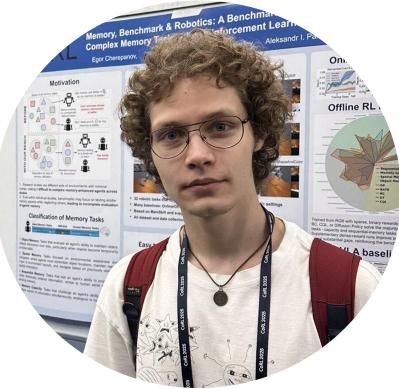


Transformers in Deep Learning

Lecture 1

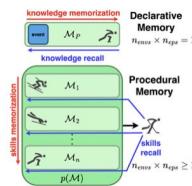
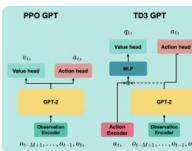
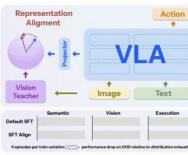
About me:



Nikita Kachaev
AIRI, Research Engineer
(ex Tinkoff, VK)

- Multimodal models, RL,
Robotics, Foundation
models

<https://tttonyalpha.github.io/>



Don't Blind Your VLA: Aligning Visual Representations for OOD Generalization

Nikita Kachaev, Mikhail Kolosov, Daniil Zelezetsky, Alexey K. Kovalev, Aleksandr I. Panov

A* conference (TBA), Oral, 2025

★ Oral ★

[Project Page](#) · [PDF](#) · [Code](#)

Memory, Benchmark & Robots: A Benchmark for Solving Complex Tasks with Reinforcement Learning

Egor Cherepanov, **Nikita Kachaev**, Alexey K. Kovalev, Aleksandr I. Panov

Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track, Spotlight, 2025

★ Spotlight ★

[Project Page](#) · [PDF](#) · [Code](#)

A New Perspective on Transformers in Online Reinforcement Learning for Continuous Control

Nikita Kachaev, Daniil Zelezetsky, Egor Cherepanov, Alexey K. Kovalev, Aleksandr I. Panov

The Thirteenth International Conference on Learning Representations (ICLR), 7th Robot Learning Workshop, Poster, 2025

[PDF](#)

Mind and Motion Aligned: A Joint Evaluation IsaacSim Benchmark for Task Planning and Low-Level Policies in Mobile Manipulation

Nikita Kachaev, Andrei Spiridonov, Andrey Gorodetsky, Kirill Muravyev, Nikita Oskolkov, Aditya

Narendra, Vlad Shakhuro, Dmitry Makarov, Aleksandr I. Panov, Polina Fedotova, Alexey K. Kovalev

Reports of the RAS: Mathematics, Informatics, and Control Processes, paper, 2025

[PDF](#)

Unraveling the Complexity of Memory in RL Agents: an Approach for Classification and Evaluation

Egor Cherepanov, **Nikita Kachaev**, Artem Zholus, Alexey K. Kovalev, Aleksandr I. Panov

Conference on Robot Learnin (CORL), RemembeRL Workshop, Poster, 2025

[PDF](#)

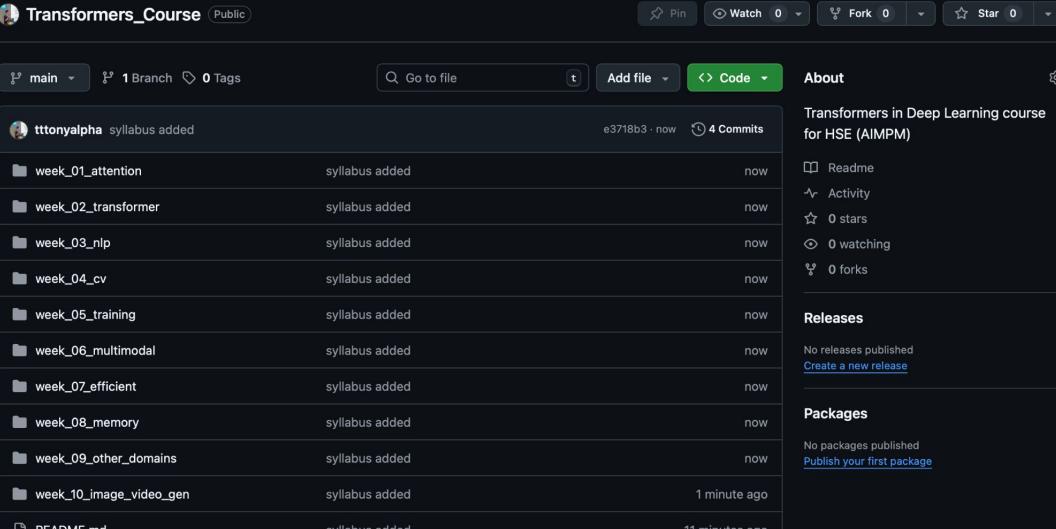
You will learn:

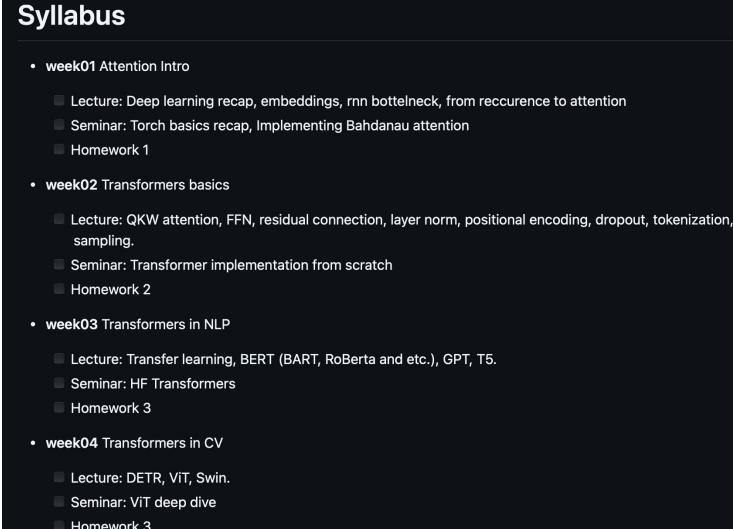
- How Transformers work
- Applications beyond NLP
- Modern Transformers best practices



Course:

Final grade =
 $0.7 \times \text{Homework} + 0.3 \times \text{Exam}$

A screenshot of a GitHub repository page for "Transformers_Course". The repository is public and has 1 branch and 0 tags. The main file list shows several folders named week_01_attention through week_10_image_video_gen, each containing a "syllabus added" file. A commit by user tttonyalpha is visible, adding the syllabus files. The repository has 4 commits and is last updated 11 minutes ago. The repository page includes sections for About, Releases, and Packages.

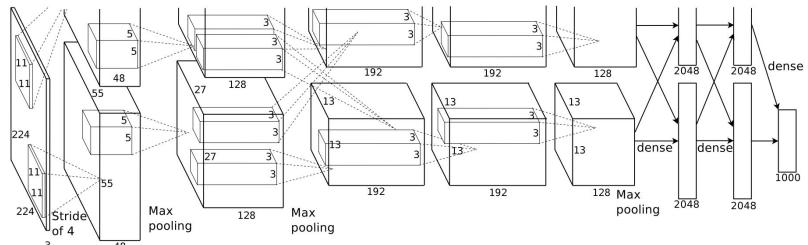
A screenshot of the syllabus page for the "Transformers_Course". The syllabus is organized into four main sections: week01 Attention Intro, week02 Transformers basics, week03 Transformers in NLP, and week04 Transformers in CV. Each section lists a lecture, a seminar, and a homework assignment.

- week01 Attention Intro
 - Lecture: Deep learning recap, embeddings, rnn bottleneck, from recurrence to attention
 - Seminar: Torch basics recap, implementing Bahdanau attention
 - Homework 1
- week02 Transformers basics
 - Lecture: QKW attention, FFN, residual connection, layer norm, positional encoding, dropout, tokenization, sampling.
 - Seminar: Transformer implementation from scratch
 - Homework 2
- week03 Transformers in NLP
 - Lecture: Transfer learning, BERT (BART, RoBERTa and etc.), GPT, T5.
 - Seminar: HF Transformers
 - Homework 3
- week04 Transformers in CV
 - Lecture: DETR, ViT, Swin.
 - Seminar: ViT deep dive
 - Homework 3

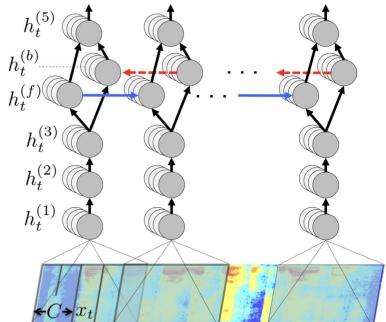
https://github.com/tttonyalpha/Transformers_Course

Deep learning before transformers:

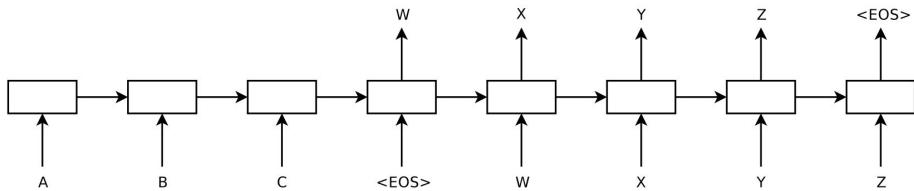
CV (AlexNet)



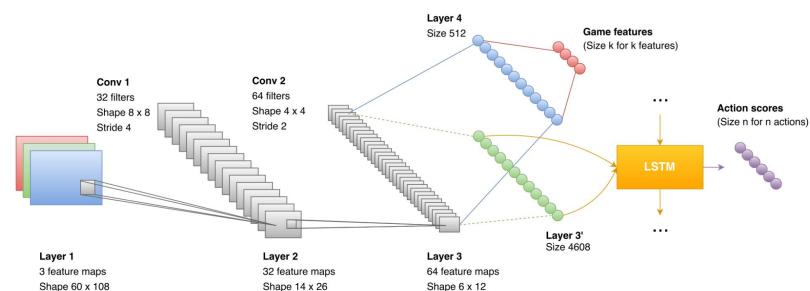
Speech (DeepSpeech)



NLP (RNNs)

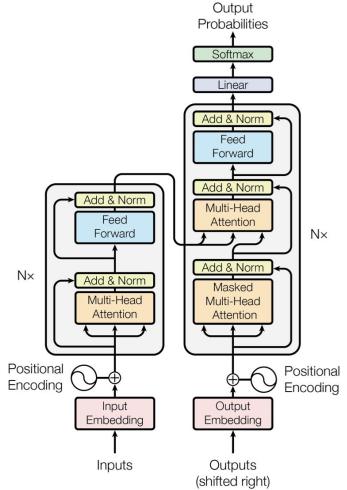


Robotics/RL (DRQN)

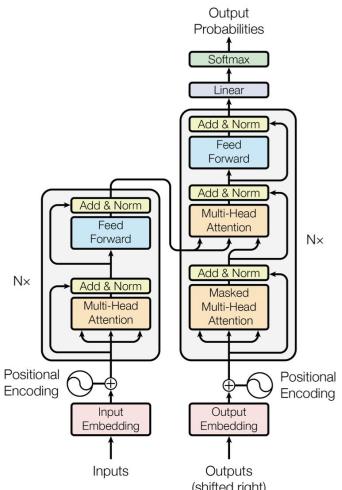


Deep learning after transformers:

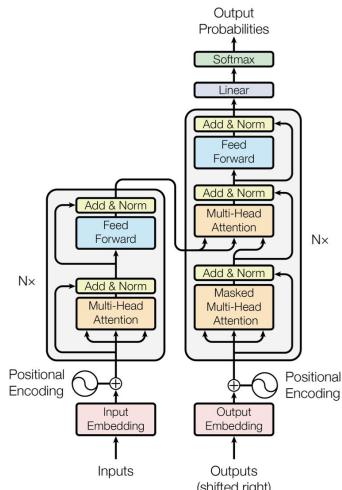
CV



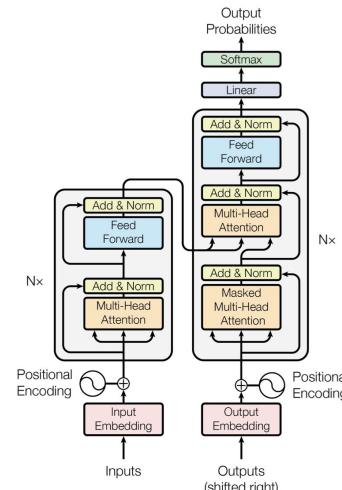
NLP



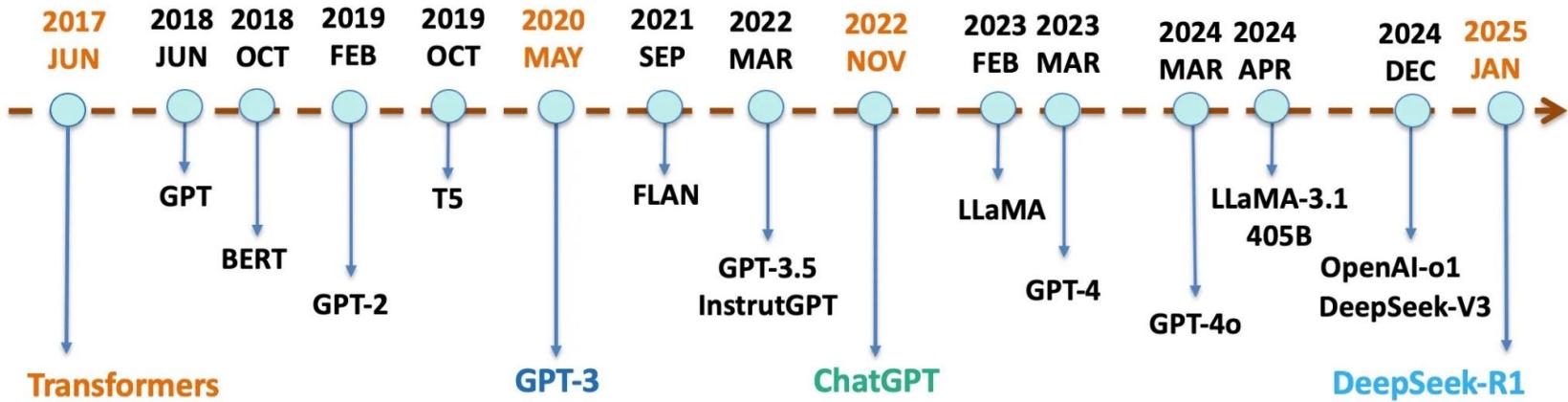
Speech



Robotics/RL



Transformers timeline:



Transformers: why do people use them so much?



Andrej Karpathy ✅

@karpathy

The Transformer is a magnificent neural network architecture because it is a general-purpose differentiable computer. It is simultaneously:

- 1) expressive (in the forward pass)
- 2) optimizable (via backpropagation+gradient descent)
- 3) efficient (high parallelism compute graph)

9:54 PM · Oct 19, 2022



...



Andrej Karpathy ✅ @karpathy · Oct 19, 2022

(3) because the compute graph is shallow and wide, mapping significantly better to our high-parallelism compute architectures (think GPUs). An earlier attempt that understood the significance and optimized for this property was the Neural GPU paper (arxiv.org/abs/1511.08228)

2

10

258



...



Andrej Karpathy ✅ @karpathy · Oct 19, 2022

Its success lies in a single architecture that simultaneously satisfies all of these properties. The original Attention Is All You Need paper is a bit haphazard and undersells the magnitude of these insights, their history and motivations. But there's a lot going on :)

10

9

317



...



Andrej Karpathy ✅ @karpathy · Oct 19, 2022

So I probably would have called the paper something like "Transformer: A general-purpose, efficient, optimizable computer" and presented it alongside the Neural Turing Machine, NeuralGPU and friends, then applied it to translation as an example. Something like that, but ok :)

20

25

443



...



