

1 We thank the reviewers for insightful comments. We have provided **code** in the supplemental for full reproducibility.

2 **Common Question: The method is for negative transfer rather than catastrophic forgetting.**

3 In transfer learning [23], “transfer” is the ability to apply knowledge learned in previous tasks to new tasks. Due to large
4 domain gap, only part of the pre-learned knowledge is useful for a new task: If the useful part is erased during transfer,
5 it is **catastrophic forgetting**; If the harmful part is preserved during transfer, it is **negative transfer** (Line 31-32).

6 Hence, catastrophic forgetting and negative transfer constitute a **dilemma** and should be mitigated jointly for optimal
7 performance. This is emphasized by the current Title and Introduction. While catastrophic forgetting has been studied
8 extensively by the community [34, 18], there has no work on mitigating negative transfer in fine-tuning. We propose a
9 novel approach to negative transfer, which is pluggable in the methods for catastrophic forgetting to tackle the dilemma.

10 **R1.1: How is Fig 1 (a)–(d) related to negative transfer? Why the new method works?**

11 **Fig 1(a)** shows that L^2 -SP performs worse than standard fine-tuning L^2 . This is a case of *negative transfer* by definition
12 [23]. In **Fig 1(b)** we delve into why negative transfer happens: the *Relative Angles* in the higher layers reveal that the
13 eigenvectors with smaller singular values are **not** transferable. This harmful part causes negative transfer in L^2 -SP since
14 it preserves all pre-trained knowledge. Similar results are observed for DELTA [18] (will be added to complete Fig 1).

15 As justified by [2], with sufficient labeled data, fine-tuning and training from scratch achieve comparably best results—
16 negative transfer does not happen in this case. Hence in **Fig 1(c)–(d)**, we analyze the singular values in this case, and
17 find that the smaller singular values are suppressed more. This hints us that the knowledge conveyed by eigenvectors
18 with smaller singular values are the causes of negative transfer and should be *shrunk*. This well motivates our approach.

19 **R2.1: BSS with continual learning & BSS in text classification with pre-trained word embeddings.**

20 In the table below: **For continual learning**, we evaluate BSS with EWC [13] on the permuted MNIST dataset. BSS
21 promotes the target task while slightly hurts the source. **For text classification**, BSS enhances the performance of
22 **BERT** [1], a state-of-the-art NLP pre-trained model. Results on Dev sets are listed, with all hyper-parameters consistent.

Method (continual learning)	task A	task B	Avg	Method (text classification)	MNLI-m	QNLI	MRPC	SST-2
fine-tuning + EWC	96.60	97.42	97.01	BERT _{base}	84.4	88.4	86.7	92.7
fine-tuning + EWC + BSS	96.46	98.04	97.25	BERT _{base} + BSS	85.0	89.6	87.9	93.2

23 **R3.1: Concern on novelty & compare with negative transfer methods in Domain Adaptation.**

24 Orthogonal to catastrophic forgetting, negative transfer is the bottleneck of transfer learning [23] and remains an *open*
25 *problem* in inductive transfer learning (a.k.a. fine-tuning in the context of deep learning). This work provides the first
26 approach to this important open problem, making a major contribution to this field.

27 Even in domain adaptation, there lacks in-depth analysis on negative transfer until [32] (CVPR’19). However, domain
28 adaptation and inductive transfer (fine-tuning) are completely different scenarios, detailed in the following table (left).
29 We are the first to address negative transfer in **fine-tuning**, to which domain adaptation methods cannot be applied.

30 The papers Reviewer #3 lists are important in domain adaptation. We will cite and discuss them in a future version.

Method	labeled source samples	source labels vs. target labels	Method	15%	30%	50%	100%
fine-tuning	unavailable	different	L^2	73.95±0.18	79.43±0.23	81.40±0.21	84.77±0.32
domain adaptation	available	identical	L^2 + re-initialize	70.32±0.32	76.36±0.29	79.98±0.28	83.35±0.33

31 **R3.2: Why not re-initialize all the high-level parameters and train again?**

32 **Fig 1** reveals that the eigenvectors with larger singular values in higher layers are transferable. If we re-initialize those
33 parameters, all pre-trained knowledge is discarded and *catastrophic forgetting* happens. Results on Stanford Dogs by
34 re-initializing Layer 4 in ResNet-50 are shown in the above table (right), which are worse than vanilla fine-tuning (L^2).

35 **R3.3: Why use feature regularization instead of parameter regularization?**

36 Parameter regularization has several disadvantages: (1) It is hard to decide weights of which layers should be regularized
37 (Line 191-192). In contrast, feature regularization can regularize each layer by taking the advantage of back-propagation.
38 (2) The parameters form high-dimensional matrix, whose SVD incurs unacceptable computational cost (Line 207-211).
39 In contrast, we can perform SVD over the feature matrix of each mini-batch, which only adds slightly more computation.

40 **R3.4: Ablation studies of the two parts (catastrophic forgetting & negative transfer).**

41 The ablation studies as requested by the reviewer have already been shown detailedly in Table 2; L^2 denotes the standard
42 fine-tuning; L^2 + BSS is the ablation study for the negative transfer part; L^2 -SP / DELTA is the ablation study for the
43 catastrophic forgetting part; And L^2 -SP + BSS / DELTA + BSS unifies the two parts to tackle the dilemma in them.

44 Through above it has been apparent that major questions raised by Reviewer #3 have been answered in the original paper.

45 References

- 46 [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language
47 understanding. In *NAACL*, 2019.
48 [2] K. He, R. Girshick, and P. Dollár. Rethinking imagenet pre-training. *arXiv preprint arXiv:1811.08883*, 2018.