

DL: Машинный перевод

План

- ❑ Статистический машинный перевод
- ❑ Эмбединги
- ❑ Нейронные сети

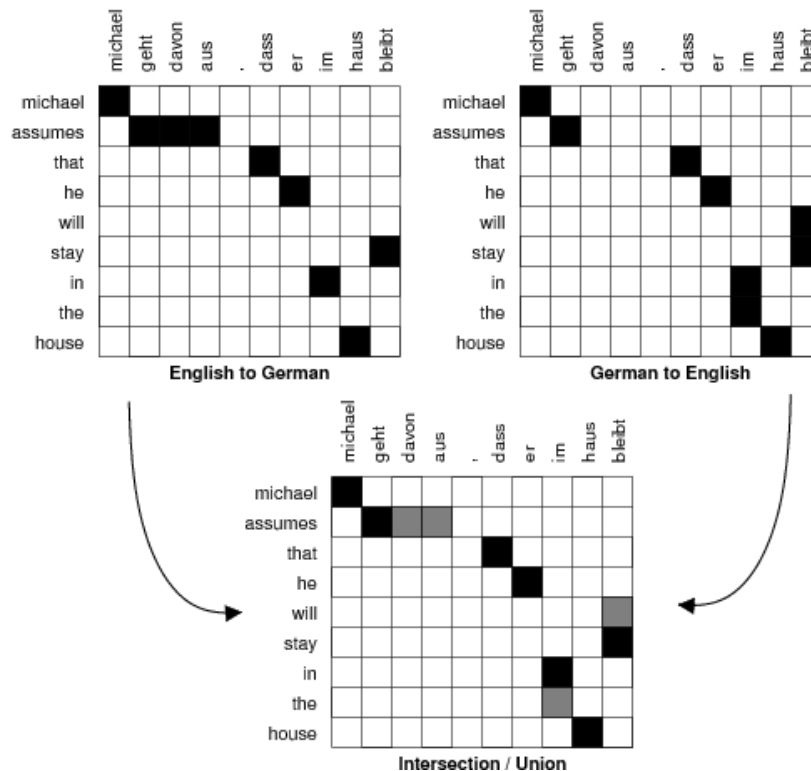
