



Deep Learning School

Механизм внимания (Attention)

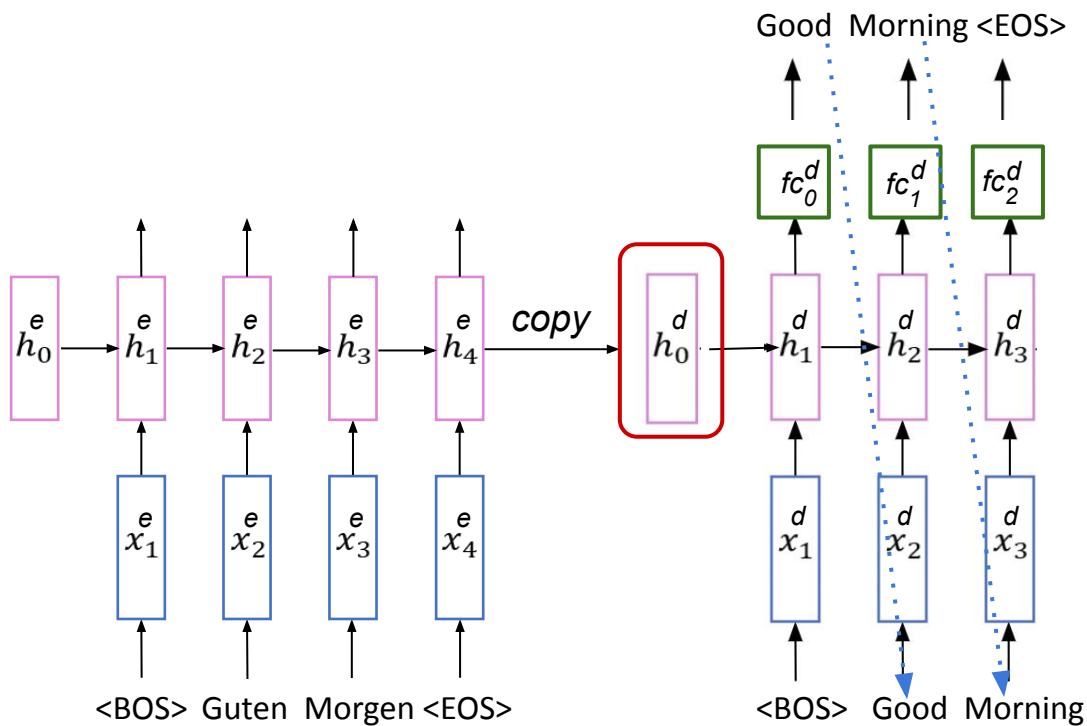
В этом видео

- Задача машинного перевода;
- Архитектура Seq2Seq для решения задачи перевода и ее обучение;
- Недостатки RNN-модели для решения задачи перевода
- **Механизм внимания (Attention)**

В этом видео

- **Механизм внимания (Attention):**
 - Механизм attention: идея;
 - Применение Attention для RNN-модели
Seq2Seq для задачи машинного перевода;
 - Разные варианты Attention

RNN для МТ: напоминание



Недостатки:

- Encoder RNN может забывать информацию из начала предложения;
- Вся информация о входном предложении содержится в одном векторе h_d^0

Attention: идея

Допустим, нам надо перевести следующее предложение на английский язык:

*Высокий статный мужчина в шляпе с прекрасным ярким
попугаем на плече только что зашел в здание на углу
Старого Косого переулка*

Attention: идея

Допустим, нам надо перевести следующее предложение на английский язык:

*Высокий статный мужчина в шляпе с прекрасным ярким
попугаем на плече только что зашел в здание на углу
Старого Косого переулка*

Как работает наша RNN-модель:

- Читает это предложение один раз и “запоминает” его в вектор скрытого состояния;
- Генерирует предложение-перевод

Attention: идея

Допустим, нам надо перевести следующее предложение на английский язык:

*Высокий статный мужчина в шляпе с прекрасным ярким
попугаем на плече только что зашел в здание на углу
Старого Косого переулка*

Что бы делали мы (люди):

- Во время генерации перевода мы бы несколько раз смотрели на разные части исходного предложения, чтобы найти те части, на которые стоит обратить внимание в данный момент времени, и лучше их запомнить

Attention: идея

Допустим, нам надо перевести следующее предложение на английский язык:

*Высокий статный мужчина в шляпе с прекрасным ярким
попугаем на плече только что зашел в здание на углу
Старого Косого переулка*

Механизм Attention — это то, что поможет декодеру модели “подглядывать” в нужные части исходного предложения во время генерации.

