

# NLP and the Web – WS 2024/2025



## Lecture 9 Neural Language Modeling

**Dr. Thomas Arnold**  
**Hovhannes Tamoyan**  
**Kexin Wang**

**Ubiquitous Knowledge Processing Lab**  
**Technische Universität Darmstadt**



# Syllabus (tentative)

<u>Nr.</u>	<u>Lecture</u>
01	Introduction / NLP basics
02	Foundations of Text Classification
03	IR – Introduction, Evaluation
04	IR – Word Representation
05	IR – Transformer/BERT
06	IR – Dense Retrieval
07	IR – Neural Re-Ranking
08	LLM – Language Modeling Foundations, Tokenization
<b>09</b>	<b>LLM – Neural LLM</b>
10	LLM – Adaptation, LoRa, Prompting
11	LLM – Alignment, Instruction Tuning
12	LLM – Long Contexts, RAG
13	LLM – Scaling, Computation Cost
14	Review & Preparation for the Exam

## Recap: LM and Sub-word tokenization

## Basic Neural Language Models

## RNN and Transformer LM

# Recap: Language Models

Language Model (LM) = Model that assigns probabilities to sequences of words



# Recap: Bigram LM

Word 1	Word 2	Count	P(W2 W1)
<S>	There	256	0.256
<S>	The	321	0.321
<S>	Oh	17	0.017
<S>	I	69	0.069
<S>	They	169	0.169
<S>	Also	123	0.123
<S>	However	45	0.045
Anything	else	...	

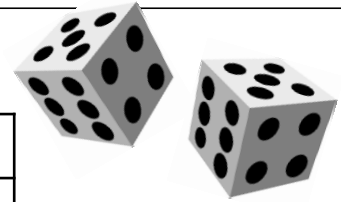
$$P(w_n|w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(w_{n-1})}$$

$$C(< S >) = 1000$$

$$\frac{C(< S > There)}{C(< S >)} = \frac{256}{1000} = 0.256$$

# Recap: Generation with a Bigram LM

Generate the next token based on conditional probabilities



There...

Word 1	Word 2	Count	$P(W2 W1)$
<S>	There	256	0.256
<S>	The	321	0.321
<S>	Oh	17	0.017
...	...	...	...
There	<b>once</b>	<b>123</b>	<b>0.0246</b>
There	<b>must</b>	<b>29</b>	<b>0.0058</b>
There	<b>has</b>	<b>69</b>	<b>0.0138</b>
...	...	...	...

There once

## Recap: OOV problem

### Big problem: zero probability n-grams

In Evaluation: Test set might contain N-Grams that did never appear in training

= the entire probability of the test set is zero! (multiplication with zero)

= Perplexity is undefined! (division by zero)

In Generation: LM has never seen the words of my prompt before

= No way to generate the next token!

Approaches: <UNK> tokens, Laplace or add-k smoothing, backoff, interpolation...

# Recap: Byte Pair Encoding (BPE)

Subword tokenization technique

Used for data compression and dealing with unknown words

Initialization:

Vocabulary = set of all individual characters

$V = \{A, B, C, \dots a, b, c, \dots 1, 2, 3, \dots !, \$, \%, \dots\}$

Repeat:

- Choose two symbols that appear as a pair most frequently (say “a” and “t”)
- Add new merged symbol (“at”)
- Replace each occurrence with the new symbol (“t”, “h”, “a”, “t” -> “t”, “h”, “at”)

Until k merges have been done



Greedy longest prefix matching. Example: “the golden snickers”

vocab: {“the”, “go”, “gold”, “##den”, “##en”, “snicker”, “##ers”, “##s”}

**the** golden snickers -> “**the**” in vocab -> [“the”]

the **golden** snickers -> “golden” not in vocab

-> prefix max match: “**gold**” -> [“the”, “gold”]

-> remaining subword: “**##en**” -> “##en” in vocab -> [“the”, “gold”, “##en”]

the golden **snickers** -> “snickers” not in vocab

-> prefix max match: “**snicker**” in vocab -> [“the”, “gold”, “##en”, “snicker”]

-> remaining subword: “**##s**” in vocab -> [“the”, “gold”, “##en”, “snicker”, “##s”]

# GPT3/4's Tokenizer



OpenAI's large language models (sometimes referred to as GPT's) process text using tokens, which are common sequences of characters found in a set of text. The models learn to understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.

