# Deep text classification

CNN, RNN

Основано на [NLP Course | For You](#)
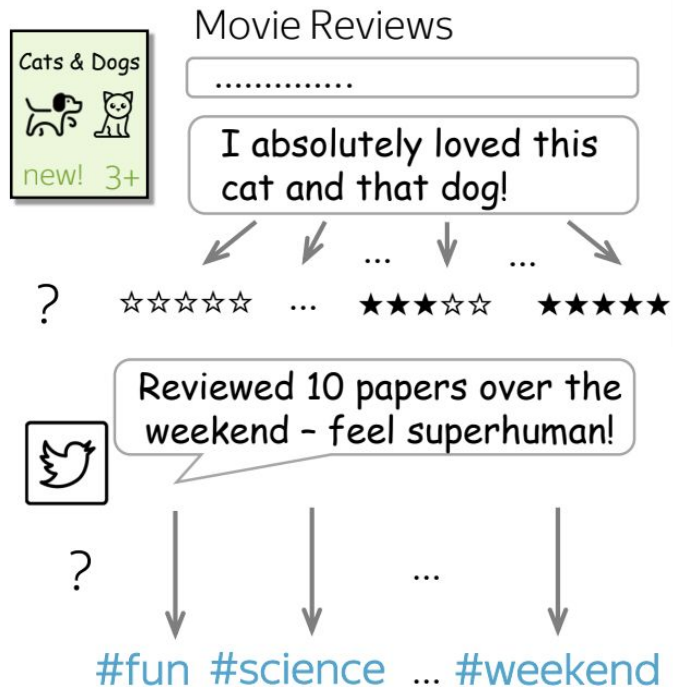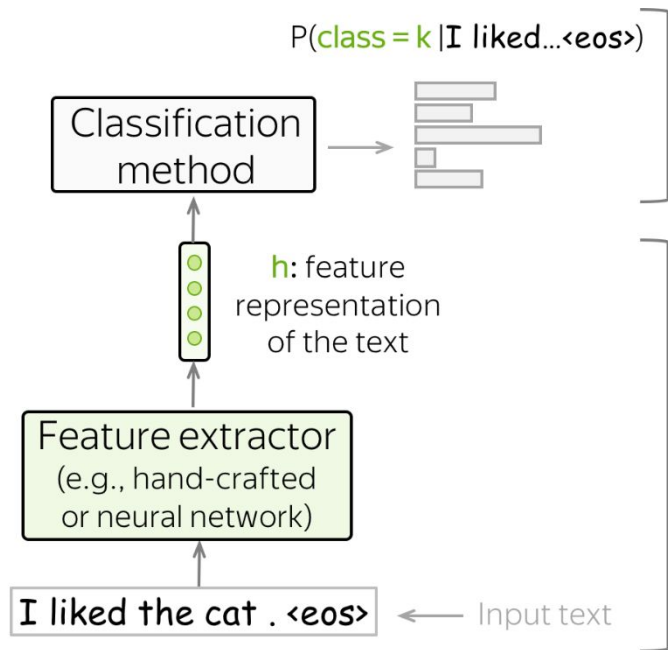
# План

- Примеры
- Последний слой
- Свёрточные архитектуры
- Recurent NN
- Embeddings

# Примеры

# Общая концепция



General Classification Pipeline

$P(\text{class} = k \mid \text{I liked...<eos>})$

Classification
method

**h**: feature representation of the text

Feature extractor
(e.g., hand-crafted or neural network)

I liked the cat . <eos> ← Input text

get probability distribution over classes

process text (document)

