

---

# GENERATING SYNTHETIC DATA VIA INTENT-PRESERVING AND INTENT-CORRUPTING AUGMENTATIONS FOR MULTI-TASK TRAINING DIALOGUE EMBEDDINGS. \*

---

**Ilya Alekseev**

Neural Networks and Deep Learning Lab  
Moscow Institute of Physics and Technology  
ilya\_alekseev\_2016@list.ru

**Denis Kuznetsov**

Neural Networks and Deep Learning Lab  
Moscow Institute of Physics and Technology  
kuznetsov.dp@phystech.edu

## ABSTRACT

Text embeddings from pre-trained language models have been proven to be extraordinarily useful for various sentence-level tasks, such as pair classification, similarity estimation, and retrieval. Corresponding models are usually trained on large amounts of clean and diverse data using contrastive loss. Unfortunately, there are no such datasets for the domain of dialogue data. In this work, we describe the process of mining a synthetic dataset of dialogues for contrastive learning with hard negative samples. We investigate various augmentation strategies for constructing dialogues with preserved or corrupted intents. To demonstrate the stated cleanliness and diversity, we train dialogue embeddings and report its performance on various tasks.

**Keywords** text embedding · dialogue · synthetic data · augmentation · natural language processing

## 1 Introduction

Obtaining embeddings is one of the key tasks in machine learning and popular one in recent years, because vector representation of an object is a convenient mathematical object. If an embedding accurately and comprehensively encodes the semantics of the original data, it opens up the possibility of using it in a wide range of tasks.

In the field of natural language processing, classical methods for text vectorization such as bag of words [1] and tf-idf [2] have long been discovered. Thanks to deep learning, we have witnessed remarkable word vectorization techniques such as word2vec [3] and GloVe [4], which convey the semantic similarity between words; CoVe [5] and ELMo [6], which encode information about the surrounding context. Recently, powerful encoder models have emerged that produce general-purpose text embeddings [7, 8, 9]. They incorporate so much semantic information about texts that it can be applied to tasks such as classification, clustering, ranking, semantic textual similarity, and more [10]. The success of these models is largely attributed to the use of contrastive learning on massive datasets.

However, the more specific the data structure, the more challenging it is to mine a large dataset. As of today, there are no encoding methods for entire dialogues. In other words, there is no way to obtain a dense vector representation that conveys universal information about a dialogue. There are language models that adapt successfully to the hierarchical and temporal nature of dialogue [11, 12], but they only address tasks at the token and utterance levels, not at the dialogue level. There are text encoders for utterances in a dialogue [13, 14], but not for the entire dialogue.

Data is almost always scarce when it comes to building dialogue models. To date, numerous methods for generating synthetic dialogues have been devised [15, 16, 17, 18], but they do not generate dialogues in pairs [19], as it is conceptually important for training SoTA text embedding [7]. The simplest way to expand a training dataset is through augmentation [20]. In this work, we will describe a method for constructing a synthetic dataset of dialogue pairs using various augmentations that preserve or alter the set of intents in the dialogue. These augmentations can be used for contrastive learning along with other tasks for pretraining powerful dialogue embeddings.

---

\**Citation:* Authors. Title. Pages.... DOI:000000/11111.

## 2 Problem Formulation

**Dialogue Data.** A dialogue is defined as the following list:

$$d = [(u_1, s_1), \dots, (u_n, s_n)],$$

where  $u_i$  represents the utterance of a participant in the dialogue at step  $s_i$ . We are interested in so-called task-oriented dialogues with two participants: the system and the user. With some approximation, they can be described as dialogues between a customer and a service worker (or a robot). During the dialogue, the customer has various intents that the worker strives to fulfill. These intents can be finding a restaurant and booking a table, calling a taxi, or purchasing a train ticket. We will consider two dialogues similar if they have a similar set of intents.