Andrej Karpathy ✅ @karpathy · Oct 19, 2022

A few people have (correctly) pointed out the hindsight here, which is fair. I don't suspect the authors would have known that 5 years later that architecture will have taken over most of AI ~unchanged, except for a re-shuffling of layernorms. Calls for a followup paper :)

12

16

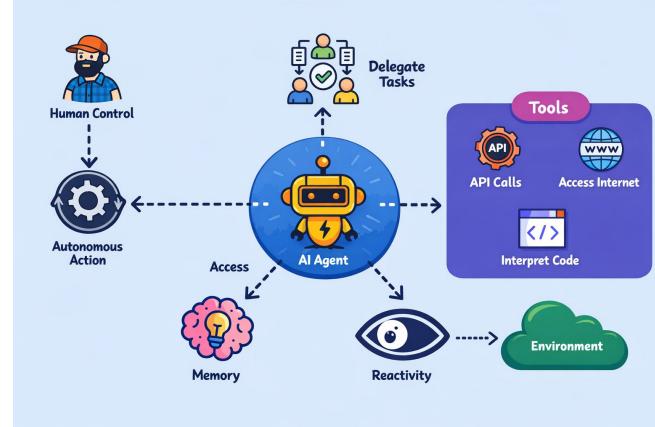
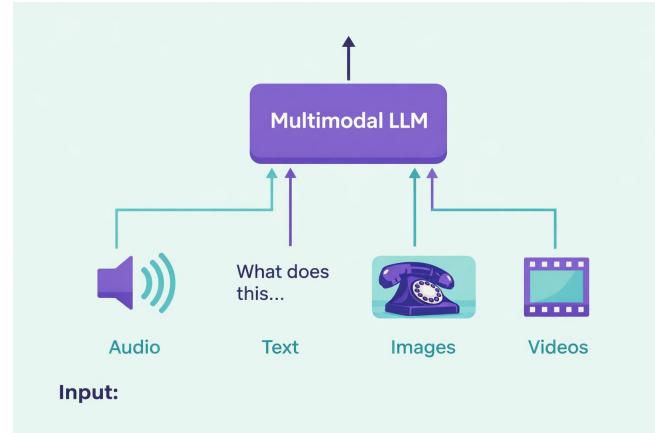
375



<https://x.com/karpathy/status/1582807367988654081?s=20>

Where we are (2026)

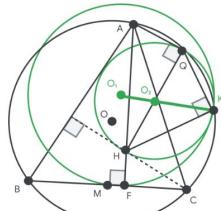
- **Multimodality** (text + image + audio + video).
- **Long-context working memory.**
- **Reasoning improvements** are increasingly driven by test-time scaling.
- **Agentic systems.**
- **Alignment & safety** moved “beyond classic RLHF”.



Alpha Geometry

IMO 2015 P3

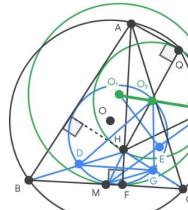
Let ABC be an acute triangle. Let (O) be its circumcircle, H its orthocenter, and F the foot of the altitude from A . Let M be the midpoint of BC . Let Q be the point on (O) such that $QH \perp QA$ and let K be the point on (O) such that $KH \perp KQ$. Prove that the circumcircles (O_1) and (O_2) of triangles FKM and KQH are tangent to each other.



AlphaGeometry

Solution

```
[...]
Construct D: midpoint BH [a]
[a], O2 midpoint HQ => BQ || O2D [20]
[...]
Construct G: midpoint HC [b]
∠GMD = ∠GO2D => M O2 G D cyclic [26]
[...]
[a],[b] => BC || DG [30]
[...]
Construct E: midpoint MK [c]
[c] => ∠KFC = ∠KO1E [104]
[...]
∠FKO1 = ∠FKO2 => KO1 || KO2 [109]
[109] => O1O2K collinear => (O1)(O2) tangent
```



Alpha Evolve

$$\max_{-1/2 \leq t \leq 1/2} \int_{\mathbb{R}} f_i(t-x) f_i(x) dx \geq \mathbb{C} \left(\int_{-1/4}^{1/4} f_i(x) dx \right)^2$$

1.5098 → 1.5053

$$\|f * f\|_2^2 \leq \mathbb{C}' \|f * f\|_1 \|f * f\|_\infty$$

0.8892 → 0.8962

$$\max_{-1/2 \leq t \leq 1/2} \left| \int_{\mathbb{R}} f(t-x) f(x) dx \right| \geq \mathbb{C}'' \left(\int_{-1/4}^{1/4} f(x) dx \right)^2$$

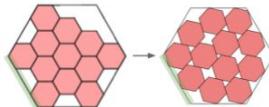
1.4581 → 1.4557

$$A(f) A(\hat{f}) \geq \mathbb{C}'''$$

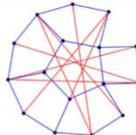
0.3523 → 0.3521

Analysis

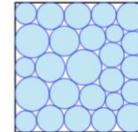
Hexagon outer edge
4.000 → 3.942



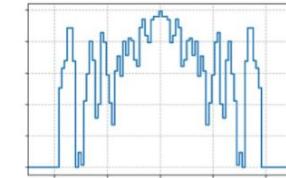
Max distance/min distance
12.890 → 12.889



Sum of radii
2.6340 → 2.6358



Geometry



$$\sup_{x \in [-2, 2]} \int_{-1}^1 f(t) g(x+t) dt \geq \mathbb{C}$$

0.380926 → 0.380924

$$|A+B| \ll |A|$$
$$|A-B| \gg |A|^{\mathbb{C}}$$

1.1446 → 1.1584

Combinatorics

Chat GPT

ChatGPT



Examples

"Explain quantum computing in simple terms"



Capabilities

"Got any creative ideas for a 10 year old's birthday?"

Allows user to provide follow-up corrections

"How do I make an HTTP request in Javascript?"

Trained to decline inappropriate requests



Limitations

May occasionally generate incorrect information

May occasionally produce harmful instructions or biased content

Limited knowledge of world and events after 2021

NVIDIA Cosmos-Reasong1

EMBODIED REASONING BENCHMARK



Based on the ego vehicle's actions in the video, what is the most likely immediate action it will take next?

- A: turn right, B: turn left, C: change to right lane, D: change to left lane

Wait, maybe the question is a trick. The video shows the road markings are double yellow lines, which means no passing, so changing lanes is prohibited. Therefore, the driver cannot change lanes, so options C and D are impossible. The only remaining options are turning right or left. But since the video doesn't show any turns, and the driver is maintaining their path, perhaps the answer is that none of the options are correct, but since we have to choose, I have no choice but to select none. Therefore, the answer is none.

Generative models

Kandinsky 5.0: A Family of Foundation Models for Image and Video Generation

Kandinsky Lab*

* A detailed list of the contributors can be found in the end of this paper.

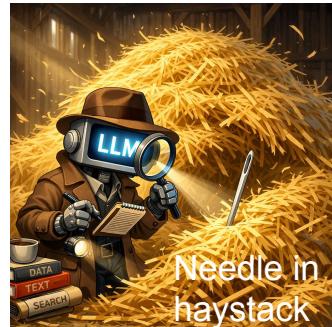


The Future

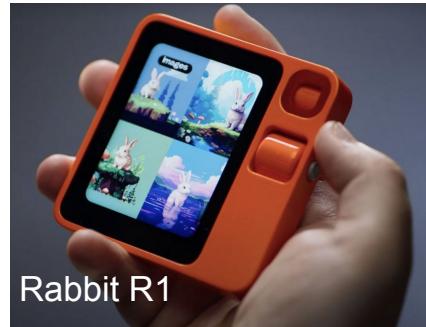
- **Reliable memory (not just long context).**
Robust retrieval, grounded citations.
- **Always-on Multimodal Assistants.** Real-time voice + vision assistants for work, education, accessibility, and robotics.
- **Verification and truthfulness at scale.** Built-in self-checking, external validators. Knowing when it doesn't know, and escalating correctly.
- **Grounding to the physical world.** More robust world models.
- **Scientific & engineering acceleration.**



1x Neo Home Robot



Needle in haystack



Rabbit R1

Key idea: embed everything

Text (token)

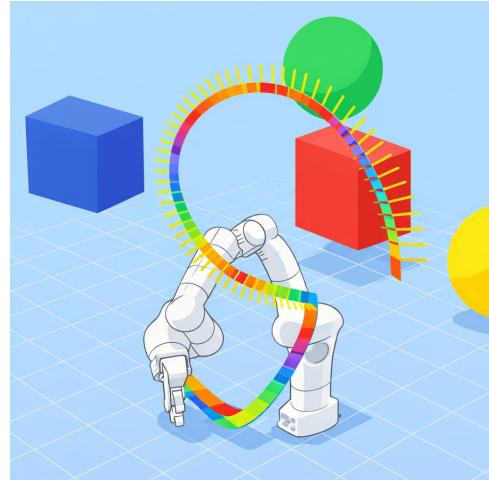


464 3290 25365 262 22514 198

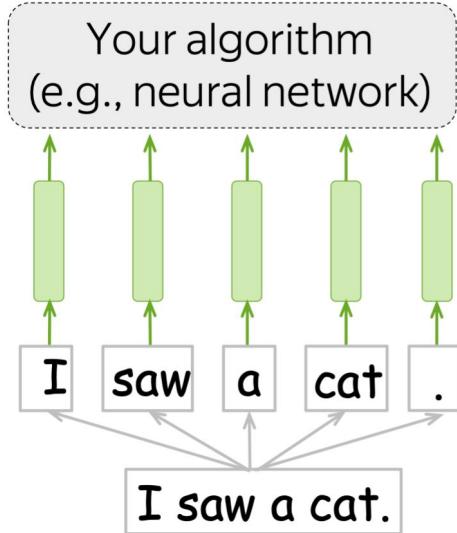
Image (patch)



Robot trajectory (timestep)



Word Embeddings



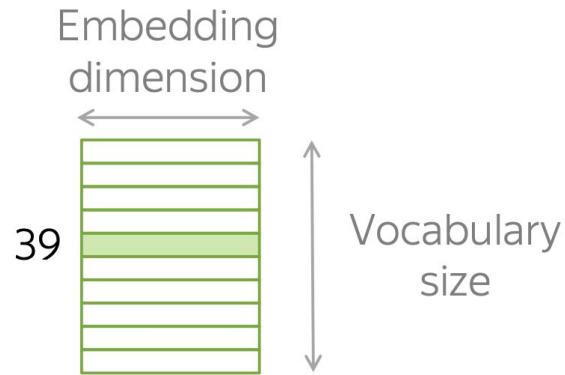
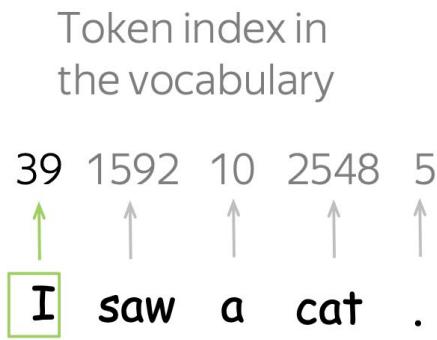
Any algorithm for solving a task

Word representation - vector
(input for your model/algorithm)

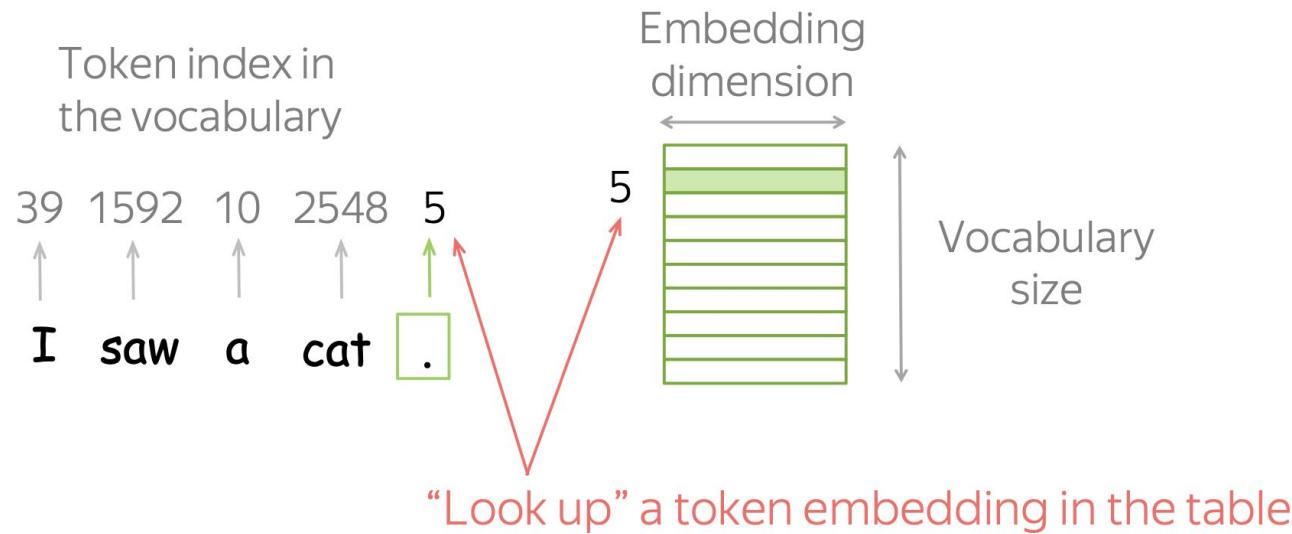
Sequence of tokens

Text (your input)

Word Embeddings



Word Embeddings



One-Hot Embeddings

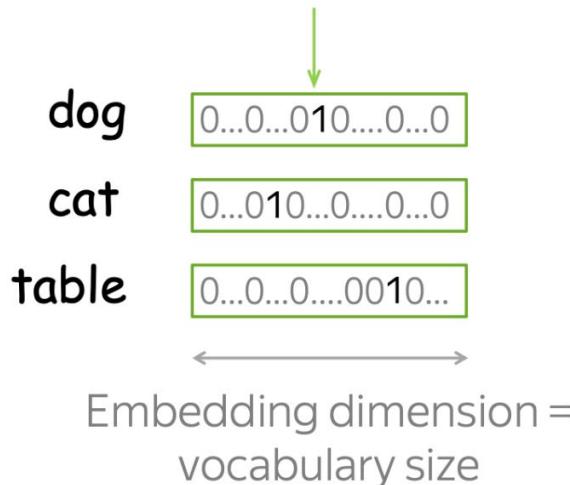
One is 1, the rest are 0



Embedding dimension =
vocabulary size

One-Hot Embeddings

One is 1, the rest are 0

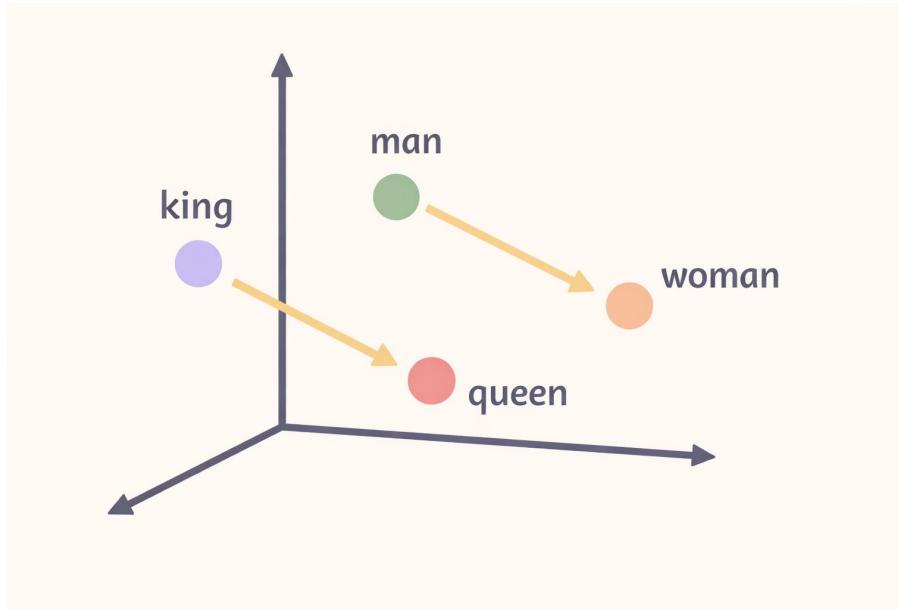


Problems:

- Vector size is too large
- Vectors know nothing about meaning

e.g., **cat** is as close to
dog as it is to **table**!

Word Embeddings: Semantics relations



Word2Vec Embeddings

We have to put information about contexts into word vectors.

How: Learn word vectors by teaching them to **predict contexts**.

Word2Vec Embeddings

... I saw a cute grey cat playing in the garden ...