You can use the tool below to understand how a piece of text might be tokenized by a language model, and the total count of tokens in that piece of text.

It's important to note that the exact tokenization process varies between models. Newer models like GPT-3.5 and GPT-4 use a different tokenizer than our legacy GPT-3 and Codex models, and will produce different tokens for the same input text.

Here is a math problem:  $234566 + 64432 / (33345) * 0.1234$

<https://platform.openai.com/tokenizer>

**Recap: LM and Sub-word tokenization**

**Basic Neural Language Models**

**RNN and Transformer LM**

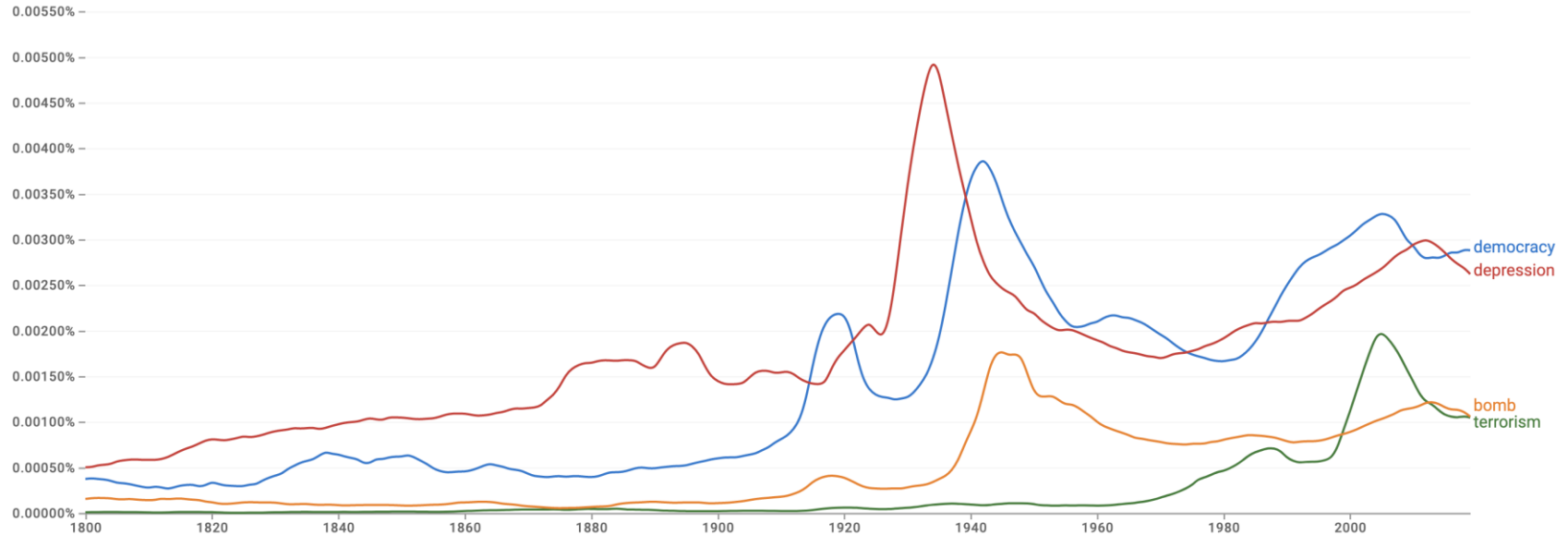
# N-Grams LMs and Long-range Dependencies

In general, count-based LMs are insufficient models of language because language has long-distance dependencies:

“The computer which I had just put into  
the machine room on the fifth floor crashed.”

