

9

Введение в NLP

word2vec

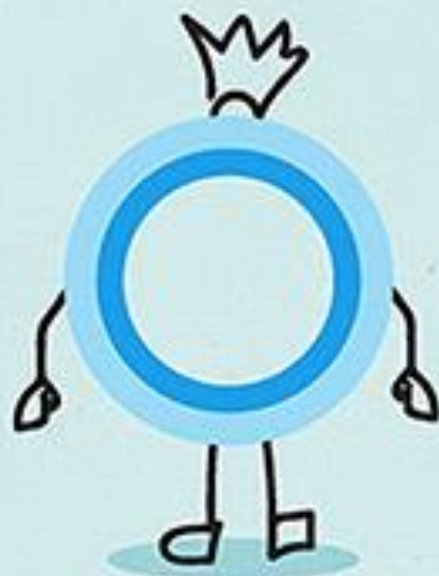
Wikipedia



東 = est = ESTE = east
西 = *ouest* = OESTE = west
北 = *nord* = NORTE = north
南 = *sud* = SUR = south



Siri



Cortana



Alexa



Google Now

Почему это вообще сложно?

Он видел их семью своими глазами

Эти типы стали есть в цехе

[Habr](#)

парень 25 лет ищет подработку на субботу и воскресенье

новые куртки есть размеры от 2000 рублей

Перед нами стол.

На столе стакан и вилка.

Что они делают?

Стакан стоит, а вилка лежит.

Если мы воткнем вилку в столешницу, вилка будет стоять.

То есть стоят вертикальные предметы, а лежат горизонтальные?

Добавляем на стол тарелку и сковороду.

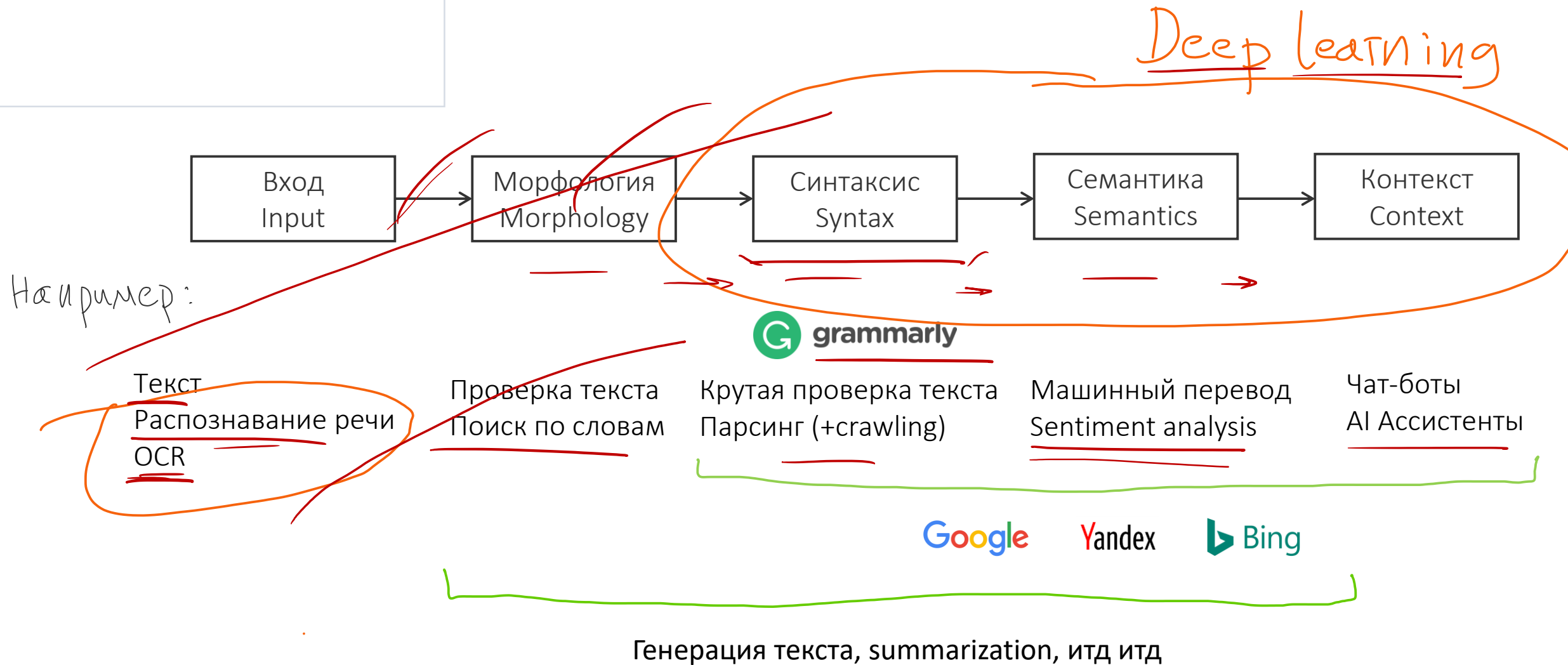
Они вроде горизонтальные, но на столе стоят.

Теперь положим тарелку в сковородку.

Там она лежит, а ведь на столе стояла.