Classification with Neural Networks

**K** classes

$P(\text{class} = k \mid ...)$

**d**-sized vector

Linear layer

softmax

**h**: feature representation of the text

Neural Network

Token embeddings

I liked the cat . <eos> ← Input text

# Последний слой

# Loss-функция

Training example: **I liked the cat on the mat <eos>**          Label: k

target ↑

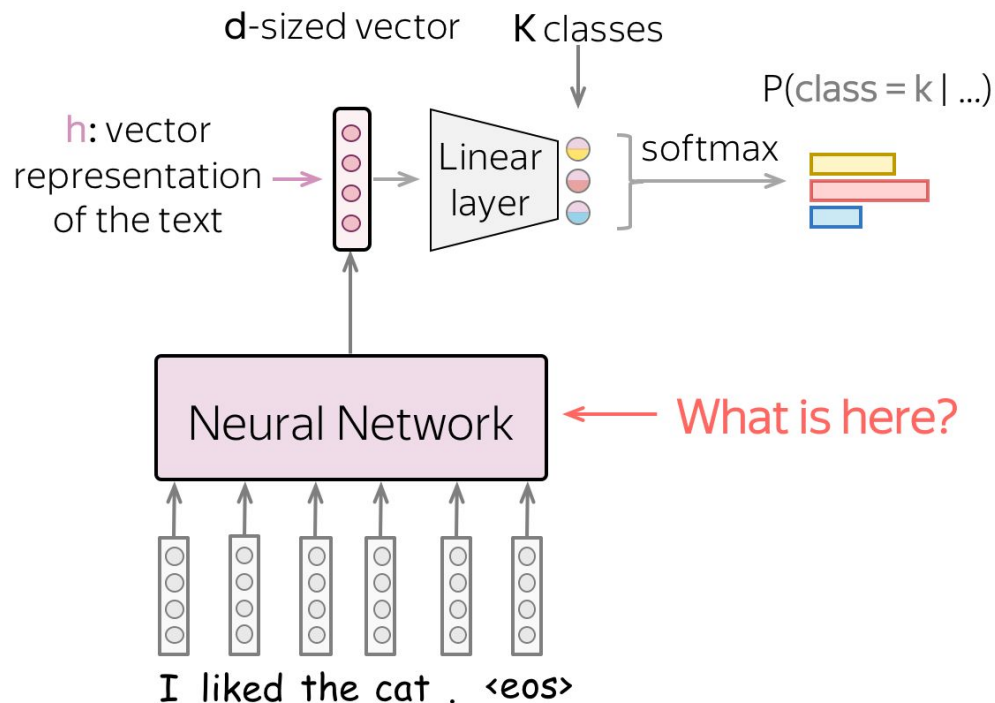Model prediction:                Target: $p^*$          Cross-entropy loss:

P(class = i |I liked...<eos>)



k

$$-\sum_{i=1}^{K} p_i^* \cdot \log P(y = i|x) \to min \quad (p_k^* = 1, p_i^* = 0, i \neq k)$$
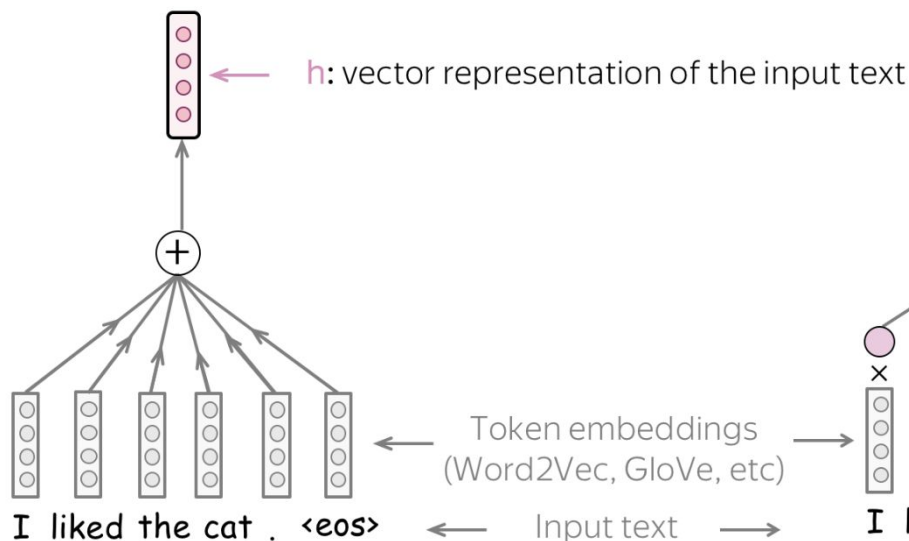
For one-hot targets, this is equivalent to
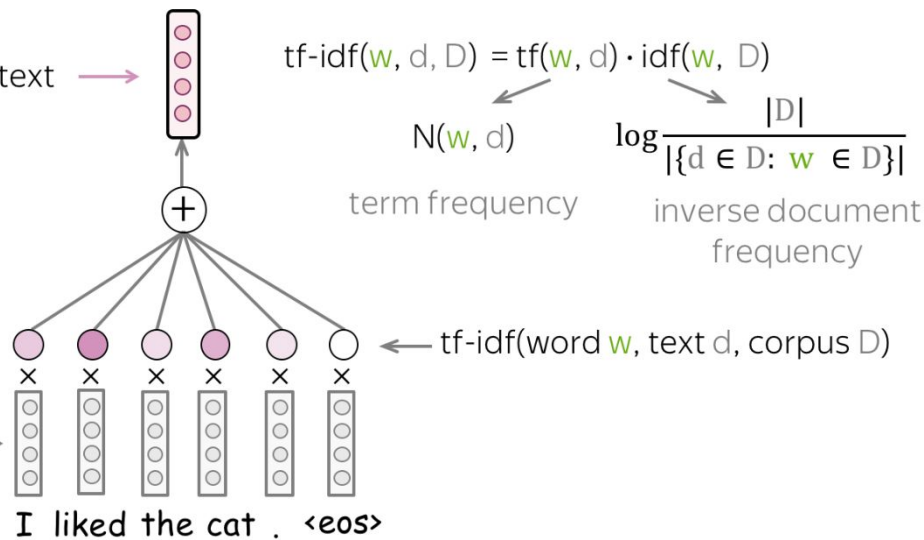
$$-\log P(y = k|x) \to min$$
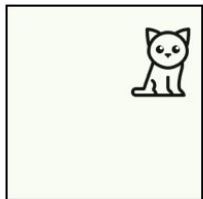
# А что внутри?