# Pre-Computed N-Grams

Google Books Ngram Viewer



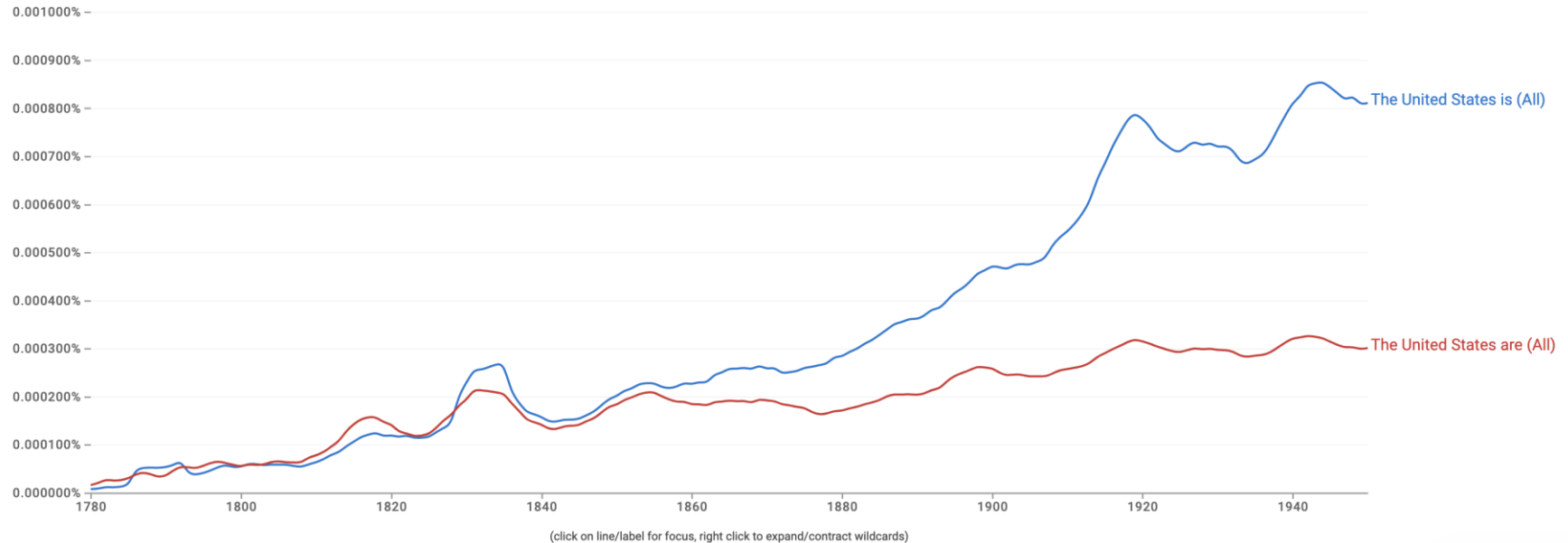
- Google n-gram viewer <https://books.google.com/ngrams/>
- Data: <http://storage.googleapis.com/books/ngrams/books/datasetv2.html>

# Pre-Computed N-Grams

Google Books Ngram Viewer



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

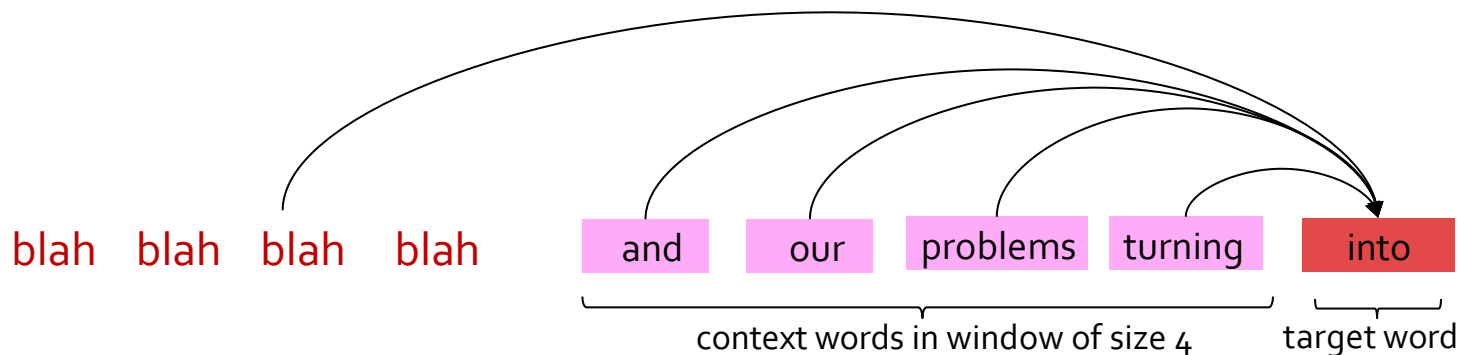


- Google n-gram viewer <https://books.google.com/ngrams/>
- Data: <http://storage.googleapis.com/books/ngrams/books/datasetsv2.html>

# LM as a Machine Learning Problem

Task: Given the embeddings of the context, predict the word on the right side.

