

# Effectiveness of Data Augmentation and Pretraining for Improving Neural Headline Generation in Low-Resource Settings

Matej Martinc<sup>1</sup>, Syrielle Montariol<sup>2</sup>, Lidia Pivovarova<sup>3</sup>, Elaine Zosa<sup>3</sup>

<sup>1</sup>Jozef Stefan Institute, <sup>2</sup>INRIA Paris, <sup>3</sup>University of Helsinki

## Headline Generation

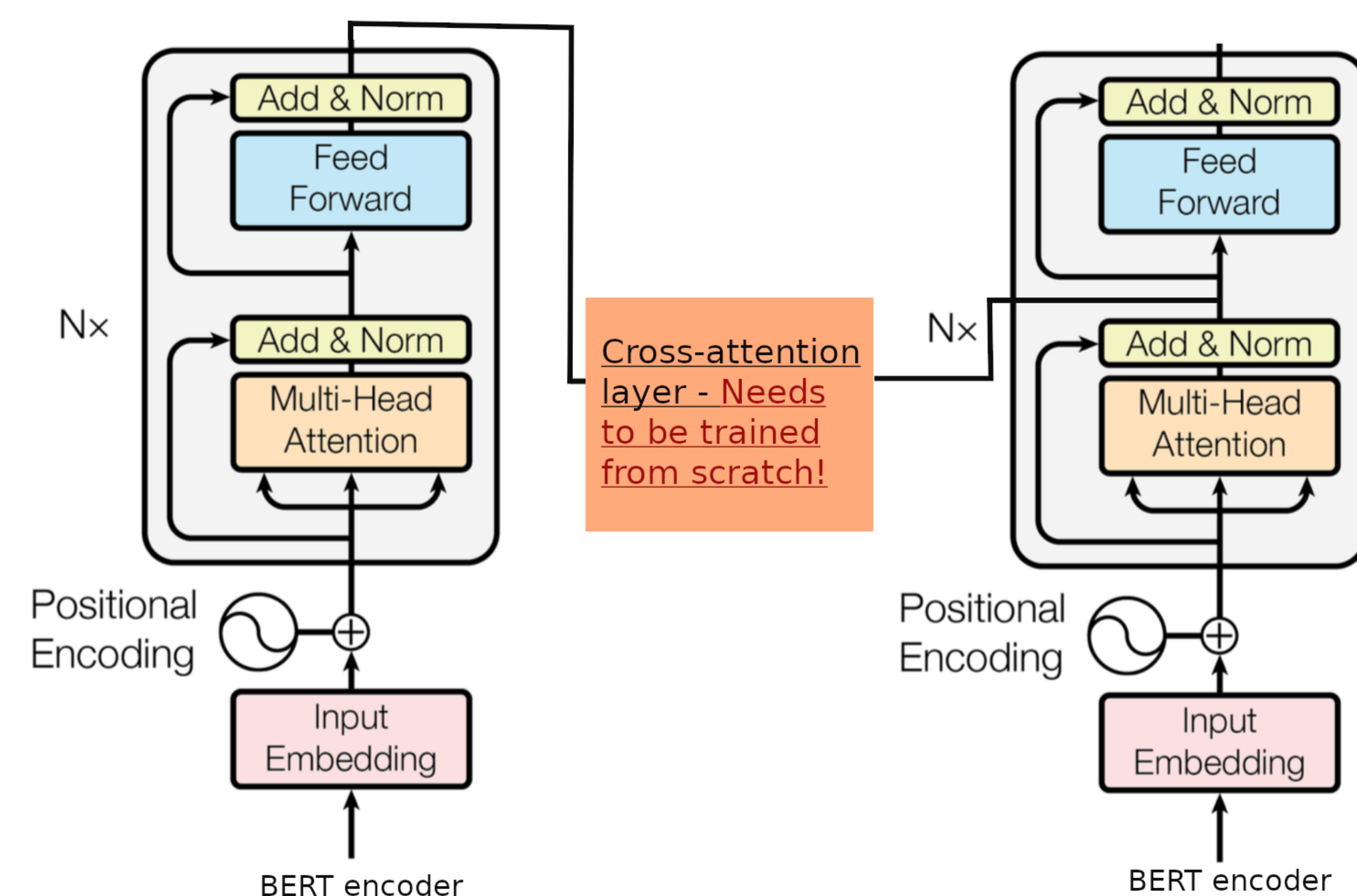
- Use Transformer Encoder-Decoder architectures, popular for **machine translation** and **summarization**.
- Headlines are a vehicle that carries the most important information about the event or topic described in the news article ⇒ **Headline generation as a summary**