**Augmentation.** By augmentation, we mean the generation of new valid examples by transforming existing ones. Valid examples are those that sufficiently resemble real-world data. Let  $D$  be the set of valid objects. Then augmentation is a non-identical mapping  $\text{aug}(d)$  that does not take objects outside the set of valid objects:

$$\text{aug} : D \rightarrow D.$$

Typically, this generation is implemented by making small changes to a valid training object. For example, image augmentation might involve slight rotations or blurring. Text augmentation is especially challenging because validity implies adherence to language rules, meaningfulness and a certain style. In the case of dialogues, it is necessary to maintain the structure and role differentiation, as indicated earlier.

**Embedding.** Embedding is a mapping of  $D$  to a vector space:

$$e : D \rightarrow E \subseteq \mathbb{R}^d.$$

For an object  $d$ , its embedding  $e(d)$  should convey some semantics of  $d$ . This is reflected in the fact that  $e(d)$  may contain lexical information or latent features useful for classification and other downstream tasks. It is especially valuable if using embeddings  $e(a), e(b)$  allows for assessing the semantic similarity of objects  $a, b$ .

## 3 Related Works

this section is incomplete, this is still a draft

AugSBERT [21] can be viewed as the similar approach to ours one, but their augmentation uses an already presented dataset of text pairs, which is absent in our case.

## 4 Method

In this section, we describe our methodology to augment dialogue data, construct dialogue encoder and train it.

### 4.1 Augmentations

**Token Insertion.** One of the simple yet effective ways to augment text is to lengthen it by inserting extra tokens. For this purpose, we added a special token ‘<mask>’ to random places in the dialogues and used transformer models trained on the MLM task to fill these masks [22]. Insertion is rejected if the token proposed by the model is only a fragment of a word [23, 24] or if the prediction probability is below a manually set threshold. To take the dialogue context into account during token insertion, multiple consecutive responses were fed into the mask-filling model at once as a compromise between feeding single utterances and entire dialogue.

**Token Replacement.** This method is identical to the previous one, except that instead of adding the ‘<mask>’ token, some tokens in the original dialogue are replaced. In this case, the mask-filling model is fed with single utterances to make replacements more diverse and random.

**Back Translation.** Translation from the original language to another language and then back to the original language. Neural machine translation models were used for this purpose [25].

**Shuffling Utterances.** Previous augmentations modify the dialogue within a single utterance, since they are methods applicable to arbitrary text data. It seems essential to learn how to change the order of utterances in a dialogue to create new valid dialogue. For this purpose, we propose using a model that measures the similarity between utterances within a dialogue. Using these similarities, it is possible to cluster utterances within each dialogue. Experiments showed that these clusters represent significant individual stages of the dialogue that can be shuffled with each other.

**Shortening Dialogue.** Individual clusters of utterances within the dialogue can be discarded, resulting in a pruned dialogue with fewer utterances.

**Lengthening Dialogue.** The special model was trained to arrange a list of given utterances. This transformer takes the text embeddings of each utterance as input sequence, without specifying information about their order in the original dialogue. It outputs ranks that can be used to "sort" the utterances to restore the original order. If some external responses are added to the original dialogue, this model generates a new, longer dialogue.

The augmentations mentioned above can be configured to either preserve or alter intent. Specifically, token replacement can be viewed as intent-corrupting augmentation, because all the keywords such as "restaurant", "taxi" etc. tend to be replaced. Pruning dialogue may remove some intents, but the result is still much similar to original dialogue, since its intents are fully encompassed by the original dialogue's intents. Shuffling utterances doesn't change any intents, it only changes their order. The rest of augmentations preserves intents because they either perform paraphrasing (back translation), or add new information (token insertion, lengthening dialogue).

To expand the set of augmentations even further, we use a composition of several ones. For more details on augmentation implementations, please refer to Appendix B.

## 4.2 Dialogue Encoder Architecture

As a method for embedding, we use RoBERTa [26] without any modifications. The input is [CLS] ut1 [SEP] ut2 [SEP] ut3, and the output is the hidden state of [CLS] token from the last layer.

For a slightly advanced dialogue language model, we use HSSA [11]. It uses BERT [22] as a backbone and modifies its attention to capture the hierarchical structure of dialogue and reach the computational trade-off between feeding a transformer separate utterances and feeding an entire dialogue.

