

1 We thank all reviewers for the constructive suggestions and the recognition on our novelty. We have carefully addressed
2 **every raised concern**. We truly hope that the reviewers and AC can reconsider the decision.

3 **Overall.** We propose a simple and novel neural similarity learning, which enables dynamic inference of CNN. The
4 factorized learning paradigm is shown effective in generic visual recognition, few-shot learning and efficient network
5 construction (in Appendix). NSL is theoretically connected with [Implicit Regularization in Matrix Factorization,
6 NeurIPS 2017] which shows that factorized learning is biased towards the minimum nuclear norm solution. Our
7 factorized learning of neural networks is empirically shown to have inductive bias that generalizes better.

8 All the reviewers recognize our novelty but concerns about the experimental evaluation and the presentation clarity.
9 We conducted all the requested experiments and will improve the presentation clarity regarding our methodology and
10 implementation. We will also release our implementation code for others to reproduce all the experimental results.

11 **Reviewer #1. [Input of $M(\cdot)$]** The input of $M(\cdot)$ in the current layer is the feature map from the bottom layer. From
12 another perspective, one can also think that the weights of the entire network depends on the input image, since the
13 weights of every layer recursively depend on the feature map from the bottom layer.

14 **[Single M for the whole network]** In the dynamic neural similarity, M is dynamic and depends on the input feature
15 map, so for different layer, M will be different (because the feature maps are different). In the shared parametrization
16 scheme, the network that generates M is shared across all the convolution kernels with the same size. Note that, the
17 size of M depends on the size of the convolution kernel.

18 **[Adaptation modules]** Adaptation modules are used to map the feature map from the bottom layer to a fixed-dimension
19 latent vector, such that all the convolution kernels with the same size can share the neural similarity generation network.

20 **[Fig. 3 and clarity]** We are very sorry about the confusion in Fig. 3 and we will modify it to reflect multiple layers. For
21 the clarity issue, we will improve the presentation of the detailed implementation in revision.

22 **[Architecture tweak]** Thanks for the suggestion. We agree with the reviewer that our method can be viewed as an
23 architectural modification which is generally useful. We use a simple way to construct a dynamic network whose
24 equivalent weights are dependent on the input. In fact, meta-learning is just one of the suitable applications for our
25 approach. Moreover, we also showcase the application of NSN in constructing efficient network in Appendix.

26 **[Additional experiments]** Thanks for
27 the constructive suggestion. We con-
28 duct all the meta-learning experiments
29 requested by the reviewer. The results
30 still show very consistent gain introduced by NSN. We will add these new results in revision.

Method	Finetune	ProtoNet	MAML
Vanilla	49.79±0.70 / 68.29±0.53	68.20±0.66 / 72.26±0.59	63.15±0.91 / 67.64±0.85
Static NSN	66.91±0.66 / 73.63±0.35	70.30±0.44 / 75.12±0.37	67.87±0.42 / 71.79±0.52
Dynamic NSN	67.82±0.71 / 78.98±0.60	71.39±0.67 / 78.72±0.79	68.90±0.65 / 72.23±0.49
Meta-learned static NSN		69.24±0.69 / 78.05±0.57	

Table 1: 5-shot acc. (%) on Mini-ImageNet. In "a / b", "a" and "b" are the accuracy with CNN-4 and CNN-9, respectively.

31 **Reviewer #2. [CNN++ in the ablation tables]** The CNN++ in all these ablation tables share the same setting with the
32 "CIFAR-10/100" subsection, so the error rate of CNN++ in all these ablation tables is the same as Table 4, *i.e.*, 7.29%.

33 **[BatchNorm]** All the baselines and our methods use batch normalization by default, which is also explicitly described
34 in Line 282 of the main paper. (Perhaps the reviewer accidentally missed this?)

35 **[Application to ResNets]** Thanks for the suggestion. We test both static NSN and
36 dynamic NSN on ImageNet-2012. To save GPU memory, we use DMS to constrain
37 the similarity matrix M and do not use NSL in 1x1 convolution. Our results show
38 very consistent gain on both ResNet-18 and ResNet-50. We will add these results and related details to revision.

Method	ResNet-18	ResNet-50
Baseline	31.69	24.77
Static NSN	30.12	23.91
Dynamic NSN	29.68	23.45

Table 2: Top-1 testing error (%) on ImageNet.

39 **[Number of parameters]** Dynamic NSN uses $\sim 25.4M$ parameters, while wide ResNet-28-10 used in LEO has more
40 than 36.5M parameters. We will add the parameter comparison for all the methods in revision.

41 **[Deformable convolution (DC)]** DC learns the continuous shape of the kernel, while ours learns both discrete shape
42 and similarity measure jointly. DC show no obvious gain in object classification. We will add discussions in revision.

43 **[Comparison to TADAM (Oreshkin et al., NIPS'18)]** This is indeed a very related work to compare with. TADAM
44 uses a ResNet-12 as the backbone network, while its additional task embedding network uses two separate fully
45 connected residual networks. The number of the total parameter used in TADAM is already much larger than our
46 method (dynamic NSN), while its 5-shot accuracy on Mini-ImageNet is still 3.2% lower than ours.

47 **[Data augmentation]** We keep our experiments fair by using the exactly same setting with the prior work. In Table 6,
48 the experimental setting exactly follows the MAML [ICML 2017] paper except for the last two rows. For the last two
49 rows, the setting is the same as the LEO [ICLR 2019] paper. For fairness, We do not perform additional augmentations.
50 Moreover, some methods like LEO perform much worse without data augmentation, so it may be unfair for LEO.

51 **Reviewer #3.** Thanks so much for the recognition on our work. Please see **Overall.** for a brief discussion about why
52 factorized learning of neural networks leads to better performance.