## Two SOTA summarization systems

### BERT-ED

#### BERT encoder-decoder

- Combination of any pretrained BERT language models.
- **Can be employed for 168 languages.**



### BART

Transformer-based model, pre-trained on denoising tasks, sentence shuffling and text infilling.

- **Multilingual version supports 50 languages.**
- **Potential problem: curse of multilinguality?**

## Pre-training

- **Sentence shuffling:** Randomly shuffle sentences, train the model to generate the original text with the correct sentence order.
- **Text infilling:** Corrupt 20% of the training corpus by an infilling scheme (spans of text are replaced with a [MASK] token), then train the encoder-decoder to generate the original text.
- **2 tasks:** Train the model with Sentence shuffling then Text infilling.

## Data Augmentation

- **BERT:** Mask 20% of words. Feed the masked article to BERT, which proposes candidate words for the masked tokens.
- **Word2Vec:** Replace random words by synonyms proposed by Word2Vec.
- **Wordnet:** Same as above but using Wordnet.
- **EDA:** Perform synonym replacement, random insertion, random swap, and random deletion.
- **Mixed:** Successively perform Word2Vec, EDA and Wordnet augmentation.

For each article in the train set, generate 5 augmented articles.

## Evaluation metrics

- **ROUGE 1, 2, L:** Current standard for evaluating summarization. **Agnostic to semantic similarity**, hence can have low correlation with human judges.
- **Semantic similarity (SS):** measures cosine distance between the sentence embedding of the true and generated headline.
- **Natural language inference (NLI):** consider the true headline as the 'premise' and the generated headline as the 'hypothesis'; use an NLI model to predict entailment, as a measure of generated headline quality. **Only used for English experiments!**

## Results

Approach		ROUGE-1	ROUGE-2	ROUGE-L	SS	NLI					
<i>English BERT-ED-based models</i>											
BASELINE	10K	10.2	1.4	9.6	24.6	15.4					
	260K	27.6	10.1	25.1	49.6	32.1					
AUGMENTATION	bert	13.2	<b>3.0</b>	2.3	0.9	12.2	2.6	30.8	6.2	15.7	<b>0.3</b>
	w2v	9.7	<b>-0.5</b>	1.6	0.2	8.9	<b>-0.7</b>	26.5	1.9	14.9	<b>-0.5</b>
	mix	10.4	<b>0.2</b>	1.7	<b>0.3</b>	9.6	<b>0.0</b>	23.9	<b>-0.7</b>	13.1	<b>-2.3</b>
PRETRAINING	2 tasks	<b>16.5</b>	<b>6.3</b>	<b>4.6</b>	<b>3.2</b>	<b>15.1</b>	<b>5.5</b>	<b>42.0</b>	<b>17.4</b>	<b>25.9</b>	<b>10.5</b>
<i>English BART-based models</i>											
BASELINE	10K	<b>29.0</b>	<b>10.9</b>	<b>26.0</b>	49.3	34.1					
	260K	31.9	13.1	28.7	<b>51.7</b>	<b>36.8</b>					
AUGMENTATION	bert	28.5	<b>-0.5</b>	10.5	<b>-0.4</b>	25.6	<b>-0.4</b>	49.1	<b>-0.2</b>	34.0	<b>-0.1</b>
	w2v	27.8	<b>-1.2</b>	10.1	<b>-0.8</b>	25.1	<b>-0.9</b>	48.2	<b>-1.1</b>	32.0	<b>-2.1</b>
	mix	27.7	<b>-1.3</b>	10.2	<b>-0.7</b>	25.0	<b>-1.0</b>	47.9	<b>-1.4</b>	32.2	<b>-1.9</b>
PRETRAINING	2 tasks	28.7	<b>-0.3</b>	10.7	<b>-0.2</b>	25.9	<b>-0.1</b>	49.2	<b>-0.1</b>	34.1	<b>0.0</b>

Table: Results on the English and Estonian datasets. Best results in a low resource setting (i.e., excluding the BART and BERT-ED models trained on English 260K dataset) per evaluation measure and language are **bolded**. For each measure, we report its absolute value (1st number) and the difference with the baseline model (2nd, colored number). Since all experiments with data augmentation and pretraining are run on the 10K dataset, differences are computed respectively to the 10K baseline, the first row for each model.