[Livejournal](#)

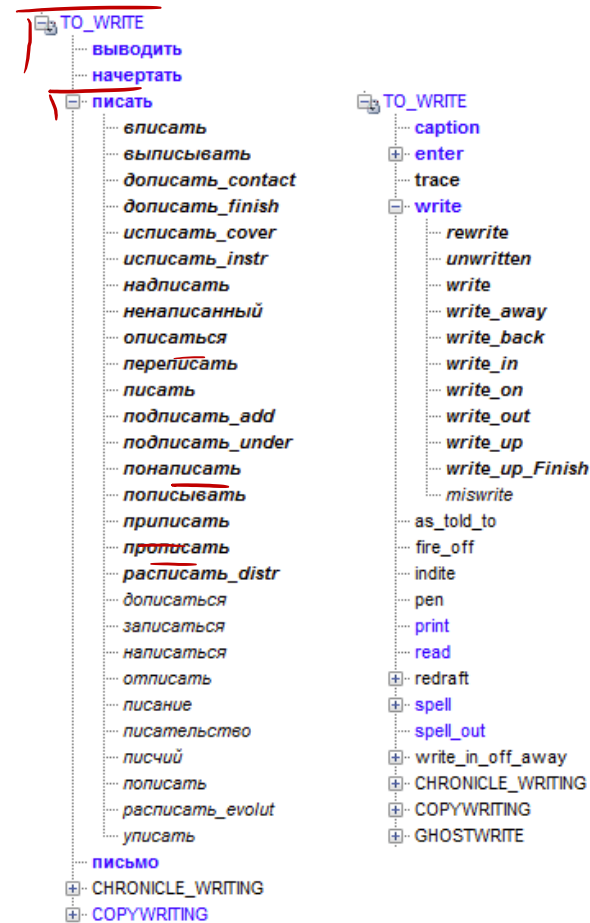
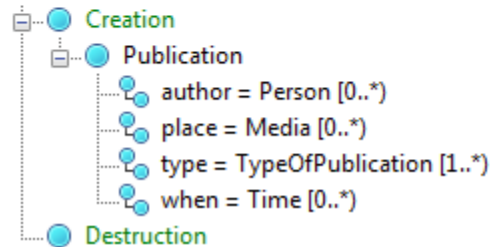
NLP pipeline



Как это делают без нейросетей

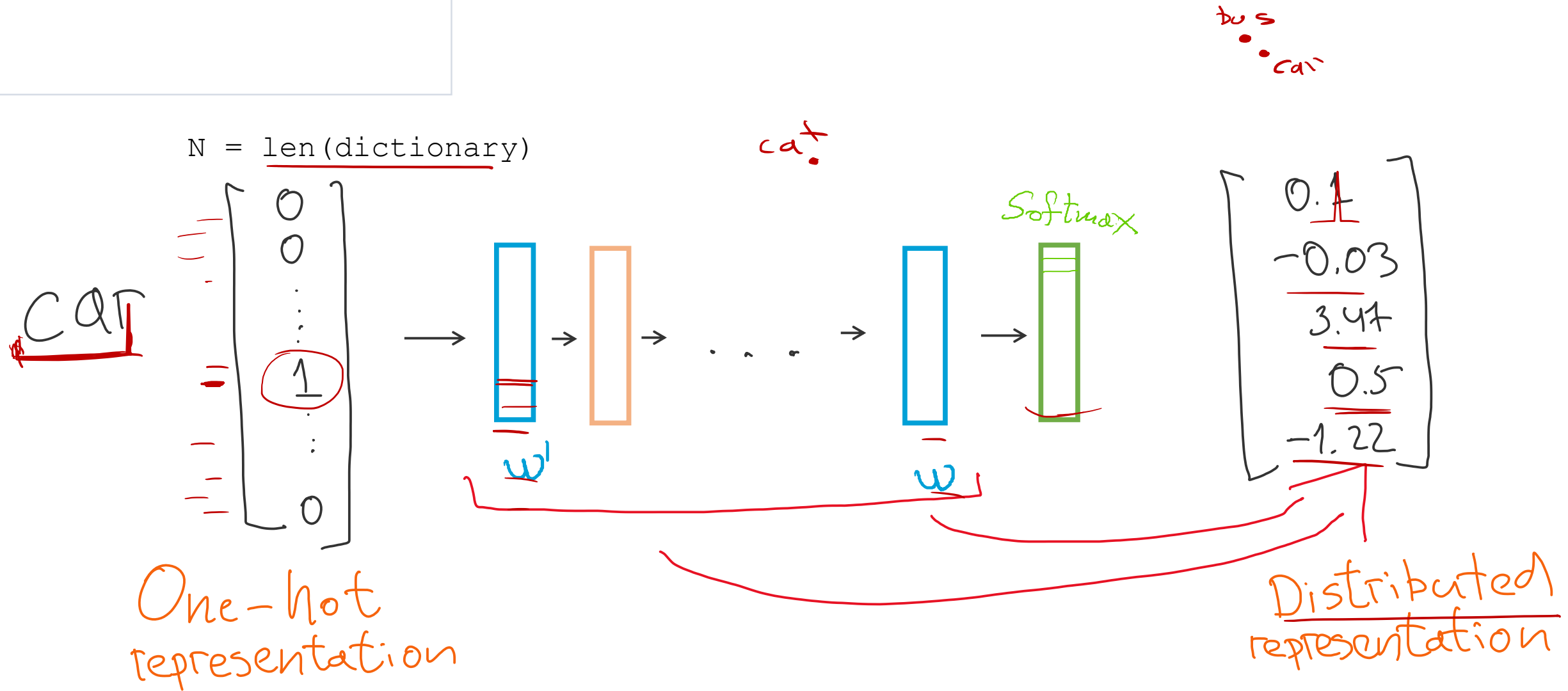
Онтоинженер Даня в 2014 году написал пост на Хабр