Word2Vec Embeddings

... I saw a cute grey cat playing in the garden ...

w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}

Word2Vec Embeddings

... **I** **saw** **a** **cute** **grey** cat playing in the garden ...

w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}

context words central word context words

Word2Vec Embeddings

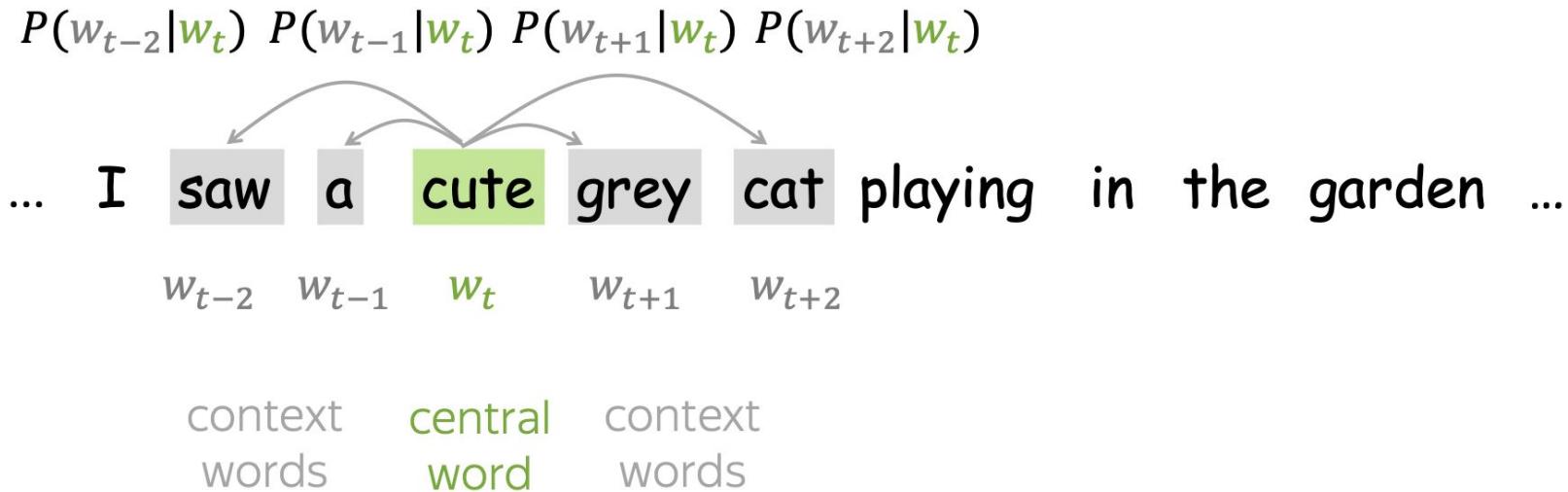
$$P(w_{t-2}|w_t) \ P(w_{t-1}|w_t) \ P(w_{t+1}|w_t) \ P(w_{t+2}|w_t)$$



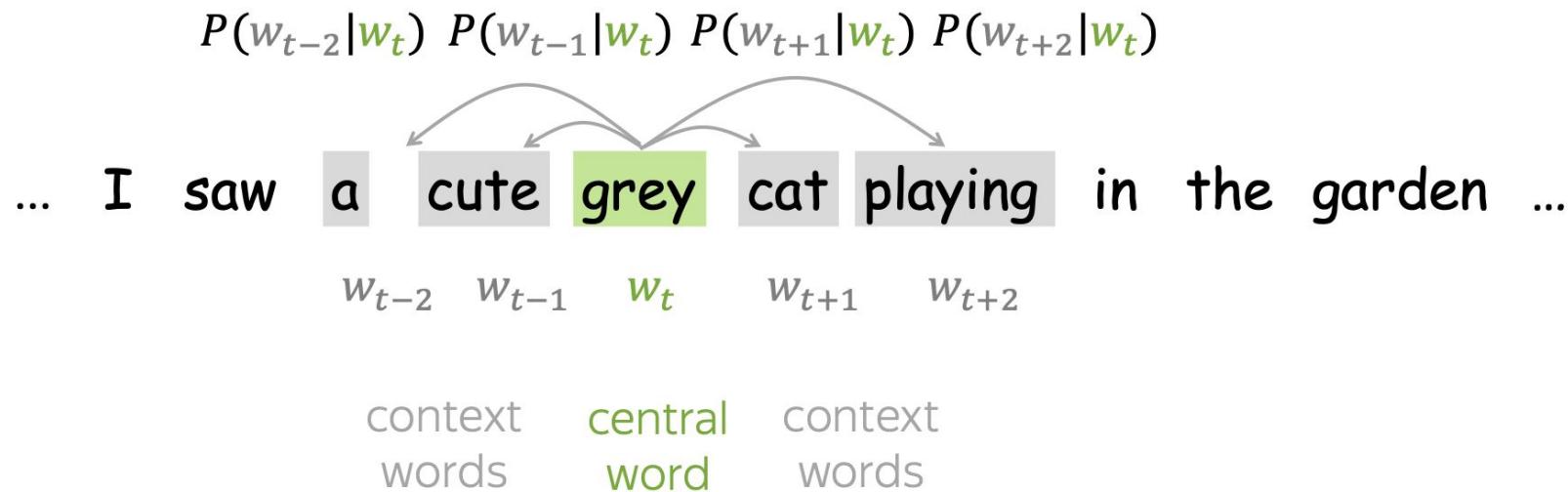
$w_{t-2} \quad w_{t-1} \quad w_t \quad w_{t+1} \quad w_{t+2}$

context central context
words word words

Word2Vec Embeddings



Word2Vec Embeddings



Likelihood of all text corpus

We want our model to think that the training data is “likely”

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m, \\ j \neq 0}} P(w_{t+j} | w_t, \theta)$$

Likelihood of all text corpus

We want our model to think that the training data is “likely”

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m, \\ j \neq 0}} P(w_{t+j} | w_t, \theta)$$

θ are all variables to optimize

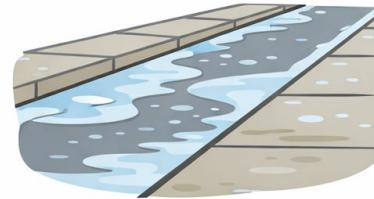
Likelihood

Likelihood

How well a guess explains observed data



Weather forecast



Wet streets

High likelihood it rained

Word2Vec Loss

$$\text{Loss} = -\frac{1}{T} \log L(\theta)$$

Word2Vec Loss

$$\text{Loss} = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

Word2Vec Loss

$$\text{Loss} = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

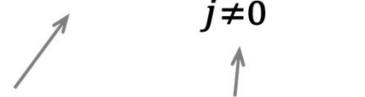


go over text

Word2Vec Loss

$$\text{Loss} = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

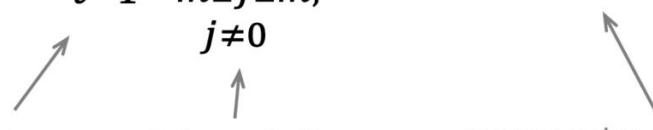
go over text with a sliding window



Word2Vec Loss

$$\text{Loss} = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | \mathbf{w}_t, \theta)$$

go over text with a sliding window compute probability of the context word given the central



Word2Vec Loss

$$\text{Loss} = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

How to compute this?



go over text

with a sliding window

compute probability of the context word given the central

Softmax function:

Softmax function $\mathbb{R}^n \rightarrow \mathbb{R}^n$:

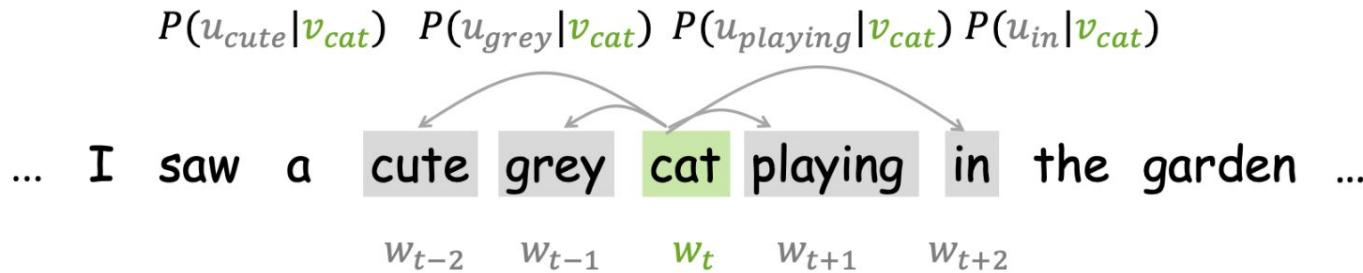
$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

- maps arbitrary values x_i to a probability distribution p_i
- “max” because amplifies probability of largest x_i
- “soft” because still assigns some probability to smaller x_i

Probability of context word

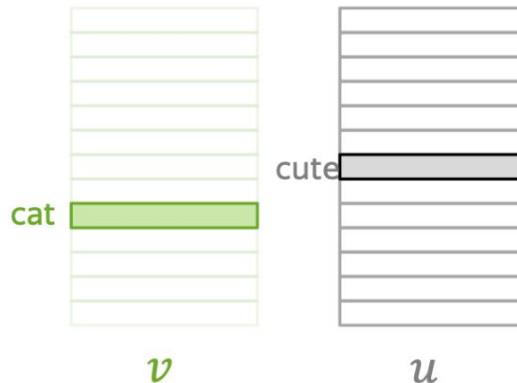
$$P(o|\textcolor{brown}{c}) = \frac{\exp(u_o^T \textcolor{brown}{v}_c)}{\sum_{w \in V} \exp(u_w^T \textcolor{brown}{v}_c)}$$

One Training step



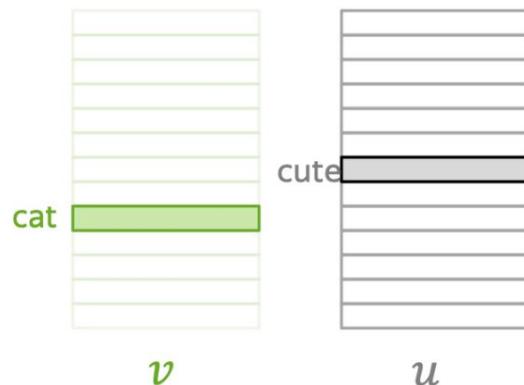
Look at the loss component for this step:

$$-\log P(cute|cat) = -\log \frac{\exp(u_{cute}^T \mathbf{v}_{cat})}{\sum_{w \in V} \exp(u_w^T \mathbf{v}_{cat})}$$



One Training step

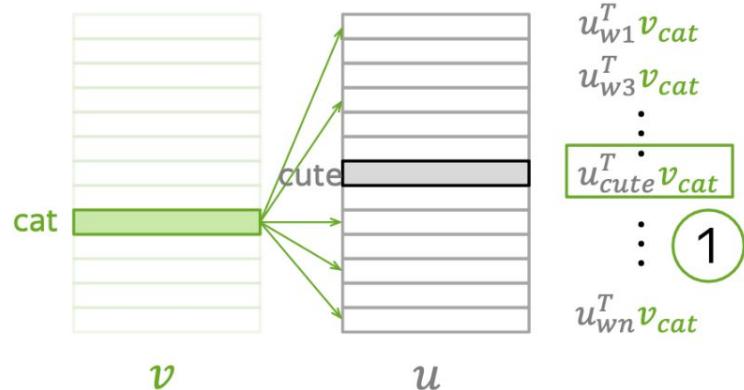
$$-\log P(\text{cute}|\text{cat}) = -\log \frac{\exp(u_{\text{cute}}^T v_{\text{cat}})}{\sum_{w \in V} \exp(u_w^T v_{\text{cat}})} = -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T v_{\text{cat}})$$



One Training step

$$-\log P(\text{cute}|\text{cat}) = -\log \frac{\exp(u_{\text{cute}}^T v_{\text{cat}})}{\sum_{w \in V} \exp(u_w^T v_{\text{cat}})} = -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T v_{\text{cat}})$$

1. Take dot product of v_{cat} with all u

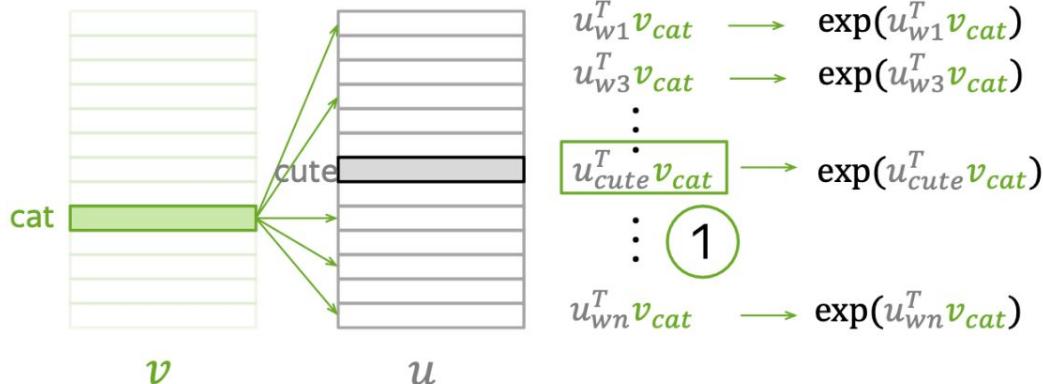


One Training step

$$-\log P(\text{cute}|\text{cat}) = -\log \frac{\exp(u_{\text{cute}}^T v_{\text{cat}})}{\sum_{w \in V} \exp(u_w^T v_{\text{cat}})} = -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T v_{\text{cat}})$$

1. Take dot product of v_{cat} with all u

2. exp



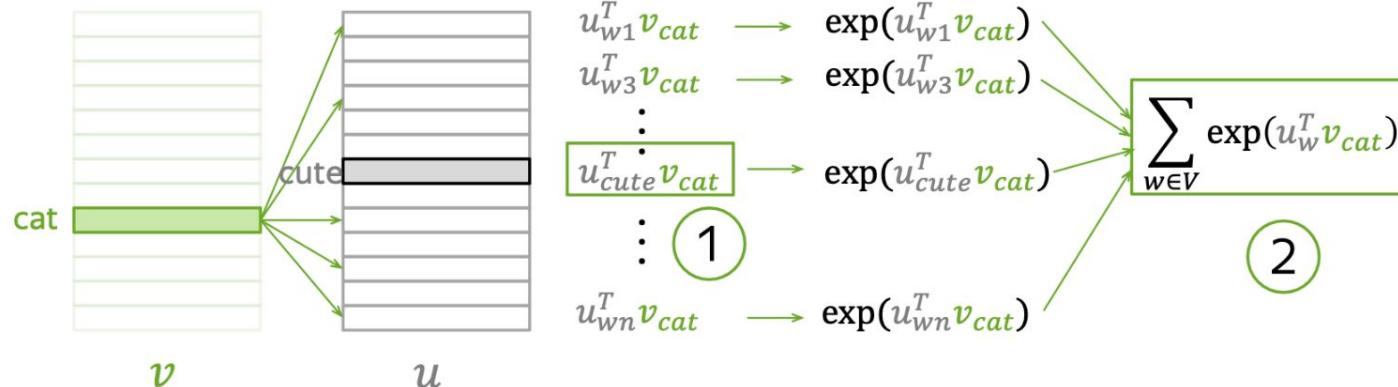
One Training step

$$-\log P(\text{cute}|\text{cat}) = -\log \frac{\exp(u_{\text{cute}}^T v_{\text{cat}})}{\sum_{w \in V} \exp(u_w^T v_{\text{cat}})} = -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T v_{\text{cat}})$$

1. Take dot product of v_{cat} with all u

2. exp

3. sum all



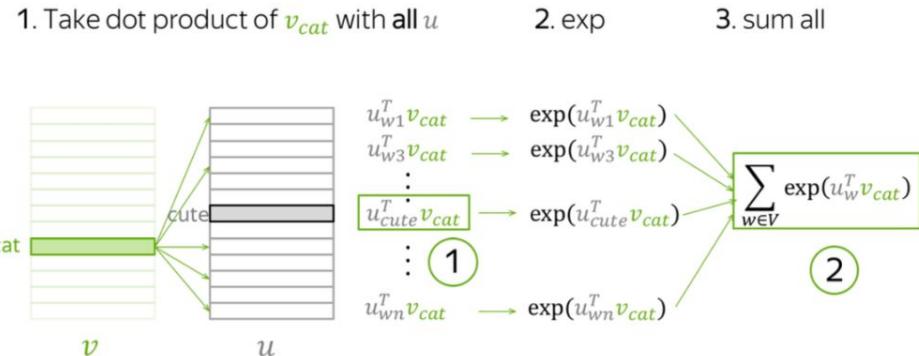
One Training step

$$-\log P(\text{cute}|\text{cat})$$

$$= -u_{\text{cute}}^T \mathbf{v}_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T \mathbf{v}_{\text{cat}})$$