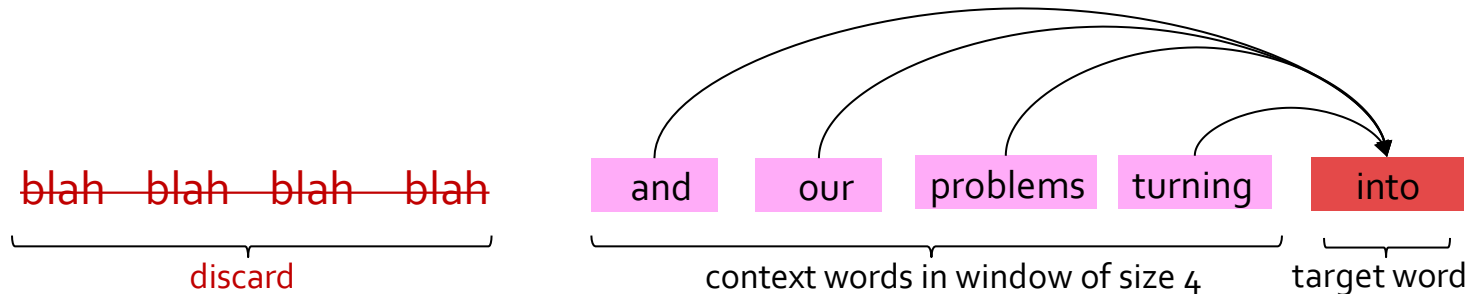
Discard anything beyond its context window



# LM as a Machine Learning Problem

Task: Given the embeddings of the context, predict the word on the right side.

