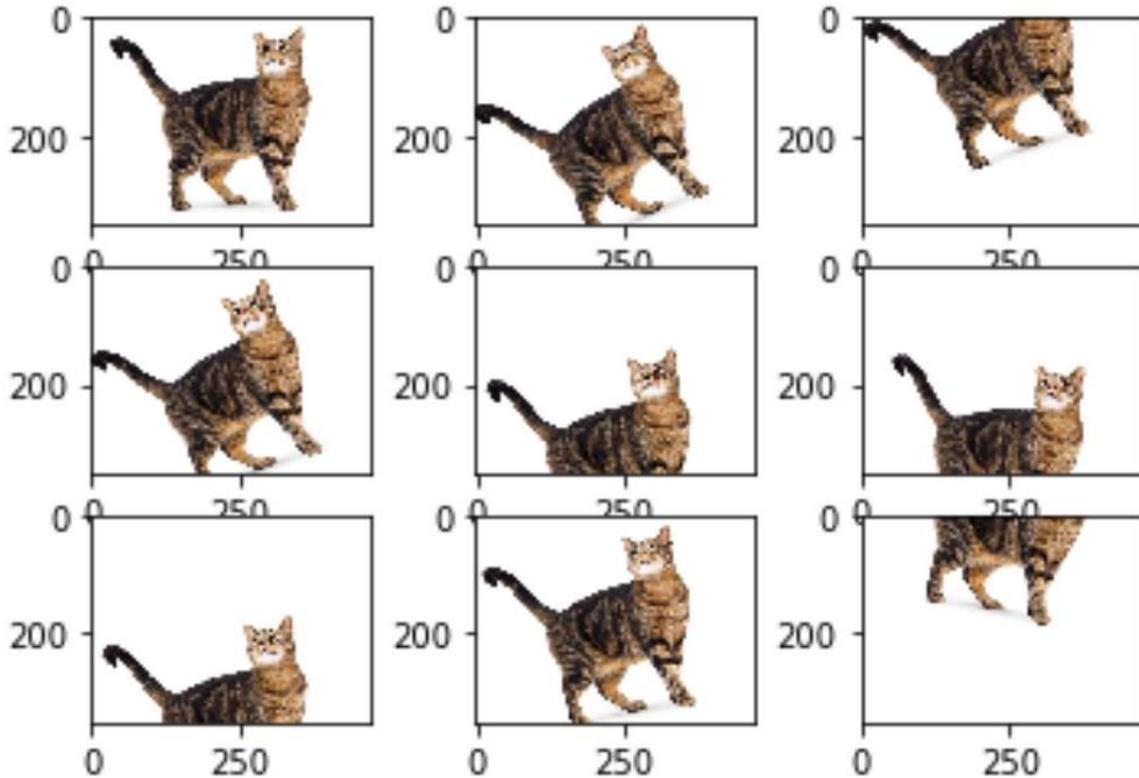


# Data Augmentation for NLP

Dominik Lewy

## What data augmentation is?

**Data augmentation (DA)** refers to strategies for increasing the diversity of training examples without explicitly collecting new data.

