

1 We thank the reviewers for their feedback. We reply to the main points below:

2 **Usefulness of the method (R1, R3) and applicability of assumptions (R4):** Actually, this work started precisely to
3 unblock the adoption of a real-world Computer Vision AutoML system: Users fine-tune models selected from a large
4 model zoo testing hundreds of combinations of different architectures, pre-training sets and hyper-parameters, but are
5 reluctant to do so without visibility of the expected ROM cost of the training. We focus on fine-tuning (R3) since it is
6 faster and typically performs better than training from scratch in vision problems. For these reasons, it is generally the
7 choice in AutoML systems, where users pay by the hour. While our work may not directly impact academic researchers,
8 it is an enabler of large-scale AutoML, which we expect will further foster academic research in the years ahead. In this
9 sense, we would say our work is useful: it enables a cost estimate that allows reducing a large search space to fit a user's
10 budget, a small step toward better accessibility, and democratization of ML. We are happy to expand the discussion in
11 this direction if it is of interest to the readers.

12 Note the technical hypotheses that the weights remain close to initialization (R4) is only used to derive the approximation,
13 which is then verified empirically to hold (albeit with a larger error margin) even when the task is quite different from
14 pre-training (e.g., fine-grained airplane classification using ImageNet pre-training, see also Fig. 10 and Section B). In
15 real-world AutoML this concern is largely moot since model selection techniques selects for further fine-tuning only
16 models pre-trained on close data.

17 **Make contributions over related literature more clear (R2):** Our main contribution is to introduce the problem of
18 predicting training time in realistic use cases (see previous point), in particular how this depends on the hyper-parameters
19 (for which some previous literature exists) and on the interaction between target task and pre-training (which, to the best
20 of our knowledge, is new). *NTK theory* [15] studies randomly initialized DNNs in the limit of infinite width and batch
21 size. We replace it with an off-the-shelf pre-trained network fine-tuned with SGD and show that its predictions can be
22 effective on practical networks. We also show how to approximate the NTK matrix efficiently while maintaining a good
23 accuracy for our end-tasks. *Gradients as feature* [25] uses a linearization-based analysis of fine-tuning which is similar
24 to ours, but it focuses on efficiently using the gradients of the network as features for a linear classifier, and is unrelated
25 to our work in both scope and methods. *Mean field approximation* [13] focuses on describing a special initialization
26 technique. However, we build on their SDE approximation in weight space and show that it can be translated to function
27 space, making a numerical solution possible for real sized DNNs.

28 **Connection to learning curve prediction methods (R3):** We thank the reviewer for pointing out this relevant area of
29 research, which we will discuss in the paper. We should note that these papers focus on predicting the effect of different
30 hyper-parameters for fixed task and architecture. However, we found the relation between target task and pre-training is
31 more difficult to model. We initially employed a similar black-box regressions techniques to predict training time, but
32 they turned out to be data-hungry and less likely to generalize on different models. This prompted us to develop a more
33 interpretable and less data expensive approach presented here.

34 **Experiment with different definition of training time (R3), different architectures (R2):** We further tested our
35 method with the suggested definition (first time training error decreases by less than ϵ in one epoch), and also increased
36 the number of iterations to 600 to ensure convergence and top accuracy. In this (more challenging) case, we achieve
37 13% avg. relative prediction error and 0.94 Pearson's correlation between predicted time and ground-truth. The method
38 also generalizes to several deeper architectures (ResNet-50, Densenet-121). However, we note that predictions for older
39 architectures (AlexNet) is less accurate (40% avg. error, 0.46 Pearson's correlation). This is expected as their loss
40 landscape is known to be very rough, so gradients at initialization are less informative.

41 **Clarifications:** We have amended the text to make all requested clarifications. More in detail:

42 • **Table 1 (R2, R4):** We report the absolute error on training-time prediction as a function of the selected threshold.
43 We show the error both using a small learning rate (our approximation holds better) and high learning rate. Different
44 thresholds reflect training time predictions at different regimes: initial (high ϵ) and final (low ϵ) convergence.

45 • **Random projections (R2):** We do not need the gradients to be sparse. We use the property of random projections to
46 preserve the expected L2 distances (and hence inner product) when applied to high dimensional vectors (see [1,5]).

47 • **Figure 1 (R4):** For each point we sample a different batch size, learning rate and a different dataset size.

48 • **Dataset details (R4):** We added a table describing for each dataset the number of images and classes.

49 • **Overestimation or underestimation? (R3):** We slightly overestimate the training time. This is due to the non-
50 linearity and high capacity of the network which is not entirely captured by the linearization approximation (see
51 Section B).

52 • **Infinite batch size assumption (R2):** Please note that this is used only in eq. (4) to give an analytical interpretation.
53 We make no such assumption in eq. (1), which is what we use for the actual training time prediction.