Статистический машинный перевод

$$\operatorname{argmax}_y P(y | x)$$

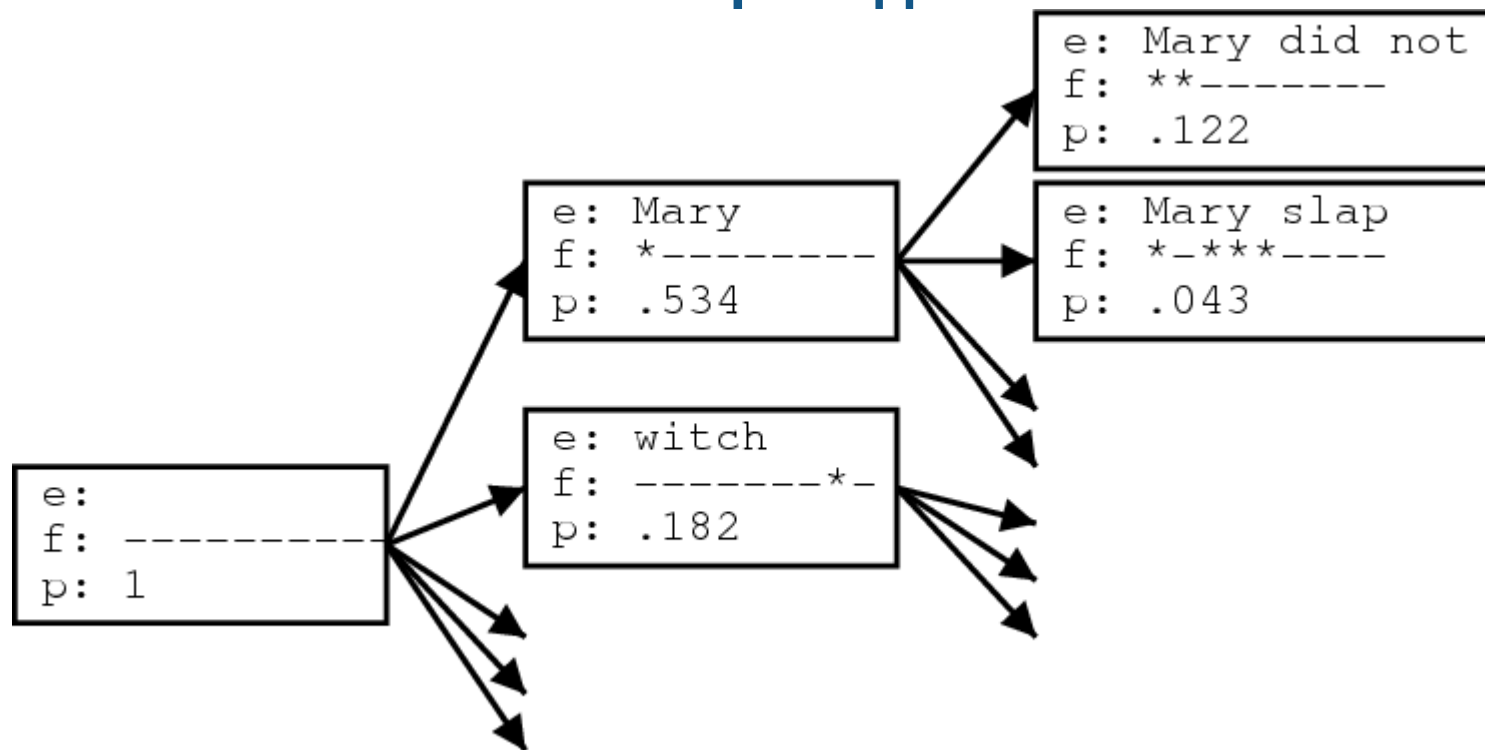
$$\operatorname{argmax}_y P(x | y) P(y)$$