Attention

Идея Attention была впервые предложена в
2014-2015 году в [этой](#) статье

Published as a conference paper at ICLR 2015

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau

Jacobs University Bremen, Germany

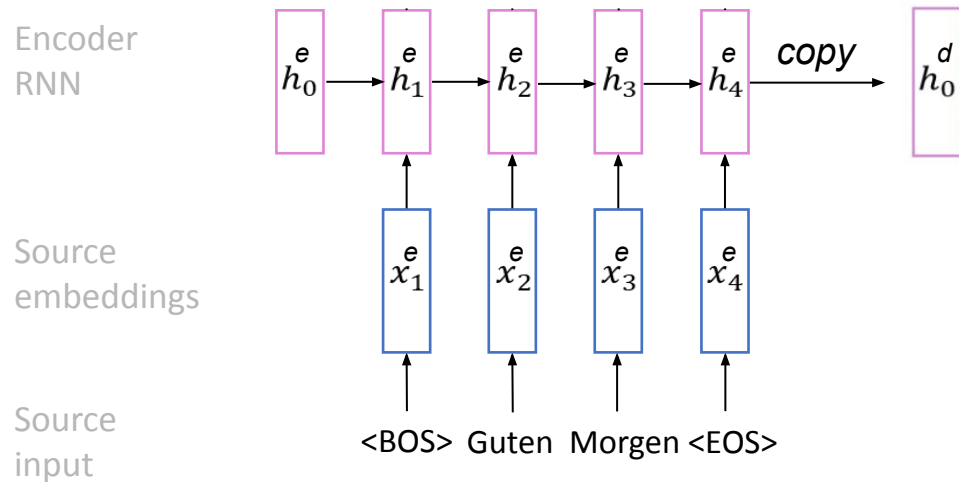
KyungHyun Cho **Yoshua Bengio***

Université de Montréal

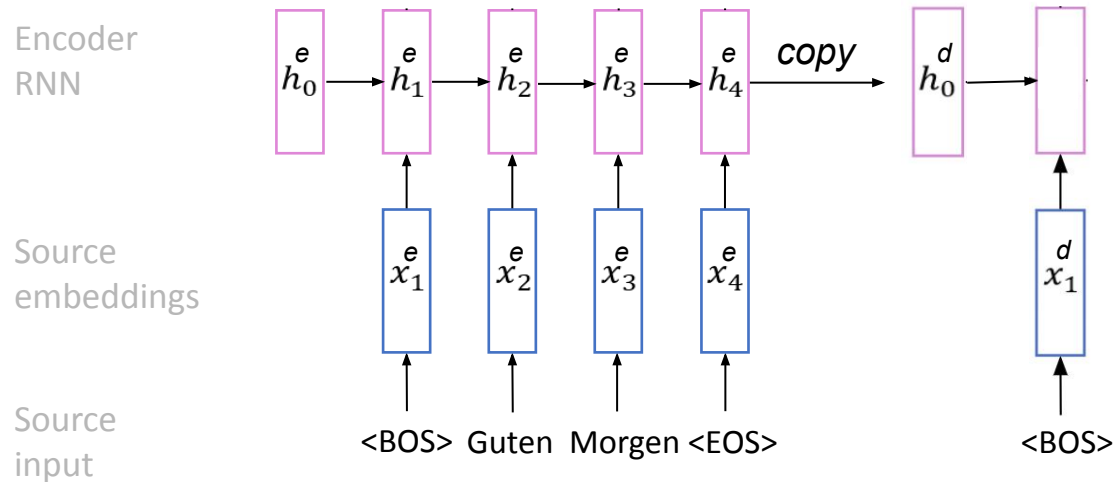
ABSTRACT

Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder-decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder-decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-)alignments found by the model agree well with our intuition.

