( $\eta = 10^{-3}$ ) at a base batch size  $B = 4096$ . The AdamW optimizer is used, and so for scaling to  
 1135 other batch sizes  $\hat{B} = \kappa B$  we use the Adam Scaling Rule (Definition C.3)<sup>12</sup> We use a pre-bn linear  
 1136 probe as discussed in Appendix H.1. Finally, the performance of BYOL can be further improved  
 1137 by employing multicrop (Caron et al., 2020) by  $\approx +2\%$  in absolute test top-1 performance on Im-  
 1138 ageNet1k compared to without multicrop, however, as this is not our focus, we omit this from the  
 1139 presented recipe.

1140 **Compute** [This section has been redacted to preserve anonymity during the peer-review process.  
 1141 If this work is accepted, the full details compute used for these experiments, including: the experi-  
 1142 ments presented, hyperparameter optimization, and the development process, will be provided.]

1143 **Additional background** Achieving large scale SSL training with ViTs to large scale SSL train-  
 1144 ing has been a long standing goal in the community. MoCo-v3 (Chen et al., 2021) enables the  
 1145 use of ViTs with contrastive learning, but achieves this through modificatinos of the ViT training

<sup>12</sup>We note that Adam (Kingma & Ba, 2015) and AdamW (Loshchilov & Hutter, 2019) are equivalent in the limit of zero weight decay, and that the Adam Scaling Rule (Definition C.3) was derived with zero weight decay (Malladi et al., 2022).

Table 6: BYOL ViT-B/16 hyperparameters.

BYOL ViT-B/16	
ImageNet1k Linear Probe Test Top-1	74.47% (Figure 19)
Weight initialization	<code>trunc_normal(.02)</code>
Backbone normalization	LayerNorm
Head normalization	BatchNorm
Synchronized BatchNorm over replicas	No
Learning rate schedule	Single Cycle Cosine
Learning rate warmup (epochs)	40
Learning rate minimum value	$1 \times 10^{-6}$
Training duration (epochs)	480
Optimizer	AdamW
Optimizer scaling rule	Adam
Base $(\beta_1, \beta_2)$	(0.9, 0.95)
Base learning rate	$1 \times 10^{-3}$
Base batch size	4096
Base teacher momentum	0.99
Weight decay	0.3
Weight decay scaling rule	None
Weight decay skip bias	Yes
Numerical precision	bf16
Augmentation stack	BYOL
Stochastic depth	0.1

1146 procedures, including gradient freezing on the image patching layer, and re-introducing Batch Normalization to post-attention MLP layers. Despite these modifications, MoCo-v3 was only trained up  
 1147 to a batch size of 6144, where model performance begins to suffer (Chen et al., 2021). In Figure 6  
 1148 we demonstrate that combining dynamic batch scaling (Appendix C.4) with the EMA Scaling Rule  
 1149 (Definition 1.1) enables BYOL to be trained using ViTs to batch sizes of 24,576 without any drop  
 1150 in performance compared to the reference batch size of 4096. We emphasize that the piecewise  
 1151 transitions in the schedules are important for preserving training dynamics.  
 1152

