

SupCL-Seq: Supervised Contrastive Learning for Downstream Optimized Sequence Representations

Hooman Sedghamiz¹ Shivam Raval¹ Enrico Santus¹ Tuka Alhanai²
Mohammad Ghassemi³

¹ DSIG - Bayer Pharmaceuticals, New Jersey, USA

² New York University, Abu Dhabi, UAE

³ Michigan State University, Michigan, USA

hooman.sedghamiz@bayer.com, sbraval@asu.edu,

enrico.santus@bayer.com, tuka.alhanai@nyu.edu, ghassem3@msu.edu

Abstract

While contrastive learning is proven to be an effective training strategy in computer vision, Natural Language Processing (NLP) is only recently adopting it as a *self-supervised* alternative to Masked Language Modeling (MLM) for improving sequence representations. This paper introduces *SupCL-Seq*, which extends the *supervised* contrastive learning from computer vision to the optimization of sequence representations in NLP. By altering the *dropout* mask probability in standard Transformer architectures (e.g. BERT_{base}), for every representation (*anchor*), we generate augmented altered *views*. A supervised contrastive loss is then utilized to maximize the system’s capability of pulling together similar samples (e.g., *anchors* and their altered *views*) and pushing apart the samples belonging to the other classes. Despite its simplicity, *SupCL-Seq* leads to large gains in many sequence classification tasks on the GLUE benchmark compared to a standard BERT_{base}, including 6% absolute improvement on CoLA, 5.4% on MRPC, 4.7% on RTE and 2.6% on STS-B. We also show consistent gains over *self-supervised* contrastively learned representations, especially in non-semantic tasks. Finally we show that these gains are not solely due to augmentation, but rather to a downstream optimized sequence representation. Code: <https://github.com/hooman650/SupCL-Seq>

1 Introduction

Sequence classification is a fundamental problem in natural language processing (NLP), as it has a wide range of applications, including but not limited to the tasks such as sentiment analysis, inference and question answering (Minaee et al., 2020). Cross-entropy loss is generally the default loss function in training neural networks for NLP downstream tasks (Zhang and Sabuncu, 2018; Sukhbaatar et al., 2015). Recently, thanks to the simplicity of augmentation methods in computer

vision (e.g., zooming, cropping, rotation, etc.), *self-supervised* and *supervised* variants of contrastive learning proved to be effective training approaches in image classification tasks (Wu et al., 2018; Hénaff et al., 2019; Khosla et al., 2020). These methods aim at optimizing the representations by minimizing the distance between similar samples and maximizing it between diverse samples (Chen et al., 2020). Gao et al. (2021) proposed to leverage the built-in *dropout* masks in attention and fully-connected layers of Transformers (Vaswani et al., 2017) to introduce noise in the embedding representations. This is obtained by simply passing *twice* the same input and using different dropout masks. In this way, for every representation (*anchor*), altered *views* are generated. Gao et al. (2021) applied this augmentation approach to improve the semantic representation of a sequence in a *self-supervised* fashion, by taking an input sentence and contrasting its similarity against its augmented version and the remaining samples in a batch. The authors further extended this approach by employing positive (i.e., entailment) and negative (i.e., contradiction) examples from natural language inference (NLI) datasets. The resulting sentence embeddings achieved large gains in semantic textual similarity (STS) tasks.

To the best of our knowledge, however, contrastive learning has not yet been applied in a supervised fashion to optimize sequence representations towards downstream tasks.¹ Inspired by the recently proposed *supervised* contrastive learning in computer vision (Khosla et al., 2020), in this paper we introduce *SupCL-Seq*, which extends the *self-supervised* contrastive method by Gao et al. (2021) to a *supervised* contrastive learning approach, in which *anchors* and altered *views*, along with their classification labels, are used to learn downstream

¹During the review process, we were made aware of a contemporaneous work by Gunel et al. (2020) on supervised contrastive learning for natural language processing. A major difference between their work and ours lays in the adopted augmentation methodology.

