



Figure 1: Projecting 50-dimensional embeddings obtained by training a simple neural network without SSE (Left), and with SSE-Graph (Center) , SSE-SE (Right) into 3D space using PCA.

Table 1: The experimental results of BERT-NER with SSE-SE and without SSE-SE when we only allow fine-tuning BERT for one epoch and two epochs. The same set of hyper-parameters are used.

Model	One Epoch				Two Epochs			
	accuracy(%)	precision(%)	recall(%)	F1 Score	accuracy(%)	precision(%)	recall(%)	F1 Score
BERT-NER	97.85	86.87	88.30	87.58	98.02	88.13	90.10	89.11
BERT-NER + SSE-SE	97.91	87.80	88.57	88.18	98.07	88.16	90.50	89.31

1 We thank the reviewers for their insightful feedback. In the following, we address their concerns and questions.
2 **R1.** It is indeed a great suggestion to examine concrete examples beyond the quantitative evaluation to get an intuition.
3 In Figure 1 (also in Figure 3 in the appendix), we use an embedding projection tool to visualize the embeddings from
4 Matrix Factorization model after applying dimension reduction to 3D space. When looking at some specific examples,
5 e.g. movies starred by the same actors, corresponding embeddings are closer to each other after using SSE-Graph
6 than the version without SSE. That is likely due to the use of item graph. Interestingly, even for SSE-SE, it does make
7 movies of the same genre closer to each other since it implicitly allows information sharing between embeddings.
8 **R2.** Thanks Reviewer 2 for raising the valuable question about the interpretation of SSE and connection to dropout. So
9 far, we have only explored the bias and variance trade-off. As shown in Theorem 1, SSE can ‘smooth’ the Rademacher
10 complexity. We have not studied whether SSE can be interpreted as Gaussian noise injection, but we feel it might
11 have similar effects. Figure 1 shows that the distributions of embeddings are definitely affected by the SSE-Graph and
12 SSE-SE and perhaps we can further study how this is related to dropout in theory.
13 **R3.** Following your suggestions and concerns, we used the BERT-NER with CRF loss to conduct the experiments in
14 Table 1. We did not have time to do pre-training, so this is only applying SSE-SE towards the fine-tuning stage of NER.
15 We fix the same set of hyper-parameters and use the same CoNLL-2003 training/test split. SSE probability of 0.01 is
16 used. Table 1 seems to imply that even for a task like NER, SSE-SE can still benefit the training procedure. Moreover,
17 for NER, it is perhaps intuitively reasonable to construct a graph over the vocabulary in the absence of Knowledge
18 Graphs, by connecting tokens of the same type: for example, an edge exists between two tokens if they are both Persons.
19 We agree that we need to dive deeper into each specific task and understand better when SSE helps. Our contention is
20 that SSE can help when the embedding space is very large, which does occur often in recommendation systems, but also
21 in NLP. It is outside of the scope of this paper to prove the usefulness of SSE to very different domains such as computer
22 vision. Currently, on all the tasks we tried (with large embedding space), SSE can always at least slightly improve
23 the performance; while for some tasks SSE can lead to significant performance gains. In general, we find SSE-Graph
24 and SSE-SE helps most when the model gets over-parameterized, either there is a large number of embeddings or
25 the dimensionality of embeddings is very large while the number of training data is limited. Moreover, we find that
26 SSE-Graph and SSE-SE can generally speed-up the training procedure of the original methods without the need to
27 touch the learning rate (as shown in Figure 2 in original submission) and that the existence of an informative knowledge
28 graph can often make this speed-up more significant.
29 We have also compared SSE-Graph with standard Graph Laplacian Regularization in Table 2 and we find that using
30 SSE-Graph can achieve much better RMSE. Moreover, it would not create a significant overhead to the training
31 procedure of Matrix Factorization as Graph Laplacian Regularization did.

Table 2: Compare SSE-Graph against Graph Laplacian Regularization.

Model	Movielens 1m	Movielens 10m
	RMSE	RMSE
SGD-MF	1.0984	1.9490
Graph Laplacian + ALS-MF	1.0464	1.9755
SSE-Graph + SGD-MF	1.0145	1.9019