### 1153 H.3 The role of Batch Normalization and Layer Normalization in BYOL with ViTs

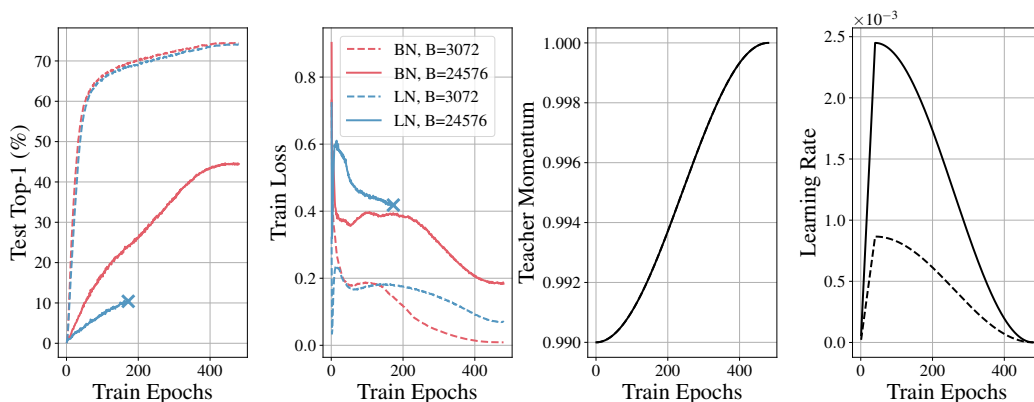


Figure 19: *BYOL ViT-B/16 on ImageNet1k* for different scalings  $\kappa$ . We present runs comparing LayerNorm (blue) to BatchNorm (red) in the projection and prediction heads of BYOL ViT models for batch size 3072 (dashed) and 24,576 (solid) *without the EMA Scaling Rule*.  $\kappa = 1$  corresponds to  $B = 4096$ . In all scenarios the transformer backbone *only* uses LayerNorm. We truncate the training of the large batch size LayerNorm variant to preserve compute (indicated by  $\times$ ).

1154 Here we compare the roles of Batch Normalization (BatchNorm, Ioffe & Szegedy (2015)) and Layer  
 1155 Normalization (LayerNorm, Ba et al. (2016)) in the projection and prediction heads of BYOL (Grill  
 1156 et al., 2020) using ViTs.



1157 It has been observed that BatchNorm plays a critical role in BYOL predictor and projector dynam-  
 1158 ics (Fetterman & Albrecht, 2020), and using either LayerNorm or *no normalization* significantly  
 1159 decrease in model performance. Subsequently, it was demonstrated (Richemond et al., 2020) that  
 1160 competitive BYOL performance could be achieved through a combination of Group Normaliza-  
 1161 tion (GroupNorm, Wu & He (2018)) and Weight Standardization (Qiao et al., 2019). Additionally,  
 1162 Richemond et al. (2020) showed that if BatchNorm is used in the backbone, one can use LayerNorm  
 1163 or *no normalization* in the predictor and projector without any performance drop.

1164 In this work, we show it is possible to train BYOL ViT using *only LayerNorm* across the back-  
 1165 bone, projector and predictor (see Figure 19), decoupling BYOL’s reliance on batch statistics, a  
 1166 desirable trait for a representation learning algorithm (Brock et al., 2021). At batch size 3072, using  
 1167 LayerNorm in the predictor and projector achieves competitive performance (74.10%), performing  
 1168 slightly worse than using BatchNorm (74.47%). At the larger batch size of 24,576, runs perform  
 1169 significantly worse as the EMA Scaling Rule was not applied.

#### 1170 H.4 Longer training duration with incremental Progressive Scaling

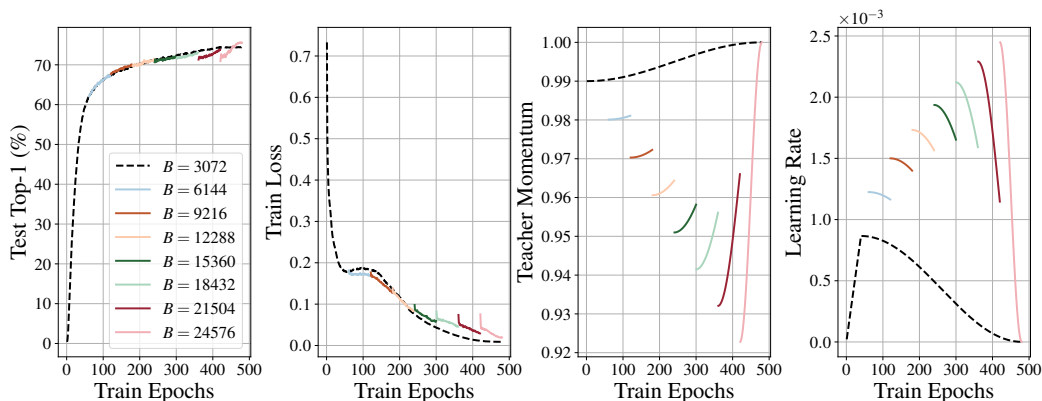


Figure 20: BYOL ViT-B/16 on ImageNet1k for different scalings  $\kappa$ . The baseline model ( $\kappa = 0.75$ , black dashed) uses batch size 3072 and teacher momentum  $\rho_B = 0.99$ . We increment the batch size by 3072 every 60 epochs to a final batch size of 24,576 using Progressive Scaling (Definition 3.2).

