

Exploring the Impact of Negative Samples of Contrastive Learning: A Case Study of Sentence Embedding

Rui Cao and Yihao Wang and Yuxin Liang and Ling Gao and Jie Zheng and Jie Ren and Zheng Wang

{ caorui, wangyihao, liangyuxin }@stumail.nwu.edu.cn, { gl, jzheng }@nwu.edu.cn, renjie@snnu.edu.cn, z.wang5@leeds.ac.uk

Background and Motivation

Contrastive learning is emerging as a powerful technique for extracting knowledge from unlabeled data. This technique requires a balanced mixture of two ingredients: positive (similar) and negative (dissimilar) samples. This is typically achieved by maintaining a queue of negative samples during training. Prior works in the area typically uses a fixed-length negative sample queue, but how the negative sample size affects the model performance remains unclear. The opaque impact of the number of negative samples on performance when employing contrastive learning aroused our in-depth exploration.

Can we apply MoCo-style contrastive structure on unsupervised textual tasks?

If could:

1. How to boost performance?
2. How to prevent model from collapsing?
3. How much negative sample do textual task needs?

This paper presents a momentum contrastive learning model with negative sample queue for sentence embedding, namely MoCoSE. We add the prediction layer to the online branch to make the model asymmetric and together with EMA update mechanism of the target branch to prevent the model from collapsing. We define a maximum traceable distance metric, through which we learn to what extent the text contrastive learning benefits from the historical information of negative samples.

Overview of Our Approach

We adapted a MoCo-style model to study the effect of the structure in text contrastive learning. The structure of the model is shown in Fig. 1.

The sentences are passed through the embedding layer and data augmentation to generate two slightly different embeddings, and then the embeddings of query and key are obtained through the online and target branches, respectively. The structure of encoder, pooler and projection of online and target branch is identical. Similar to MoCo V2, a prediction layer composed of MLP is added to the upper branch to prevent model from collapsing. We use the same InfoNCE loss as MoCo.

The PyTorch style pseudo-code for training MoCoSE with the negative sample queue is shown in Algorithm 1.

Experiment and Analysis

1. Main Results

We train with a corpus of 1 million sentences randomly selected from the English Wikipedia and experiment in seven standard semantic text similarity (STS) tasks using [CLS] token embeddings from the online encoder output as sentence embeddings.

Learning rate:	
MoCoSE-BERT-base	3e-5
MoCoSE-BERT-large	1e-5
Batch size:	
MoCoSE-BERT-base	64
MoCoSE-BERT-large	32
Validate per step	
100	
EMA decay weight	0.75-0.95
Weight decay	1e-6
Negative queue size	512

In addition to the semantic similarity task, we evaluate seven migration learning tasks.

The proposed base model surpasses SimCSE in most of unsupervised tasks, and the large model obtains better results in most migration tasks.

2. Ablation Study

Symmetric Two-branch Structure

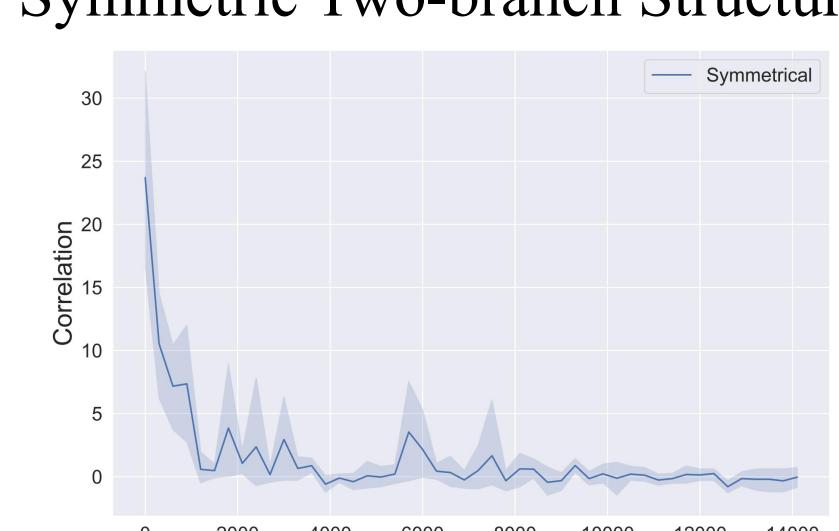


Figure 5: Experiment on a symmetric two-branch structure with EMA decay weight set to 0.