## 4.3 Pre-train Tasks

To make the encoder produce rich vector representations, it is necessary to train it to perform tasks where semantic features are engaged. A popular task is contrastive learning:

$$\mathcal{L} = -\log \frac{\exp(\text{sim}(x, y))}{\sum_z \exp(\text{sim}(x, z))}.$$

Here,  $x = e_\theta(a)$ ,  $y = e_\theta(b)$  are embeddings of a pair of semantically similar objects and  $z = e_\theta(c)$  is semantically distant from  $x$ ,  $\text{sim}$  is a metric similarity function. This task trains embeddings to convey semantic similarity as metric similarity. We propose two pre-train tasks with contrastive loss: dialogue retrieval and dialogue variant of inverse cloze task.

**Dialogue Retrieval.** In this case,  $a, b$  are two intent-preserving augmentations of the same dialogue, and the remaining dialogues from the training batch are used as negative examples. This task simulates the second stage of training text embeddings, which uses a large batch and weakly supervised text pairs. In addition to this, in order to reduce required batch size, we utilize intent-corrupting augmentations as hard negative samples. We believe that this pre-train task plays a major role in training a rich vector representation of the dialogue. Notably, the technique with two augmentations and siamese networks resembles BYOL [27], well-known method for self-supervised learning (maybe move this sentence to Related Works).

**Inverse Cloze Task.** For this task, a dialogue and some span of utterances from it are taken as  $a, b$ , and negative examples are the remaining dialogues and spans in the batch. This is a dialogue variant of ICT [28]. This task partially simulates the first stage of training text embeddings, in which a retrieval-oriented language model is pretrained. Recall that the goal of this stage is to train the CLS token to aggregate information about the entire input sequence. Therefore, this task can be easily discarded by taking the pre-trained RetroMAE [29] instead of the usual BERT.

**Dialog-post Tasks.** We employ tasks from Dialog-post paper [11]. They aim to train a versatile language model that extends the domain from plain texts to the dialogue data. it's a draft of the paragraph

## 4.4 Summary of Method

Our method can be divided into two big stages. At the first stage, we develop and apply augmentations for dialogue data. In result, the dataset consists of samples with following schema: original dialogue and two sets of augmented versions with preserved and corrupted intents. At the second stage, we apply the multitask learning with all the tasks listed in section 4.3. Tasks of dialogue retrieval and inverse cloze task

## 5 Experiments

**Dataset.** In all experiments, we used the same dataset of dialogues. This dataset is a combination of several task-oriented dialogue datasets from DialogStudio [30]. All the dialogues were filtered based on their length, resulting in a dataset comprising 501K dialogues. For more details on the dataset used, please refer to Appendix A.

**Evaluation.** We evaluate dialogue embeddings in transfer learning setting. We adopt methods from the work [31] and complement them with one more: multi-label classification of services in the MultiWOZ 2.2 dataset. All these methods utilize dialogue embeddings as features, without modifying the encoder model itself.

**Baseline Solution.** As our baseline solution, we employ BERT and BERT-RetroMAE without the pretraining method proposed by us. The results are presented in Table 1.

	Service Clf	Categorization	Relatedness	Retrieval
BERT	0.49			
BERT-RetroMAE				
BERT + Contr loss	0.68			

Table 1: Dialogue embedding evaluation results.

## Acknowledgments

This was supported in part by.....

## References

- [1] Wisam A. Qader, Musa M. Ameen, and Bilal I. Ahmed. An overview of bag of words;importance, implementation, applications, and challenges. In *2019 International Engineering Conference (IEC)*, pages 200–204, 2019.
- [2] Karen Spärck Jones. A statistical interpretation of term specificity and its application in retrieval. *J. Documentation*, 60:493–502, 2021.
- [3] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013.
- [4] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics.
- [5] Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. Learned in translation: Contextualized word vectors, 2018.
- [6] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- [7] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighof. C-pack: Packaged resources to advance general chinese embedding, 2023.
- [8] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training, 2022.
- [9] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards general text embeddings with multi-stage contrastive learning, 2023.
- [10] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embedding benchmark, 2023.
- [11] Zhenyu Zhang, Lei Shen, Yuming Zhao, Meng Chen, and Xiaodong He. Dialog-post: Multi-level self-supervised objectives and hierarchical model for dialogue post-training. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10134–10148, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [12] Yiyang Li, Hai Zhao, and Zhuosheng Zhang. Back to the future: Bidirectional information decoupling network for multi-turn dialogue modeling, 2022.