4. get loss (for this one step)

$$J_{t,j}(\theta) = \underbrace{-u_{\text{cute}}^T \mathbf{v}_{\text{cat}}}_{1} + \underbrace{\log \sum_{w \in V} \exp(u_w^T \mathbf{v}_{\text{cat}})}_{2}$$



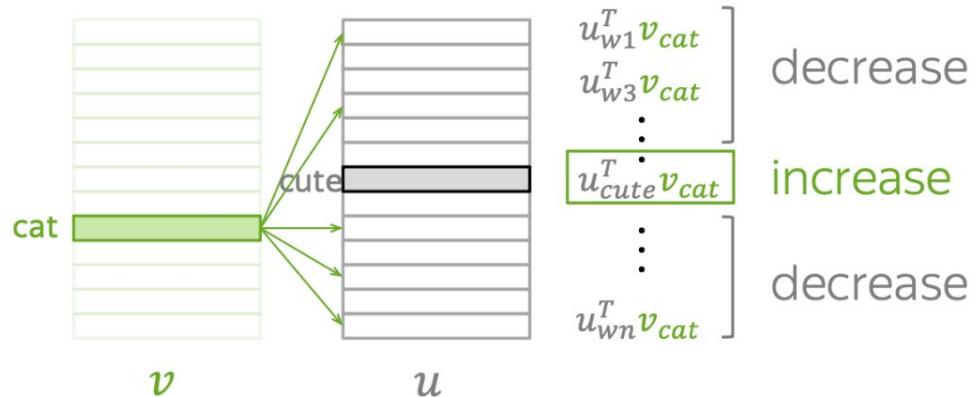
5. evaluate the gradient, make an update

$$\mathbf{v}_{\text{cat}} := \mathbf{v}_{\text{cat}} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial \mathbf{v}_{\text{cat}}}$$

$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \quad \forall w \in V$$

One Training step

$$-\log P(\text{cute}|\text{cat}) = -u_{\text{cute}}^T \mathbf{v}_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T \mathbf{v}_{\text{cat}})$$

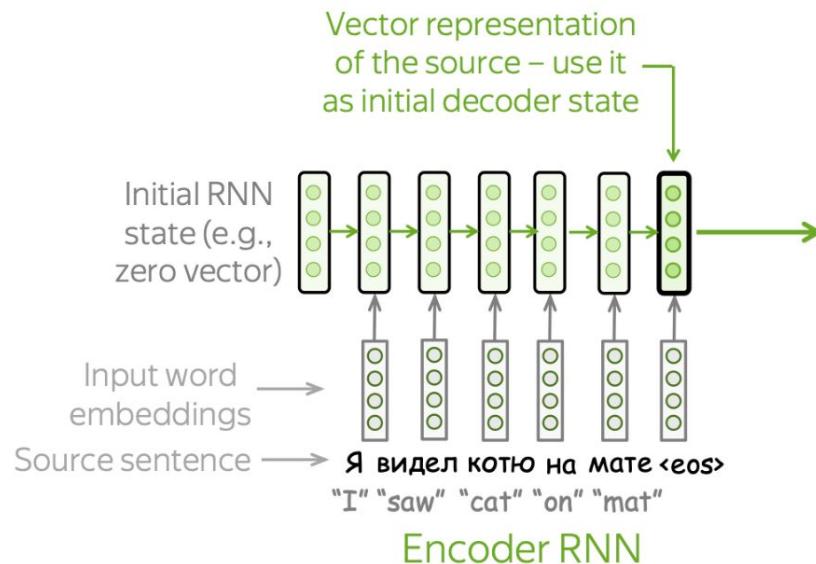


Machine Translation before transformers

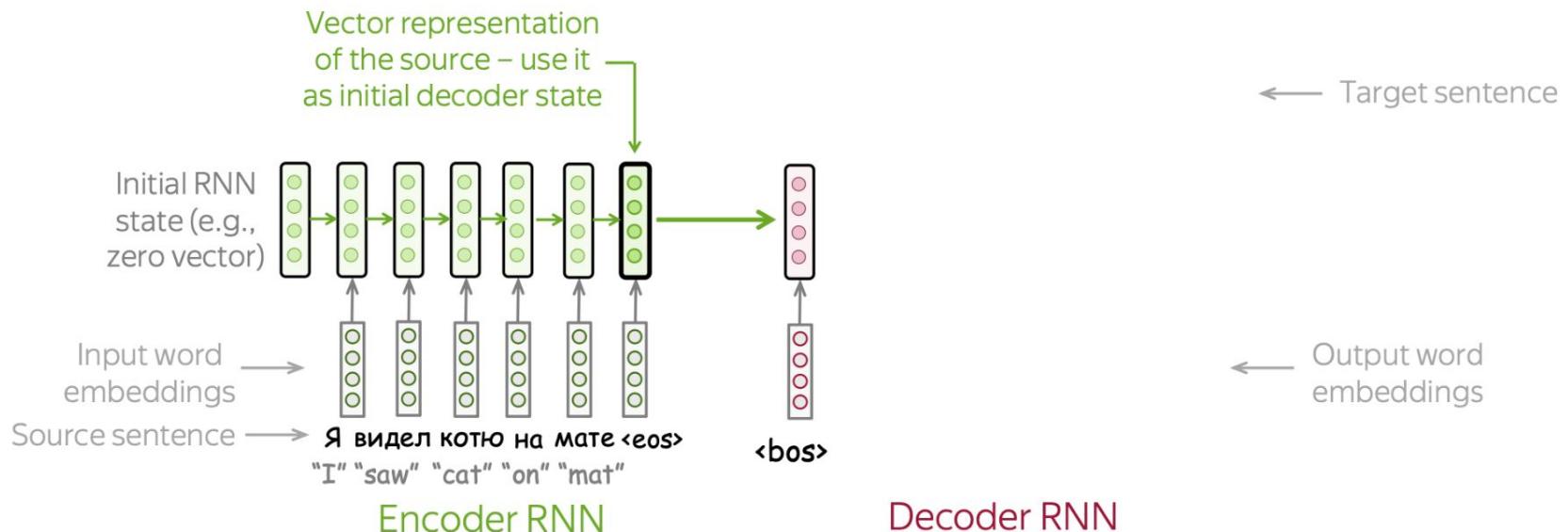


GTA San Andreas (2005)