Для этого потребуется модель предметной области (онтология) и правила извлечения информации. Созданием онтологий и правил занимается специальный отдел компьютерных лингвистов, которых мы называем онтоинженерами. Пример онтологии, моделирующей факт публикации:



Deep NLP

Из символьного в непрерывное



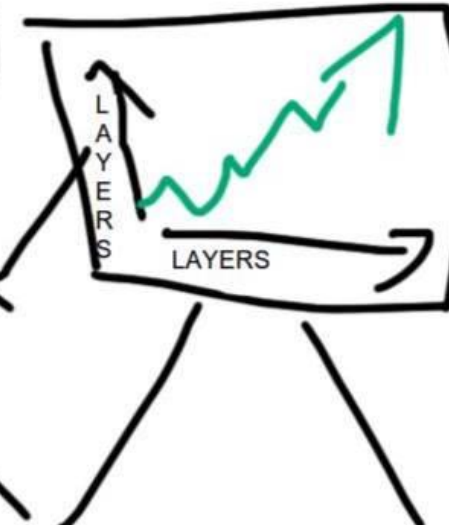
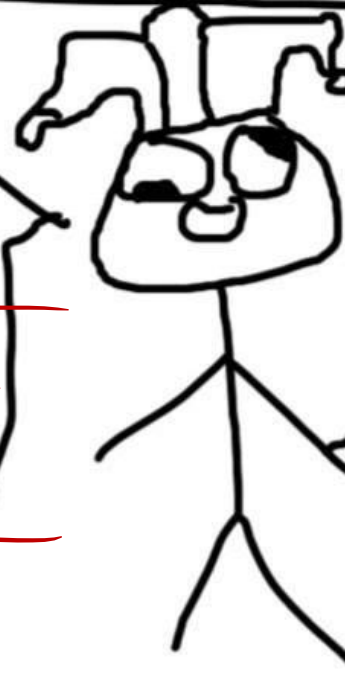
STATISTICAL LEARNING

Gentlemen, our learner overgeneralizes because the VC-Dimension of our Kernel is too high, Get some experts and minimize the structural risk in a new one. Rework our loss function, make the next kernel stable, unbiased and consider using a soft margin



NEURAL NETWORKS

STACK
MORE
LAYERS



word2vec

The quick brown fox jumps over the lazy dog

Source Text

Training Samples

<u>The</u> quick brown fox jumps over the lazy dog.	→	(<u>the</u> , <u>quick</u>) (<u>the</u> , <u>brown</u>)
The <u>quick</u> brown fox jumps over the lazy dog.	→	(quick, <u>the</u>) (quick, <u>brown</u>) (quick, fox)
The quick brown fox jumps over the lazy dog.	→	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown <u>fox</u> jumps over the lazy dog.	→	(<u>fox</u> , quick) (<u>fox</u> , brown) (fox, jumps) (fox, over)

Skip-gram

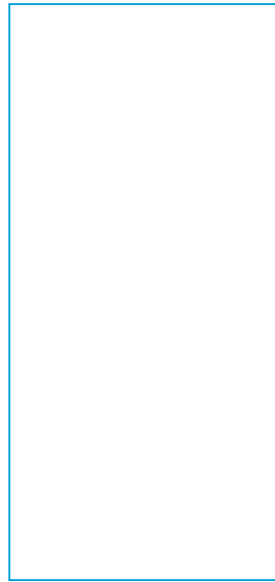
word2vec

The quick brown fox jumps over the lazy dog

fox -> quick
fox -> brown
fox -> jumps
fox -> over

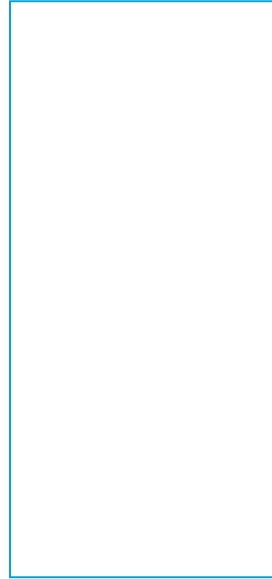
len(dict)