# №0. Усреднение эмбедингов
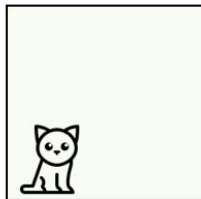
## Sum of embeddings
(Bag of Words, Bag of Embeddings)
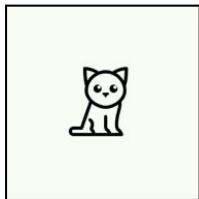
## Weighted sum of embeddings
(e.g., using tf-idf weights)

$h$: vector representation of the input text

$\text{tf-idf}(w, d, D) = \text{tf}(w, d) \cdot \text{idf}(w, D)$

$N(w, d)$ — term frequency

$\log \dfrac{|D|}{|\{d \in D:\ w \in D\}|}$ — inverse document frequency

$\text{tf-idf}(\text{word } w, \text{ text } d, \text{ corpus } D)$

Token embeddings
(Word2Vec, GloVe, etc)

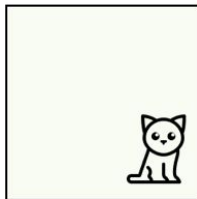I liked the cat . <eos>

Input text

I liked the cat . <eos>

# CNN

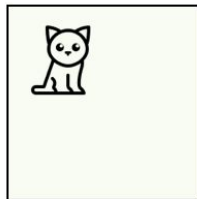Label: cat    Label: cat    Label: cat    Label: cat    Label: cat

We don't care where the cat is, we care that it is somewhere.
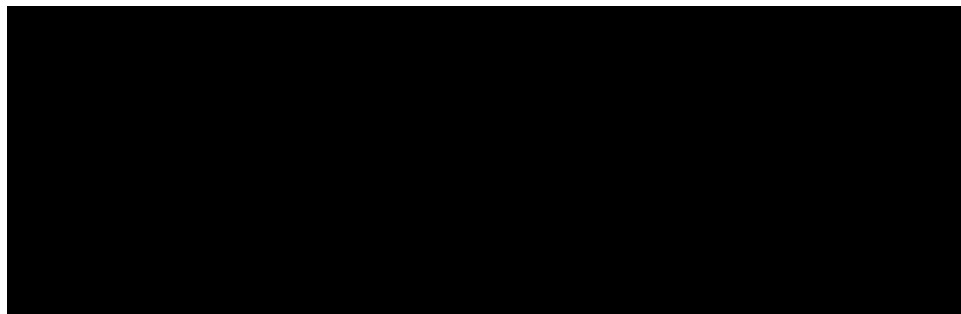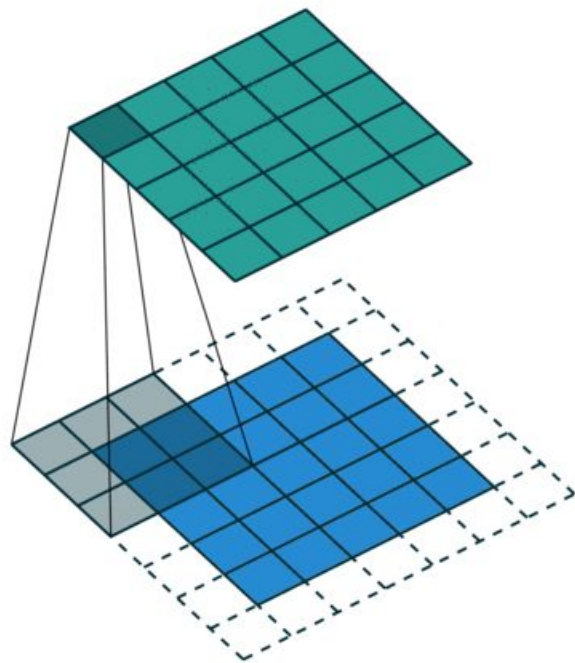
Then why don't we process all these cats similarly?

An absolutely great movie! I watched the premiere with my friends.