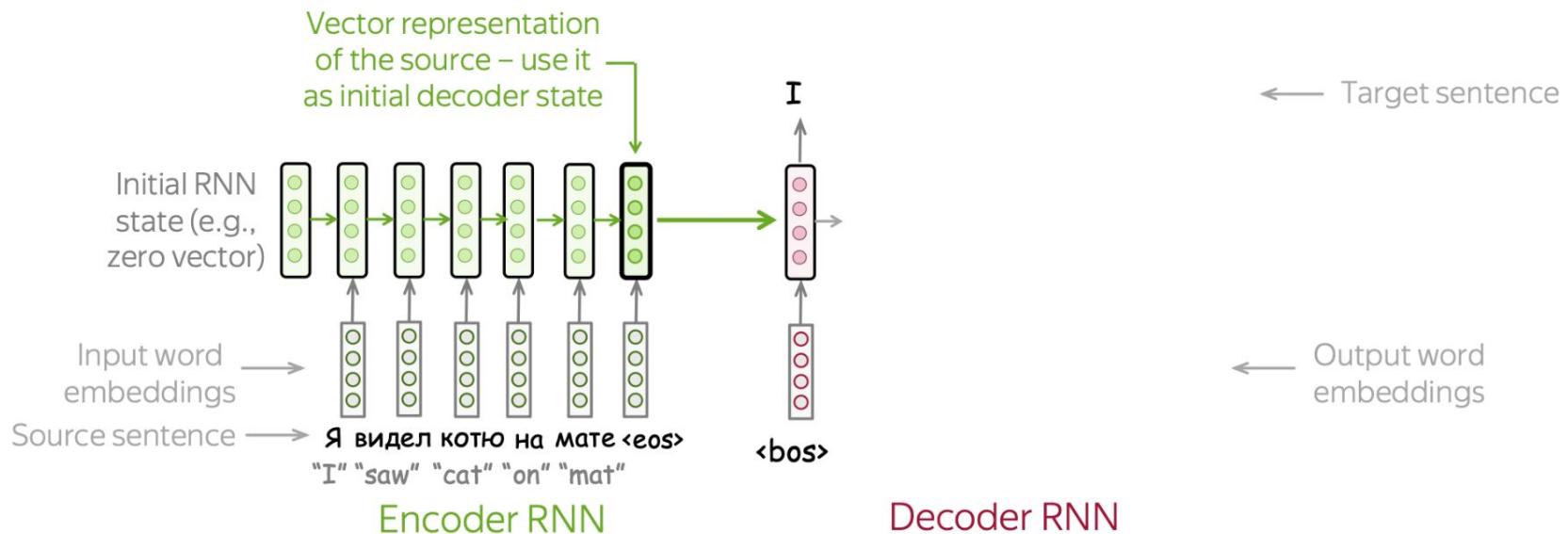
Seq2Seq models: RNN Encoder-Decoder



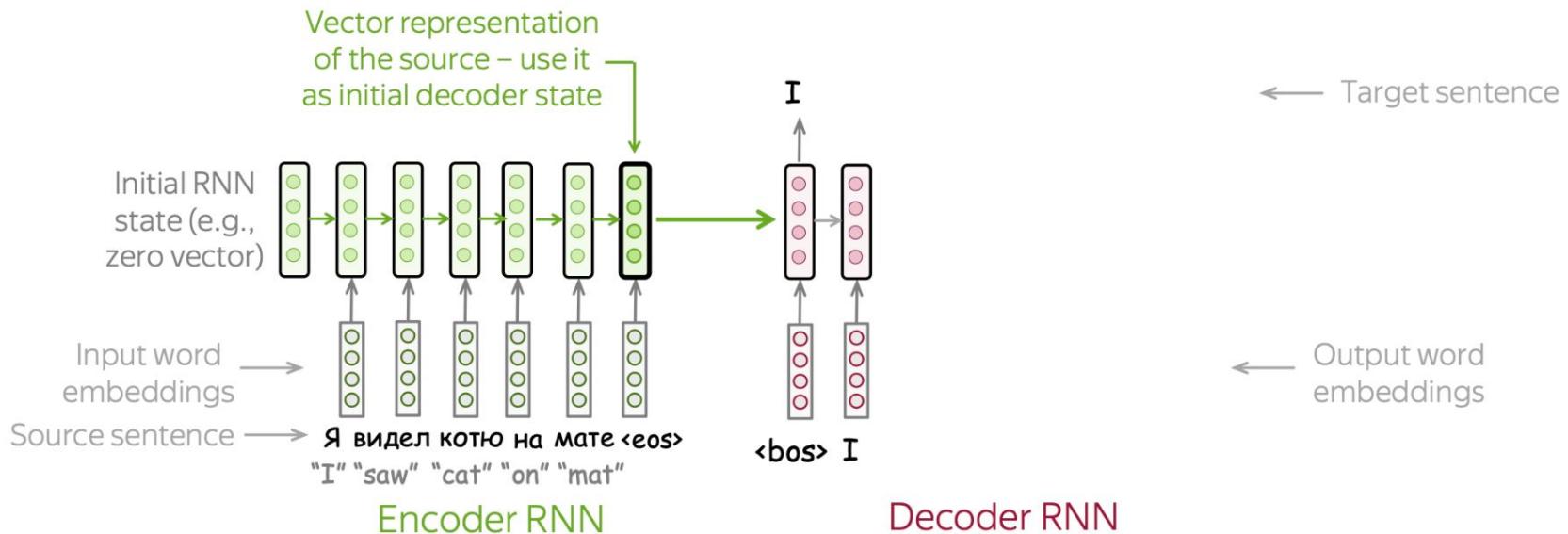
Seq2Seq models: RNN Encoder-Decoder



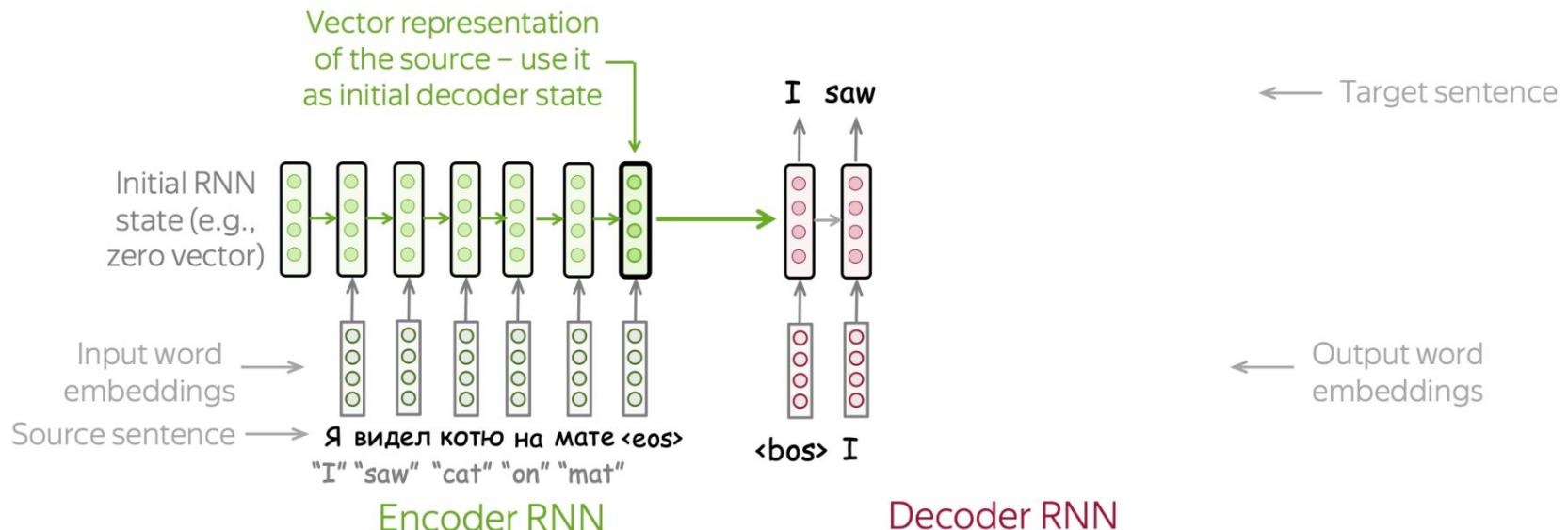
Seq2Seq models: RNN Encoder-Decoder



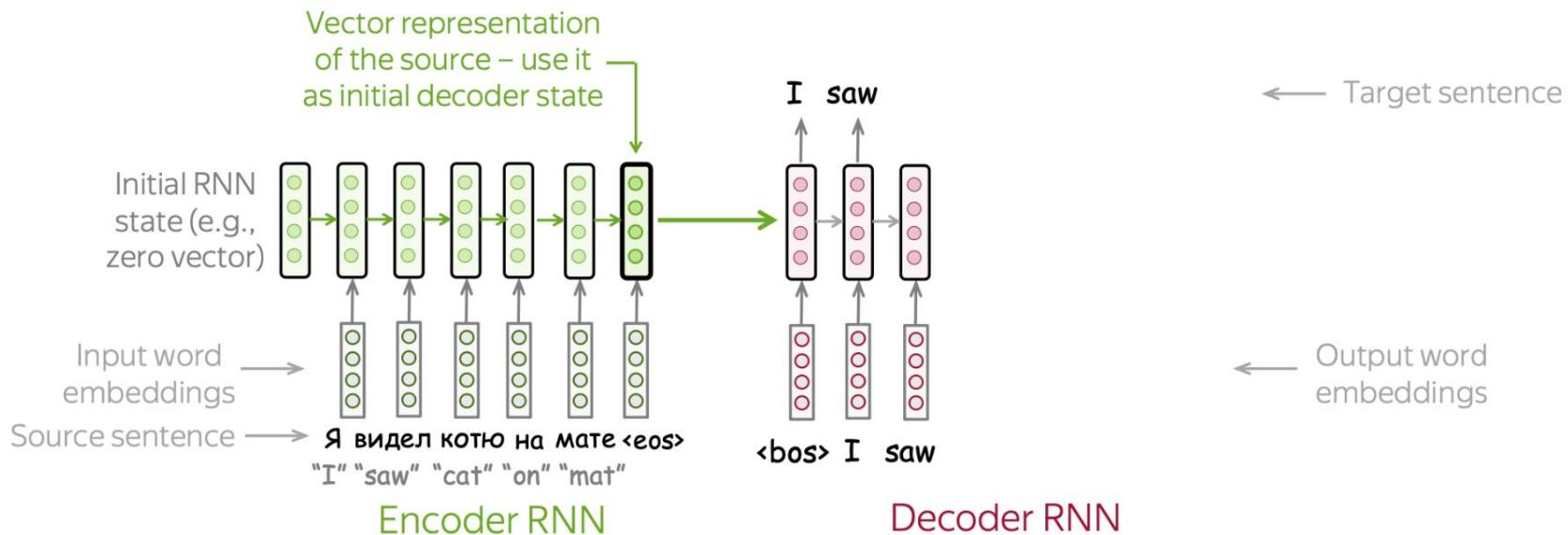
Seq2Seq models: RNN Encoder-Decoder



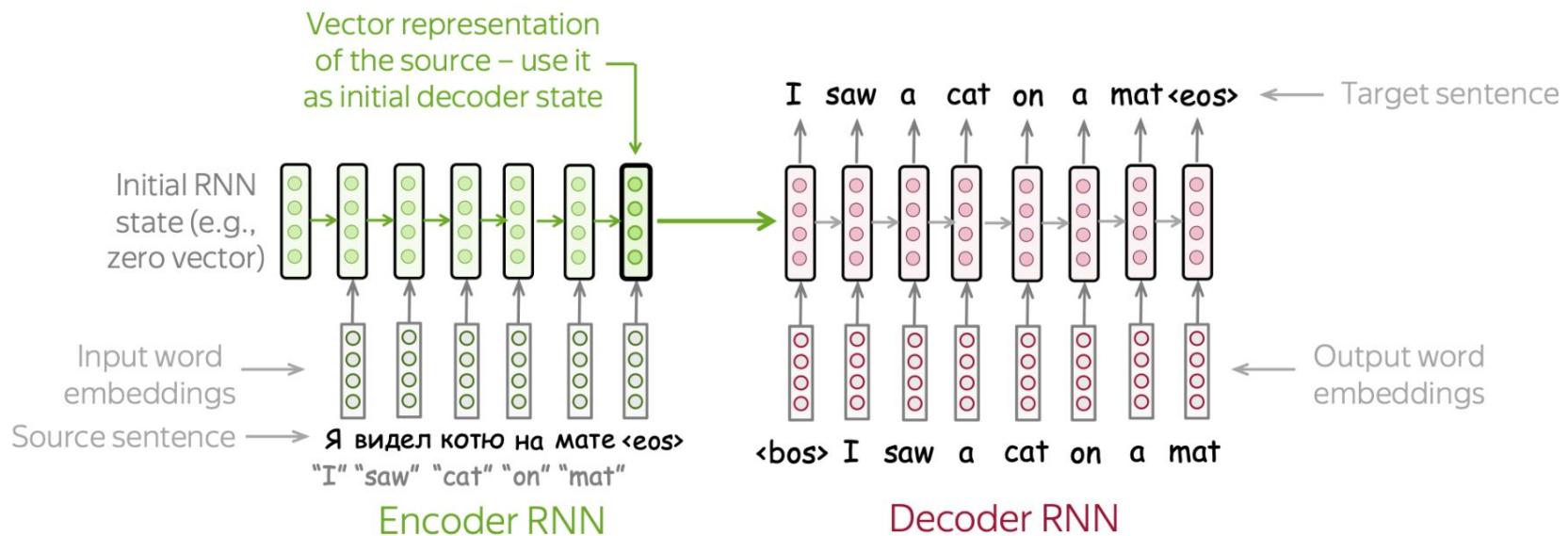
Seq2Seq models: RNN Encoder-Decoder



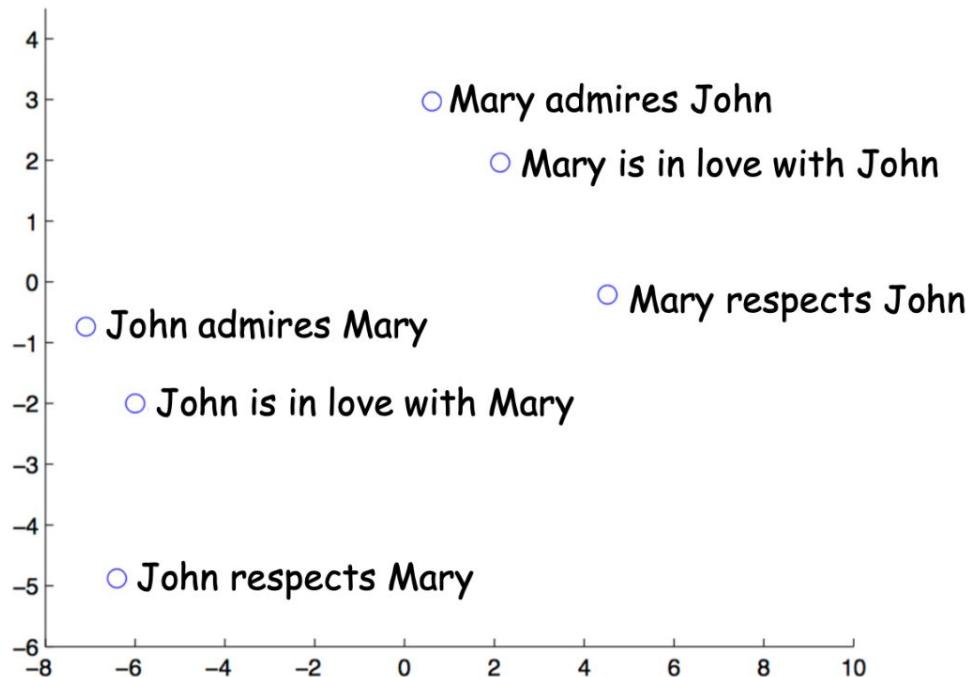
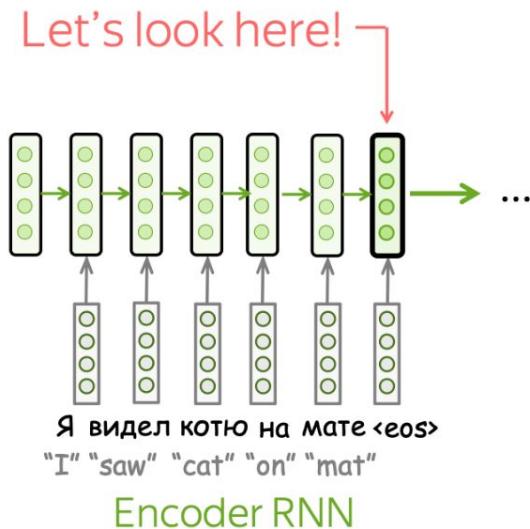
Seq2Seq models: RNN Encoder-Decoder



Seq2Seq models: RNN Encoder-Decoder



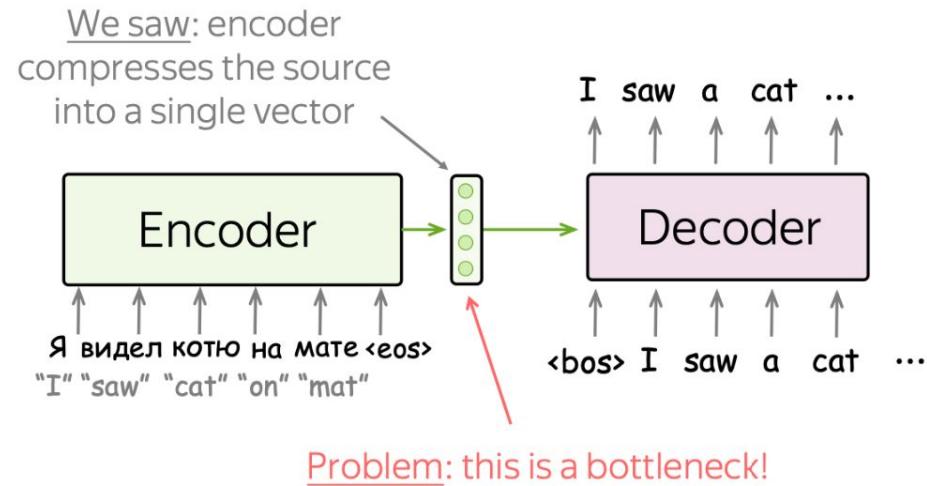
What does final encoder state represent?



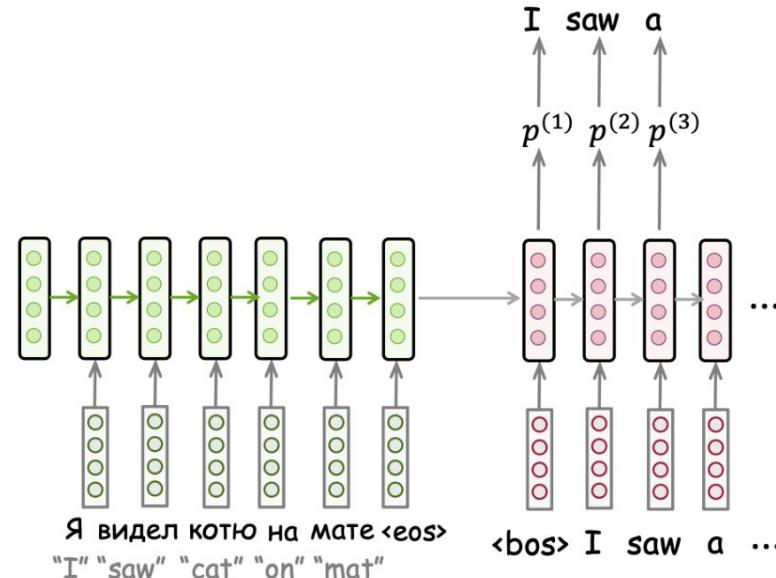
Seq2Seq models: RNN bottleneck

Fixed source representation is bad:

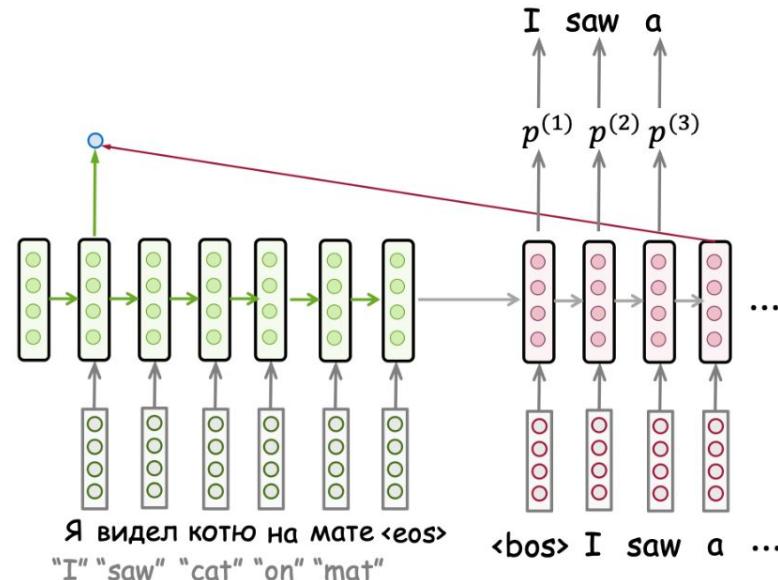
- for **encoder**, it is hard to compress a sentence
- for **decoder**, different information may be needed at different steps



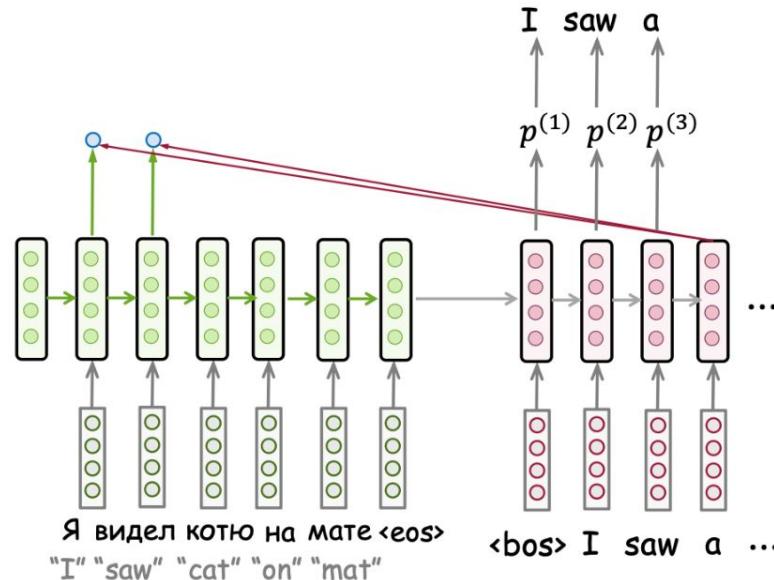
Seq2Seq models: Encoder-Decoder Attention



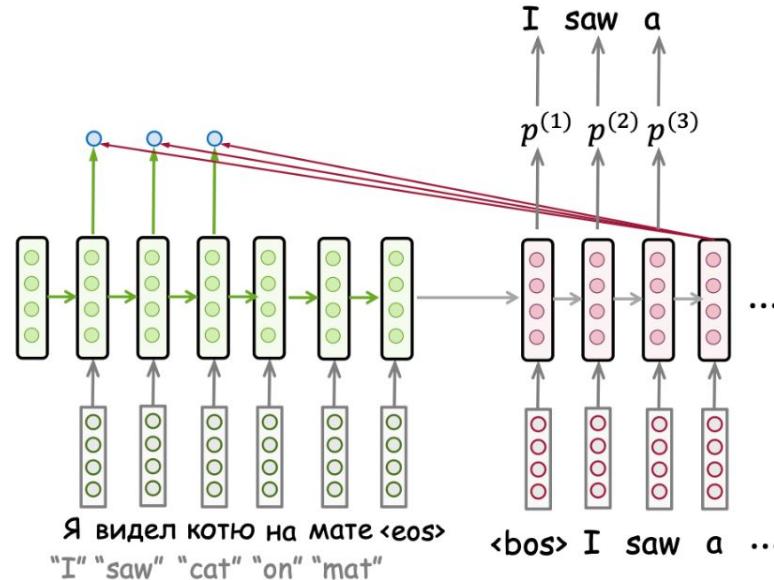
Seq2Seq models: Encoder-Decoder Attention



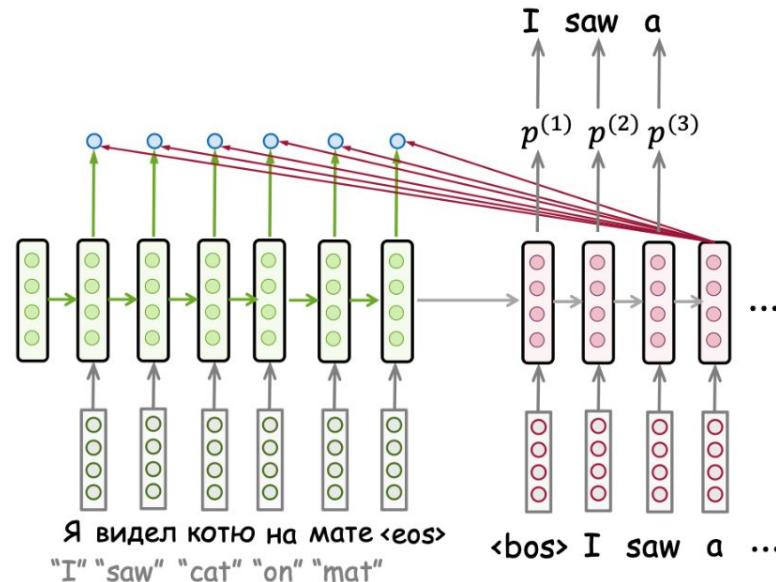
Seq2Seq models: Encoder-Decoder Attention