1171 Here we use the same base hyperparameters as Table 6, except that we train for 480 instead of 300  
 1172 epochs. To mitigate the student impulse phenomena discussed in Section 3.4, in Figure 20 we in-  
 1173 vestigate increasing the batch size every 60 epochs using Progressive Scaling (Definition 3.2). We  
 1174 observe that this more gradual procedure enables closer tracking of the baseline train loss trajec-  
 1175 tory. Additionally, this procedure results in a scaled linear probe performance that outperforms the  
 1176 baseline (75.64% compared to the baseline performance of 74.47%). The same procedure can be ap-  
 1177 plied to the LayerNorm variant discussed in Appendix H.3, which produces a similar result (75.09%  
 1178 compared to the baseline performance of 74.10%).

#### 1179 H.5 Building intuition around Progressive Scaling and momentum sensitivity

1180 Our final BYOL ViT results are to help build intuition around Progressive Scaling (Definition 3.2),  
 1181 as well as when the EMA Scaling Rule is most important. In Figure 21 we explore transition-  
 1182 ing from the baseline batch size 4096 model to batch size 24,576 in a *single transition* after  
 1183 60 epochs. After this transition, we continue training for 240 epochs for a range of momenta:  
 1184  $\rho \in \{0.8, 0.9, 0.95, 0.97, 0.9867, 0.994, 0.999\}$  *without* the EMA Scaling Rule.

1185 We observe that after the transition, any  $0.9 \leq \rho \leq 0.994$  produces a linear probe performance that  
 1186 matches or outperforms the baseline at the end of training. This indicates that after the initial training  
 1187 period, BYOL becomes less sensitive to the choice of teacher momentum. Note that without the  
 1188 initial 60 epochs of training with batch size 4096, *all models*, including those employing the EMA  
 1189 Scaling Rule diverge (see  $B = 24, 576$  in Figure 6).

1190 We present an illustration for why this might happen in Figure 22. First, we see that using the EMA  
 1191 Scaling Rule *always* keeps the model within the acceptable momentum region. We also see that

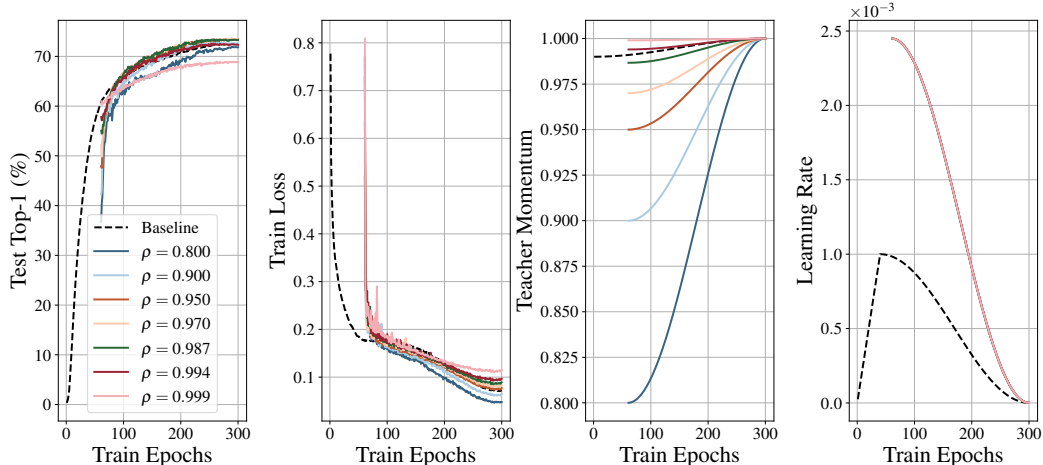


Figure 21: *BYOL ViT-B/16 on ImageNet1k* for different momenta  $\rho$ . The baseline model ( $\rho = 0.99$ , black dashed) uses batch size 4096. At the 60th epoch we apply Progressive Scaling (Definition 3.2) and transition to batch size 24576. We train for a further 240 epochs without EMA scaling for a range of momenta:  $\rho \in \{0.9, 0.95, 0.97, 0.9867, 0.994\}$ .

