

# Введение в принципы модели BERT

---

или по верхам о модели, взорвавшей NLP

Маша Шеянова, [masha.shejanova@gmail.com](mailto:masha.shejanova@gmail.com)

# Byte Pair Encodings (BPE)

---

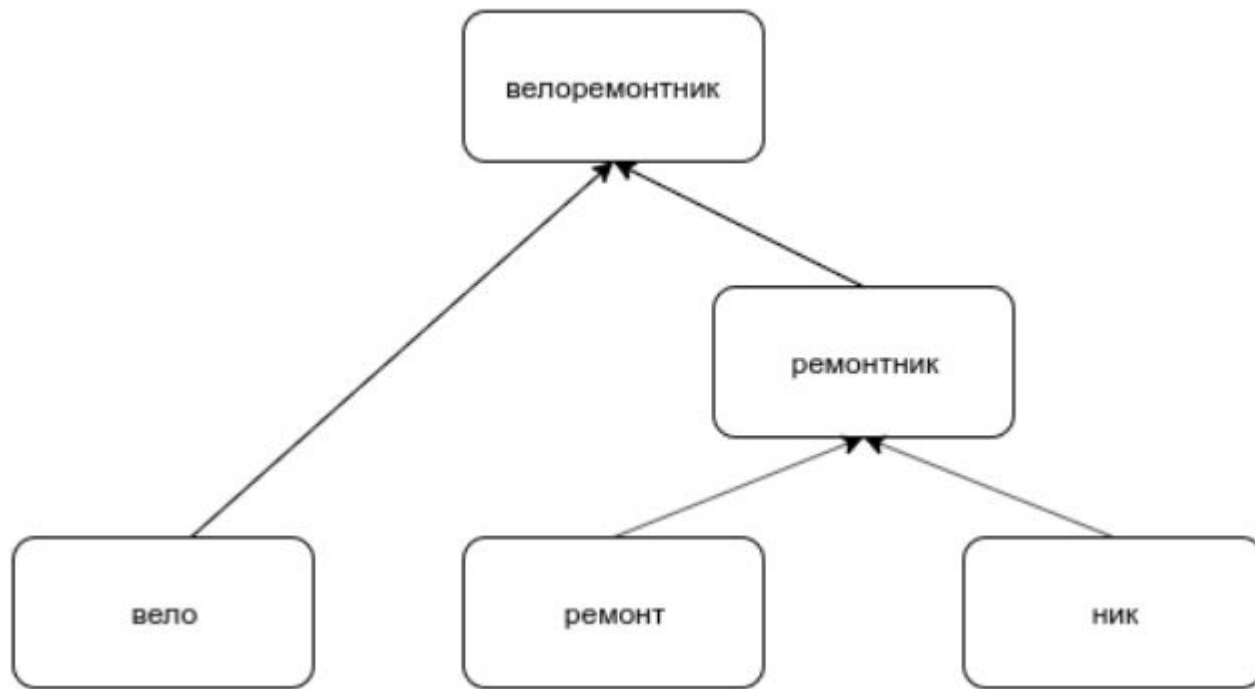
# Что такое BPE

Альтернативный способ токенизации для обучения эмбедингов слов. Чем-то похож на fasttext, но лучше.

Алгоритм:

- вначале у нас столько “токенов”, сколько не-пробельных символов
- пока мы не получили столько токенов, сколько мы хотим в итоге получить
  - находим два самых часто встречаемых друг с другом элемента
  - сливаем, их образуя новый токен
  - повторяем

# Пример



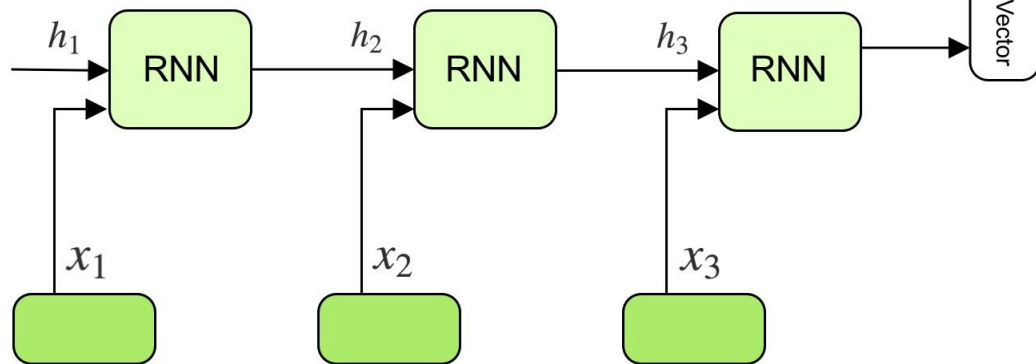
Если в обучающем корпусе не было слова *велоремонтник*, то получится *(вело, ремонтник)* или *(вело, ремонт, ник)*.