Энкодер будет работать точно так же,
как и раньше



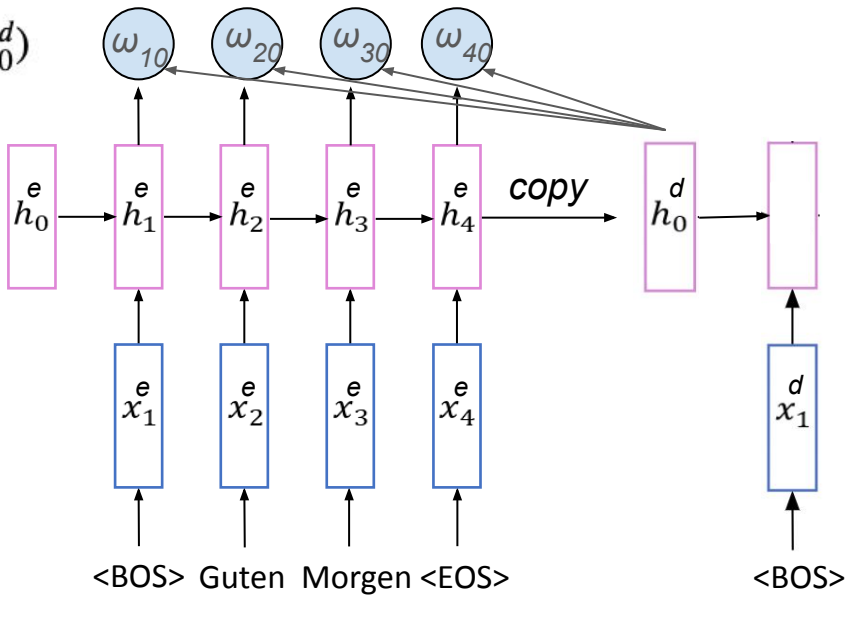
Перед тем, как обновить скрытое состояние декодера, нам нужно “обратить внимание” на входную последовательность



0.56 1.42 -0.45 -0.08

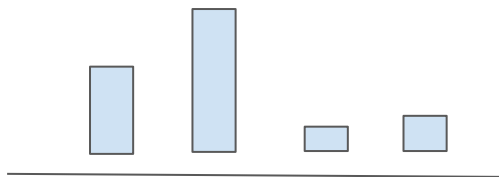
w_{i0} — это вес (важность) i -го входного токена для текущего состояния декодера