Discard anything beyond its context window

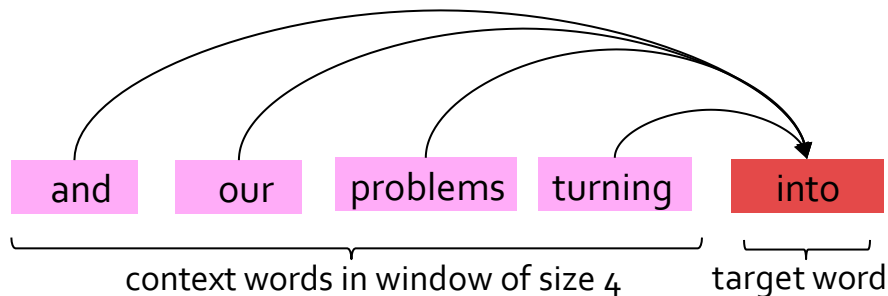




# LM as a Machine Learning Problem

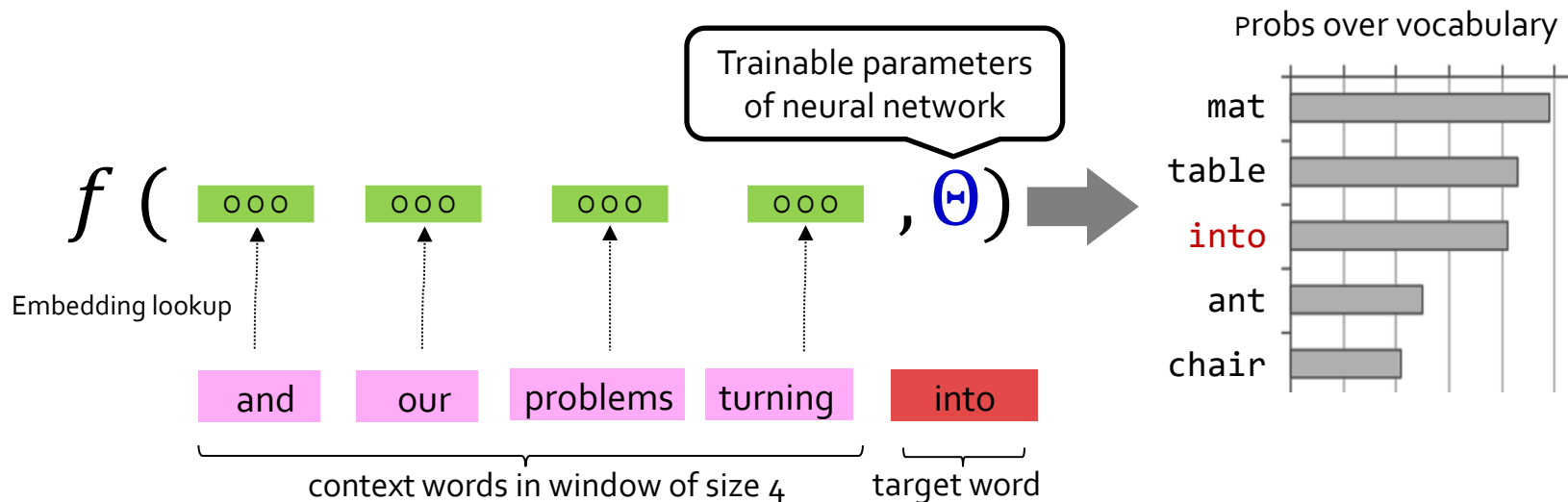
Task: Given the embeddings of the context, predict the word on the right side.

Discard anything beyond its context window



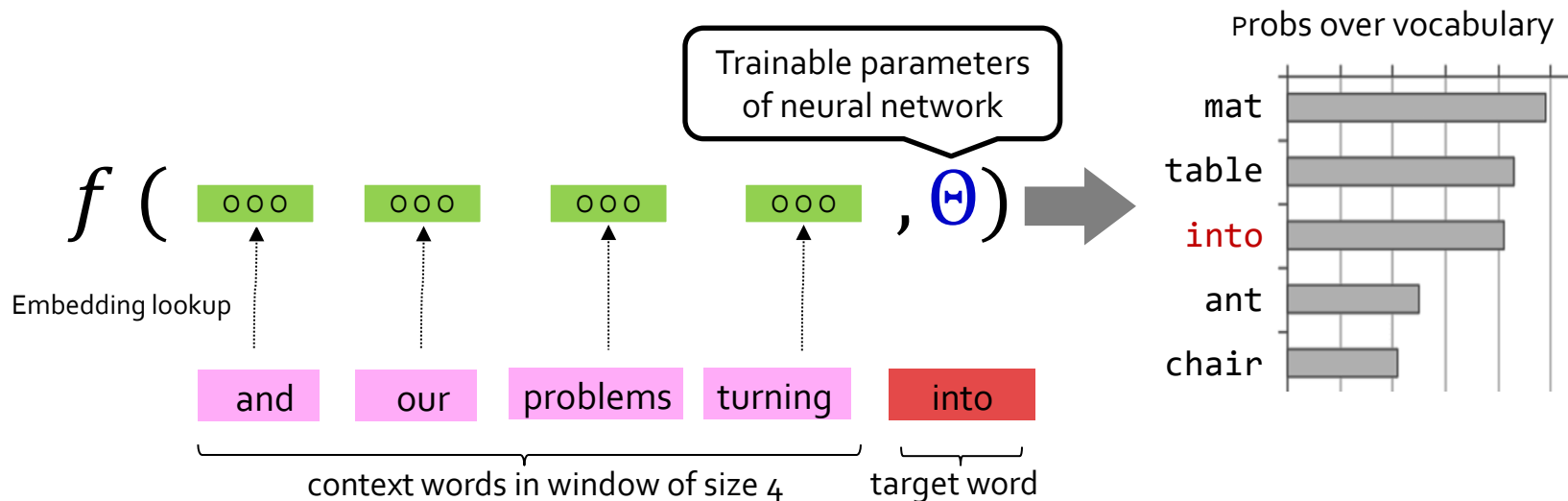
# A Fixed-Window Neural LM

Training this model is basically optimizing its parameters  $\Theta$  such that it assigns high probability to the **target word**.