$$P(x, a | y) P(y)$$

Статистический машинный перевод: alignment



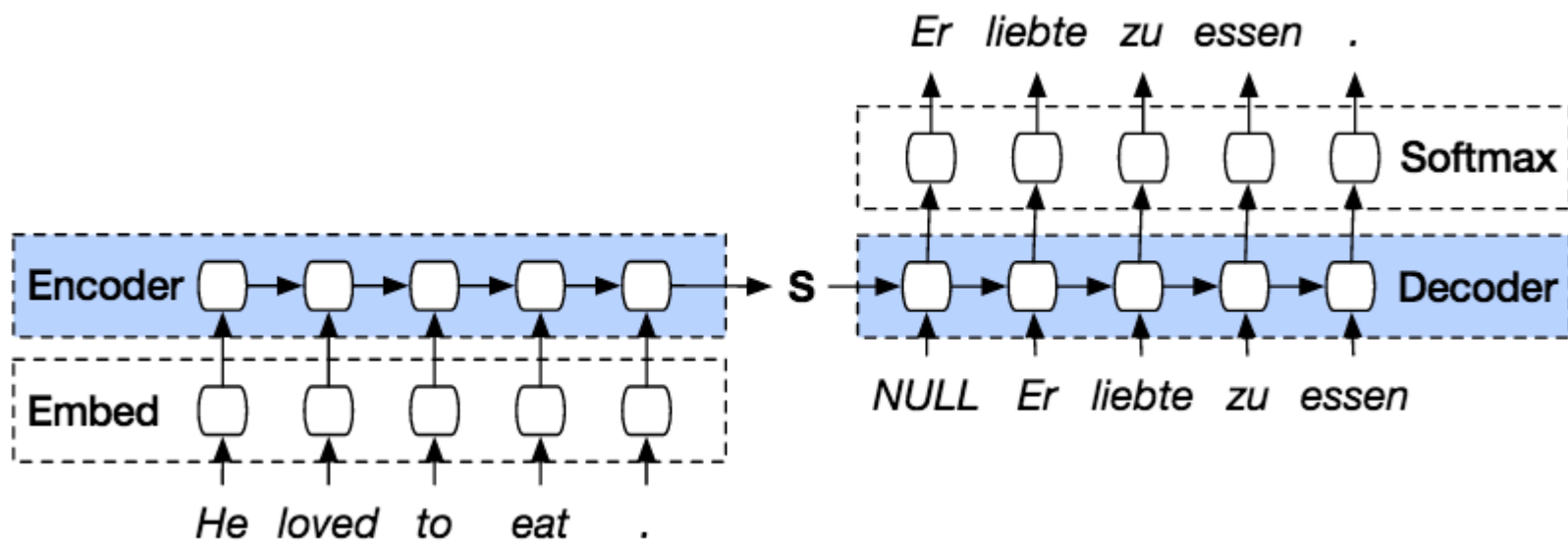
Статистический машинный перевод



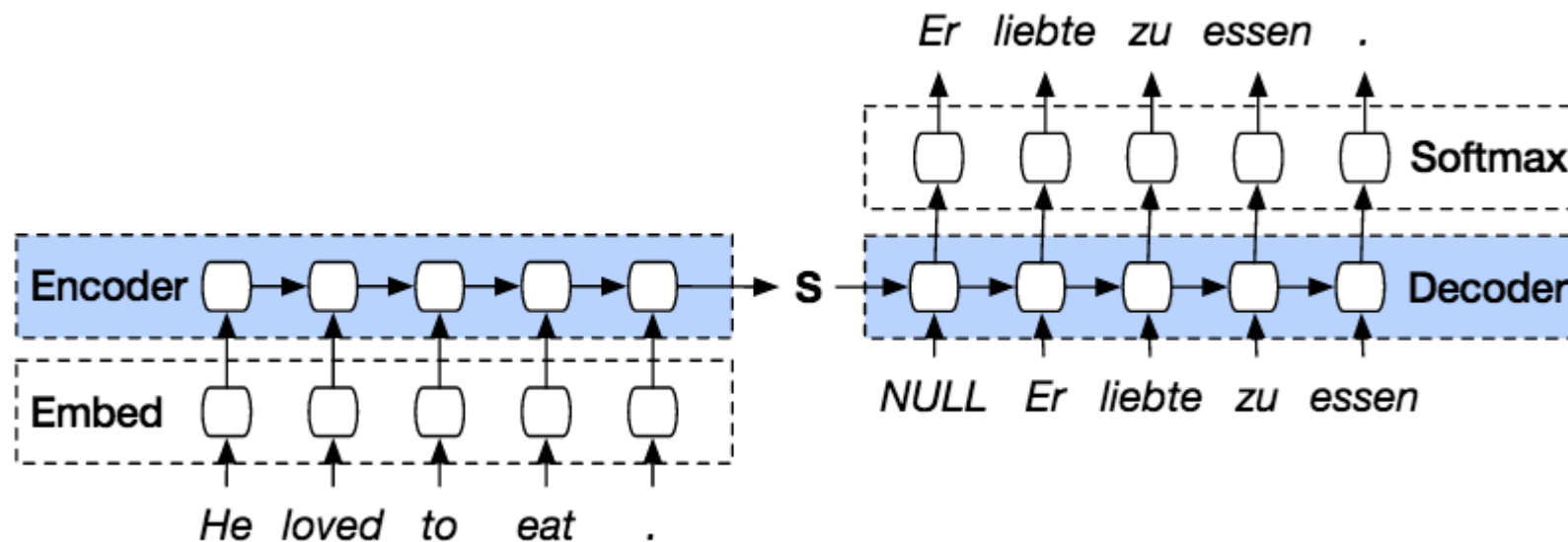
Статистический машинный перевод