1192 *not* using the EMA Scaling Rule can keep the model within the acceptable momentum region for a  
 1193 range of batch sizes, depending on how large wide in momenta the acceptable region is at the base  
 1194 batch size. Finally, we see that the momentum value matters much more at low values of momenta  
 1195 (the acceptable momentum region shrinks), whereas at large momenta, this region of acceptability  
 1196 widens.

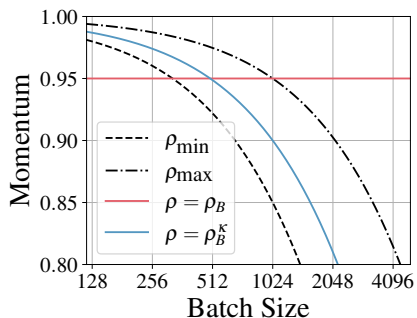


Figure 22: A hypothetical scenario where there is an upper and lower limit for momenta qualitatively leading to the same result.. We assume at base batch size  $B = 1024$  there is an upper ( $\rho_{\max}$ , black dashdot) and lower ( $\rho_{\min}$ , black dashed) limit for valid momenta. We show what happens if we start with  $\rho_B = 0.95$  at a batch size of 4096, and scale with ( $\rho = \rho_B^\kappa$ , blue) and without ( $\rho = \rho_B$ , red) the EMA Scaling Rule.

## 1197 H.6 Scaling a ResNet-50 BYOL using LARS and Progressive Scaling

1198 Here we investigate whether Progressive Scaling and the EMA Scaling Rule can be used in practice  
 1199 where there is no known optimizer SDE approximation. We use the default 300 epoch configuration  
 1200 for BYOL (Grill et al., 2020) in Figure 23. We see that although trajectories during training do not  
 1201 match, we are able to match or surpass the linear probe performance of the BYOL baseline at the  
 1202 larger batch size if 32,768. This indicates that the contributions of our work have practical utility  
 1203 beyond the theoretical constraints.

1204 **Compute** [This section has been redacted to preserve anonymity during the peer-review process.  
 1205 If this work is accepted, the full details compute used for these experiments, including: the experi-  
 1206 ments presented, hyperparameter optimization, and the development process, will be provided.]