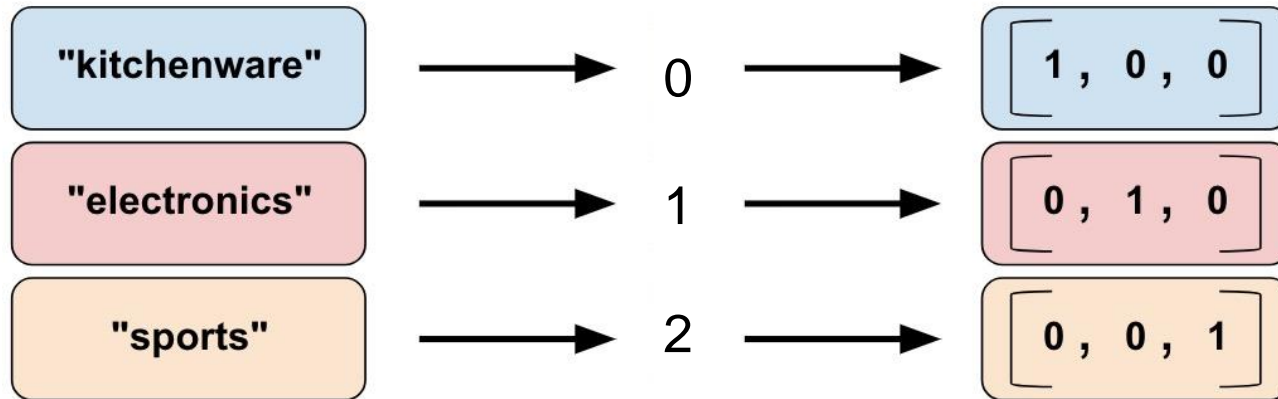
# A Fixed-Window Neural LM

- It will also lay the foundation for the future models (RNN, transformers, ...)
- But first we need to figure out how to train neural networks!



# Feeding Text to Neural Net

- In practice this is implemented in this way:
  1. Turn each word into a unique index
  2. Map each index into a one-hot vector