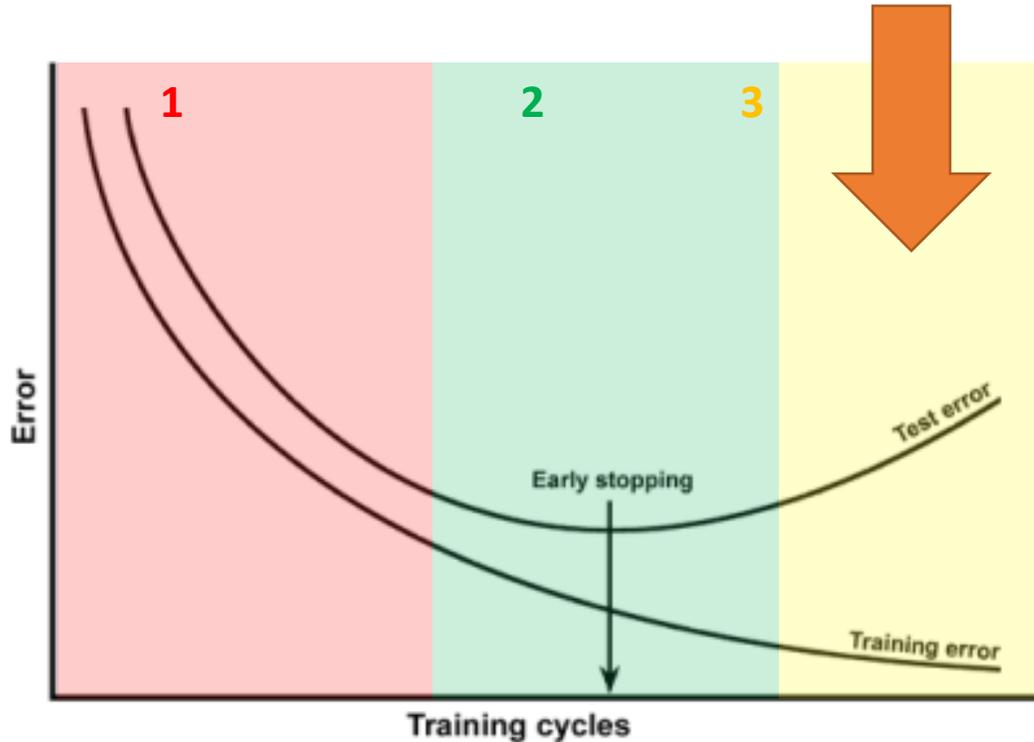


Nearest neighbors in word2vec



## Why do we use it?

### Train/Test error curves



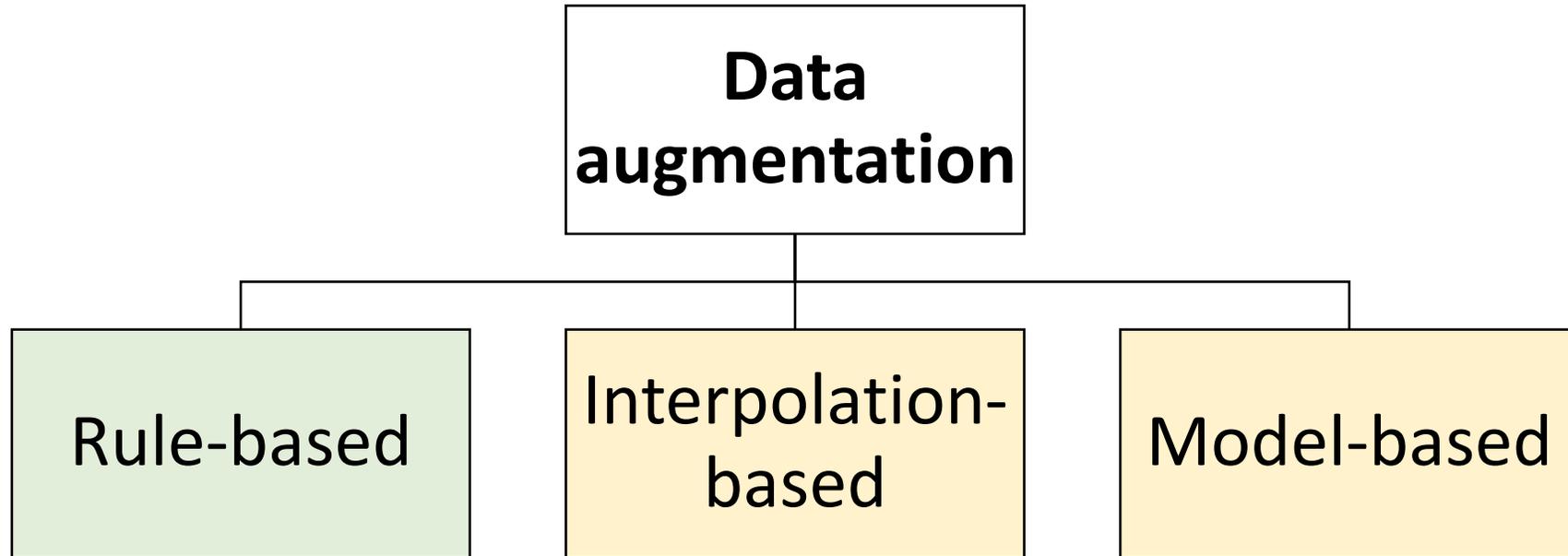
Train Error	Test Error	Strategy
High	High	Increase the capacity/complexity of the model
Low	Optimal	OK
Low	Much worse than train	Data augmentation, Regularization (weight decay, dropout)

When discussing generalization and overfitting three scenarios arise:

1. Model trained excludes the true data-generating process – this corresponds to underfitting and induces bias
2. Model matched the true data-generating process
3. Model included the data-generating process but also many other possible generating processes – this corresponds to overfitting, in such scenario variance rather than bias dominates the error

The goal of regularization is to take the model from third scenario to the second.

How can we categorize data augmentation methods in NLP space?



Copy is modified

Synthetic sample is created

**Example of a rule-based method**

**EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks**

**Jason Wei<sup>1,2</sup> Kai Zou<sup>3</sup>**

<sup>1</sup>Protago Labs Research, Tysons Corner, Virginia, USA

<sup>2</sup>Department of Computer Science, Dartmouth College

<sup>3</sup>Department of Mathematics and Statistics, Georgetown University

`jason.20@dartmouth.edu`    `kz56@georgetown.edu`

## EDA – Easy Data Augmentation

Operation	Sentence
None	A sad, superior human comedy played out on the back roads of life.
SR	A <i>lamentable</i> , superior human comedy played out on the <i>backward</i> road of life.
RI	A sad, superior human comedy played out on <i>funniness</i> the back roads of life.
RS	A sad, superior human comedy played out on <i>roads</i> back <i>the</i> of life.
RD	A sad, superior human out on the roads of life.

Table 1: Sentences generated using EDA. SR: synonym replacement. RI: random insertion. RS: random swap. RD: random deletion.