- [13] Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić, Pei-Hao Su, Tsung-Hsien Wen, and Ivan Vulić. Convert: Efficient and accurate conversational representations from transformers, 2020.
- [14] Zhihan Zhou, Dejiao Zhang, Wei Xiao, Nicholas Dingwall, Xiaofei Ma, Andrew O. Arnold, and Bing Xiang. Learning dialogue representations from consecutive utterances. In *NAACL 2022*, 2022.
- [15] Sungdong Kim, Minsuk Chang, and Sang-Woo Lee. Neuralwoz: Learning to collect task-oriented dialogue via model-based simulation, 2021.
- [16] Biswesh Mohapatra, Gaurav Pandey, Danish Contractor, and Sachindra Joshi. Simulated chats for building dialog systems: Learning to generate conversations from instructions, 2021.
- [17] Dazhen Wan, Zheng Zhang, Qi Zhu, Lizi Liao, and Minlie Huang. A unified dialogue user simulator for few-shot data augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3788–3799, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [18] Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng Zhang, and Minlie Huang. Augesc: Dialogue augmentation with large language models for emotional support conversation, 2023.
- [19] Timo Schick and Hinrich Schütze. Generating datasets with pretrained language models, 2021.
- [20] Heydar Soudani, Evangelos Kanoulas, and Faegheh Hasibi. Data augmentation for conversational ai. *arXiv preprint arXiv:2309.04739*, 2023.
- [21] Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. Augmented sbert: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks, 2021.
- [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [23] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation, 2016.
- [24] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [25] Jörg Tiedemann and Santhosh Thottingal. OPUS-MT — Building open translation services for the World. In *Proceedings of the 22nd Annual Conferenec of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal, 2020.
- [26] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- [27] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning, 2020.
- [28] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering, 2019.
- [29] Shitao Xiao, Zheng Liu, Yingxia Shao, and Zhao Cao. Retromae: Pre-training retrieval-oriented language models via masked auto-encoder, 2022.
- [30] Jianguo Zhang, Kun Qian, Zhiwei Liu, Shelby Heinecke, Rui Meng, Ye Liu, Zhou Yu, , Huan Wang, Silvio Savarese, and Caiming Xiong. Dialogstudio: Towards richest and most diverse unified dataset collection for conversational ai, 2023.
- [31] Che Liu, Rui Wang, Junfeng Jiang, Yongbin Li, and Fei Huang. Dial2vec: Self-guided contrastive learning of unsupervised dialogue embeddings, 2022.
- [32] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art natural language processing, 2020.

## A Dialogue Dataset

Large dataset is important for contrastive learning. For training our models, we took a merge of some task-oriented datasets from DialogStudio collection [30]. All of them are listed in the table 2.

Name	# dialogues	# utterances	# tokens
AirDialogue	321K	4086K	49.4M
SimJointGEN	100K	1584K	22.5M
MS-DC			
MetaLWOZ			
MULTIWOZ2_2			
SGD			
KETOD			
FRAMES			
Disambiguation			
ABCD			
AirDialogue			
BiTOD			
Taskmaster1			
<b>Total</b>			
<b>Filtered</b>	501K	6320K	83.7M

Table 2: All the datasets are taken from DialogStudio collection [30]. BERT tokenizer is used to count # tokens column.

## B Auxiliary Models

As mentioned in Section 4.1, we trained special models to perform dialogue-level augmentations.