$$w_{i0} = s(h_i^e, h_0^d)$$



0.20 0.69 0.03 0.08

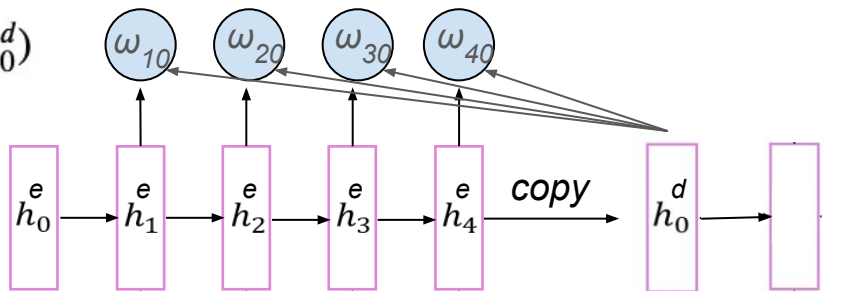
SoftMax



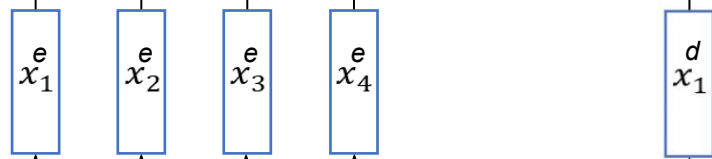
$$w_{i0} = s(h_i^e, h_0^d)$$

ω_{10} ω_{20} ω_{30} ω_{40}

Encoder
RNN



Source
embeddings



Source
input

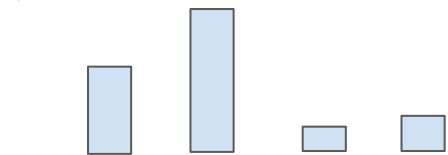
<BOS> Guten Morgen <EOS>

<BOS>

$$a_0 = \sum_{i=1}^4 h_i^e \cdot w_{i0}$$

a_0

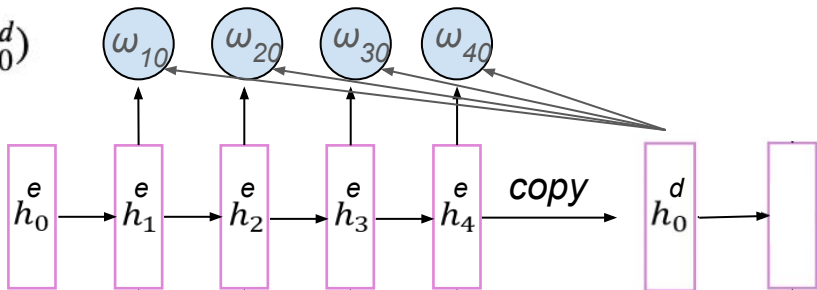
SoftMax



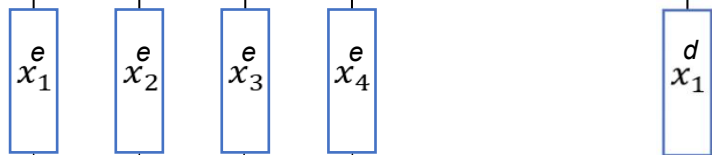
$$w_{i0} = s(h_i^e, h_0^d)$$

w_{10} w_{20} w_{30} w_{40}

Encoder
RNN



Source
embeddings



Source
input

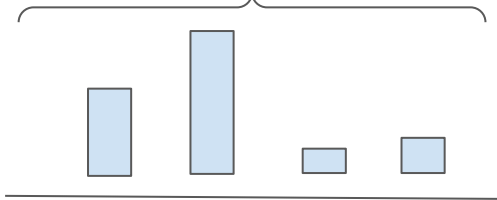
<BOS> Guten Morgen <EOS>

<BOS>

$$a_0 = \sum_{i=1}^4 h_i^e \cdot w_{i0}$$

a_0

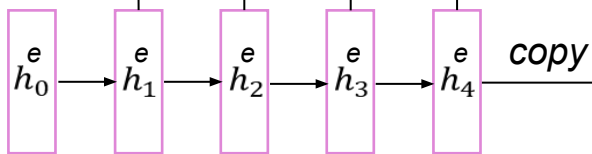
SoftMax



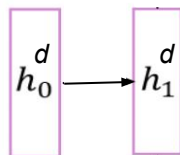
$$w_{i0} = s(h_i^e, h_0^d)$$

w_{10} w_{20} w_{30} w_{40}

Encoder
RNN

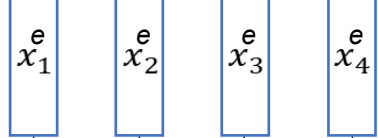


copy



$$h_1^d = RNN([x_1^d, a_0], h_0^d)$$

Source
embeddings

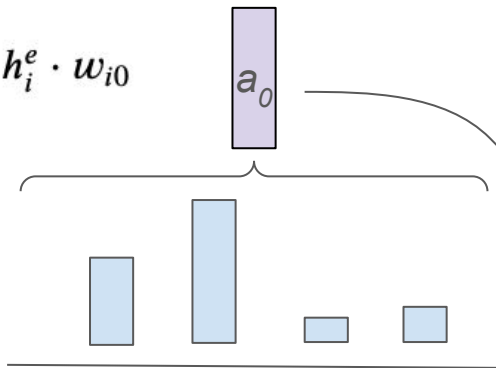


Source
input

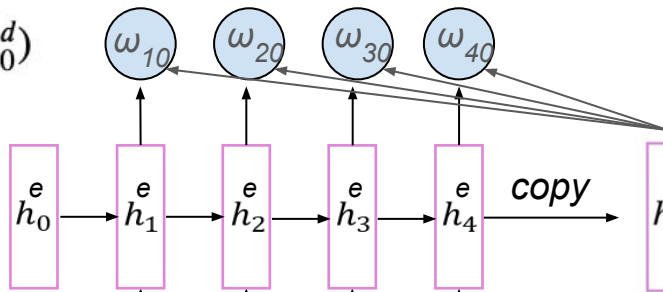
<BOS> Guten Morgen <EOS>

<BOS>

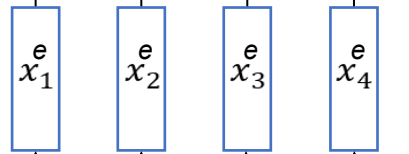
$$a_0 = \sum_{i=1}^4 h_i^e \cdot w_{i0}$$