EMA Hyperparameters

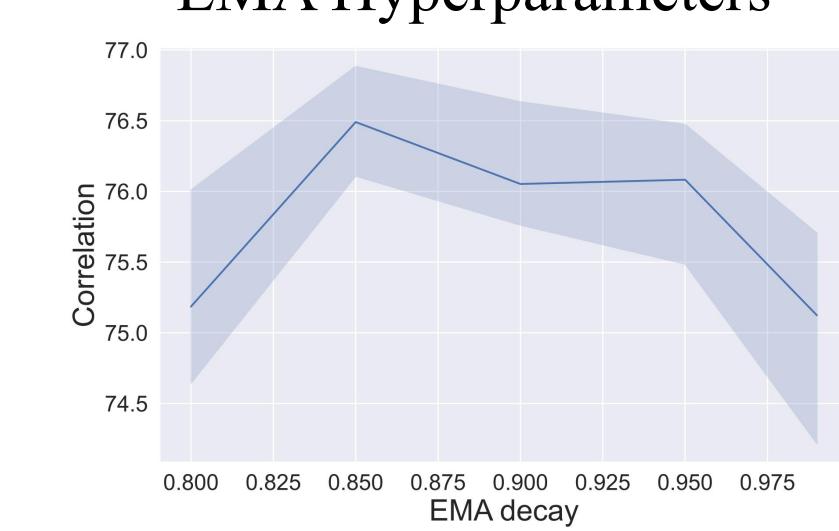


Figure 8: Effect of EMA decay weight on model performance.

Different Data Augmentations

Augmentation Methods		Avg
Dropout only		76.76
+ FGSM ($\epsilon=5\text{-}9$)		77.04
+ Position_shuffle (True)		78.00
+ Token dropout (prob=0.1)		41.32
+ Feature dropout (prob=0.01)		76.33
+ Feature dropout (prob=0.1)		71.62
+ Typo		22.32
+ Synonym replace (roberta-base)		28.70
+ Paraphrasing (xlnet-base-cased)		60.45
+ Backtranslation (en-de-en)		69.35
Epsilon 1e-9	5e-9	1e-8
Avg.	75.61	76.64
Avg.	75.39	76.62
Avg.	76.62	76.26

Table 9: Different parameters of FGSM in data augmentation affect the model results.

FGSM first calculates the derivative of model with respect to the input, and use a sign function to obtain its specific gradient direction. Then, after multiplying it by a step size, the result x_{0002_ing} perturbation is added to the original input to obtain the sample under the FGSM attack.

We use the alpang toolkit in our data augmentation experiments. All the parameter listing is default value given by official. <https://github.com/makeeward/alpang>

Table 4: The impact of different combinations of projection and predictor on the model.