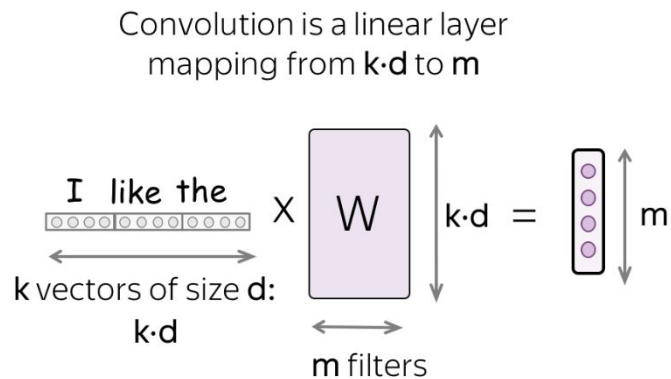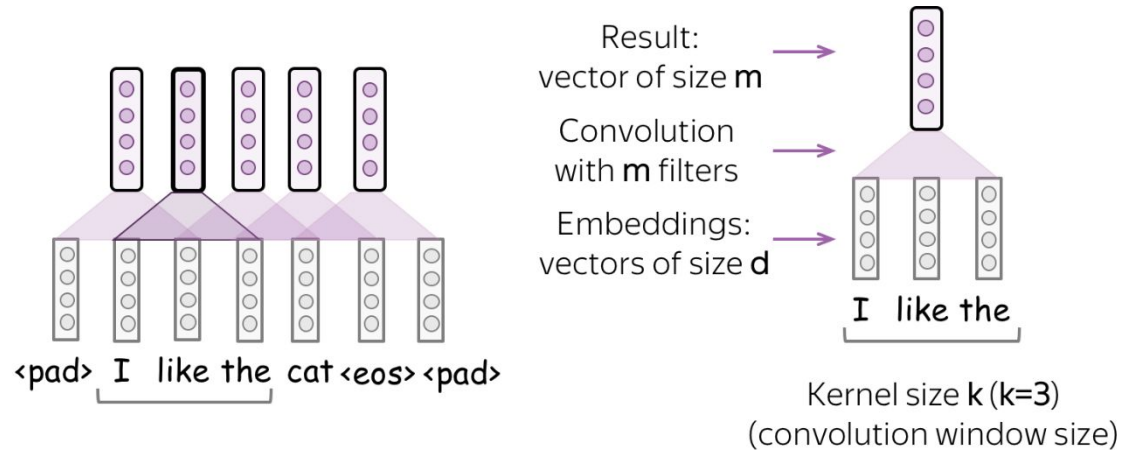The movie about cats was absolutely great, and the cats were cute.