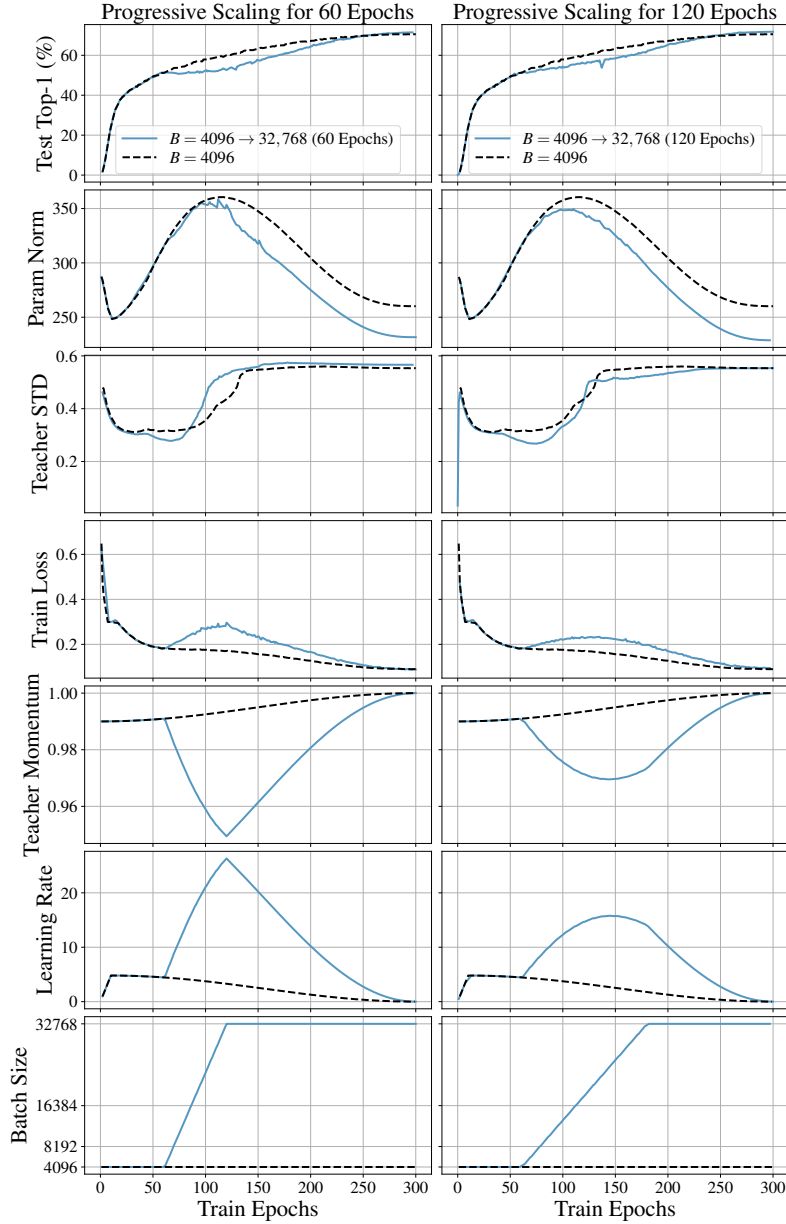


Figure 23: *ResNet50 BYOL on ImageNet1k using LARS* for different configurations of progressive scaling. The baseline (black dashed) uses batch size 4096 and momentum  $\rho_B = 0.99$ . We consider progressive scaling (blue) smoothly from epoch 60 for 60 epochs (left) and 120 epochs (right) up until batch size 32,768, scaling the learning rate linearly, and applying the EMA Scaling Rule.

## 1207 H.7 Preventing collapse phenomena in DINO at scale

1208 Until now, our representative SSL method has been BYOL for reasons discussed in Section 3.4.  
 1209 Here, we will turn our attention to Distillation with NO labels (DINO) (Caron et al., 2021), which  
 1210 has the update rule presented in Definition H.4.

1211 **Definition H.4** (DINO Update). *DINO learns unsupervised features by matching predictions over*  
 1212 *emergent pseudo-labels of a student backbone and head  $f(\cdot; \theta)$  to those of an EMA teacher  $f(\cdot; \zeta)$*   
 1213 *through a cross-entropy guided distillation procedure. DINO has an additional centering procedure,*  
 1214 *which is a form of batch normalization with momentum  $\rho_c = 0.9$  which we do not scale using the*

Table 7: DINO ViT-B/16 Training hyperparameters.

	DINO ViT-B/16
CIFAR10 Linear Probe Top-1 ( $\rho_B = 0.996$ )	85.38%
CIFAR10 Linear Probe Top-1 ( $\rho_B = 0.992$ )	86.96%
Weight initialization	<code>trunc_normal (.02)</code>
Normalization	Layer Norm
Learning rate schedule	Single Cycle Cosine
Learning rate warmup (epochs)	50
Learning rate minimum value	$1 \times 10^{-6}$
Training duration (epochs)	280
Optimizer	AdamW
Optimizer scaling rule	Adam
Base ( $\beta_1, \beta_2$ )	(0.9, 0.95)
Base learning rate	$3 \times 10^{-4}$
Base batch size ( $B$ )	1024
Base teacher momentum ( $\rho_B$ )	0.992 or 0.996
Base weight decay	0.04
Weight decay scaling rule	Linear
Weight decay skip bias	Yes
Center Momentum	0.9
Center Momentum Scaling Rule	None
Precision	bf16
Augmentation stack	DINO Multi-crop