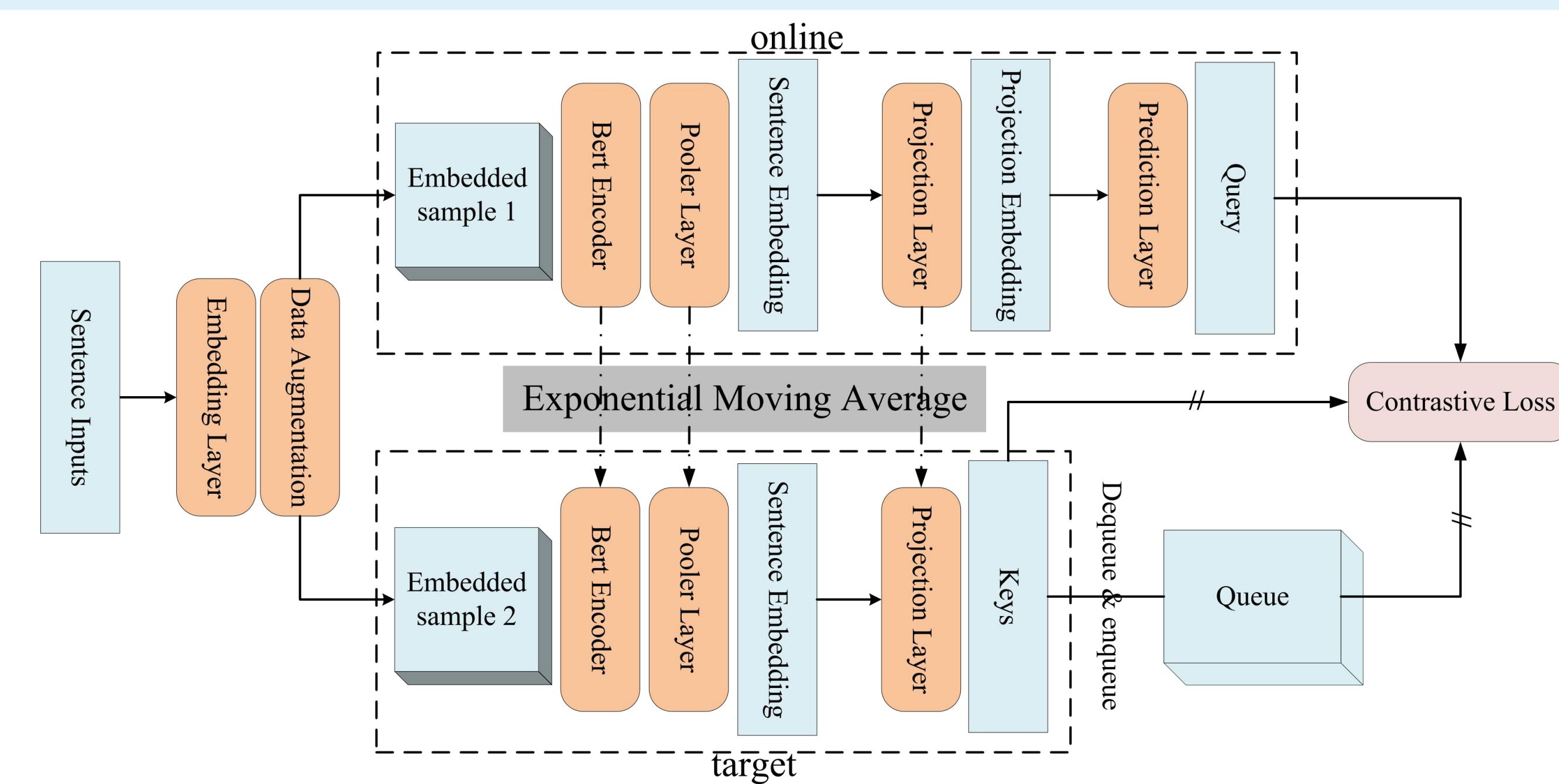


Fig.1 The model structure of MoCoSE. The embedding layer consists of a BERT embedding layer with additional data augmentation. The pooler, projection, and predictor layers all keep the same dimensions with the encoder layer. The MoCoSE minimizes contrastive loss between query, queue and keys (i.e. InfoNCE loss).

Predictor Mapping Dimension

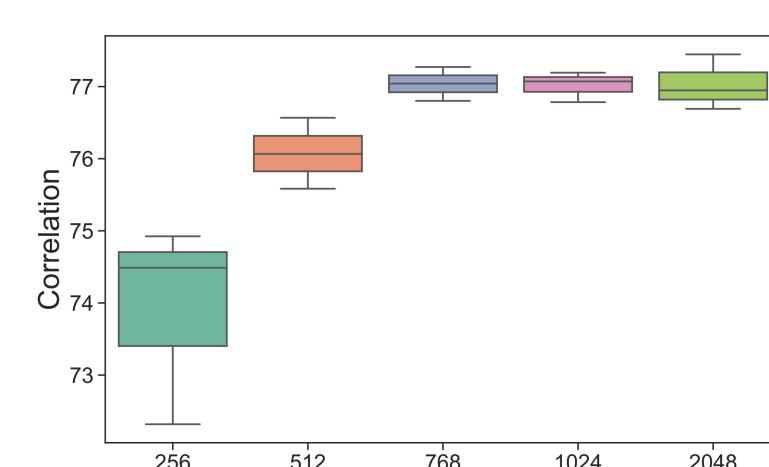


Figure 10: Impact of dimensions of the sentence embedding on the model with fixed queue size of 512.