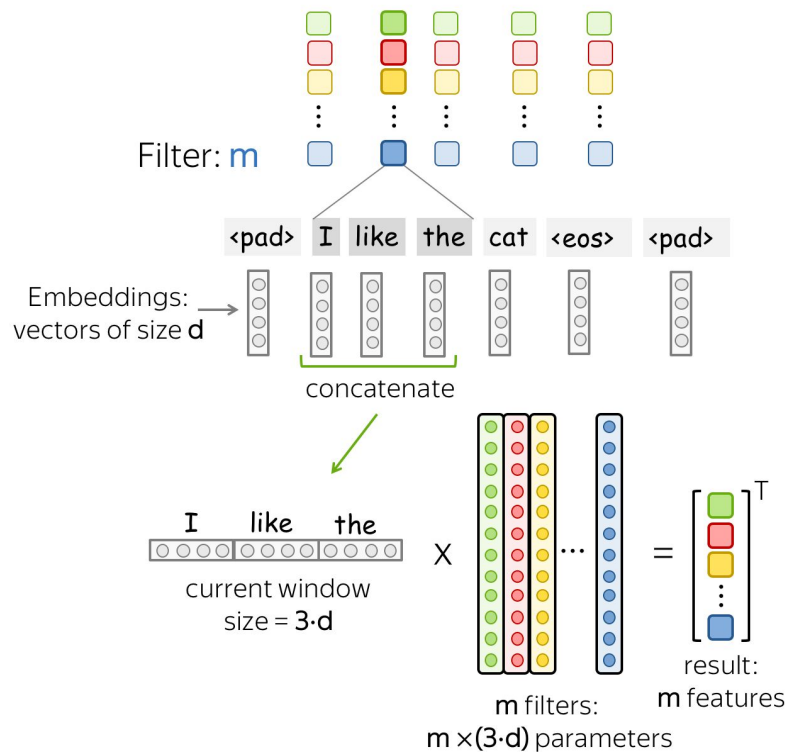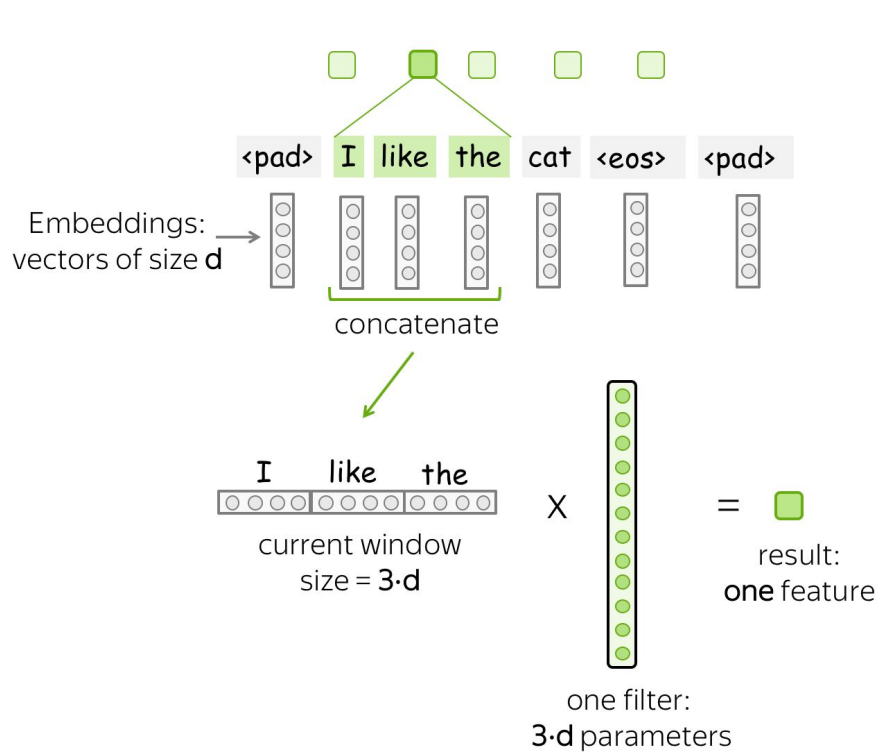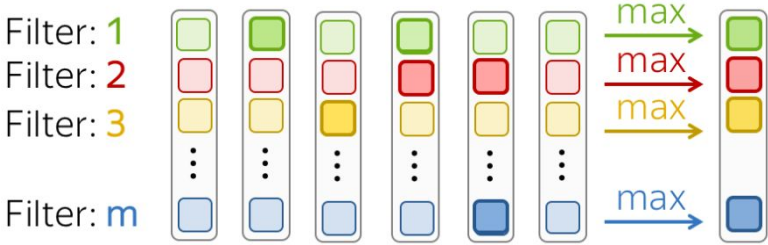The movie is about cats running around, and it is absolutely great.