## EDA – Easy Data Augmentation

### Data sets include:

- SST-2 – Stanford Sentiment Treebank
- CR – customer reviews
- SUBJ – subjective/objective dataset
- TREC – question type dataset
- PC – Pro-Con dataset

<b>Model</b>	<b>Training Set Size</b>			
	500	2,000	5,000	full set
RNN	75.3	83.7	86.1	87.4
+EDA	79.1	84.4	87.3	88.3
CNN	78.6	85.6	87.7	88.3
+EDA	80.7	86.4	88.3	88.8
<i>Average</i>	76.9	84.6	86.9	87.8
+EDA	79.9	85.4	87.8	<b>88.6</b>

### Architectures include:

- RNN – LSTM-RNN
- CNN

Table 2: Average performances (%) across five text classification tasks for models with and without EDA on different training set sizes.

## EDA – Easy Data Augmentation

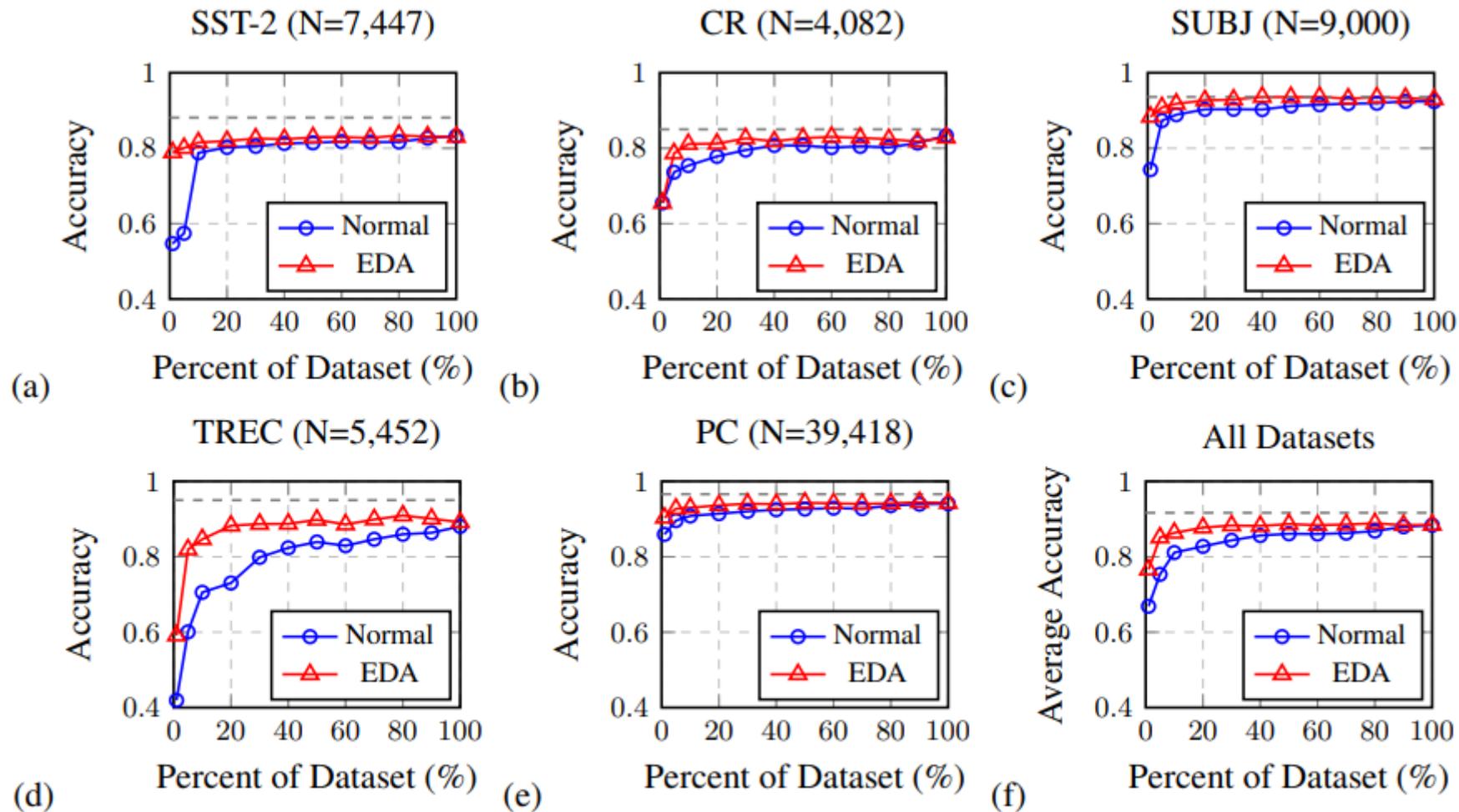


Figure 1: Performance on benchmark text classification tasks with and without EDA, for various dataset sizes used for training. For reference, the dotted grey line indicates best performances from Kim (2014) for SST-2, CR, SUBJ, and TREC, and Ganapathibhotla (2008) for PC.