fox $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$



u

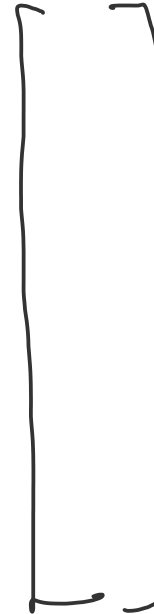
256-1024



v



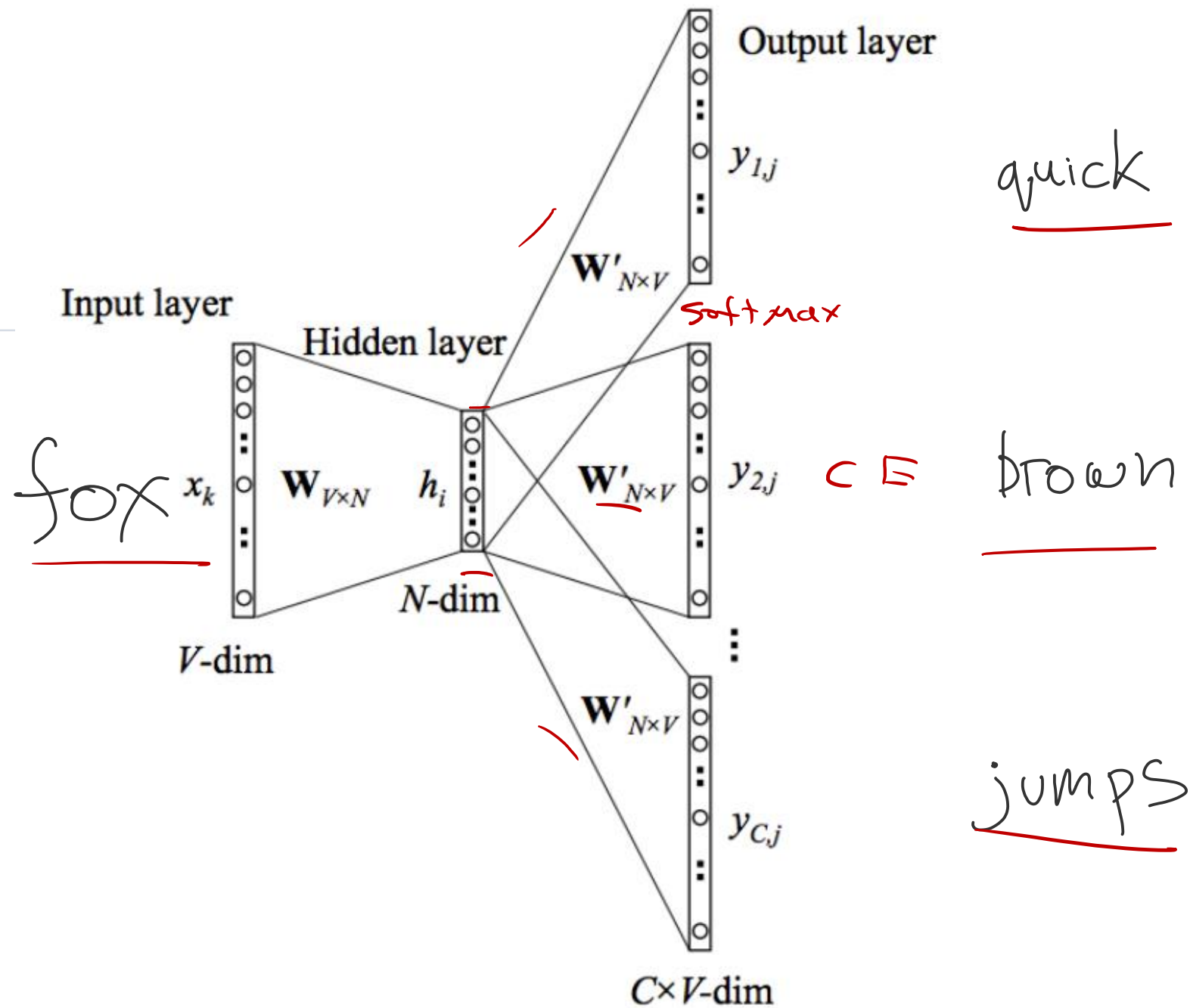
len(dict)