# CNN

# CNN



Result: vector of size $m$

Convolution with $m$ filters

Embeddings: vectors of size $d$

`<pad>` I like the cat `<eos>` `<pad>`

Kernel size $k$ ($k=3$)
(convolution window size)

I like the

Convolution is a linear layer mapping from $k \cdot d$ to $m$

I like the

$\times$  W  $k \cdot d$  $=$  $m$

$k$ vectors of size $d$: $k \cdot d$

$m$ filters

# CNN

# CNN

# CNN

# CNN

# CNN



This is our **h**: vector representation of the text

concatenate representations from different convolutions

max-over-time pooling

max    max    max

Convolution

&lt;pad&gt; I like the cat &lt;eos&gt; &lt;pad&gt;    &lt;pad&gt; I like the cat &lt;eos&gt; &lt;pad&gt;    &lt;pad&gt; I like the cat &lt;eos&gt; &lt;pad&gt;

kernel size=2    kernel size=3    kernel size=4

# CNN



Global Pooling (max-over-time)

Convolution

Pooling

Convolution

The final is always global pooling: you need it to get a single vector

Block convolution+pooling
You can stack several of them!

<pad> I like the cat on the red mat . <eos> <pad>

# RNN
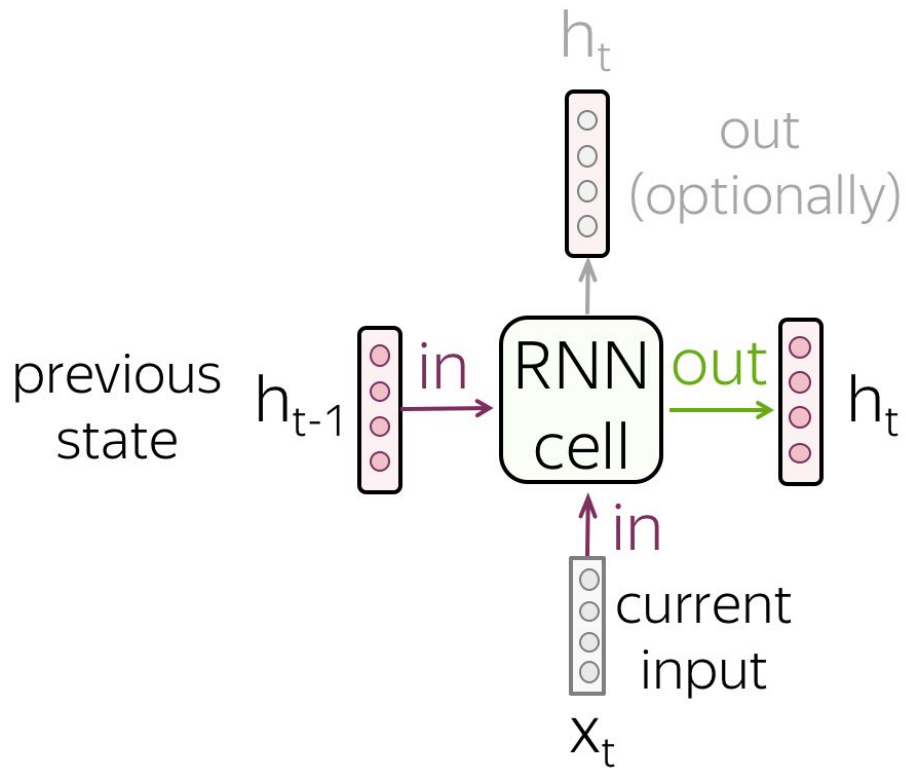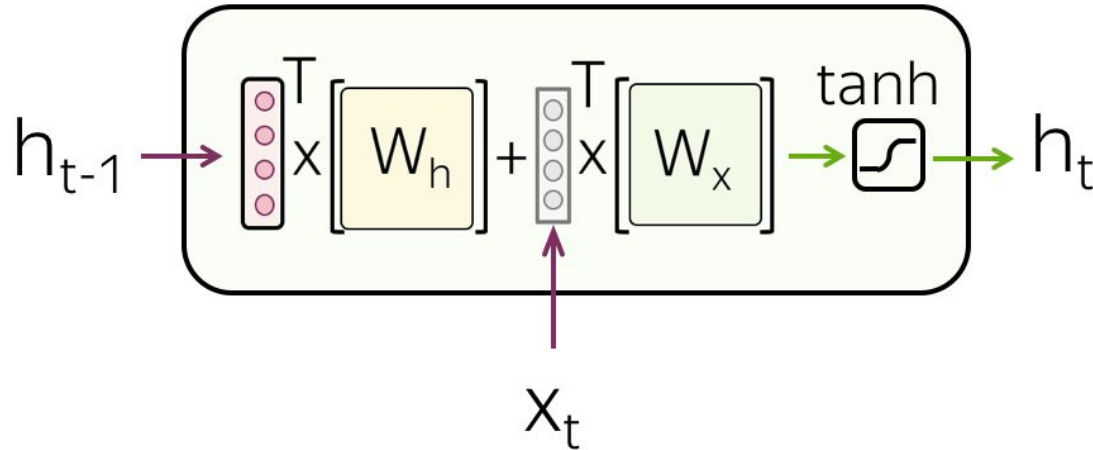
# RNN



## Vanilla RNN

$$h_t = \tanh(h_{t-1} W_h + x_t W_x)$$

# Twp layer RNN



This is our h: vector representation of the text

2nd layer RNN

$h_0^{(2)}$ $h_1^{(2)}$ $h_2^{(2)}$ $h_3^{(2)}$ $h_4^{(2)}$ $h_5^{(2)}$

Final state of the 2nd-layer RNN

$h_1^{(1)}$ $h_2^{(1)}$ $h_3^{(1)}$ $h_4^{(1)}$ $h_5^{(1)}$

1st layer RNN

$h_0^{(1)}$ $h_1^{(1)}$ $h_2^{(1)}$ $h_3^{(1)}$ $h_4^{(1)}$

Final state of the 1st-layer RNN

I   like   the   cat   <eos>

# Bidirectional RNN

This is our h: vector representation of the text

Concatenate representations from forward and backward RNN

Initial RNN state (e.g., zero vector)