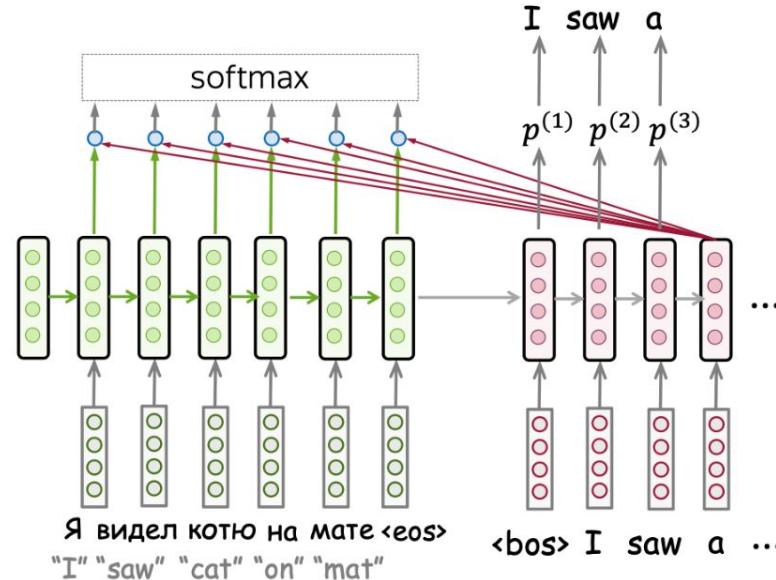
Seq2Seq models: Encoder-Decoder Attention



Seq2Seq models: Encoder-Decoder Attention

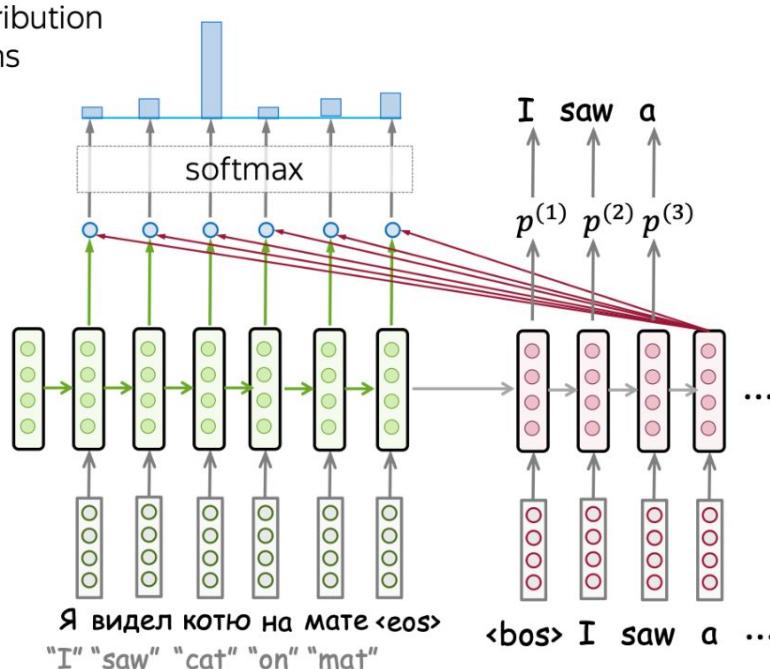


Seq2Seq models: Encoder-Decoder Attention



Seq2Seq models: Encoder-Decoder Attention

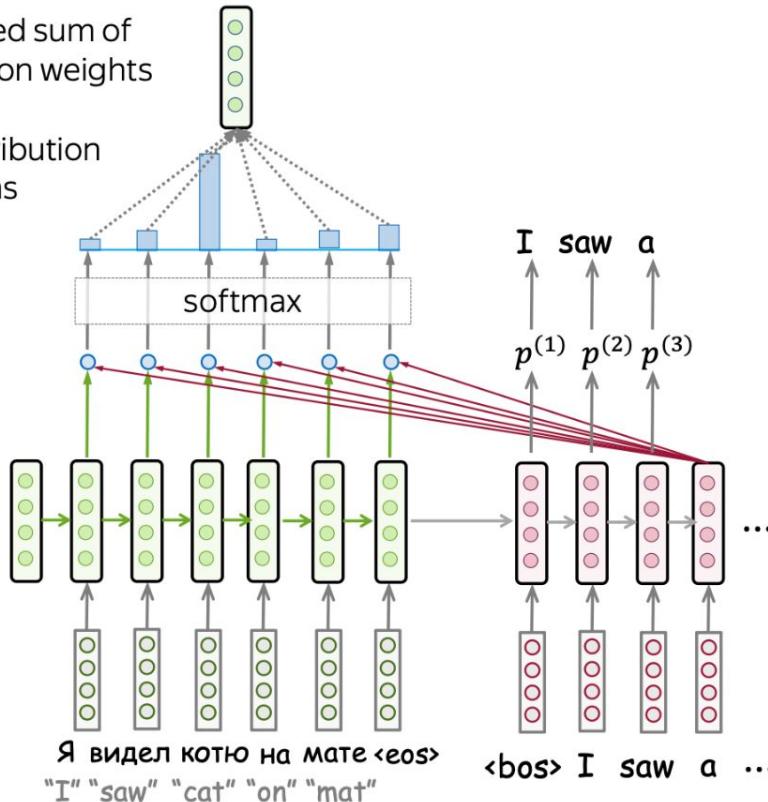
Attention weights: distribution
over source tokens



Seq2Seq models: Encoder-Decoder Attention

Attention output: weighted sum of encoder states with attention weights

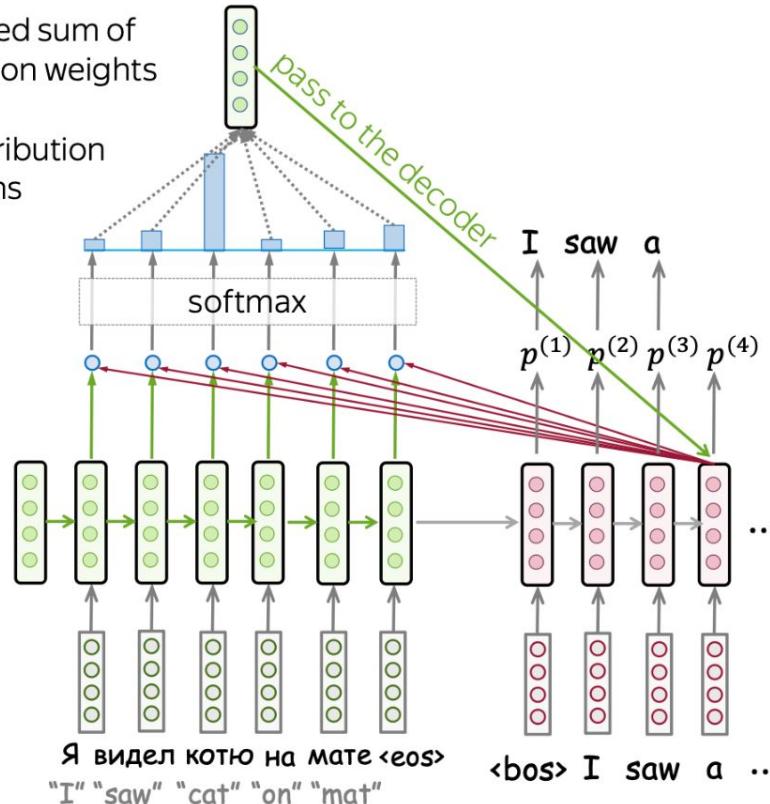
Attention weights: distribution over source tokens



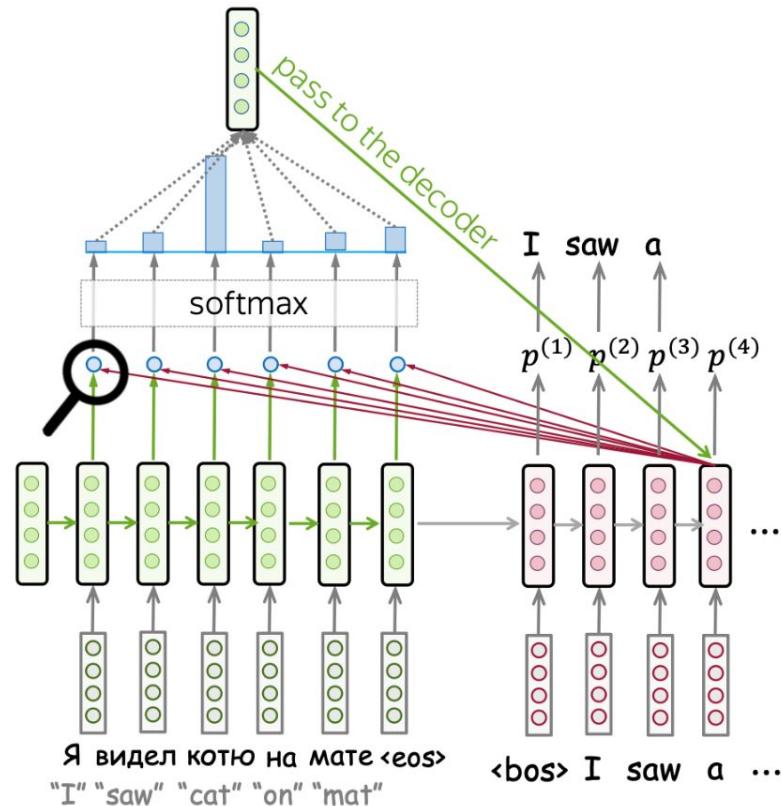
Seq2Seq models: Encoder-Decoder Attention

Attention output: weighted sum of encoder states with attention weights

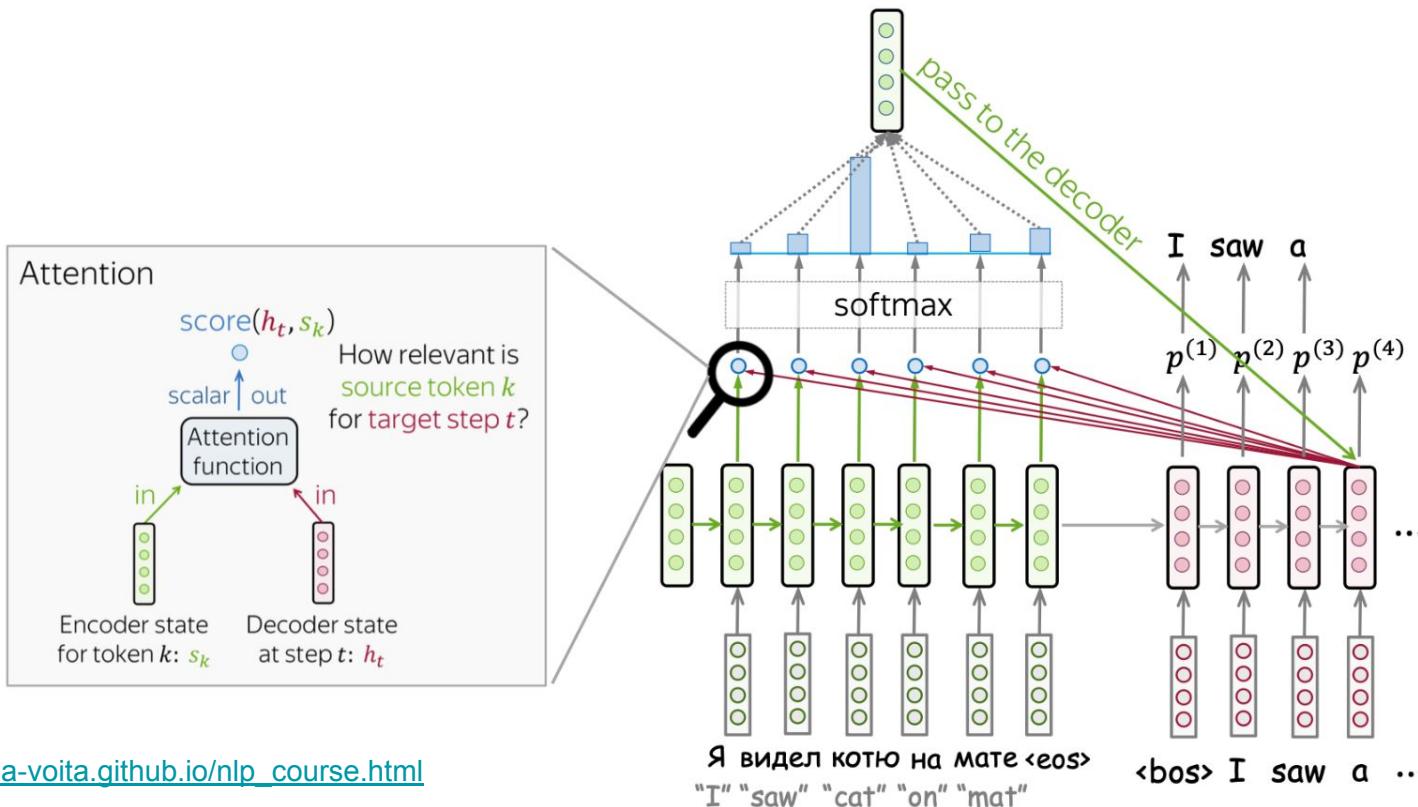
Attention weights: distribution over source tokens