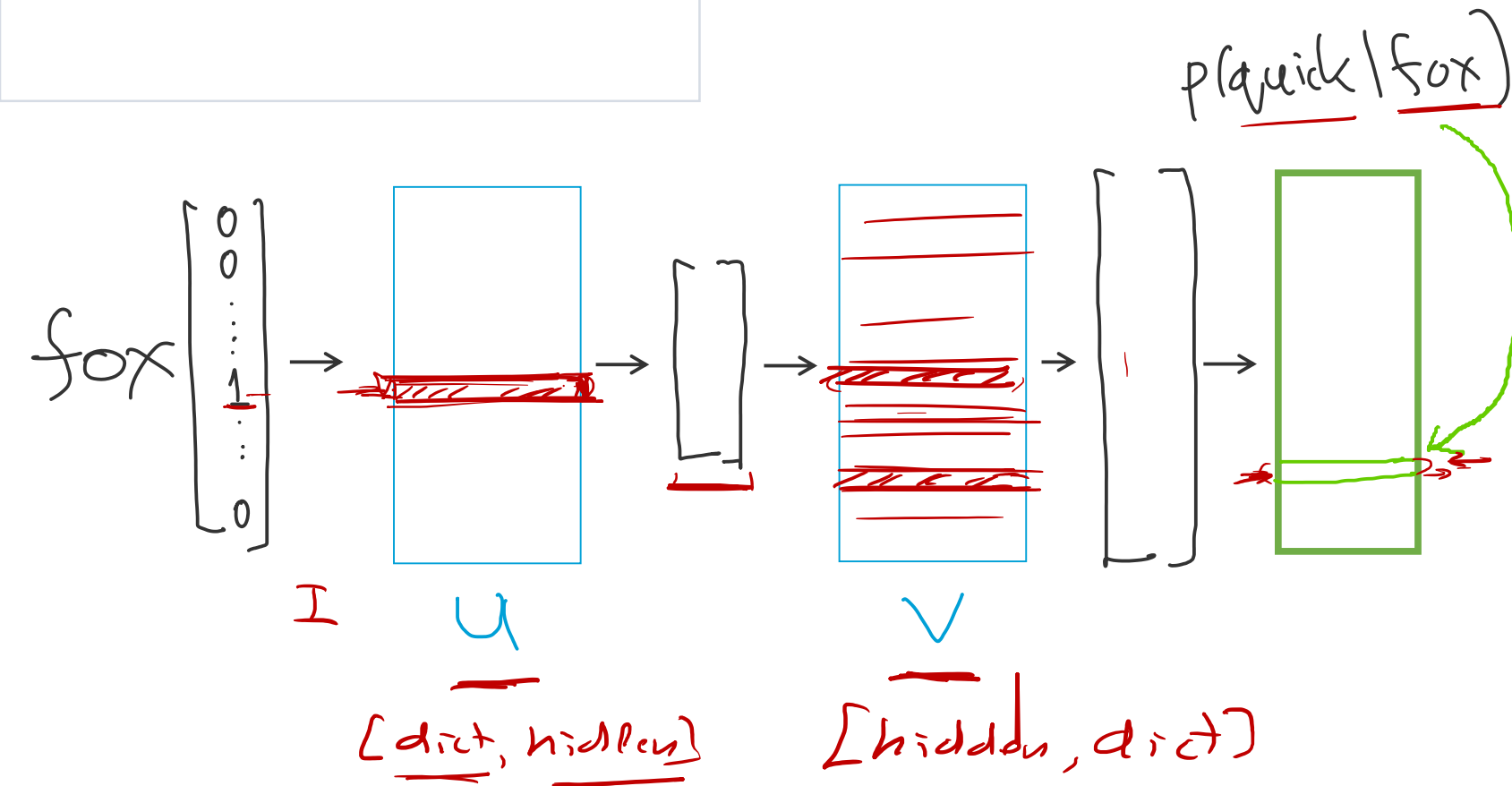
Softmax

quick

$$L = - \sum_j \ln p(c = y_j | x_j)$$



Word vector



$$p(o|i) = \frac{e^{u_i \cdot v_o}}{\sum_k e^{u_i \cdot v_k}}$$

$$L = - \sum_j \ln p(o|i)$$

$$= - \sum_j \ln \frac{e^{u_i \cdot v_o}}{\sum_k e^{u_i \cdot v_k}}$$

$$w(\text{fox}) = \underline{u(\text{fox})} + \underline{v(\text{fox})}$$

$$\text{fox}$$

$$[1 \ 0.5 \ -0.3 \ \dots \ 0.02]$$

Word vector



$$p(o|i) = \frac{e^{u_i \cdot v_o}}{\sum_k e^{u_i \cdot v_k}}$$

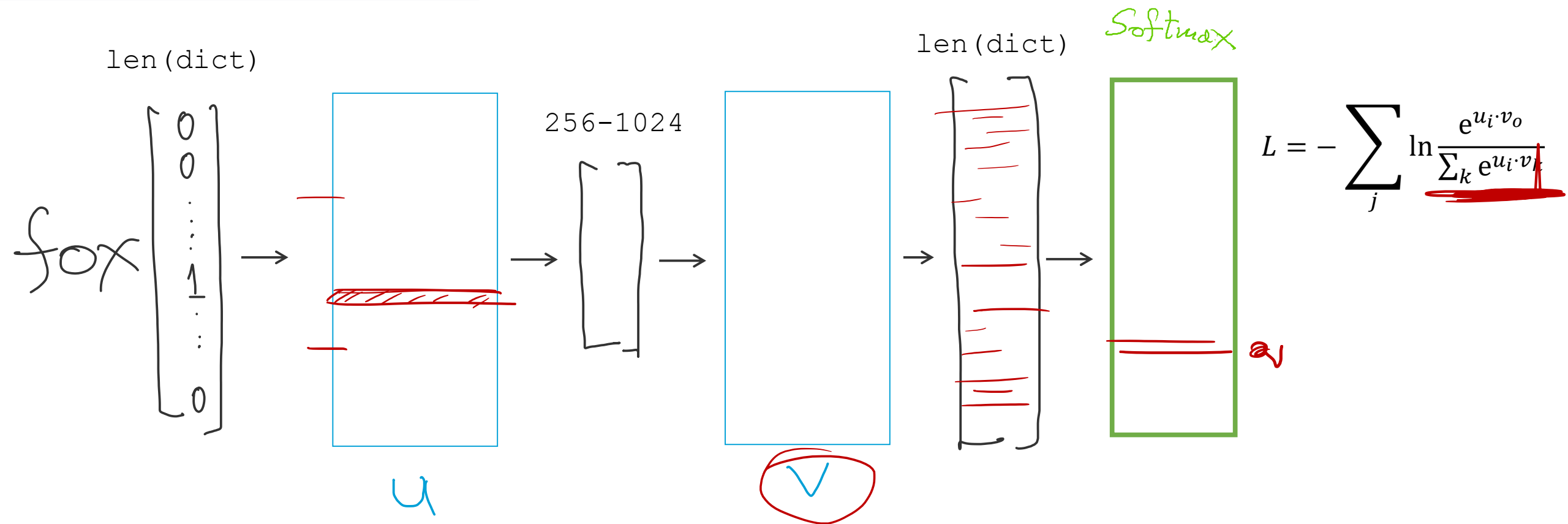
$$\begin{aligned} L &= - \sum_j \ln p(o|i) \\ &= - \sum_j \ln \frac{e^{u_i \cdot v_o}}{\sum_k e^{u_i \cdot v_k}} \end{aligned}$$