Backward RNN

Final state of the backward RNN

Forward RNN

Final state of the forward RNN

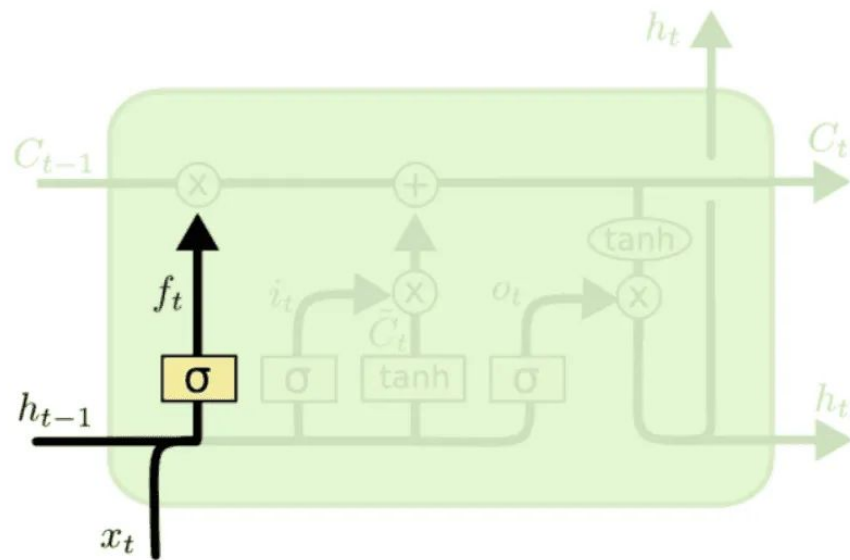Initial RNN state (e.g., zero vector)
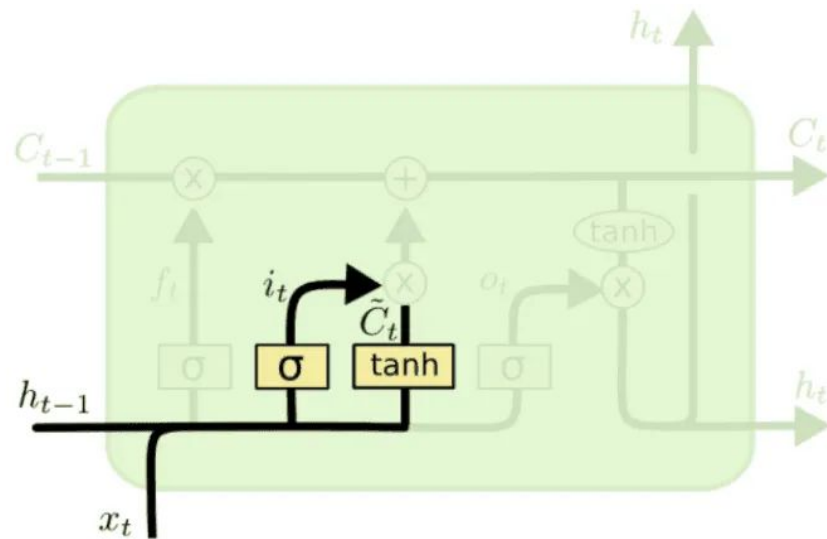
<bos>   I   like   the   cat   <eos>

# LSTM



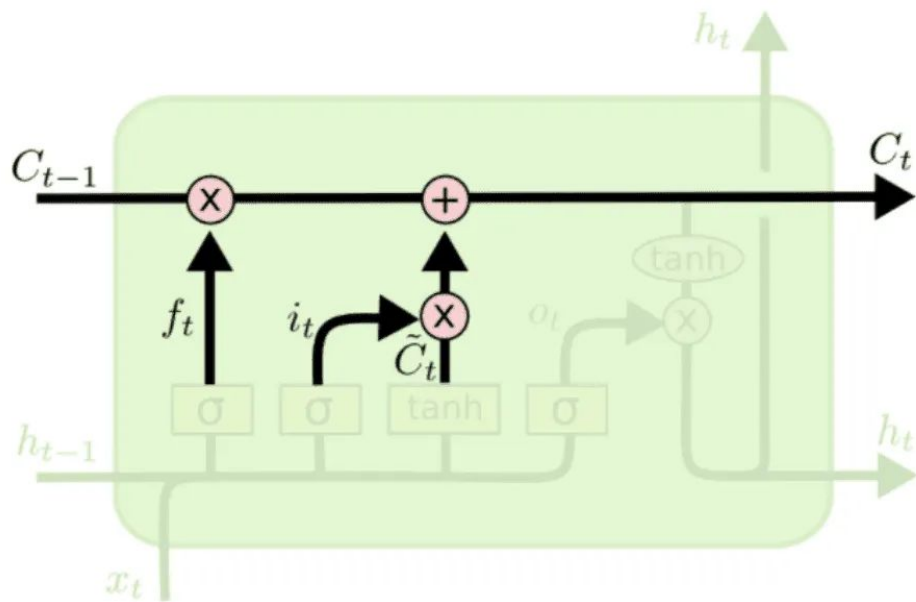$$f_t = \sigma(h_{t-1}W_1^f + x_t W_2^f + b_f)$$