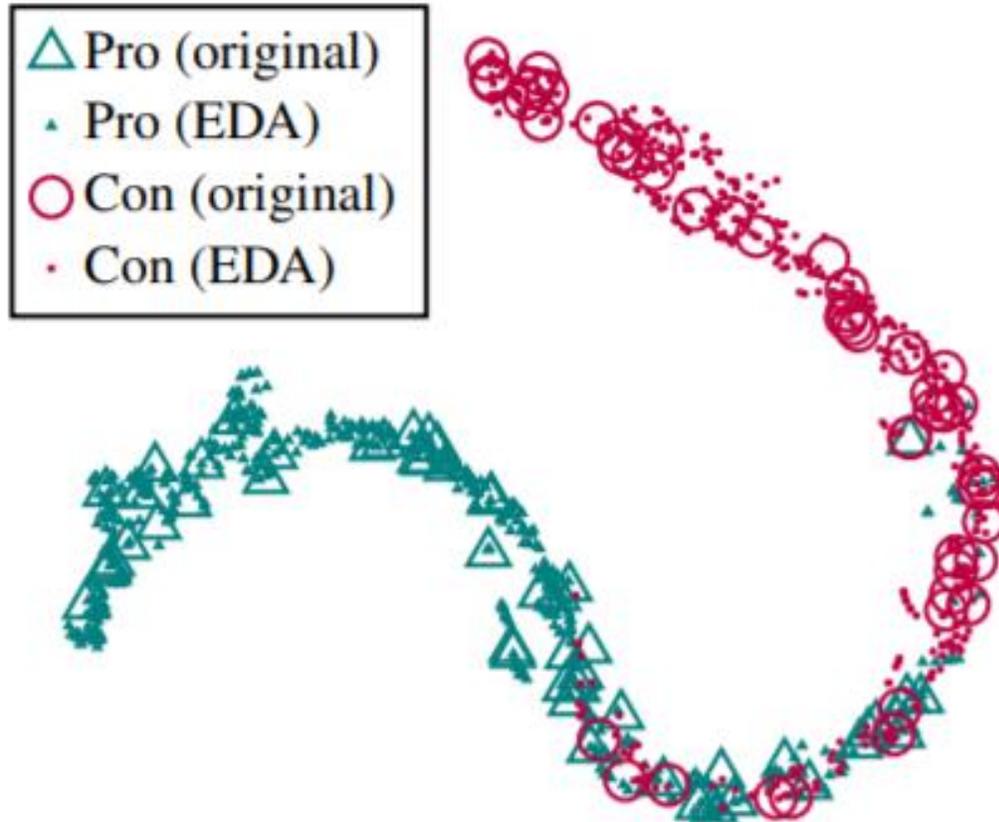
**EDA – Easy Data Augmentation**

Figure 2: Latent space visualization of original and augmented sentences in the Pro-Con dataset. Augmented sentences (small triangles and circles) closely surround original sentences (big triangles and circles) of the same color, suggesting that augmented sentences maintained their true class labels.

## EDA – Easy Data Augmentation

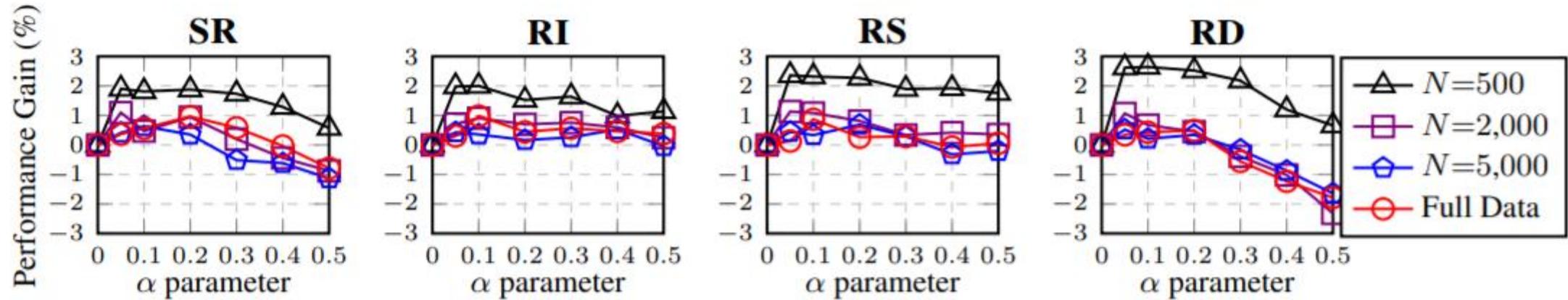


Figure 3: Average performance gain of EDA operations over five text classification tasks for different training set sizes. The  $\alpha$  parameter roughly means “percent of words in sentence changed by each augmentation.” SR: synonym replacement. RI: random insertion. RS: random swap. RD: random deletion.