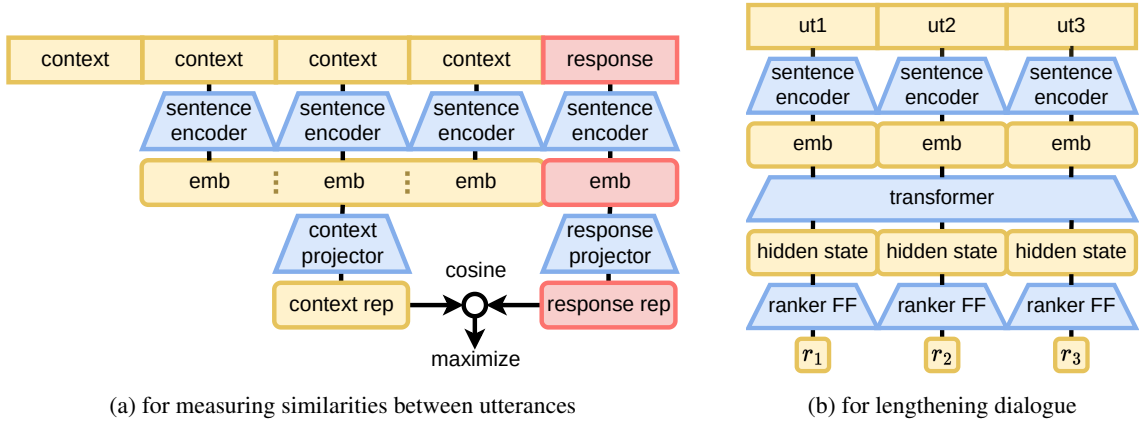


Figure 1: Auxiliary models for performing augmentations.

### B.1 Pairwise Model

We utilized a model for measuring the similarity between dialog utterances. In its basic form, it can be implemented as follows: take sentence embeddings of the utterances and compare the cosine similarities between them. Training such a model on sequential utterances using contrastive learning yields commendable utterance embeddings [14].

However, this approach has a significant drawback; it does not consider the context from several preceding utterances. In a dialogue, it is crucial to compare not just pairs of utterances but pairs of context-response. Therefore, we employed the following model (fig. 1a):

1. First, embeddings for all context utterances  $c = [u_1, \dots, u_k]$  and the response  $r$  are obtained using a pretrained sentence encoder.

2. The embeddings of the context are concatenated and passed to a projector that outputs a vector representing the context.
3. The response embedding is fed into a second encoder, resulting in a vector representing the response.
4. The cosine similarity between the obtained vectors is computed as a measure of the context and response similarity.

aws-ai/dse-bert-large model from hugging face [32] was used as the sentence encoder. The model was trained with a contrastive loss using in-batch negative sampling, with the following parameters: batch size is 128, temperature is 0.05, context size is 3, projection size is 256. Only 3 last layers of sentence encoder were fine-tuned in order to decrease computational cost. Trained model reaches 0.955 retrieval accuracy@5.

Batches were formed from "context-response" pairs from the entire dialog dataset, where negative examples were not samples from the same dialog but entirely random examples from the dataset. This allows batching of arbitrary sizes, not limited to the dialog size, making the pre-training task more challenging.

The resulting model closely resembles the ConveRT model [13] for obtaining utterance embeddings. The drawbacks of the latter model are, firstly, that it is proprietary, and secondly, its architecture is highly specific and does not utilize the familiar BERT-like backbone.

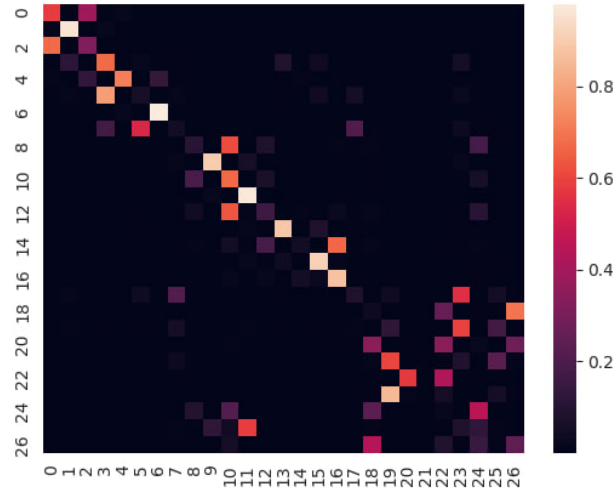


Figure 2: All similarities between contexts and responses within a dialogue. It is easy to see, that consecutive utterances form clusters.