$$w_{i0} = s(h_i^e, h_0^d)$$



Source embeddings

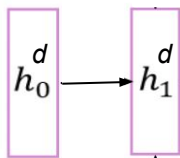


Source input

<BOS> Guten Morgen <EOS>



fc_0^d



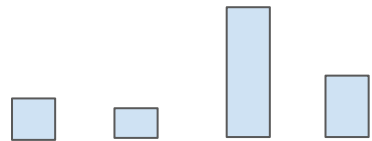
$$h_1^d = RNN([x_1^d, a_0], h_0^d)$$

<BOS>

$$a_1 = \sum_{i=1}^4 h_i^e \cdot w_{i1}$$

a_1

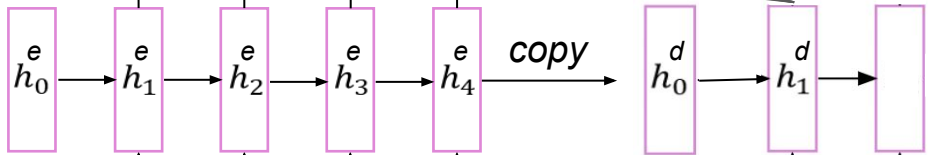
SoftMax



$$w_{i1} = s(h_i^e, h_1^d)$$

w_{11} w_{21} w_{31} w_{41}

Encoder
hidden
states



Source token
embeddings

x_1 x_2 x_3 x_4

Source
input

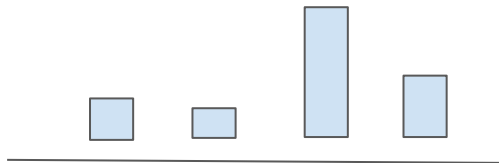
<BOS> Guten Morgen <EOS>

<BOS> Good

$$a_1 = \sum_{i=1}^4 h_i^e \cdot w_{i1}$$

a_1

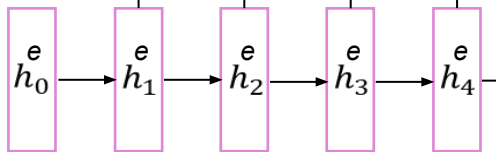
SoftMax



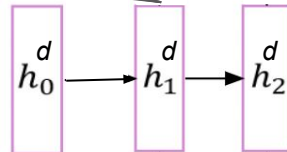
$$w_{i1} = s(h_i^e, h_1^d)$$

ω_{11} ω_{21} ω_{31} ω_{41}

Encoder
hidden
states



copy



$$h_2^d = RNN([x_2^d, a_1], h_1^d)$$

Source token
embeddings

Source
input

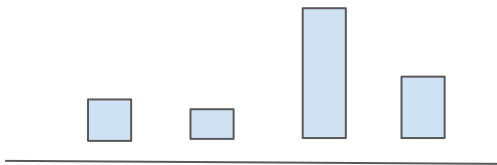
<BOS> Guten Morgen <EOS>

<BOS> Good

$$a_1 = \sum_{i=1}^4 h_i^e \cdot w_{i1}$$

a_1

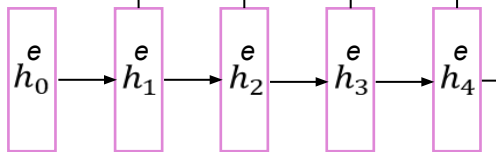
SoftMax



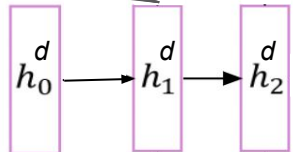
$$w_{i1} = s(h_i^e, h_1^d)$$

ω_{11} ω_{21} ω_{31} ω_{41}

Encoder hidden states

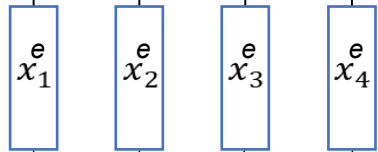


copy



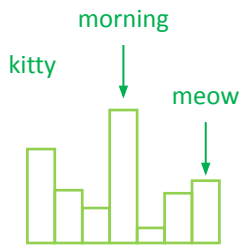
$$h_2^d = RNN([x_2^d, a_1], h_1^d)$$

Source token embeddings



Source input

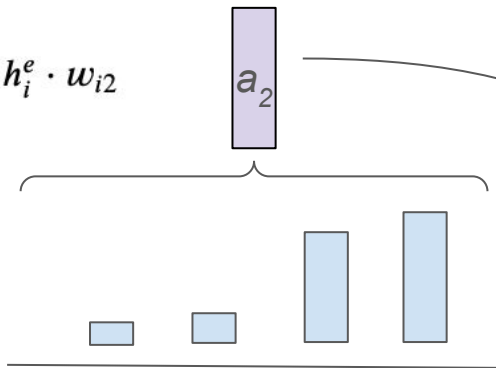
<BOS> Guten Morgen <EOS>



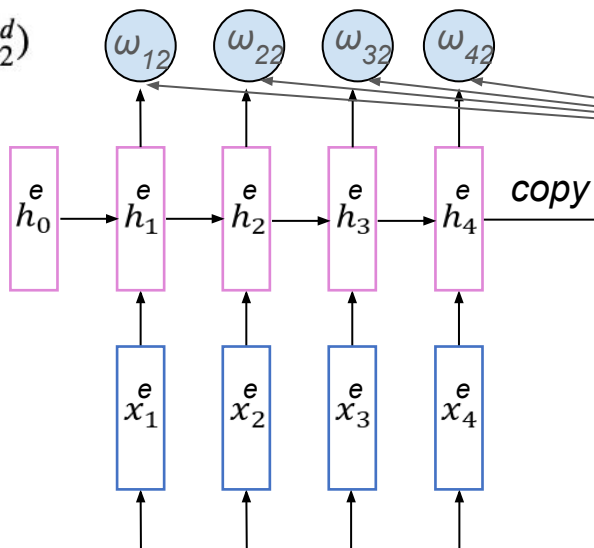
fc_2^d

<BOS> Good

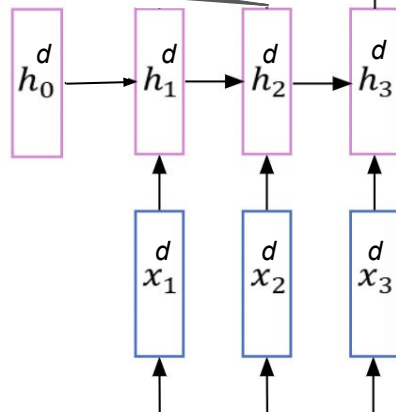
$$a_2 = \sum_{i=1}^4 h_i^e \cdot w_{i2}$$