$$w(\text{fox}) = u(\text{fox}) + v(\text{fox})$$

fox

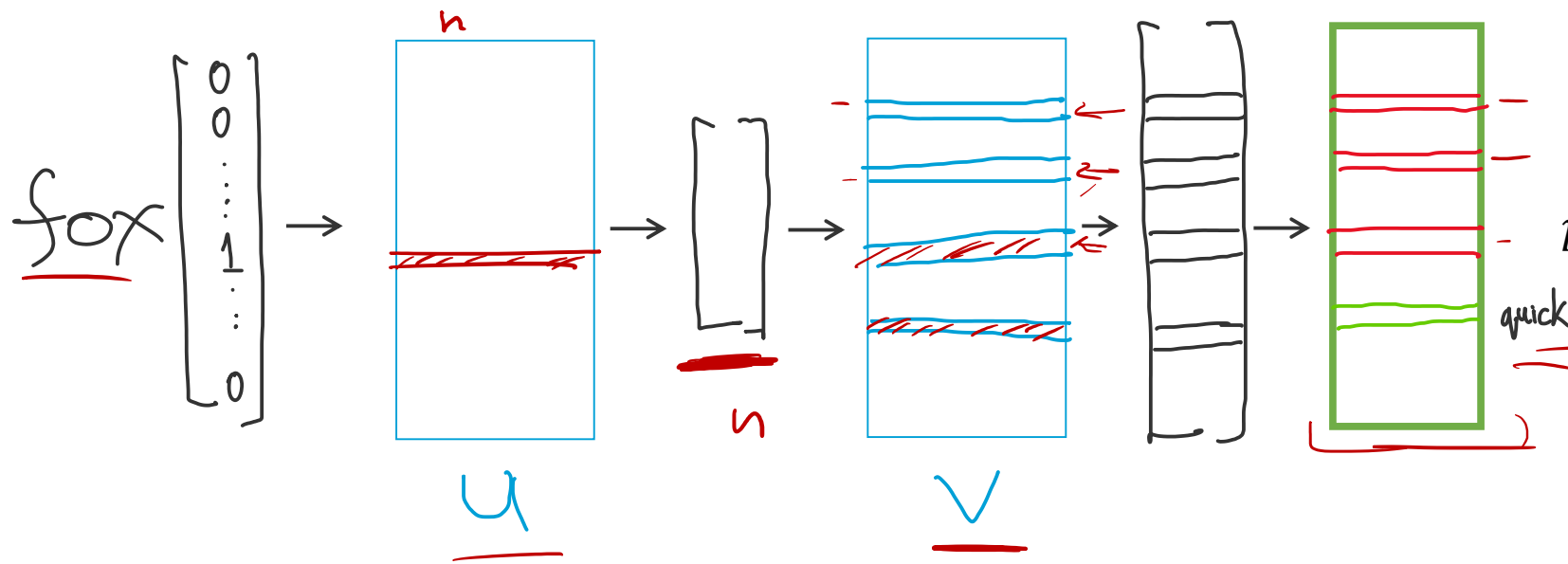
$$[1 \ 0.5 \ -0.3 \ \dots \ 0.02]$$

Небольшая проблемка



Negative sampling

brood, jump, quick



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$p(D = 1 | i, o) = \sigma(u_i \cdot v_o) = \frac{1}{1 + e^{-u_i \cdot v_o}}$$

$$p(D = 0 | i, o) = 1 - \sigma(u_i \cdot v_o) = \frac{1}{1 + e^{u_i \cdot v_o}}$$