1215 *EMA Scaling Rule. The update for the parameters of DINO is*

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \times \frac{1}{B} \sum_{x \in \mathbb{B}} \nabla_{\boldsymbol{\theta}} \mathcal{L}(x; \boldsymbol{\theta}_t, \boldsymbol{\zeta}_t, \mathbf{c}_t) \quad (78)$$

$$\boldsymbol{\zeta}_{t+1} = \rho \boldsymbol{\zeta}_t + (1 - \rho) \boldsymbol{\theta}_{t+1} \quad (79)$$

$$\mathbf{c}_{t+1} = \rho_c \mathbf{c}_t + (1 - \rho_c) \mathbb{E}_{x'} \boldsymbol{\zeta}(x') \quad (80)$$

$$\text{with } \mathcal{L}(x; \boldsymbol{\theta}_t, \boldsymbol{\zeta}_t, \mathbf{c}_t) = H(f(x_1, \boldsymbol{\theta}_t), f(x_2, \boldsymbol{\zeta}_t) - \mathbf{c}_t) + (x_1 \leftrightarrow x_2), \quad (81)$$

1216 where  $H(\mathbf{a}, \mathbf{b}) \equiv -\sum_{m=1}^M p_m(\mathbf{a}) \log p_m(\mathbf{b})$  is the cross-entropy between categorical distributions  
 1217 over  $M$  (emergent pseudo-)classes given logits  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^M$ ,  $x_1$  and  $x_2$  are two views of a single  
 1218 variate  $x$ , often produced by augmentations, and  $x_1 \leftrightarrow x_2$  denotes symmetrization over  $x_1$  and  $x_2$ .

1219 In practice, DINO employs multi-crop (Caron et al., 2021). We omit this detail for clarity of presen-  
 1220 tation, although we *do* use multi-crop in the experiments that follow.

1221 Our interest DINO is due to the difficulty in its optimization<sup>13</sup>, and in particular, preventing collapse  
 1222 phenomena in DINO at batch sizes above 1024, which is an open research problem. In this section,  
 1223 we will show that a combination of the EMA Scaling Rule (Definition 1.1) and Progressive Scaling  
 1224 (Definition 3.2) enable training of DINO beyond batch size 1024 without sacrificing performance.

1225 **Hyperparameters** Base hyperparameters are presented in Table 7.

1226 **Compute** [This section has been redacted to preserve anonymity during the peer-review process.  
 1227 If this work is accepted, the full details compute used for these experiments, including: the experi-  
 1228 ments presented, hyperparameter optimization, and the development process, will be provided.]

1229 **Results** In Figures 24 and 25 we show the results obtained training DINO on CIFAR-10 with  
 1230  $\rho_B = 0.996$  and  $\rho_B = 0.992$  respectively at the reference batch size of 1024. We employ smooth  
 1231 Progressive Scaling (Definition 3.2) between epochs 120 and 180.

1232 At batch size 2048, the training loss matches the reference *only* when the EMA Scaling Rule is  
 1233 applied, whereas the run *without* the scaling rule diverges from the reference. The impact of this

<sup>13</sup>For an example, see <https://github.com/facebookresearch/dino/issues/43#issuecomment-881453515>.

1234 divergence is emphasized as we consider the larger batch size of 4096. Here, there is also a gap *with*  
 1235 the EMA Scaling Rule, however is approximately three times smaller than the gap *without* the EMA  
 1236 Scaling Rule.

1237 Additionally, we observe that using  $\rho_B = 0.992$  yields higher Top-1 accuracy over  $\rho_B = 0.996$ , and  
 1238 in our experiments, using the EMA Scaling Rule *always* performs better in terms of linear probe  
 1239 performance than not using the scaling rule.