Seq2Seq models: Encoder-Decoder Attention



Seq2Seq models: Encoder-Decoder Attention



Computation Pipeline

Attention input

s_1, s_2, \dots, s_m
all encoder states

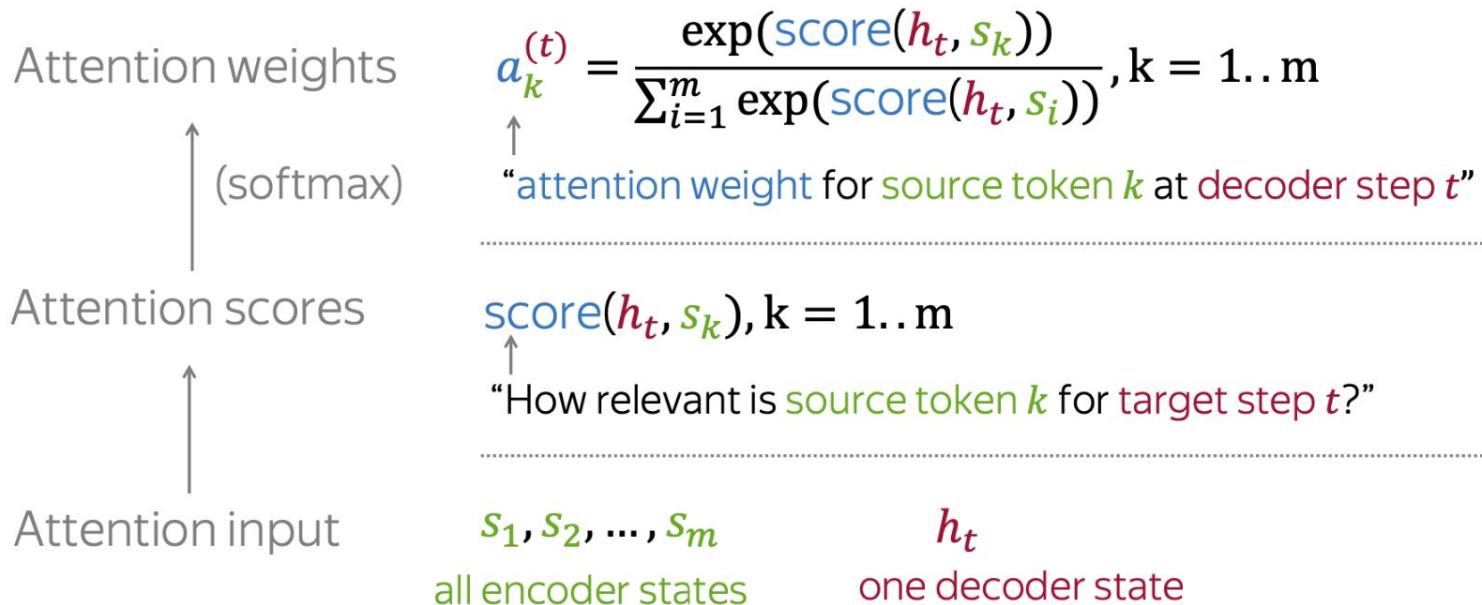
h_t
one decoder state

Computation Pipeline

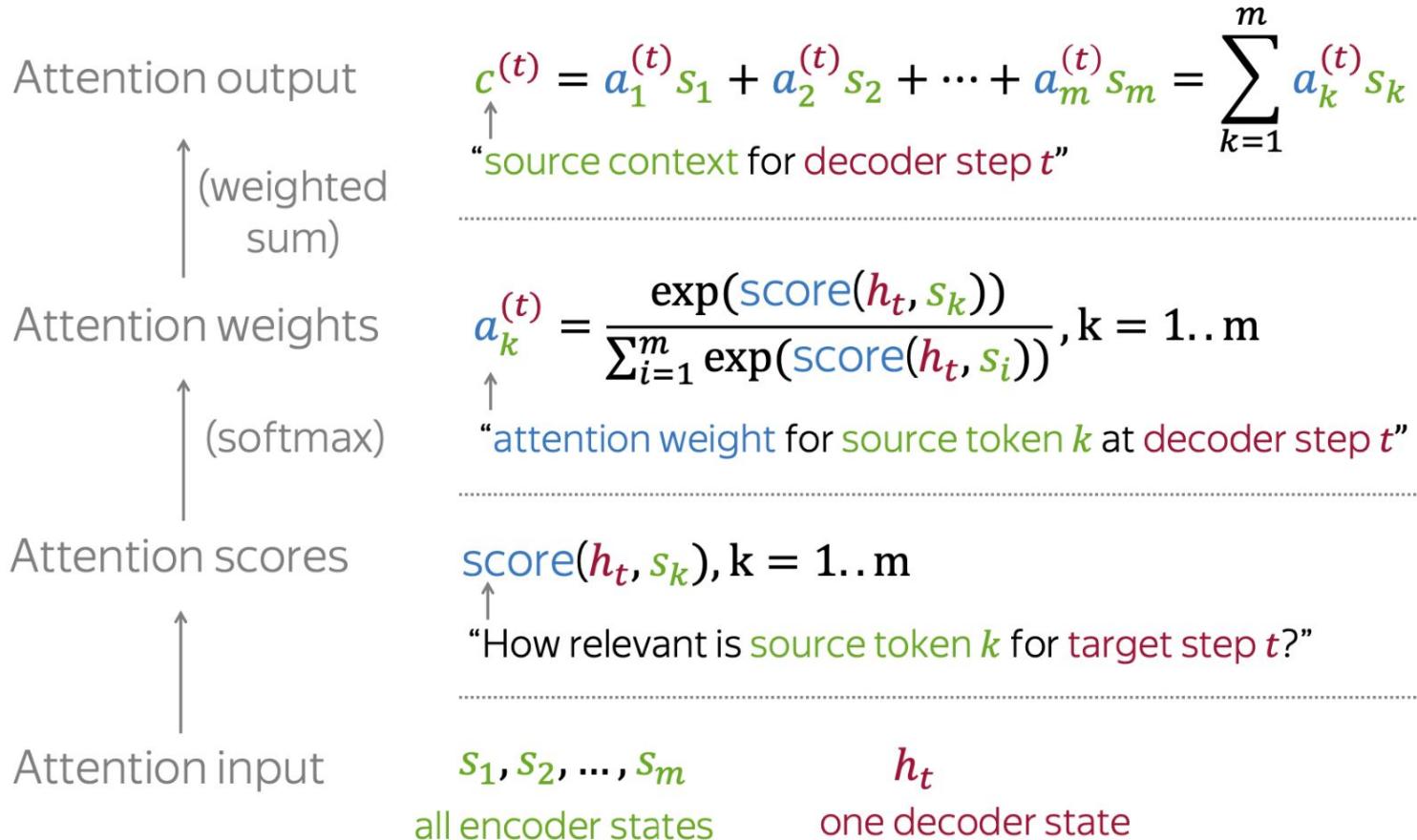
Attention scores
↑
Attention input

$\text{score}(h_t, s_k), k = 1..m$
“How relevant is source token k for target step t ? ”
.....
 s_1, s_2, \dots, s_m h_t
all encoder states one decoder state