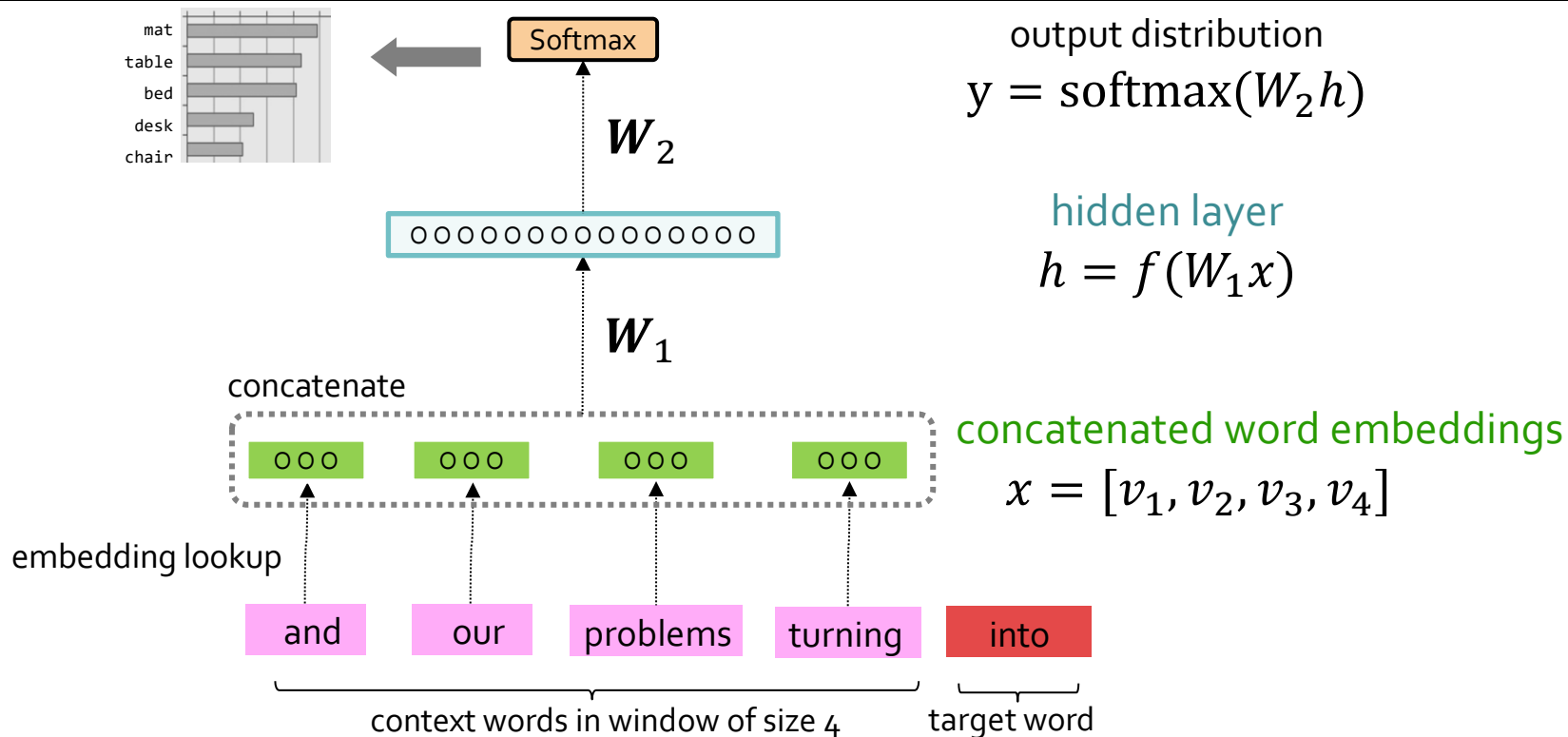


# Feeding Text to Neural Net

- In practice this is implemented in this way:
  1. Turn each word into a unique index
  2. Map each index into a one-hot vector
  3. Lookup the corresponding word embedding via matrix multiplication

$$\begin{array}{c} [0 \quad 0 \quad 0 \quad 1 \quad 0] \\ \text{One-hot vector} \end{array} \times \begin{array}{c} \begin{bmatrix} 8 & 2 & 1 & 9 \\ 6 & 5 & 4 & 0 \\ 7 & 1 & 6 & 2 \\ 1 & 3 & 5 & 8 \\ 0 & 4 & 9 & 1 \end{bmatrix} \\ \text{Embedding Weight Matrix} \end{array} = \begin{array}{c} [1 \quad 3 \quad 5 \quad 8] \\ \text{Hidden layer output} \end{array}$$