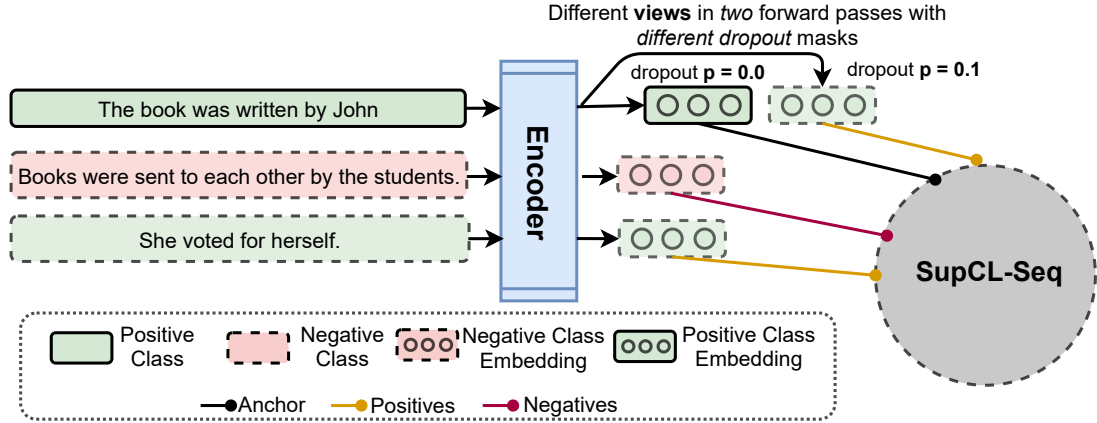


Figure 1: *SupCL-Seq* applied to COLA (Warstadt et al., 2018). *SupCL-Seq* first forward propagates the input N times (in this example $N = 2$) through the same encoder (e.g. $\text{BERT}_{\text{base}}$) with N different dropout masks (e.g. $p = 0$ and $p = 0.1$ respectively) and obtains their corresponding noisy embedding views. The noisy embeddings that belong to the same class are then employed as the positive pairs for the original input (*anchor* with dropout mask of $p = 0$). In this way, the samples belonging to the *negative* class effectively are used as negatives.

optimized sequence representations by means of contrastive learning (see Figure 1 for details).

SupCL-Seq is simple and can be applied to any sequence classification task, on any arbitrary number of classes. In our experiments, *SupCL-Seq* leads to large gains in several tasks of the GLUE benchmark (Wang et al., 2018), including CoLA (6% Matthew’s correlation coefficient absolute improvement), MRPC (5.4% accuracy score absolute improvement), RTE (4.3% accuracy score absolute improvement) and STS-B (2.6% Spearman’s rank correlation coefficient absolute improvement).

The main contributions of this paper are:

- The adaptation from computer vision to NLP of a *supervised* contrastive learning approach for sequence classification tasks (*SupCL-Seq*), which extends Gao et al. (2021)’s approach by optimizing the sequence representations for any downstream task, independently on the number of classes.
- Empirical demonstration that *SupCL-Seq* leads to significant gains in many text classification tasks in GLUE (Wang et al., 2018) using a standard transformer such as $\text{BERT}_{\text{base}}$ (Devlin et al., 2018).

2 Method

SupCL-Seq extends the self-supervised contrastive learning (Gao et al., 2021) for improving semantic representations to a supervised setting, in which representations are optimized towards the downstream task, independently on the number of labels.

The augmentation step is obtained by forward propagating the input batch N times in the same encoder with N distinct *dropout* masks (i.e., different dropout probabilities). The generated altered *views*, along with their *anchor*’s label, are then used to optimize the sequence representations through a supervised contrastive loss function (Khosla et al., 2020). Figure 1 details our training approach.

Formally, our pipeline consists of a single *Encoder Transformer*, $\text{Enc}(\cdot)$ (i.e., $\text{BERT}_{\text{base}}$ with $\approx 110\text{M}$ parameters (Devlin et al., 2018)). This encoder generates N altered embeddings, $\tilde{\mathbf{x}}_n = \text{Enc}(\mathbf{x}, p_n)$, for each input \mathbf{x} and *dropout* probability p_n .² A contrastive loss function is then applied in a supervised fashion to maximize the encoder’s capability of building downstream optimized sequence representations (see Section 2.1). After this contrastive training, the encoder parameters are frozen and a linear classification layer is then trained with cross-entropy. In the remainder of this section, we review the self-supervised contrastive function (Gao et al., 2021) and its extended *supervised* counterpart inspired by Khosla et al. (2020).