**Example of an interpolation-based method**

# **Leveraging BERT with Mixup for Sentence Classification (Student Abstract)**

**Amit Jindal,<sup>1</sup> Dwaraknath Gnaneshwar,<sup>1</sup> Ramit Sawhney,<sup>2</sup> Rajiv Ratn Shah<sup>3</sup>**

<sup>1</sup>Manipal Institute of Technology  
{amitj646, dwarakasharma} @gmail.com

<sup>2</sup>Netaji Subhas Institute of Technology  
ramits.co@nsit.net.in

<sup>3</sup>Indraprastha Institute of Information Technology, Delhi  
rajivratn@iiitd.ac.in

## Data Augmentation via Mixing Images

Image 1 -  
label: dog



Image 2 -  
label: cow



Mixup -  
label: (dog:0.3, cow:0.7)



SamplePairing -  
label: dog



CutMix -  
label: (dog:0.7, cow:0.3)



SmoothMix -  
label: (dog:0.3, cow:0.7)



AttentiveMix -  
label: (dog:0.14, cow:0.86)



RandomSquare -  
label: (dog:0.42, cow:0.58)

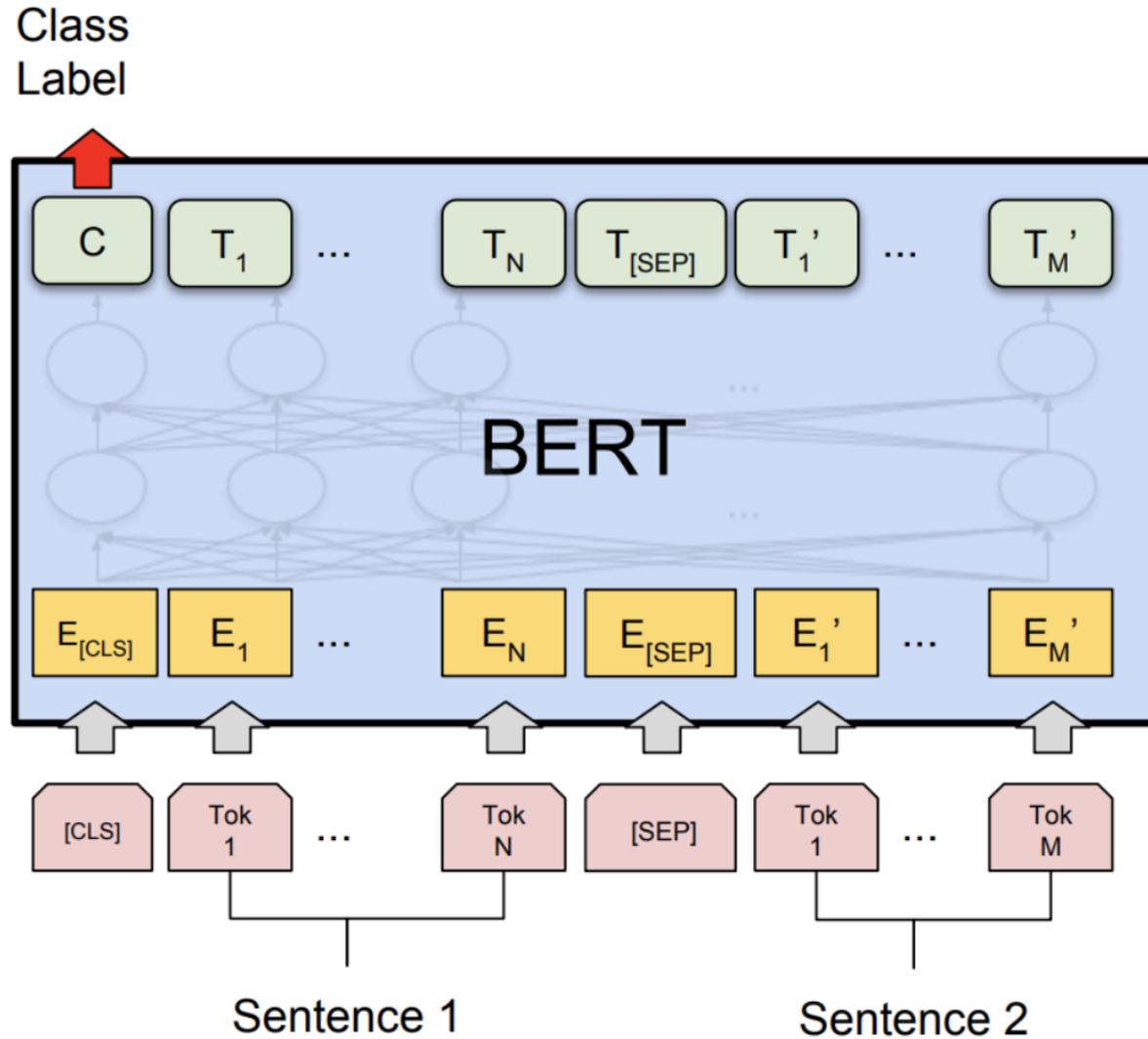


## Mixup for Sentence Classification

**Sentence A:**      *The quick brown fox jumps over the lazy dog.*

**Sentence B:**      *Pack my box with five dozen liquor jugs.*

## Mixup for Sentence Classification



## Mixup for Sentence Classification

Experiments error rates					
Model	IMDB	SST-1	MR	TREC	SUBJ
CNN	7.67	55.0	21.61	8.8	7.59
LSTM	12.85	53.6	21.67	13.5	8.60
BERT	6.46	46.75	<b>10.83</b>	2.62	2.20
BERT + Input Mixup	6.17	45.55	12.81	2.41	1.505
BERT + Manifold Mixup	<b>6.02</b>	<b>44.20</b>	11.94	<b>2.20</b>	<b>1.501</b>

Table 1: Test error (%) of the testing methods using BERT. Best results highlighted in Bold.