$$w_{i2} = s(h_i^e, h_2^d)$$



copy



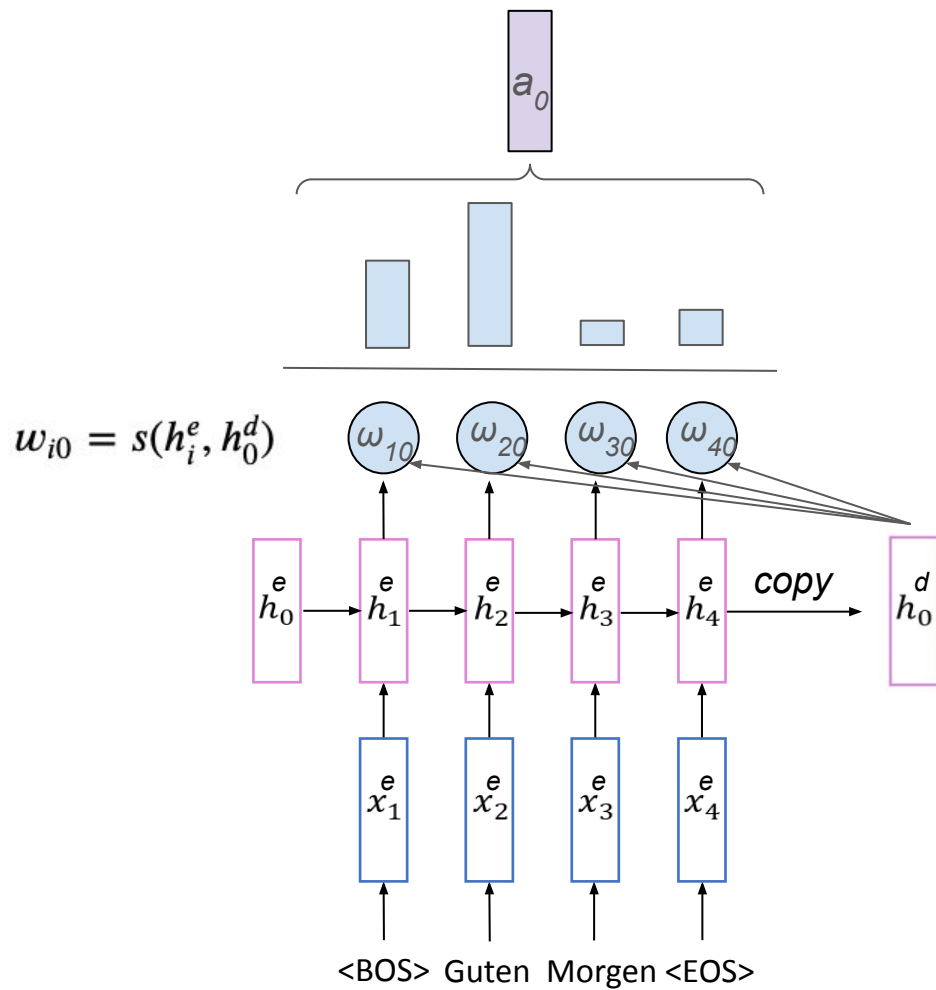
$$h_3^d = RNN([x_3^d, a_2], h_2^d)$$



fc_3^d

<BOS> Good morning

Source input



Функция s может быть обучаемой или не обучаемой.

Если она обучаемая, то ее веса обучаются совместно со всеми параметрами сети.

Варианты функций Attention

$$w = s(h_1, h_2)$$

- Dot product attention:
(без обучаемых параметров)

$$s(h_1, h_2) = h_1^T h_2$$

- Однослойная нейросеть
(с обучаемыми параметрами)