# LSTM



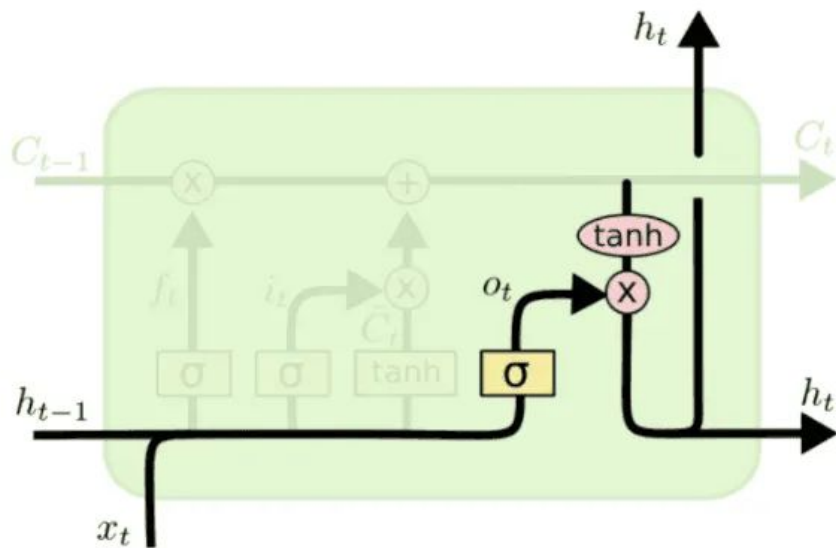$$\tilde{C}_t = \tanh(h_{t-1}W_1^C + x_t W_2^C + b_c)$$

$$i_t = \sigma(h_{t-1}W_1^i + x_t W_2^i + b_i)$$

# LSTM



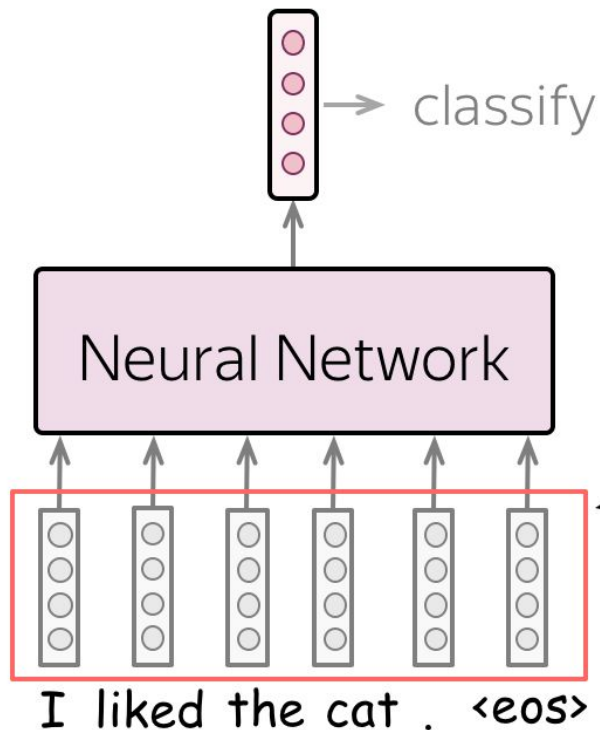$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

# LSTM



$$o_t = \sigma(h_{t-1}W_1^o + x_t W_2^o + b_o)$$
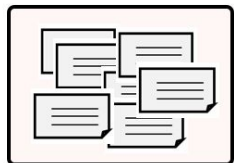
$$h_t = o_t \odot \tanh(c_t)$$

# Embeddings



classify

Neural Network

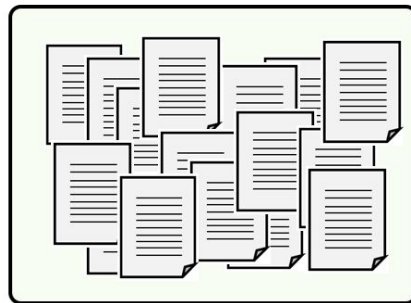I liked the cat . <eos>

Input word embeddings:

- Train from scratch
- Take pretrained (Word2Vec, GloVe)
- Initialize with pretrained, then fine-tune

# Embeddings

Training data for text classification (labeled)

- Not huge, or not diverse, or both
- Domain: task-specific

Training data for word embeddings (unlabeled)

- Huge diverse corpus (e.g., Wikipedia)
- Domain: general
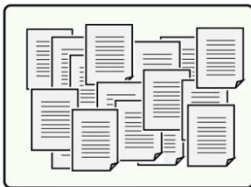
# Embeddings

- Train from scratch

  What they will know:

  May be not enough to learn relationships between words
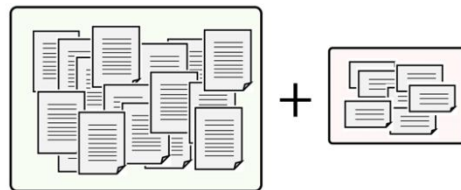
- Take pretrained (Word2Vec, GloVe)

  What they will know:

  Know relationships between words, but are **not** specific to the task

- Initialize with pretrained, then fine-tune

  What they will know:

  Know relationships between words and adapted for the task

"**Transfer**" knowledge from a huge unlabeled corpus to your task-specific model