# Attention and Transformers

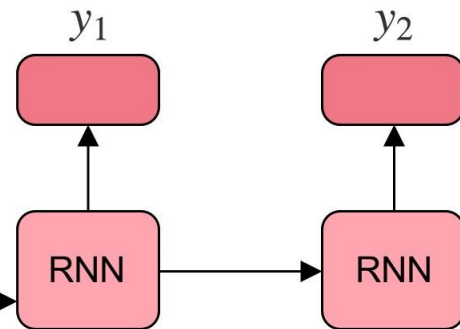
---

# seq2seq (стандартная)

Encoder



Decoder

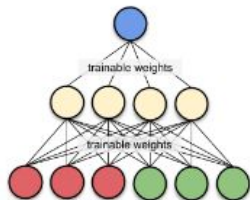


seq2seq + attention

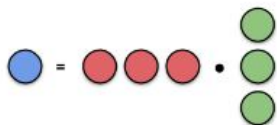
# Виды attention (ИСТОЧНИК)



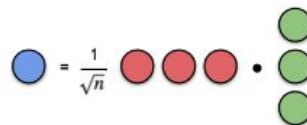
Additive / Concat



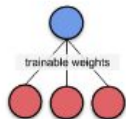
Dot product



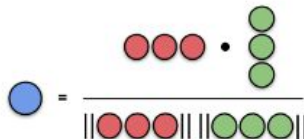
Scaled dot product



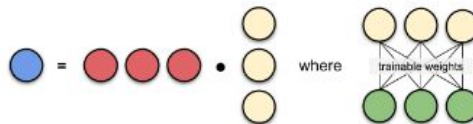
Location-based



Cosine similarity



General





# Виды attention (ИСТОЧНИК)

Name	Alignment score function	Citation
Content-base attention	$\text{score}(s_t, h_i) = \text{cosine}[s_t, h_i]$	<a href="#">Graves2014</a>
Additive(*)	$\text{score}(s_t, h_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a [s_t; h_i])$	<a href="#">Bahdanau2015</a>
Location-Base	$\alpha_{t,i} = \text{softmax}(\mathbf{W}_a s_t)$ Note: This simplifies the softmax alignment to only depend on the target position.	<a href="#">Luong2015</a>
General	$\text{score}(s_t, h_i) = s_t^\top \mathbf{W}_a h_i$ where $\mathbf{W}_a$ is a trainable weight matrix in the attention layer.	<a href="#">Luong2015</a>
Dot-Product	$\text{score}(s_t, h_i) = s_t^\top h_i$	<a href="#">Luong2015</a>
Scaled Dot-Product(^)	$\text{score}(s_t, h_i) = \frac{s_t^\top h_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	<a href="#">Vaswani2017</a>

## Трансформер (*Attention is all you need*, 2017)

Это очень большая и сложная модель. Если хотите глубоко понять устройство — переходите по ссылкам на слайде. В общих чертах:

- как и seq2seq, состоит из энкодера и декодера
- но не использует RNN, полностью заменив передачу вектора состояния на attention (и в энкодере это тоже происходит, называется *self-attention*)
- и энкодер, и декодер многослойные (в оригинальной версии, 6 слоёв), attention применяется на каждом из них

BERT — это большой энкодер из трансформера, обученный предсказывать пропущенные слова и угадывать, идёт ли одно предложение за другим.

# о применениях BERT

---

# BERT как контекстные эмбединги

BERT использует BPE и в процессе обучения создаёт векторное представление для BPE-токенов.

За счёт self-attention эмбединги каждого слова знают о его контексте.

```
print ("Similarity of 'bank' as in 'bank robber' to 'bank' as in 'bank vault':", same_bank)
```

```
Similarity of 'bank' as in 'bank robber' to 'bank' as in 'bank vault': 0.9456751
```

```
print ("Similarity of 'bank' as in 'bank robber' to 'bank' as in 'river bank':", different_bank)
```

```
Similarity of 'bank' as in 'bank robber' to 'bank' as in 'river bank': 0.6797334
```

# BERT для классификации

Кроме эмбедингов BPE-токенов, BERT обучает представление добавленного токена [CLS], который “отвечает” за весь текст.

Дальше:

- можно использовать эмбединг этого символа как эмбединг текста, и сверху добавлять свою модель
- (advanced) можно дообучать его

# Практика

---

# Использование модели BERT

- для классификации
- как эмбединги