As a result, the obtained model is capable of recognizing individual stages in dialogues (fig. 2). This behavior is achieved due to two factors. Firstly, when a dialogue furthers to a new topic, the similarity between two consecutive utterances drops substantially. This is clearly visible, for example, in dialogues in which ordering a taxi replaces booking a table in a restaurant. Secondly, within one topic, there are also small drops in similarity in cases where there is a transition from one question to another. For example, the question “how many people should I book for?” replaces the question “which restaurant do you prefer?”. Moreover, since these questions relate to one topic, they remain close to themselves and distant to questions on other topics. Therefore, clusters are obtained.

(!provide the dialogue and change pic to vector instead of raster!)

## B.2 Listwise Model

We trained the special model to merge utterances of two different dialogues (fig. 1b). It is a transformer over text embeddings of the utterances. Sentence encoder is sentence-transformers/all-mpnet-base-v2 from hugging face. We used 4-layer transformer with 4 attention heads and hidden dimension twice smaller than sentence encoder’s one, i.e. 384. Only 3 last layers of sentence encoder were fine-tuned in order to decrease computational costs.

Output ranks are transformed with softmax function. Then KL-divergence between output and target probabilities are minimized. Target probabilities are defined as softmax over true ranks of utterances, i.e.  $-i$  for  $i$ -th utterance in dialogue.



Resulting model trains to "sort" given utterances. Thanks to the attention mechanism of transformers, this can be viewed as asking the children at physical education class to look at each other and line up by height.

To measure the sorting quality, we need to utilize appropriate metric. All traditional ranking metrics such as nDCG are designed to compare with gold ranks, not just sorting quality. So during validation of our model, we were converting the ranks to a permutation over the original sequence of  $n$  elements. Then, we calculated the number of transpositions. It is easy to implement and can be normalized by maximum possible number of transpositions  $n(n-1)/2$ , resulting in  $[0, 1]$ -ranged metric. Our trained model reaches 0.96 value.

## C Composition of Augmentations

In order to maximize diversity of training data, we use not only 5 basic augmentations described in section 4.1, but also 4 extra compositions of augmentations. All the resulting pipelines are defined in fig. 3.

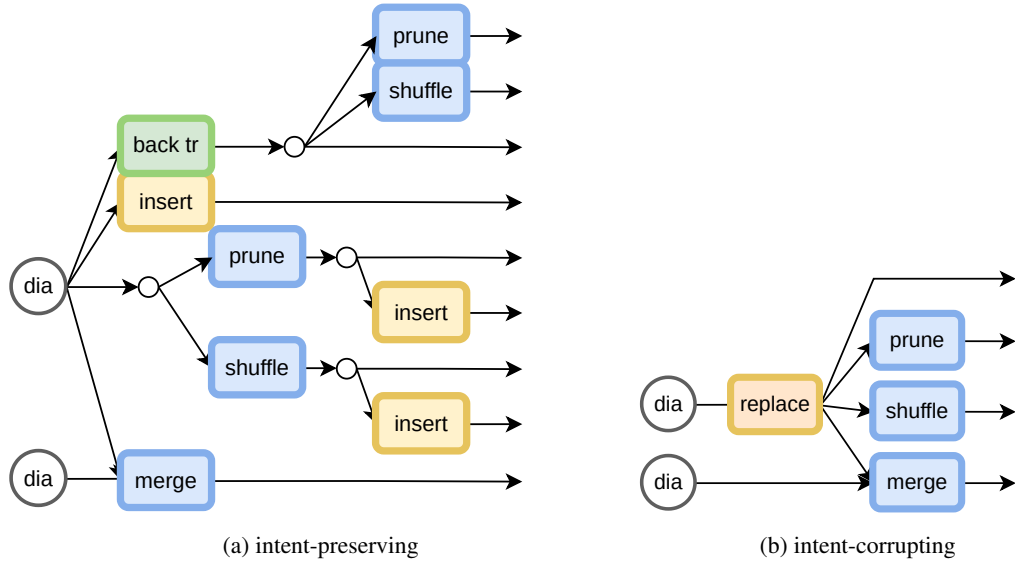


Figure 3: Compositions of augmentations.