### Data sets include:

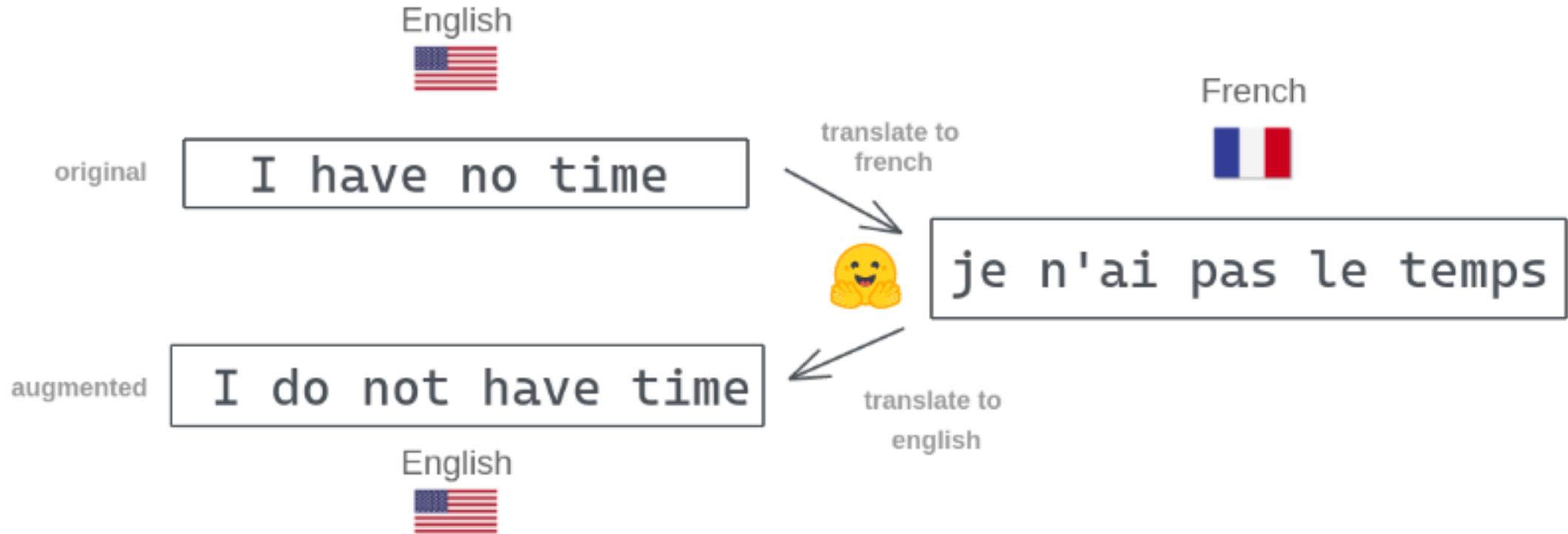
- IMDB – binary sentiment classification dataset
- MR – movie review data set with binary sentiment
- SST-1 – Stanford Sentiment Treebank
- TREC – question type dataset
- SUBJ – subjective/objective dataset

## Mixup for Sentence Classification

Table 2: Test error (%) Manifold Mixup for different sets of eligible layers  $S$  on IMDB

$S$	IMDB
$\{0\}$	6.17
$\{0, 1\}$	6.10
$\{0, 1, 2\}$	6.27
$\{0, 1, 2, 3\}$	6.21
$\{0, 1, 2, 3, 4\}$	6.15
$\{0, 1, 2, 3, 4, 5\}$	6.22
$\{0, 1, 2, 3, 4, 5, 6\}$	6.09
$\{0, 1, 2, 3, 4, 5, 6, 7\}$	<b>6.02</b>
$\{0, 1, 2, 3, 4, 5, 6, 7, 8\}$	6.23
$\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$	6.25
$\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$	6.28
$\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11\}$	6.40

Example of a model-based method



## Where is augmentation mostly used now?

- Low-Resource Languages
- Mitigating Bias
- Fixing Class Imbalance
- Few-Shot Learning

## What is the array of tasks in NLP to be considered?

- Sentence Classification
- Summarization
- Question Answering
- Sequence Tagging
- Parsing
- Grammatical Error Correction
- Neural Machine Translation
- Data to text (NLG)
- ....

**Sequence Tagging – rule-based equivalent**

**An Analysis of Simple Data Augmentation for Named Entity Recognition**

**Xiang Dai<sup>1,2,3</sup> Heike Adel<sup>1</sup>**

<sup>1</sup>Bosch Center for Artificial Intelligence, Renningen, Germany

<sup>2</sup>University of Sydney, Sydney, Australia

<sup>3</sup>CSIRO Data61, Sydney, Australia

dai.dai@csiro.au heike.adel@de.bosch.com

## Sequence Tagging – rule-based equivalent

	Instance												
None	She O	did O	not O	complain O	of O	headache B-problem	or O	any B-problem	other I-problem	neurological I-problem	symptoms I-problem	. O	
LwTR	L. O	One O	not O	complain O	of O	headache B-problem	he O	any B-problem	interatrial I-problem	neurological I-problem	current I-problem	. O	
SR	She O	did O	non O	complain O	of O	headache B-problem	or O	whatsoever B-problem	former I-problem	neurologic I-problem	symptom I-problem	. O	
MR	She O	did O	not O	complain O	of O	neuropathic B-problem	pain I-problem	syndrome I-problem	or O	acute B-problem	pulmonary I-problem	disease I-problem	. O
SiS	not O	complain O	She O	did O	of O	headache B-problem	or O	neurological B-problem	any I-problem	symptoms I-problem	other I-problem	. O	