When the dimension of embedding is low, this causes considerable damage to the performance of the model, while the dimension rises to certain range, the performance of the model stays steady.

Batch Size

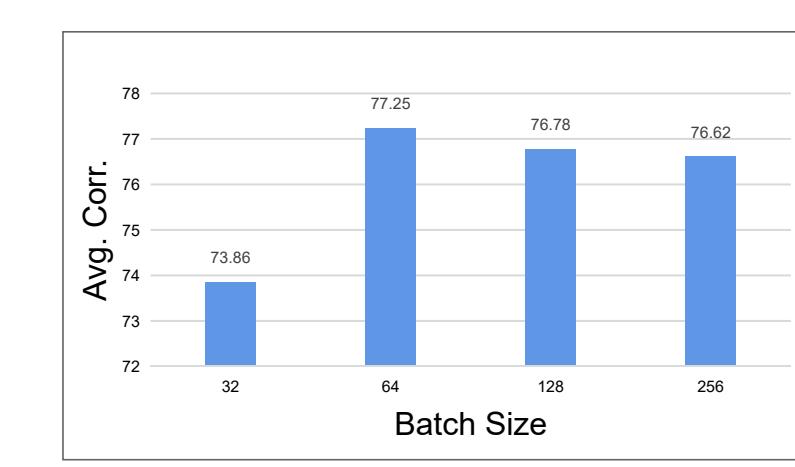


Table 6(b): Impact of batch size on the model with fixed queue size of 512.

The model performance does not improve with increasing batch size, which contradicts the general experience in image contrastive learning.

Distribution of Singular Values

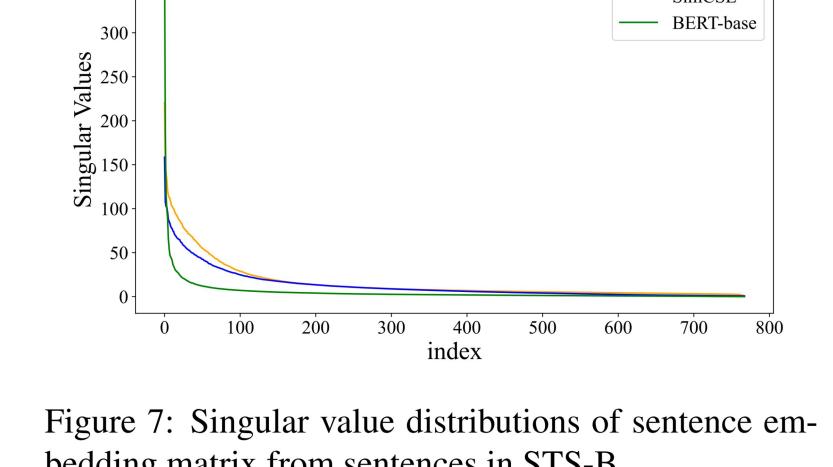


Figure 7: Singular value distributions of sentence embedding matrix from sentences in STS-B.

The model is able to alleviate the rapid decline of singular values compared to other methods, making the curve smoother, i.e., our model is able to make the sentence embedding more isotropic. The rapid decrease of the singular value means that the distribution is wider in the first few directions (singular value represents the variance) and smaller in the later directions, the decrease of the singular value slows down means the variance of the distribution becomes more consistent and uniform in all directions.

3. The study of Negative Sample Queue

We test the size of negative sample queue to the model performance. With queue size longer than 1024, the results get unstable and worse. We suppose this may be due to the random interference introduced to the training by filling the initial negative sample queue.

Initial Size	128	256	512	1024	4096	Queue Size
w.o. init.	76.40	76.19	75.38	76.63	50.17	
init. 1/4 queue	75.92	76.34	77.30	76.20	50.42	
init. 1/2 queue	76.16	76.39	76.94	76.57	38.74	
init. all (normal)	76.87	75.81	76.29	76.45	45.80	

Table 7: Correlation performance of initializing different proportion of negative queue with different negative queue size.