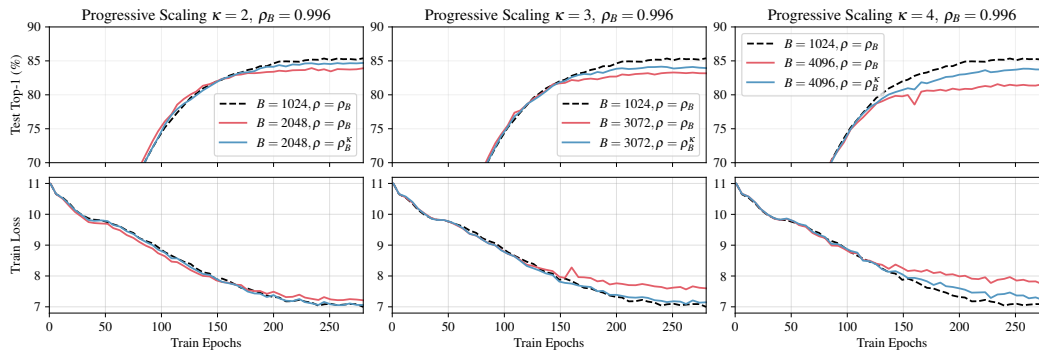


Figure 24: *DINO ViT-B/16 on CIFAR-10* for different scalings  $\kappa$  and base teacher momentum  $\rho_B = 0.996$ . The baseline model ( $\kappa = 1$ , black dashed) uses batch size 1024 and center momentum  $\rho_c = 0.9$ , and is scaled up from batch size 2048 (left) to 4096 (right) with (blue,  $\rho = \rho_B^\kappa$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule. Between epochs 100 and 180 we scale the batch size using progressive scaling (Definition 3.2).

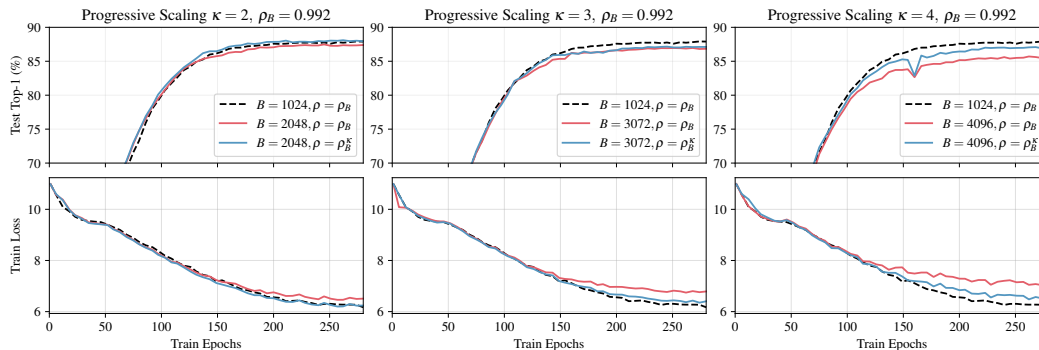


Figure 25: *DINO ViT-B/16 on CIFAR-10* for different scalings  $\kappa$  and base teacher momentum  $\rho_B = 0.992$ . The baseline model ( $\kappa = 1$ , black dashed) uses batch size 1024 and center momentum  $\rho_c = 0.9$ , and is scaled up from batch size 2048 (left) to 4096 (right) with (blue,  $\rho = \rho_B^\kappa$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule. Between epochs 100 and 180 we scale the batch size using progressive scaling (Definition 3.2).

1240 In Figure 26 we show how the hyperparameters  $\rho$ ,  $B$  and learning rate change with the progressive  
 1241 scaling in Definition 3.2.

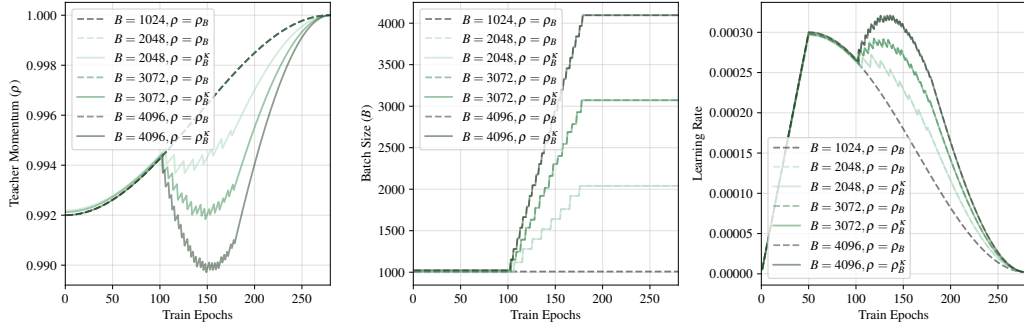


Figure 26: *DINO ViT-B/16 on CIFAR-10* for different scalings  $\kappa$  and base teacher momentum  $\rho_B = 0.992$ . We show how the hyperparameters  $\rho$ ,  $B$  and learning rate change with the Progressive Scaling in Definition 3.2. These hyperparameters correspond to the training runs in Figure 25. Those for Figure 24 are identical, with the exception of  $\rho$  that starts at 0.996 instead of 0.992.