2.1 Contrastive Learning

Let $i \in I \equiv \{1 \cdots MN\}$ be the index of all the encoded sequence embeddings $\tilde{\mathbf{X}} \equiv \{\tilde{\mathbf{x}}_1 \cdots \tilde{\mathbf{x}}_{MN}\}$ in an input batch. Each sample i is forward propagated N time using distinct drop-out masks, gener-

²We employ the $\text{BERT}_{\text{base}}$ ’s last layer’s hidden-state of the first token of the sequence (i.e., pooled CLS embeddings) as $\tilde{\mathbf{x}}_n$, which is then L_2 normalized.

ating altered views denoted as $\tilde{\mathbf{x}}_{j(i)} = \text{Enc}(\mathbf{x}_i, p_n)$, where $j(i)$ refers to the indices of the altered *view(s)* for the sample i (also called *positive pairs*). In self-supervised contrastive learning (Gao et al., 2021; Khosla et al., 2020), the cost function is formulated as:

$$\mathcal{L}_i^{\text{self-sup}} = - \sum_{i \in I} \log \frac{e^{\text{sim}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_{j(i)})/\tau}}{\sum_{b \in B(i)} e^{\text{sim}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_b)/\tau}}, \quad (1)$$

Where, $B(i) \equiv I \setminus \{i\}$ is the set of so called *negative* pairs for the anchor $\tilde{\mathbf{x}}_i$. τ is a temperature scaling parameter. $\text{sim}(\cdot)$ stands for any similarity function such as cosine similarity or the inner product.

The main shortcoming of self-supervised contrastive learning is that since the class *labels* of the inputs are ignored, the samples from the same class might end up being employed as the negative pairs (e.g. $B(I)$) and therefore hurt the training performance. For instance, in CoLA (Warstadt et al., 2018) the aim is to determine whether the input is grammatical or ungrammatical. An unsupervised contrastive learning in this case might employ a grammatically correct sentence as the negative pair for the input anchor (see Figure 1).

In order to avoid this limitation and make the system able to learn in a *supervised* fashion, Khosla et al. (2020) extended Equation 1 to account for input labels. Given M annotated samples $\{\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i\}_{i=1 \dots M}$ passed N times through the *dropout* masks, the *supervised* contrastive learning loss is defined as:

$$\mathcal{L}_i^{\text{sup}} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{e^{\text{sim}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_p)/\tau}}{\sum_{b \in B(i)} e^{\text{sim}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_b)/\tau}}, \quad (2)$$

Where $P(i) \equiv \{p \in B(i) : \tilde{\mathbf{y}}_p = \tilde{\mathbf{y}}_i\}$ is the positive pair set distinct from sample i and $|\cdot|$ stands for cardinality (for details on derivation of Equation 2 see Khosla et al. (2020)). *SupCL-Seq* employs $\mathcal{L}_i^{\text{sup}}$ as contrastive loss function.

3 Experiments

We performed a set of experiments to i) evaluate the effect of number and level of dropout

Task	Drop-out	Batch size	Score
CoLA	[0.0,0.1,0.2,0.3,0.4]	800	61.2
	[0.0,0.1,0.2,0.3]	640	57.9
	[0.0,0.1,0.2]	480	58.9
	[0.0,0.1]	320	57.9
	[0.1,0.1]	256	60.7
RTE	[0.0,0.1,0.2,0.3,0.4]	800	63.5
	[0.0,0.1,0.2,0.3]	640	62.4
	[0.0,0.1,0.2]	480	69.3
	[0.0,0.1]	320	63.8
	[0.1,0.1]	256	65.3

Table 1: Effects of different dropout masks and number of views on CoLA and RTE tasks. Score denotes Matthews Correlation Coefficient.

passes on two challenging datasets (see 3.1); ii) compare the performance of a standard BERT_{base} (Devlin et al., 2018) architecture with a *SupCL-Seq*-empowered BERT_{base} model on several benchmarks in GLUE (Wang et al., 2018) (see 3.2); iii) compare the performance of *SupCL-Seq* with the *self-supervised* contrastive approach introduced by Gao et al. (2021) in a subset of tasks (see 3.3); and, finally, iv) assess whether the improvements achieved with our approach are solely due to augmentation (i.e., *dropout* masks) and to which extend contrastive loss helped (see 3.2.1).