Add random initial process to First-in-First-out queue

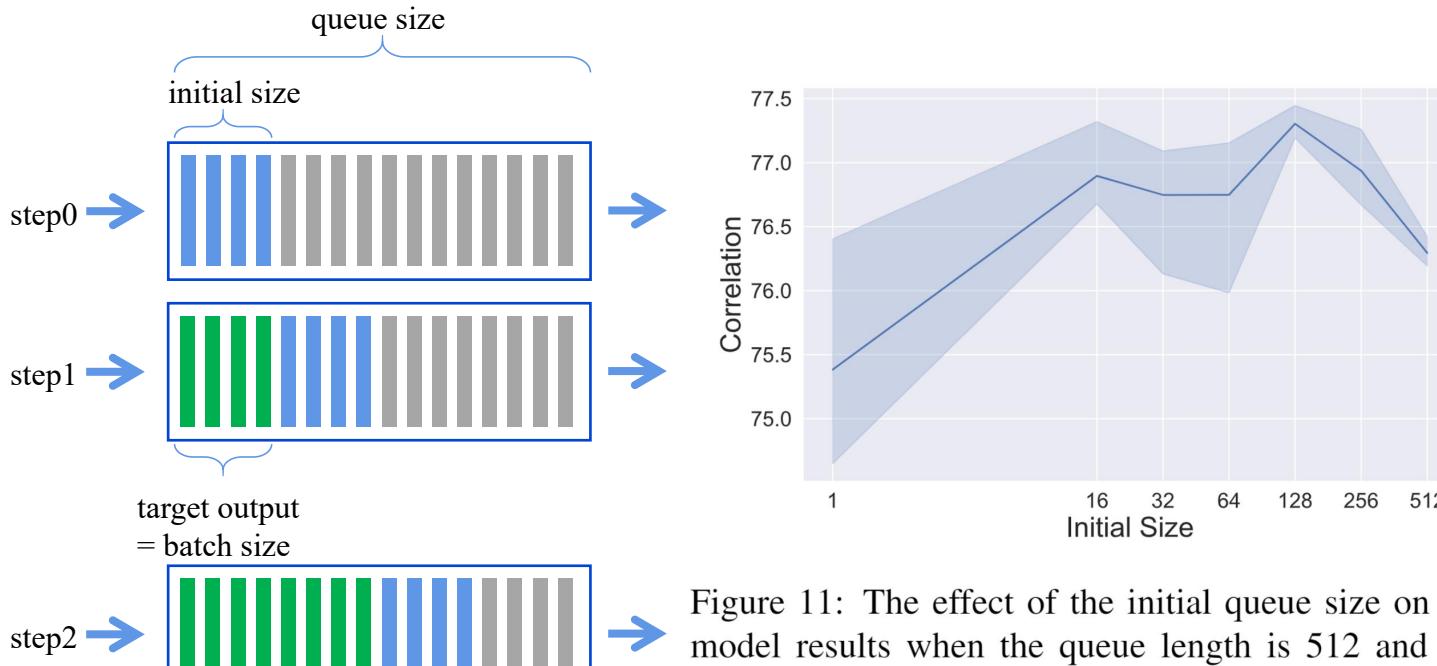


Figure 11: The effect of the initial queue size on the model results when the queue length is 512 and the batch size is 64.

We initialize a smaller negative queue, then fill the queue to its set length in the first few updates, and then update normally.

queue size	initial size	target output = batch size
step0	1	1
step1	16	16
step2	32	32

Figure 11: The effect of the initial queue size on the model results when the queue length is 512 and the batch size is 64.

;

We initialize a smaller negative queue, then fill the queue to its set length in the first few updates, and then update normally.

set queue size = 1024	0	selected	1024
Corr.	0~	256~	512~
512	768	1024	256~768
Avg.	76.10	77.02	75.71
		76.18	76.86

Table 8: The impact of negative samples at different locations in the queue on the model performance.

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.

;

We find that the increase in queue length affects the model performance not only because of the increased number of negative samples, but more because it provides historical information within a certain range.