1242 We also attempted to use a sharp batch size transition (Figures 27 and 28), which leads to the  
 1243 collapse phenomena observed in prior work. This collapse happens with and without the EMA  
 1244 Scaling Rule. We suspect this is due to dynamics specific to DINO’s early phase that are even more  
 1245 challenging to replicate under discretization than those of BYOL.

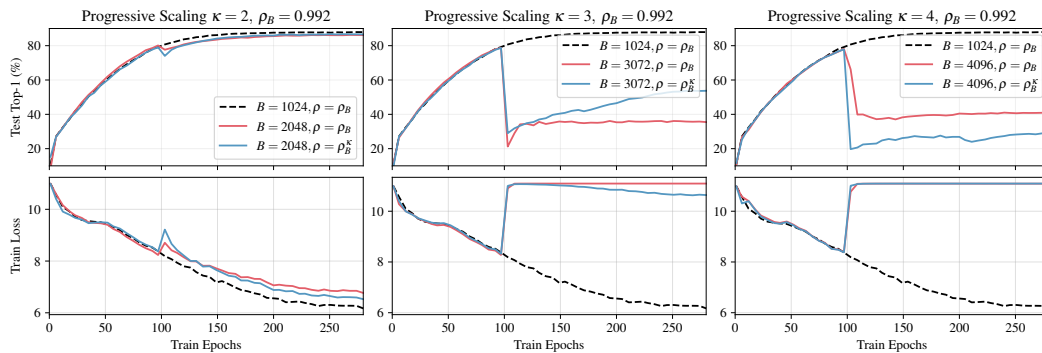


Figure 27: *DINO ViT-B/16 on CIFAR-10* for different scalings  $\kappa$  and base teacher momentum  $\rho_B = 0.992$ . The baseline model ( $\kappa = 1$ , black dashed) uses batch size 1024 and center momentum  $\rho_c = 0.9$ , and is scaled up from batch size 2048 (left) to 4096 (right) with (blue,  $\rho = \rho_B^\kappa$ ) and without (red,  $\rho = \rho_B$ ) the EMA Scaling Rule. Progressive Scaling is employed with a sharp transition at epoch 100, leading to a collapse phenomenon.

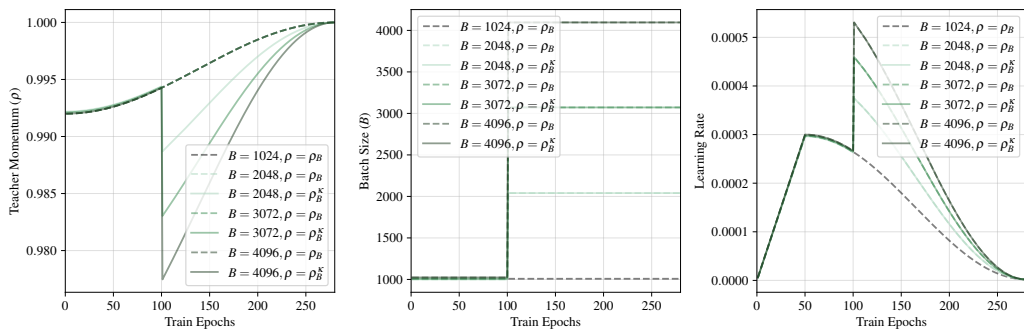


Figure 28: *DINO ViT-B/16 on CIFAR-10* with  $\rho_B = 0.992$  and a sharp transition in batch size at epoch 100. We show how the hyperparameters  $\rho$ ,  $B$  and learning rate change with sudden scaling. These hyperparameters correspond to the training runs in Figure 27.

1246 Our results in this section show it is possible to scale DINO to large batch sizes *without* sacrificing  
 1247 performance by using *both* the EMA Scaling Rule and Progressive Scaling, providing the batch size  
 1248 schedule of Progressive Scaling is not sudden. This resolves an open problem in SSL research.