Seq2seq 2014

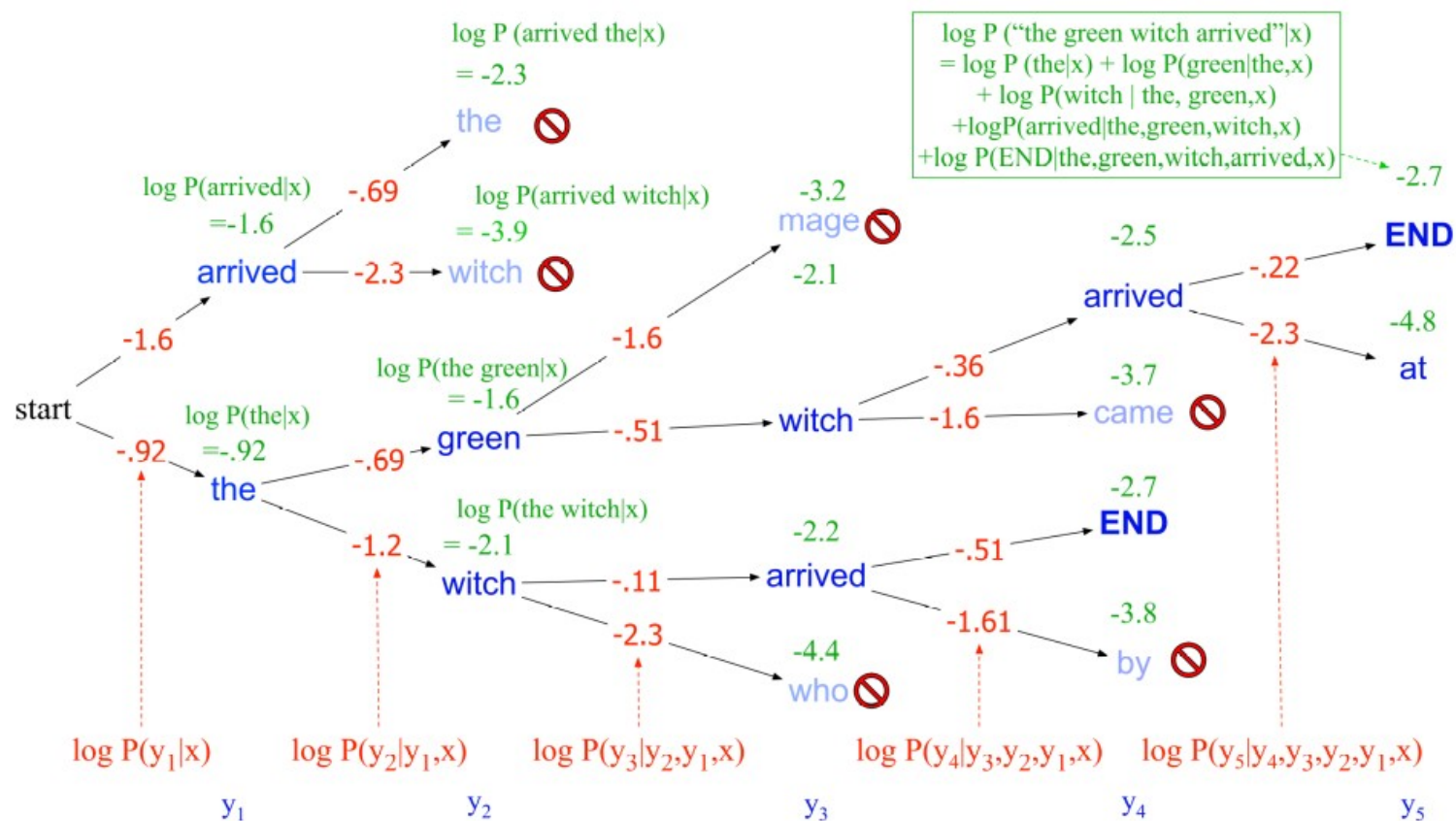


Seq2seq greedy vs beam search



$$\text{score}(y_1, y_2, \dots, y_n) = \log P(y_1, y_2, \dots, y_n | x) = \sum_{i=1}^n \log P(y_i | y_1, \dots, y_{i-1} | x)$$