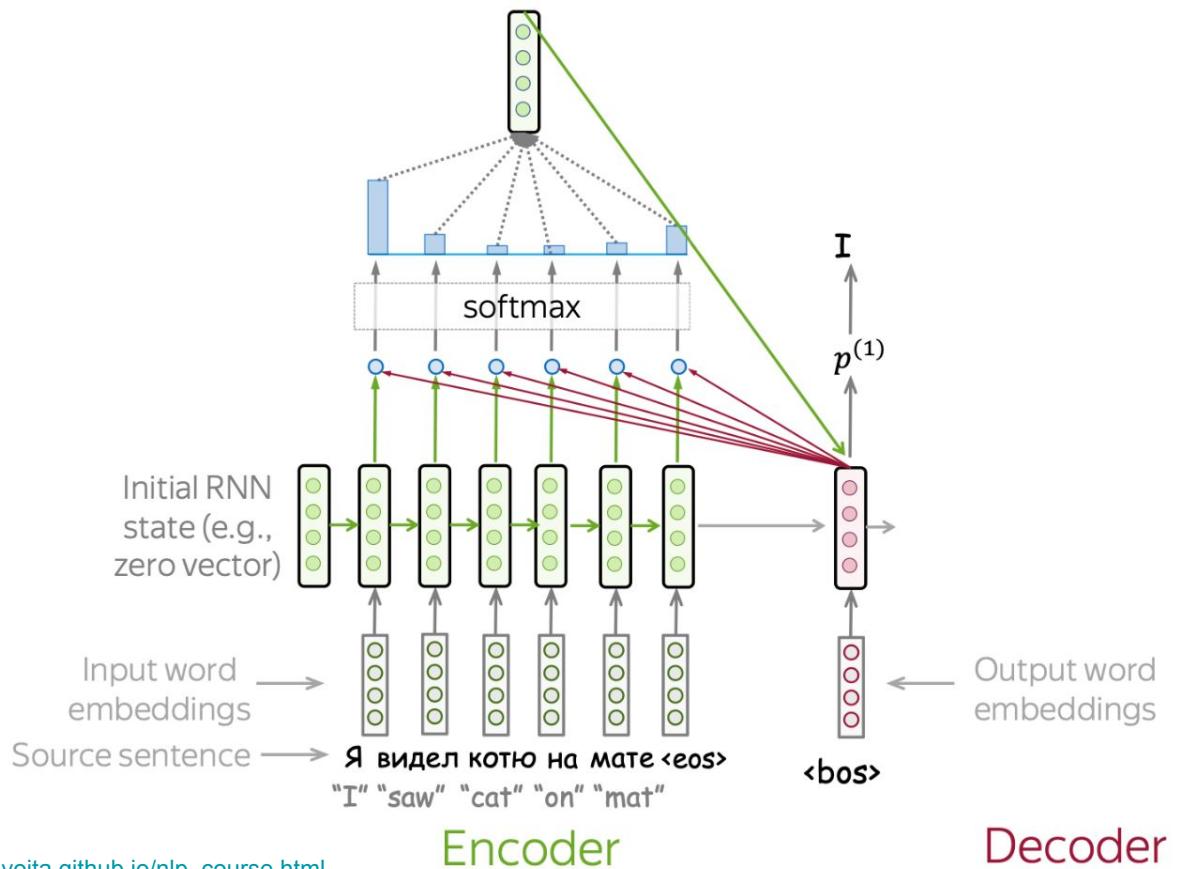
Computation Pipeline



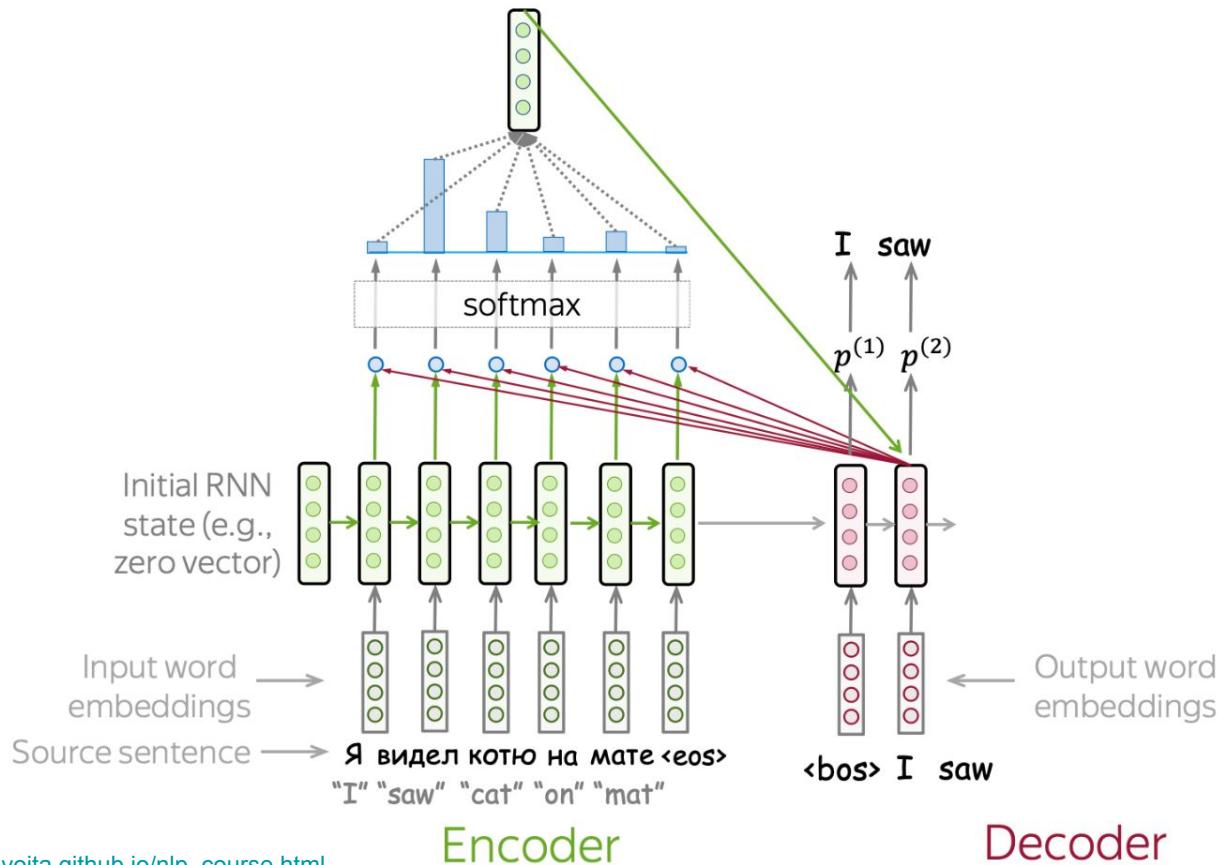
Computation Pipeline



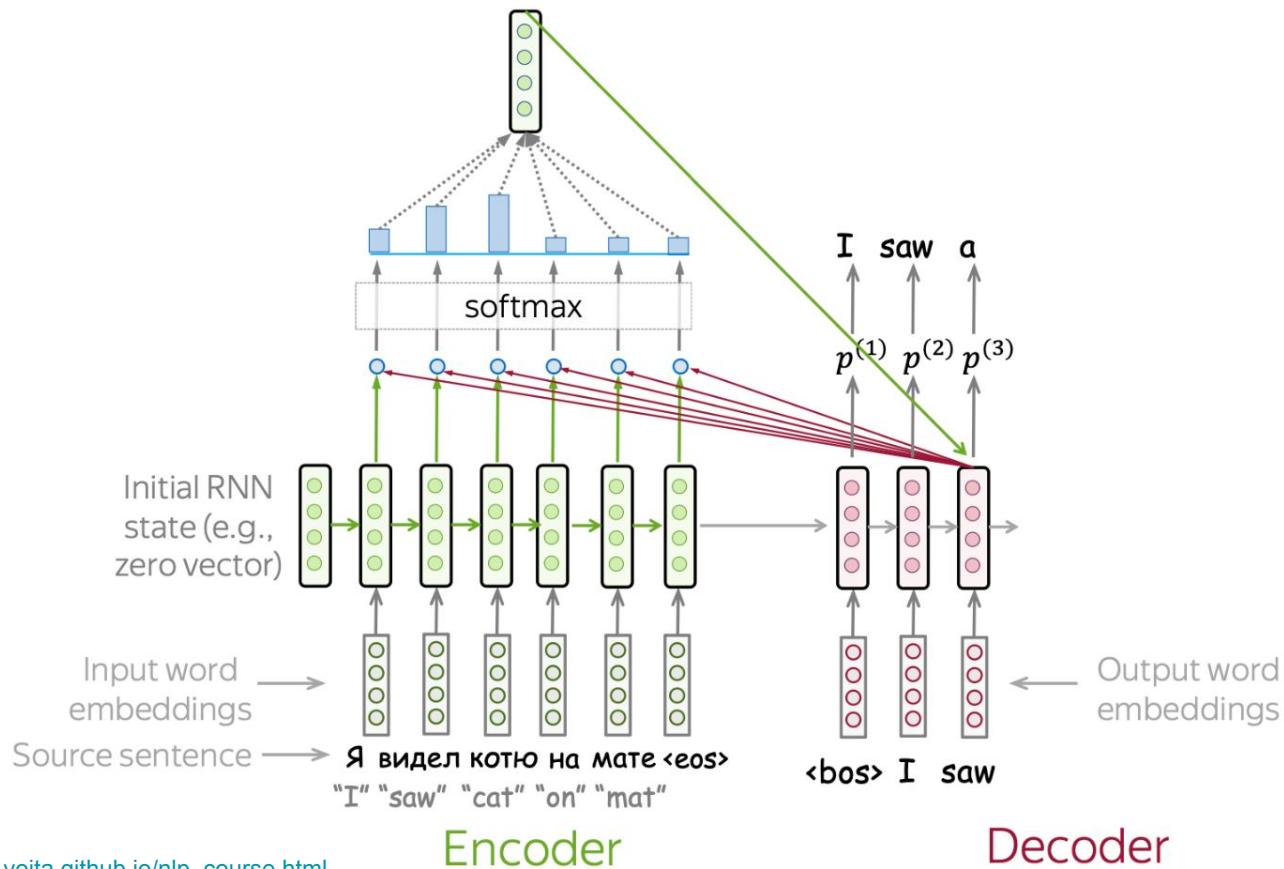
Model Learns to Pick Relevant Tokens



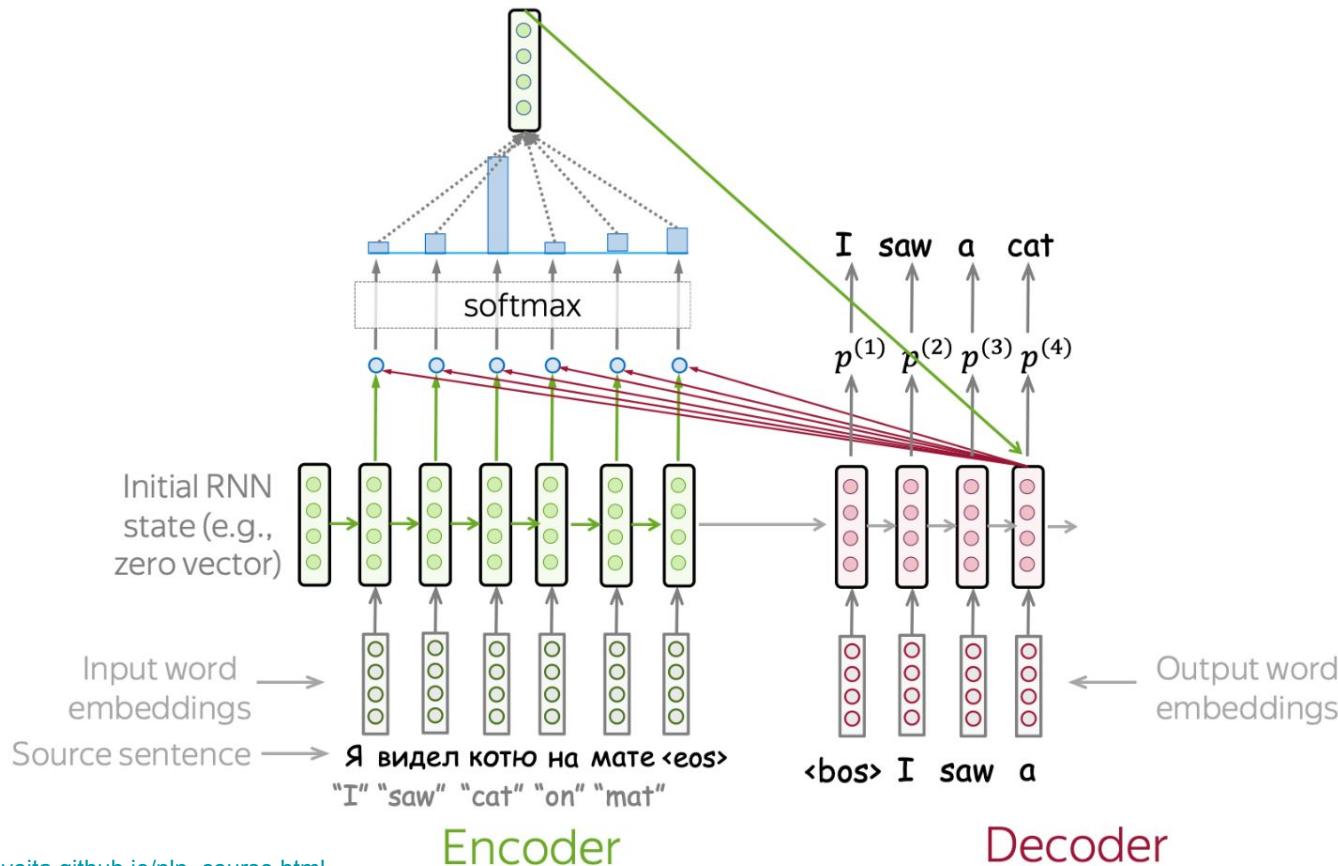
Model Learns to Pick Relevant Tokens



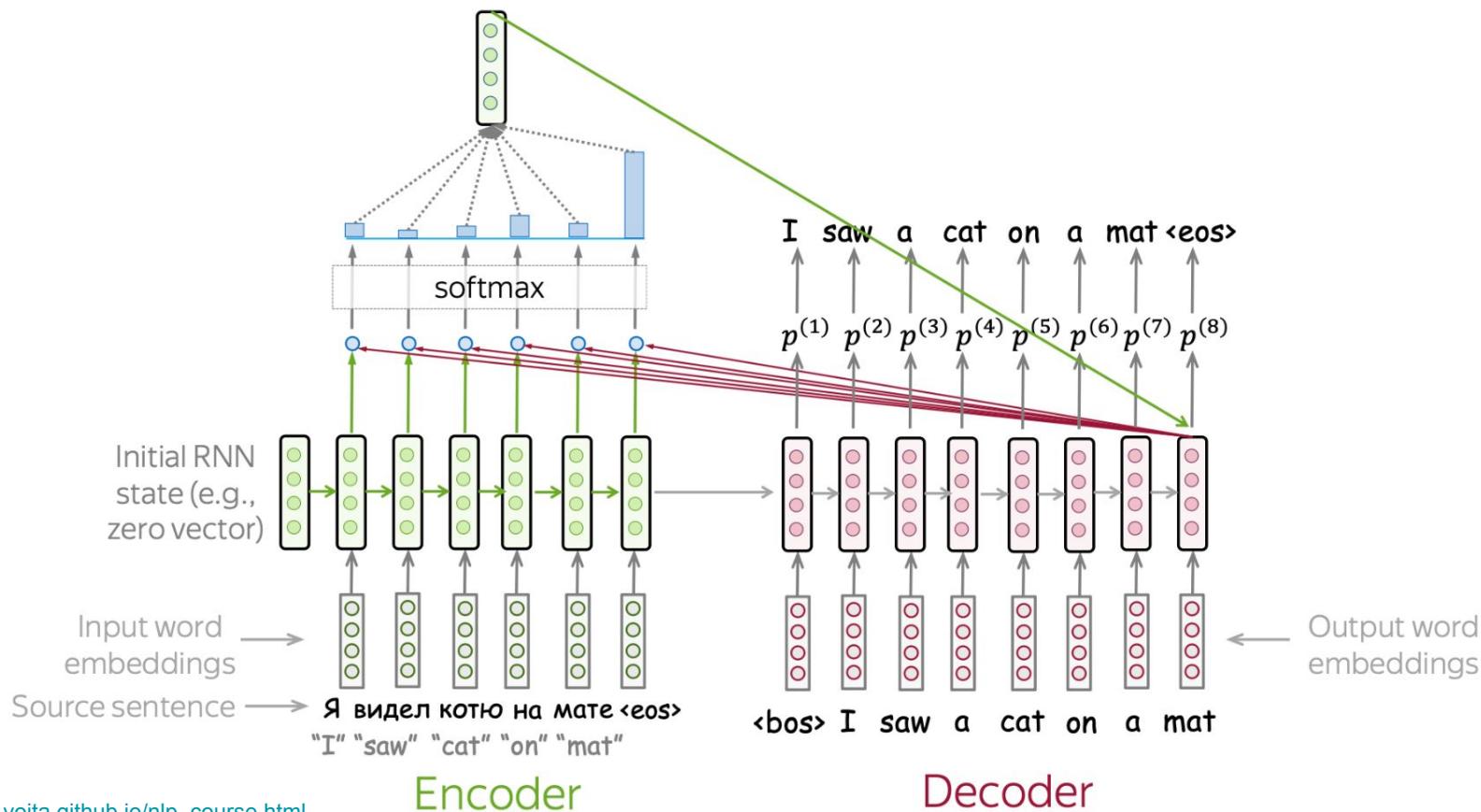
Model Learns to Pick Relevant Tokens



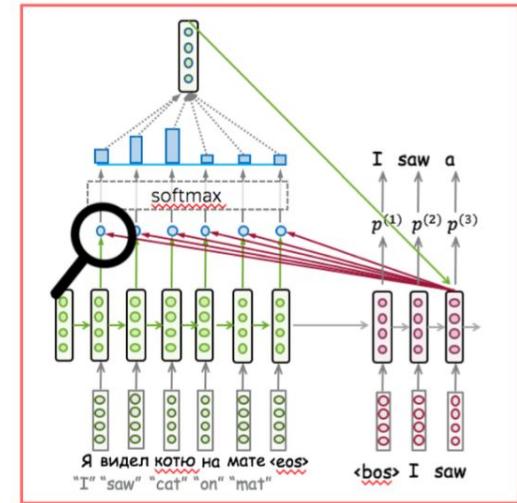
Model Learns to Pick Relevant Tokens



Model Learns to Pick Relevant Tokens

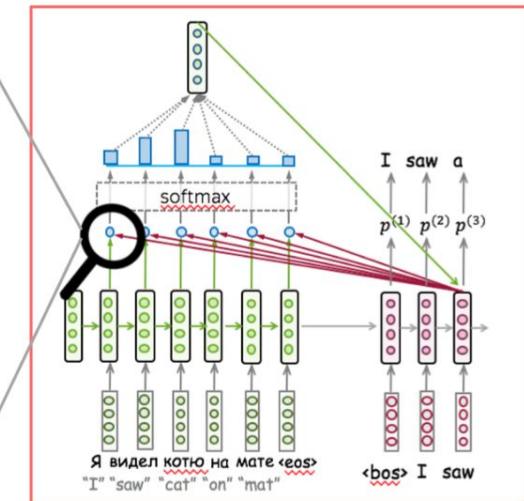
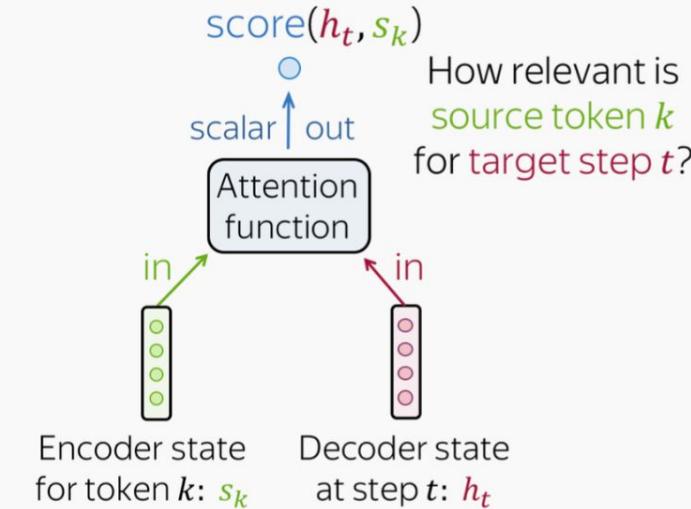


Attention Score Functions



Attention Score Functions

Attention



Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

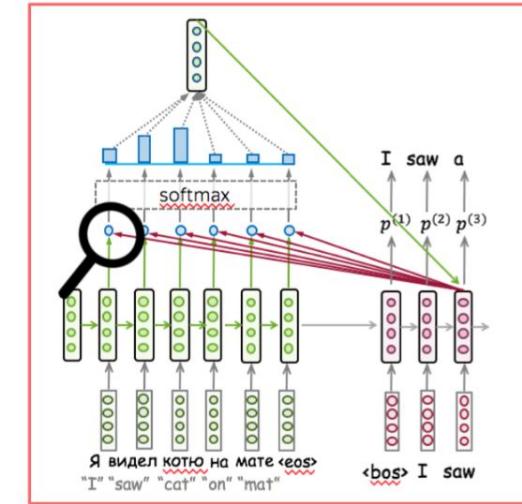
$$\begin{matrix} h_t^T \\ \text{---} \\ \text{pink dots} \end{matrix} \times \begin{matrix} \times \\ \text{---} \\ \text{green dots} \end{matrix} s_k$$

- Bilinear: $\text{score}(h_t, s_k) = h_t^T W s_k$

$$\begin{matrix} h_t^T \\ \text{---} \\ \text{pink dots} \end{matrix} \times \begin{matrix} W \\ \text{---} \\ \text{---} \end{matrix} \times \begin{matrix} \times \\ \text{---} \\ \text{green dots} \end{matrix} s_k$$

- Multi-Layer Perceptron: $\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$

$$\begin{matrix} w_2^T \\ \text{---} \\ \text{blue dots} \end{matrix} \times \tanh \left[\begin{matrix} W_1 \\ \text{---} \end{matrix} \times \begin{matrix} h_t \\ \text{---} \\ \text{pink dots} \end{matrix} \right] \begin{matrix} \times \\ \text{---} \\ \text{green dots} \end{matrix} s_k$$



Bahdanau model (sep 2014)

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.

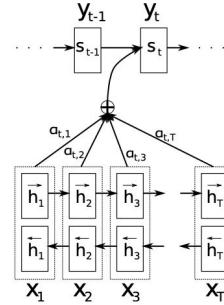
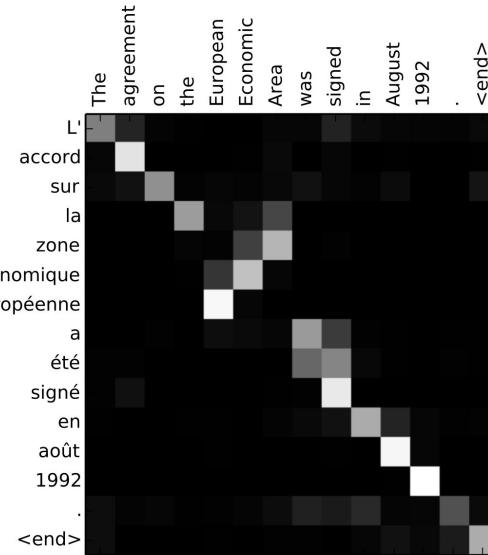


Figure 1: The graphical illustration of the proposed model trying to generate the t -th target word y_t given a source sentence (x_1, x_2, \dots, x_T) .

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}, \quad (6)$$

where

$$e_{ij} = a(s_{i-1}, h_j)$$

The probability α_{ij} , or its associated energy e_{ij} , reflects the importance of the annotation h_j with respect to the previous hidden state s_{i-1} in deciding the next state s_i and generating y_i . Intuitively, this implements a mechanism of attention in the decoder. The decoder decides parts of the source sentence to pay attention to. By letting the decoder have an attention mechanism, we relieve the encoder from the burden of having to encode all information in the source sentence into a fixed-length vector. With this new approach the information can be spread throughout the sequence of annotations, which can be selectively retrieved by the decoder accordingly.

Bahdanau model (sep 2014)



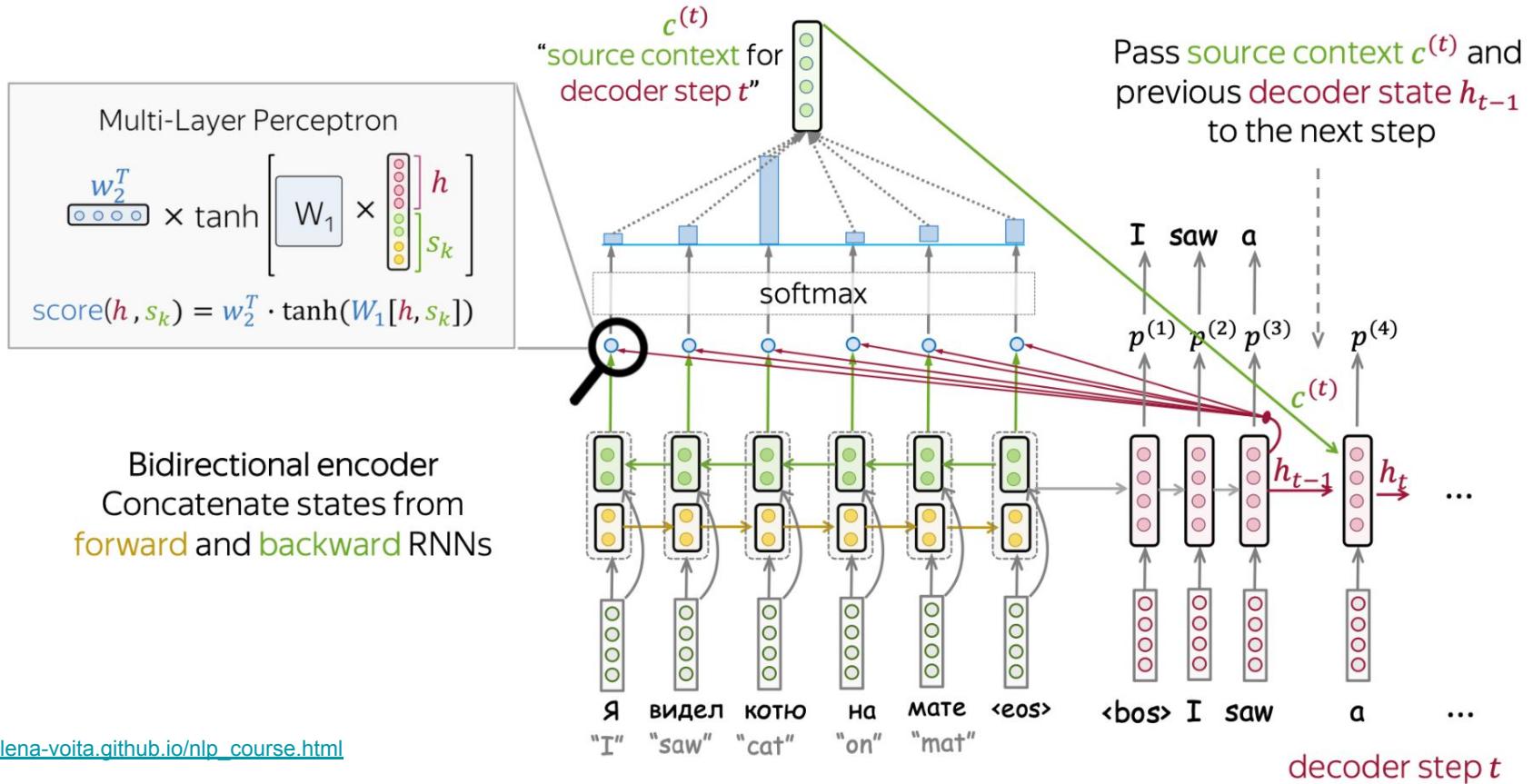
So I started thinking about how to avoid the bottleneck between encoder and decoder RNN. My first idea was to have a model with two “cursors”, one moving through the source sequence (encoded by a BiRNN) and another one moving through the target sequence. The cursor trajectories would be marginalized out using dynamic programming. KyungHyun Cho recognized this as an equivalent to Alex Graves’ RNN Transducer model. Following that, I may have also read Graves’ hand-writing recognition paper. The approach looked inappropriate for machine translation though.

The above approach with cursors would be too hard to implement in the remaining 5 weeks of my internship. So I tried instead something simpler - two cursors moving at the same time synchronously (effectively hard-coded diagonal attention). That sort of worked, but the approach lacked elegance.

So one day I had this thought that it would be nice to enable the decoder RNN to learn to search where to put the cursor in the source sequence. This was sort of inspired by translation exercises that learning English in my middle school involved. Your gaze shifts back and forth between source and target sequence as you translate. I expressed the soft search as softmax and then weighted averaging of BiRNN states. It worked great from the very first try to my great excitement. I called the architecture RNNSearch, and we rushed to publish an ArXiV paper as we knew that Ilya and co at Google are somewhat ahead of us with their giant 8 GPU LSTM model (RNN Search still ran on 1 GPU).

As it later turned out, the name was not great. The better name (attention) was only added by Yoshua to the conclusion in one of the final passes.

Bahdanau model (sep 2014)



Luong model (aug 2015)