Table 1: Original training instance and different types of augmented instances. We highlight changes using blue color. Note that LwTR (Label-wise token replacement) and SiS (Shuffle within segments) change token sequence only, whereas SR (Synonym replacement) and MR (Mention replacement) may also change the label sequence.

**Sequence Tagging – interpolation-based equivalent**

**SeqMix: Augmenting Active Sequence Labeling via Sequence Mixup**

**Rongzhi Zhang**  
Georgia Tech

`rongzhi.zhang@gatech.edu`

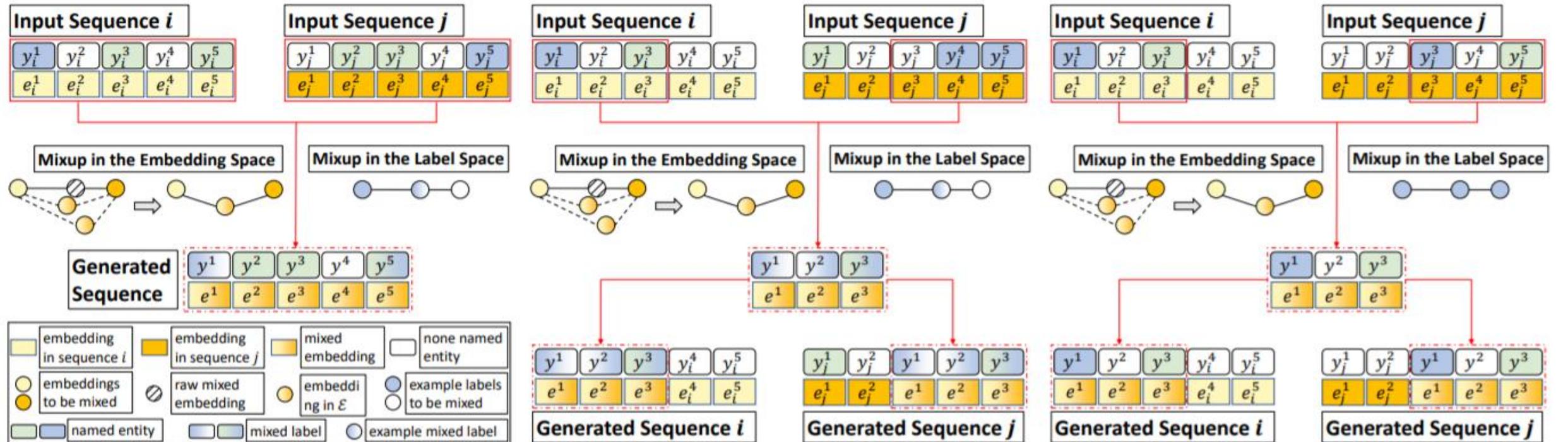
**Yue Yu**  
Georgia Tech

`yueyu@gatech.edu`

**Chao Zhang**  
Georgia Tech

`chaozhang@gatech.edu`

Sequence Tagging – interpolation-based equivalent



(a) Whole sequence mixup

(b) Sub-sequence mixup

(c) Label-constrained sub-sequence mixup

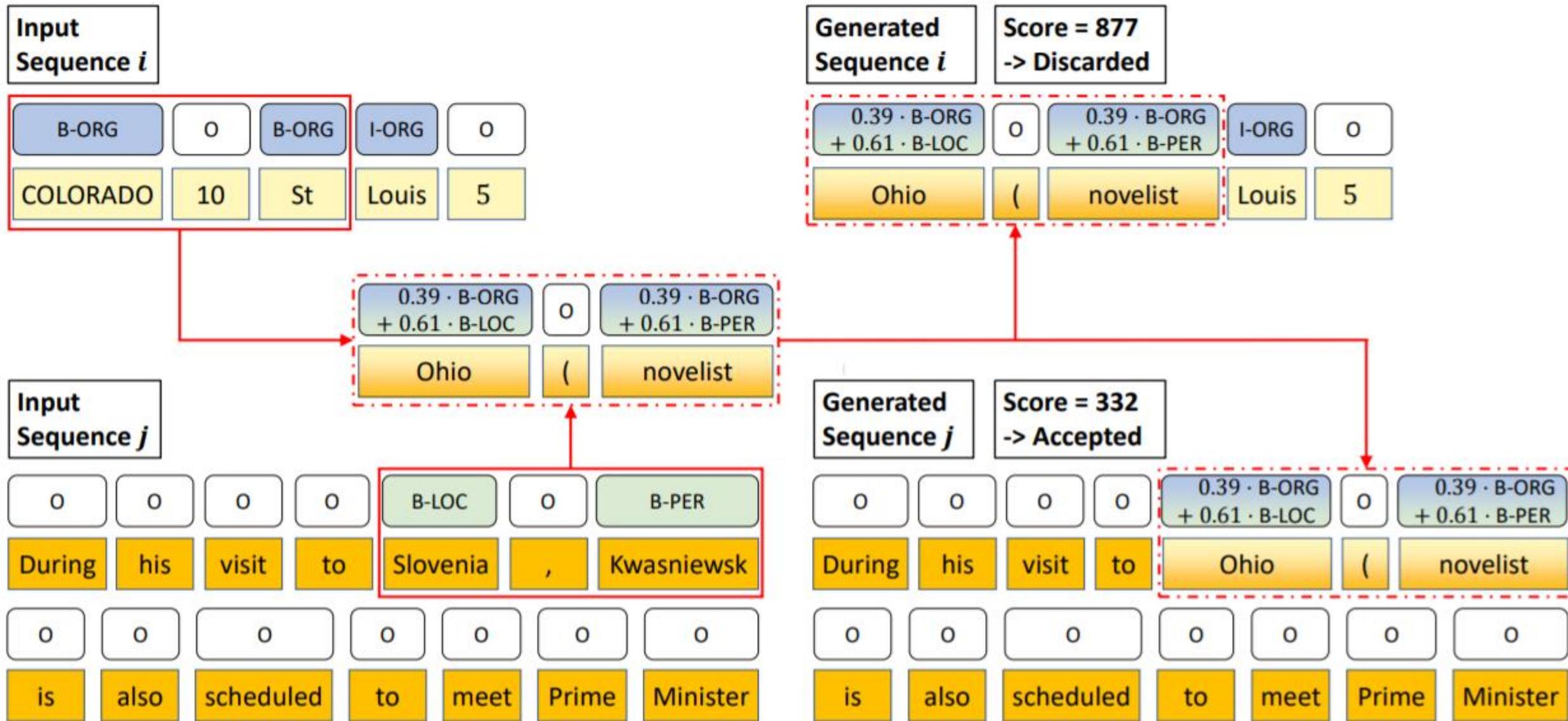


Figure 5: A generation case of sub-sequence mixup.

## Other methods

DA Method	Ext.Know	Pretrained	Preprocess	Level	Task-Agnostic
SYNONYM REPLACEMENT (Zhang et al., 2015)	✓	×	tok	Input	✓
RANDOM DELETION (Wei and Zou, 2019)	×	×	tok	Input	✓
RANDOM SWAP (Wei and Zou, 2019)	×	×	tok	Input	✓
BACKTRANSLATION (Sennrich et al., 2016)	×	✓	Depends	Input	✓
SCPN (Wieting and Gimpel, 2017)	×	✓	const	Input	✓
SEMANTIC TEXT EXCHANGE (Feng et al., 2019)	×	✓	const	Input	✓