$$s(h_1, h_2) = \sigma(W[h_1, h_2])$$

- Многие другие

Эффект применения Attention

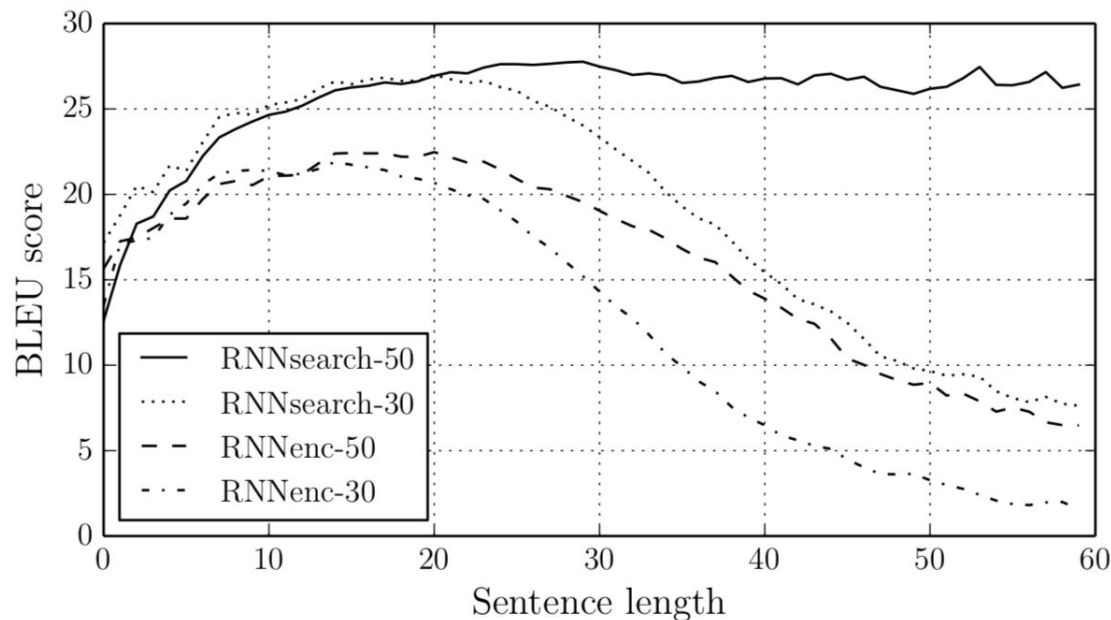
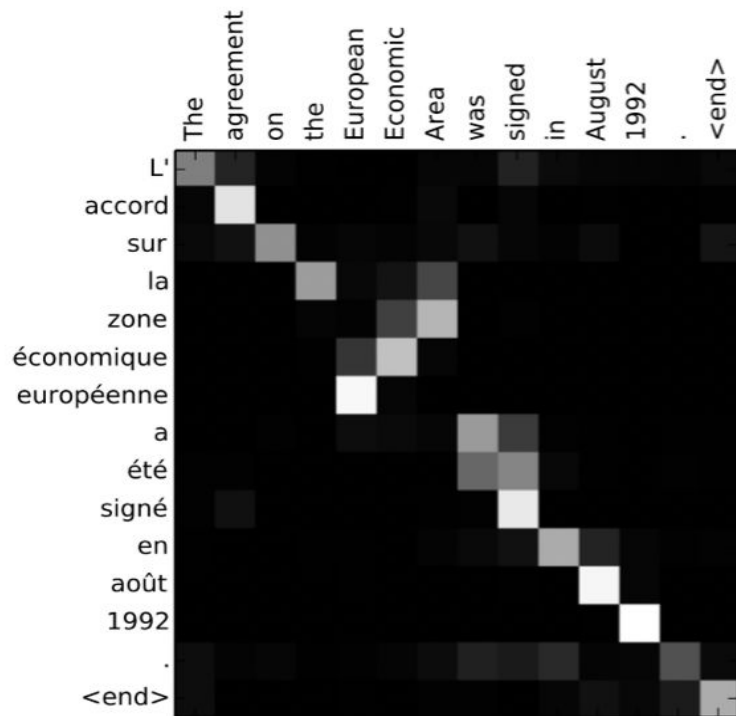
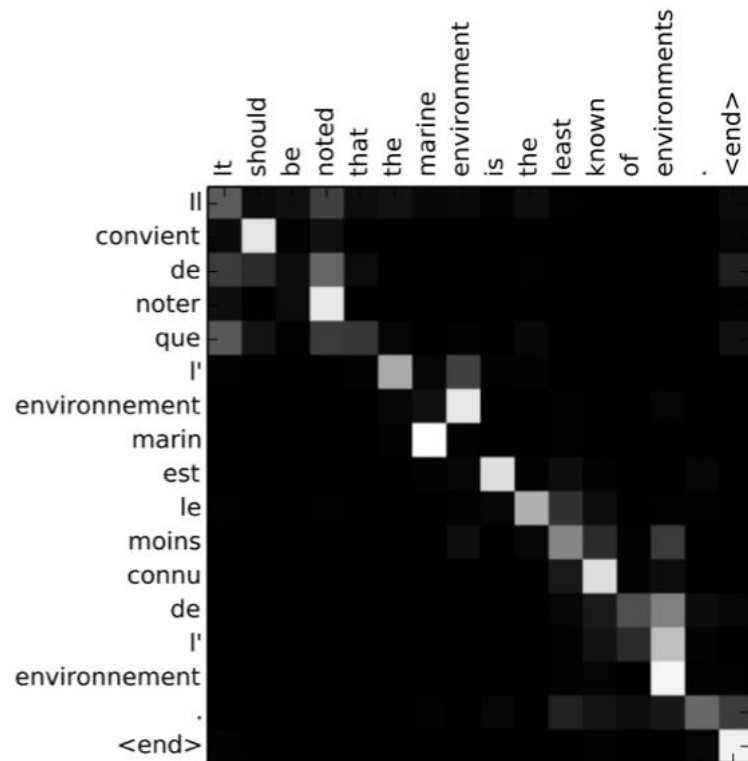


Figure 2: The BLEU scores of the generated translations on the test set with respect to the lengths of the sentences. The results are on the full test set which includes sentences having unknown words to the models.

Attention visualization



Attention visualization



Революция после Attention

Изобретение Attention было поворотным моментом в задаче машинного перевода и всей области NLP

В 2017 году Google представили новую архитектуру нейросети, которую называли Transformer. Основой ее был механизм Attention. Transformer сотворили революцию в NLP, и потом были успешно адаптированы и к другим доменам.

Итоги видео

В этом видео:

- Идея механизма Attention;
- Attention для RNN-модели Seq2Seq;
- Разные варианты Attention;
- Attention может быть интерпретируемым.