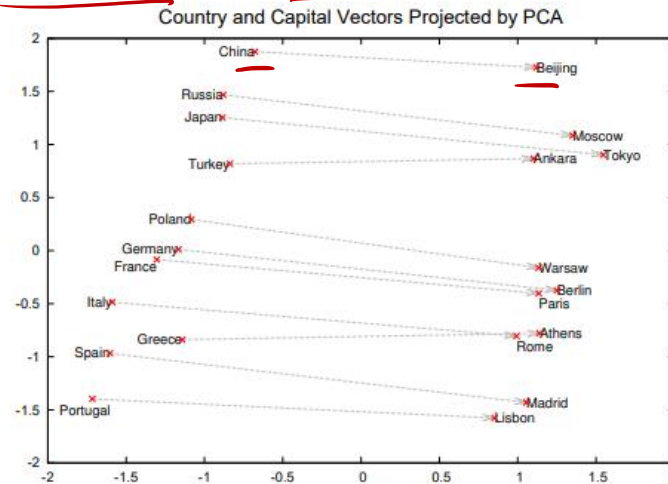
$$L = -\log \sigma(u_i \cdot v_o) - \sum_{k \sim P(\omega)} \log \sigma(-u_i \cdot v_k)$$

$$P(\omega) \sim U(\omega)^{3/4}$$

Свойства word2vec

Аналогии

$$\underline{w(Paris)} - \underline{w(France)} + \underline{w(Russia)} = ? \quad \underline{w(Moscow)}$$



[Mikolov'13](#)

king:queen::man:[woman, Attempted abduction, teenager, girl]

house:roof::castle:[dome, bell_tower, spire, crenellations, turrets]

new york times:sulzberger::fox:[Murdoch, Chernin, Bancroft, Ailes]

love:indifference::fear:[apathy, callousness, timidity, helplessness, inaction]

donald trump:republican::barack obama:[Democratic, GOP, Democrats, McCain]

building:architect::software:[programmer, SecurityCenter, WinPcap]

[Chris Nicholson on Quora](#)

T-SNE



Ок, а из полезного в хозяйстве?

CAT

$$\begin{bmatrix} 0.1 \\ -0.03 \\ 3.47 \\ 0.5 \\ -1.22 \end{bmatrix}$$

The film is powerful, accessible and funny

The movie is well done, but slow

Ah yes, and then there's the music

[Stanford Sentiment Treebank](#)

$$\underline{v(sentence)} = \sum_w \underline{P(word)} \underline{v(word)}$$

fastText от Facebook

DOV

Character N-grams

Each word w is represented as a bag of character n -gram. We add special boundary symbols $<$ and $>$ at the beginning and end of words, allowing to distinguish prefixes and suffixes from other character sequences. We also include the word w itself in the set of its n -grams, to learn a representation for each word (in addition to character n -grams). Taking the word *where* and $n = 3$ as an example, it will be represented by the character n -grams:

<wh, whe, her, ere, re>
and the special sequence

<where>.

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c.$$

хсреднить
все
n-grams

- Лучше word2vec
- Решает проблему OOV (out of the vocabulary)
- Есть pretrained вектора
- Работает из коробки

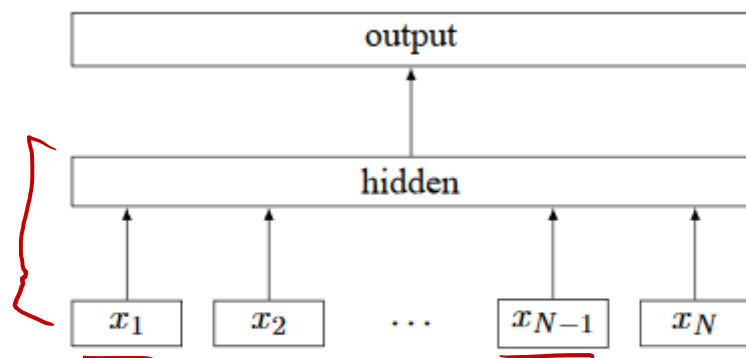
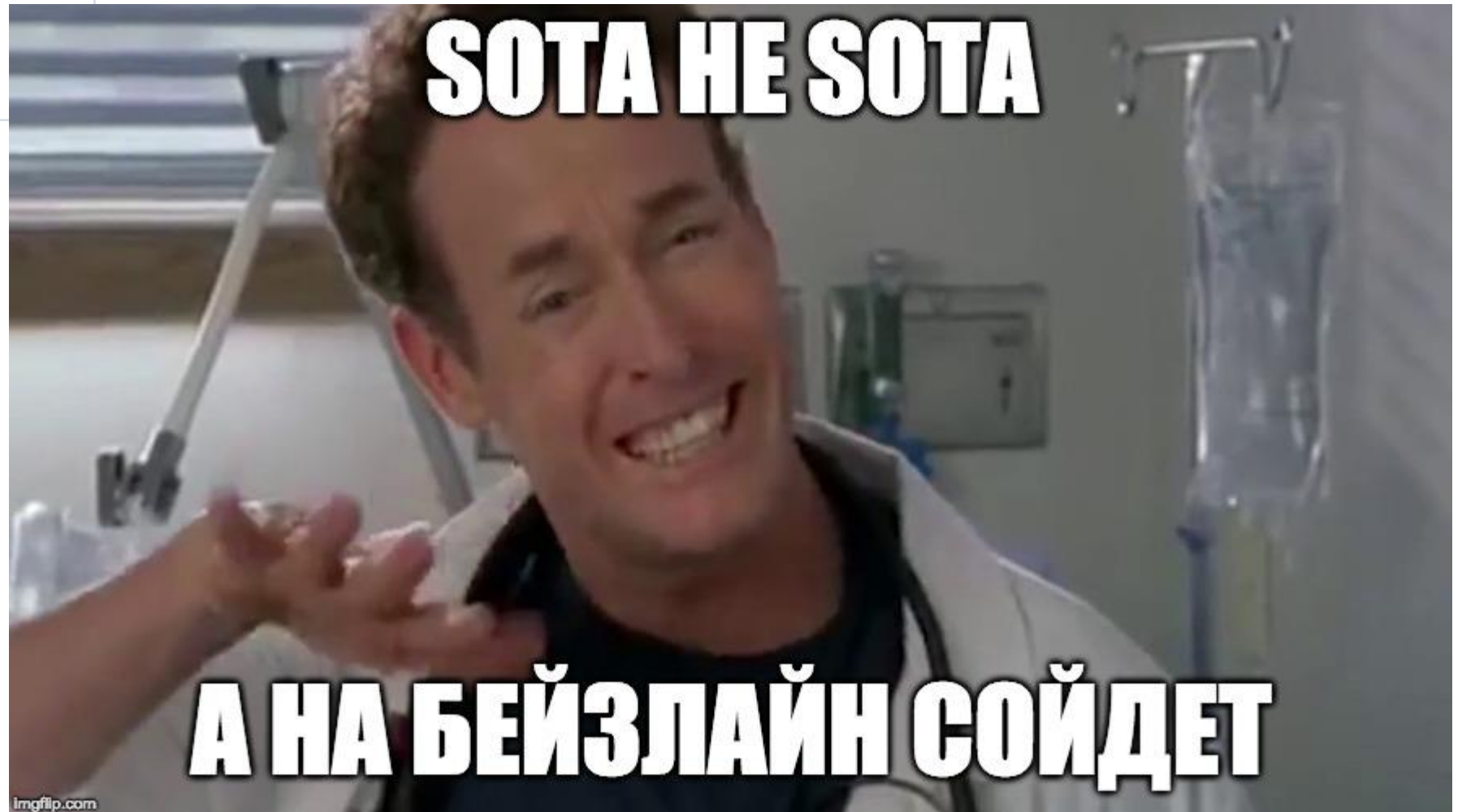