CONTEXTUALAUG (Kobayashi, 2018)	×	✓	-	Input	✓
LAMBADA (Anaby-Tavor et al., 2020)	×	✓	-	Input	×
GECA (Andreas, 2020)	×	×	tok	Input	×
SEQMIXUP (Guo et al., 2020)	×	×	tok	Input	×
SWITCHOUT (Wang et al., 2018b)	×	×	tok	Input	×
EMIX (Jindal et al., 2020a)	×	×	-	Emb/Hidden	✓
SPEECHMIX (Jindal et al., 2020b)	×	×	-	Emb/Hidden	Speech/Audio
MIXTEXT (Chen et al., 2020c)	×	×	-	Emb/Hidden	✓
SIGNEDGRAPH (Chen et al., 2020b)	×	×	-	Input	×
DTREEMORPH (Şahin and Steedman, 2018)	×	×	dep	Input	✓
$Sub^2$ (Shi et al., 2021)	×	×	dep	Input	Substructural
DAGA (Ding et al., 2020)	×	×	tok	Input+Label	×
WN-HYPERS (Feng et al., 2020)	✓	×	const+KWE	Input	✓
SYNTHETIC NOISE (Feng et al., 2020)	×	×	tok	Input	✓
UEDIN-MS (DA part) (Grundkiewicz et al., 2019)	✓	×	tok	Input	✓
NONCE (Gulordava et al., 2018)	✓	×	const	Input	✓
XLDA (Singh et al., 2019)	×	✓	Depends	Input	✓
SEQMIX (Zhang et al., 2020)	×	✓	tok	Input+Label	×
SLOT-SUB-LM (Louvan and Magnini, 2020)	×	✓	tok	Input	✓
UBT & TBT (Vaibhav et al., 2019)	×	✓	Depends	Input	✓
SOFT CONTEXTUAL DA (Gao et al., 2019)	×	✓	tok	Emb/Hidden	✓
DATA DIVERSIFICATION (Nguyen et al., 2020)	×	✓	Depends	Input	✓
DIPS (Kumar et al., 2019a)	×	✓	tok	Input	✓
AUGMENTED SBERT (Thakur et al., 2021)	×	✓	-	Input+Label	Sentence Pairs

The end.  
Thank you!

# Discussion