beam search



beam search: нормировка

$$\text{score}(y_1, y_2, \dots, y_n) = \log P(y_1, y_2, \dots, y_n | x) = \sum_{i=1}^n \log P(y_i | y_1, \dots, y_{i-1} | x)$$

$$\frac{1}{n} \sum_{i=1}^n \log P(y_i | y_1, \dots, y_{i-1} | x)$$

- N-gram overlap between machine translation output and reference translation
- Compute precision for n-grams of size 1 to 4
- Add brevity penalty (for too short translations)

$$\text{BLEU} = \min \left(1, \frac{\text{output-length}}{\text{reference-length}} \right) \left(\prod_{i=1}^4 \text{precision}_i \right)^{\frac{1}{4}}$$

- Typically computed over the entire corpus, not single sentences

SYSTEM A: Israeli officials responsibility of airport safety
2-GRAM MATCH 1-GRAM MATCH

REFERENCE: Israeli officials are responsible for airport security

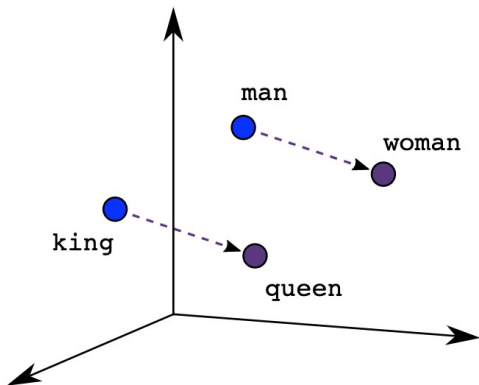
SYSTEM B: airport security Israeli officials are responsible
2-GRAM MATCH 4-GRAM MATCH

Metric	System A	System B
precision (1gram)	3/6	6/6
precision (2gram)	1/5	4/5
precision (3gram)	0/4	2/4
precision (4gram)	0/3	1/3
brevity penalty	6/7	6/7
BLEU	0%	52%

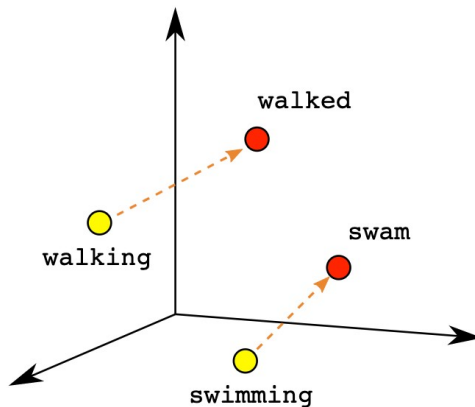
Word Embedding: one hot

motel [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0] AND
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0] = 0

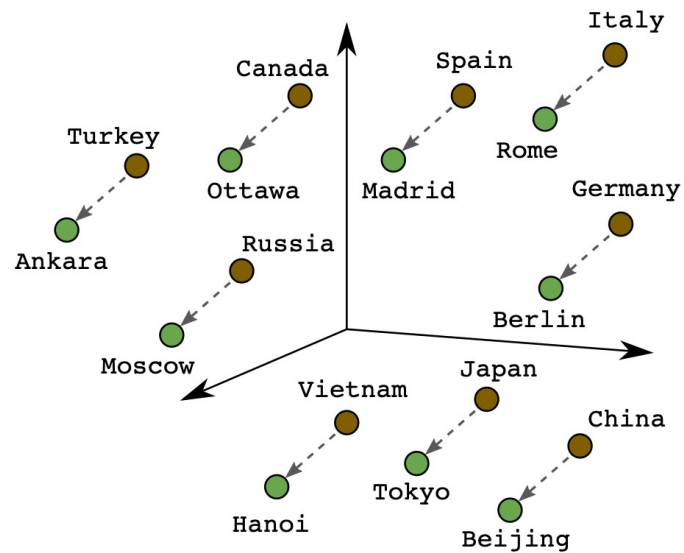
Word Embedding: word2vec



Male-Female



Verb Tense



Country-Capital

Word Embedding: word2vec