## Examples

True headline	EXAMPLE
BART 260K	fighting n. y. c. soda ban, industry focuses on personal choice
BART 10K	soda industry fights new york city's soda ban
BERT-ED 260K	soft-drink industry takes aim at sugary drinks
BERT-ED 10K	soft - drink industry seeks to fight sugary drinks ban on sugary drinks
BERT-ED 10K + BERT aug	in new york's york taxes's taxes s
BERT-ED 10K + shuffling	in new york city, new york city's new york city's bans law
BERT-ED 10k + infilling	u. s. and new york's new new york city mayor' campaign campaign moves new york's mayor's campaign campaign philippines : s. o. p. to be suspended s. a. lawmakers s. ban s. a.'s
BERT-ED 10K + 2 tasks	new yorkers face a challenge to soda industry in new yorkers in new yorkers' campaign campaign in new york city's

Table: Examples of generated English headlines.

- BART 10k produces results very similar to BART 260k but can hallucinate.
- Performance of BERT-ED-based model trained on 10K drops compared to the 260K dataset.
- Data augmentation only slightly improves the performance on English.
- Pretraining on 2 tasks works the best, resulting in longer headlines with meaningful beginnings.

Approach		ROUGE-1	ROUGE-2	ROUGE-L	SS				
<i>Estonian BERT-ED-based models</i>									
BASELINE		3.9	0.3	3.8	17.9				
AUGMENTATION	bert	9.8	<b>5.9</b>	2.5	2.2	9.4	5.6	36.9	19.0
	w2v	8.5	<b>4.6</b>	2.1	1.8	8.1	4.3	34.4	16.5
PRETRAINING	infilling	13.9	<b>0.1</b>	4.3	4.0	13.2	9.4	44.0	26.1
	shuffling	11.3	<b>7.4</b>	2.8	2.5	10.7	6.9	40.7	22.8
	2 tasks	<b>17.6</b>	<b>13.7</b>	<b>6.5</b>	<b>6.2</b>	<b>16.3</b>	<b>12.5</b>	<b>49.8</b>	<b>31.9</b>
<i>Estonian BART-based models</i>									
BASELINE		26.2	12.3	24.4	56.7				
AUGMENTATION	bert	25.4	<b>-0.8</b>	11.6	<b>-0.7</b>	23.8	<b>-0.6</b>	55.9	<b>-0.8</b>
	w2v	23.0	<b>-3.2</b>	9.8	<b>-2.5</b>	21.5	<b>-2.9</b>	53.5	<b>-3.2</b>
PRETRAINING	infilling	<b>27.1</b>	<b>0.9</b>	<b>12.9</b>	<b>0.6</b>	<b>25.2</b>	<b>0.8</b>	<b>57.2</b>	<b>0.5</b>
	shuffling	26.6	<b>0.4</b>	12.6	<b>0.3</b>	24.8	<b>0.4</b>	56.9	<b>0.2</b>
	2 tasks	26.6	<b>0.4</b>	12.3	<b>0.0</b>	24.6	<b>0.2</b>	56.6	<b>-0.1</b>

## Datasets

News datasets for evaluation of headline generation:

Language	train set	test set
English KPTimes 260K - <b>BASELINE</b>	259,923	10,000
English KPTimes 10K	10,000	10,000
Croatian	32,223	3,582
Estonian	10,750	7,747

## Conclusion

- Pretrained multilingual NLG model for a specific low-resource language should be picked over training the cross-attention layer from scratch.
- Pretraining and data augmentation has little effect on BART-based models.
- Pretraining > augmentation.
- SS and NLI are highly correlated w/ ROUGE.
- The paper repository:



• More details: MATEJ.MARTINC@IJS.SI