3.1 Dropout Levels

In order to study the effect of the number and the level of dropout passes, we assessed the performance of several configurations of BERT_{base} on CoLA (Warstadt et al., 2018) and RTE (Dagan et al., 2006) datasets. Gao et al. (2021) empirically showed that using two distinct dropout masks with the same probability of $p = 0.1$ lead to the highest performance in their settings. In our supervised experiments, instead, we can generate *views* with different levels of noise, as the system can always rely on their labels. Therefore we choose different parameters, using intervals of 0.1 for the dropout probabilities. Table 1 reports the results for both datasets. While clear improvements are visible on **CoLA** when more masks are applied, experiments on **RTE** show that this is not always the case. In the latter dataset, in fact, performance fluctuates largely across the settings, achieving the highest score when three passes are used. This suggests that the number and level of dropout passes is a task-dependent hyper-parameter.

3.2 GLUE Tasks

In order to assess the benefit of *SupCL-Seq*, we compared the performance of a standard

System Nr. Training Samples	QQP	CoLA	MRPC	RTE	STS-B	SST-2	QNLI	WNLI	MNLI
BERT _{base} - <i>Standard</i>	71.2	55.2	86.6/80.3	64.6	88.0/87.7	92.6	90.5	56.6	84.1
BERT _{base} - <i>SupCL-Seq</i>	85.9	61.2	89.7/85.7	69.3	89.3/89.1	93.2	91.0	56.6	84.5
BERT _{base} - <i>Dropout Augmented</i>	-	60.9	88.9/84.5	63.1	80.4/81.6	93.4	-	-	-

Table 2: GLUE Test results. BERT_{base} - *Standard* is our implementation using the reported hyper-parameters in Devlin et al. (2018) for each task. BERT_{base} - *Dropout Augmented* is the standard version trained also on augmented samples. Matthews Correlation Coefficient is reported for CoLA, Pearson/Spearman correlations for STS-B, F1/Accuracy for MRPC, F1 score for QQP, and accuracy scores are reported for the other tasks.

BERT_{base} architecture with the one of a *SupCL-Seq*-empowered BERT_{base} model on numerous tasks from the GLUE benchmark (for a detailed description of each task see Wang et al. (2018)). GLUE also includes a regression task (i.e., **STS-B**), which requires no architecture modifications³. For the classification experiments, we deploy the hyper-parameters reported in Devlin et al. (2018). Appendix A details our grid-search-based training details for *SupCL-Seq*. Results are described in Table 2, rows one and two. As it can be noticed, in all cases the *SupCL-Seq*-empowered BERT_{base} model obtains equal or higher performance compared to the standard implementation.

3.2.1 Is it the dropout augmentation or the loss-function?

In order to study whether the performance gain observed in the previous experiments is solely due to the *dropout* augmentation, we ran a new set of experiments on the smaller datasets (i.e., MRPC, RTE, STS-B and SST-2) in which the standard BERT_{base} is trained also on the augmented samples (for the training parameters, see Appendix A). Table 2, third row, shows the results. While we notice performance gains compared to the BERT_{base} - *Standard* in a few tasks, augmentation does not always help. For example, the score for the augmented row is lower in the **RTE** dataset. Interestingly, dropout augmentation significantly hurts the performance (≈ 8 points) in **STS-B** dataset, where *MSE* loss is employed. We also observe that for all CoLA, MRPC and STS-B, *SupCL-Seq* outperforms the augmented variant, suggesting that its gains are due to the combination of augmentation and contrastive learning, rather than from only the former.

³To employ *SupCL-Seq*, we rounded to first decimal digit and grouped by similar labels, employing the Mean Squared Error (MSE) Loss for the baselines.

System	RTE	CoLA	MRPC
BERT _{base} - <i>Self-supervised-CL</i>	55.6	35	79/68.3
BERT _{base} - <i>SupCL-Seq</i>	69.3	61.2	89.7/85.7

Table 3: Comparison of unsupervised and supervised contrastive loss.

3.3 Supervised Versus Unsupervised contrastive Learning

Since, to the best of our knowledge, the only previous attempt of using contrastive learning for improving sequence representation in NLP was performed by Gao et al. (2021) – they used a *self-supervised* approach to improve the semantic representation, adopting a loss similar to Equation 1 –, in Table 3.3 we compare the performance of a linear layer trained on top of their representations with the one of a linear layer trained on top of our representations, which are instead optimized in a *supervised* fashion while the parameters of the base model are kept frozen. *SupCL-Seq* significantly outperforms the re-implementation of Gao et al. (2021), with larger gains in non-semantic tasks (e.g. CoLA), suggesting that our representations are optimized for the given downstream tasks.