Figure 1: Model architecture of fastText for a sentence with N ngram features x_1, \dots, x_N . The features are embedded and averaged to form the hidden variable.

BACK 
TO THE PRESENT



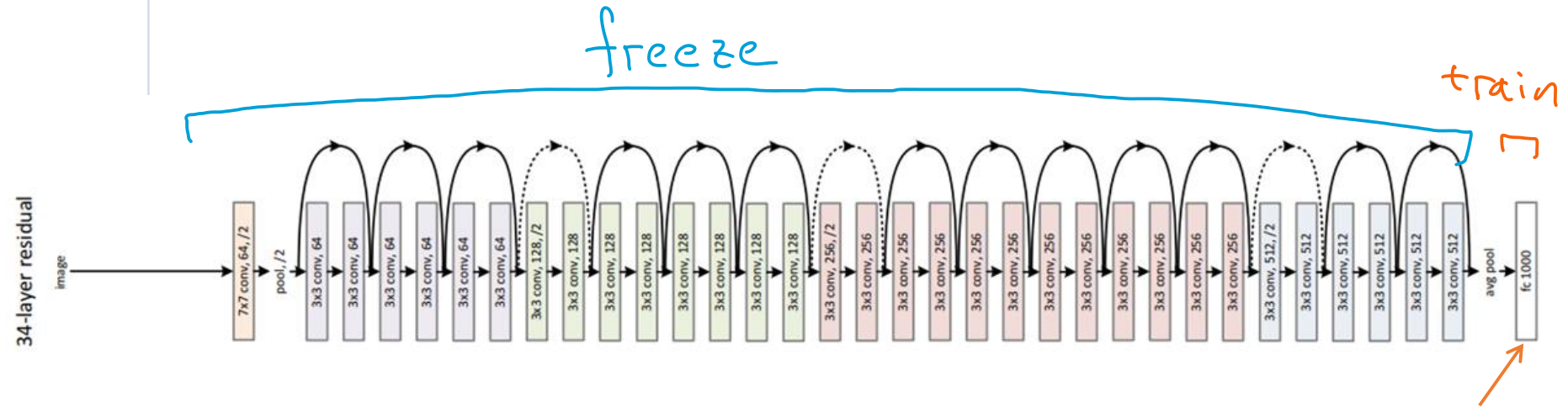
BERT

ELMO



BACK 
TO THE PRESENT

Pretrained vectors



CAT $\rightarrow \begin{bmatrix} 0.1 \\ -0.03 \\ 3.47 \\ 0.5 \\ -1.22 \end{bmatrix}$

```
CLASS torch.nn.Embedding(
```

```
CLASSMETHOD from_pretrained(
```

Тоже не нужно тренировать самим!

<https://code.google.com/archive/p/word2vec/>

<https://nlp.stanford.edu/projects/glove/>

<https://github.com/facebookresearch/fastText>

**ЧТО ЖЕ СЛУЧИТСЯ В СЛЕДУЮЩЕЙ
СЕРИИ???**

УЗНАЕМ В СЛЕДУЮЩЕЙ СЕРИИ!!!