4 Discussion and Conclusion

In this paper, we introduced *SupCL-Seq* a *supervised* contrastive learning framework for optimizing sequence representations for downstream tasks. In a series of experiments, we showed that *SupCL-Seq* leads to large performance gains in almost all GLUE tasks when compared to both a standard BERT_{base} architecture and an augmented BERT_{base} (i.e., improvements are not only due to augmentation). We also investigated the effect of number and level of *dropout* passes, finding that this has to be treated as a task-dependent hyper-parameter, to be fine tuned in a validation set. Finally, we compared our *supervised* approach to the *self-supervised* method by Gao et al. (2021), showing

consistent performance improvements, especially in non-semantic tasks, where the *self-supervised* approach is weaker. These encouraging results open the door to multi-task learning applications of *SupCL-Seq*, where the optimization needs to be constrained towards multiple objectives.

Acknowledgments

We would like to thank the reviewers and the chairs for their insightful reviews and suggestions.

References

- Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. [Big self-supervised models are strong semi-supervised learners](#). *CoRR*, abs/2006.10029.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The pascal recognising textual entailment challenge. In *Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment*, pages 177–190, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. [Simcse: Simple contrastive learning of sentence embeddings](#). *CoRR*, abs/2104.08821.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Veselin Stoyanov. 2020. Supervised contrastive learning for pre-trained language model fine-tuning. In *International Conference on Learning Representations*.
- Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aäron van den Oord. 2019. [Data-efficient image recognition with contrastive predictive coding](#). *CoRR*, abs/1905.09272.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. [Supervised contrastive learning](#). *CoRR*, abs/2004.11362.
- Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Narjes Nikzad, Meysam Chenaghlu, and Jianfeng Gao. 2020. [Deep learning based text classification: A comprehensive review](#). *CoRR*, abs/2004.03705.
- Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. 2015. Training convolutional networks with noisy labels. 3rd International Conference on Learning Representations, ICLR 2015 ; Conference date: 07-05-2015 Through 09-05-2015.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). *CoRR*, abs/1706.03762.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2018. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. [Huggingface’s transformers: State-of-the-art natural language processing](#). *CoRR*, abs/1910.03771.
- Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. 2018. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhilu Zhang and Mert R. Sabuncu. 2018. [Generalized cross entropy loss for training deep neural networks with noisy labels](#). *CoRR*, abs/1805.07836.

A Training Details

Task	Learning Rate	Batch Size	dropout
CoLA	$5e-05$	128	[0.0, 0.1, 0.2]
MRPC	$1e-4$	128	[0.0, 0.05, 0.1, 0.2]
RTE	$1e-4$	48	[0.0, 0.1, 0.2]
STS-B	$1e-4$	64	[0.0, 0.05, 0.1, 0.2]
SST-2	$5e-05$	320	[0.0, 0.1, 0.2]
WNLI	$1e-04$	320	[0.0, 0.1, 0.2]
QNLI	$5e-05$	48	[0.0, 0.2]
QQP	$5e-05$	16	[0.0, 0.2, 0.3, 0.4, 0.5]
MNLI	$5e-05$	8	[0.1, 0.1]

Table 4: Contrastive learning training details per GLUE task. All of the tasks were trained for 5 epochs (except *QNLI*, *QQP* and *MNLI* that were trained for 2, 1 and 3 epochs respectively) and $\tau = 0.05$.

SupCL-Seq is implemented on top of the Huggingface’s trainer python package (Wolf et al., 2019)⁴. For the *sim(.)* (similarity) function, we employed *inner dot product*. For supervised contrastive learning, we employed the hyperparameters detailed in Table 4. We used a grid search strategy for our hyperparameter optimization, where the number of dropouts and their corresponding probability were set to two (i.e. [0.1, 0.1]) and five respectively ([0.0, 0.1, 0.2, 0.3, 0.4]). For the learning rate we employed a range of [$5e-05$, $1e-4$].

⁴<https://github.com/huggingface/transformers>