# A Fixed-Window Neural LM



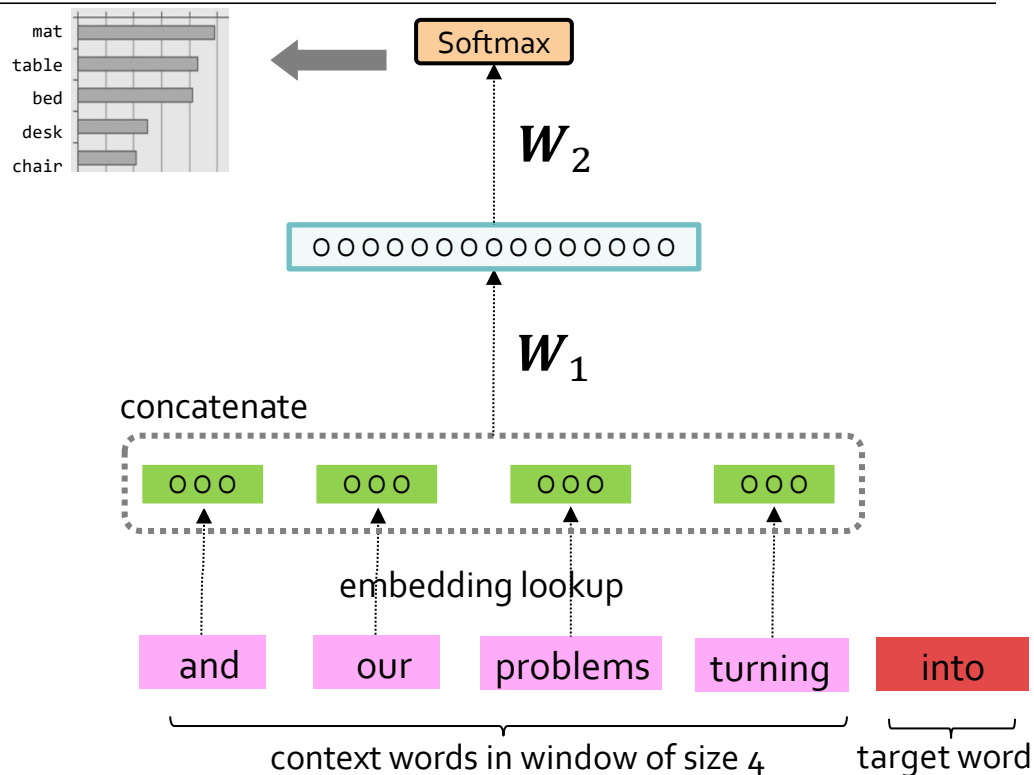
# A Fixed-Window Neural LM

Improvements over n-gram LM:

- Tackles the sparsity problem
- Model size is  $O(n)$  not  $O(\exp(n))$  —  $n$  being the window size.

Remaining problems:

- Fixed window is too small
- Enlarging window enlarges  $W$  — Window can never be large enough!
- It's not deep enough to capture nuanced contextual meanings



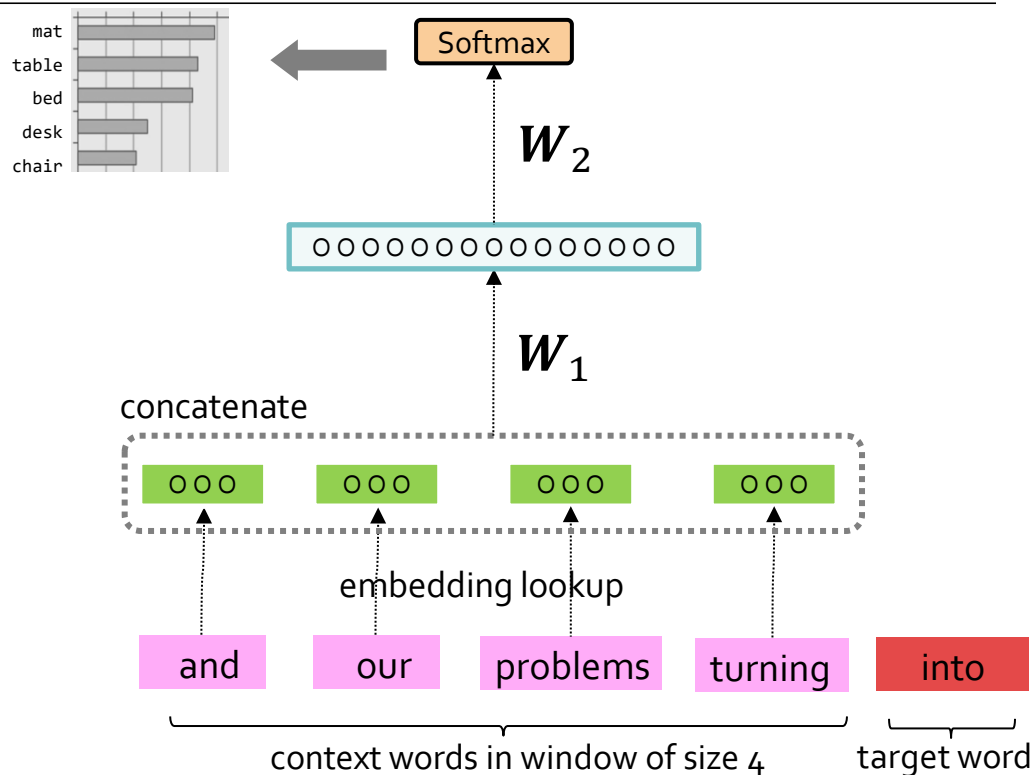
# A Fixed-Window Neural LM

Improvements over n-gram LM:

- Tackles the sparsity problem
- Model size is  $O(n)$  not  $O(\exp(n))$  —  $n$  being the window size.

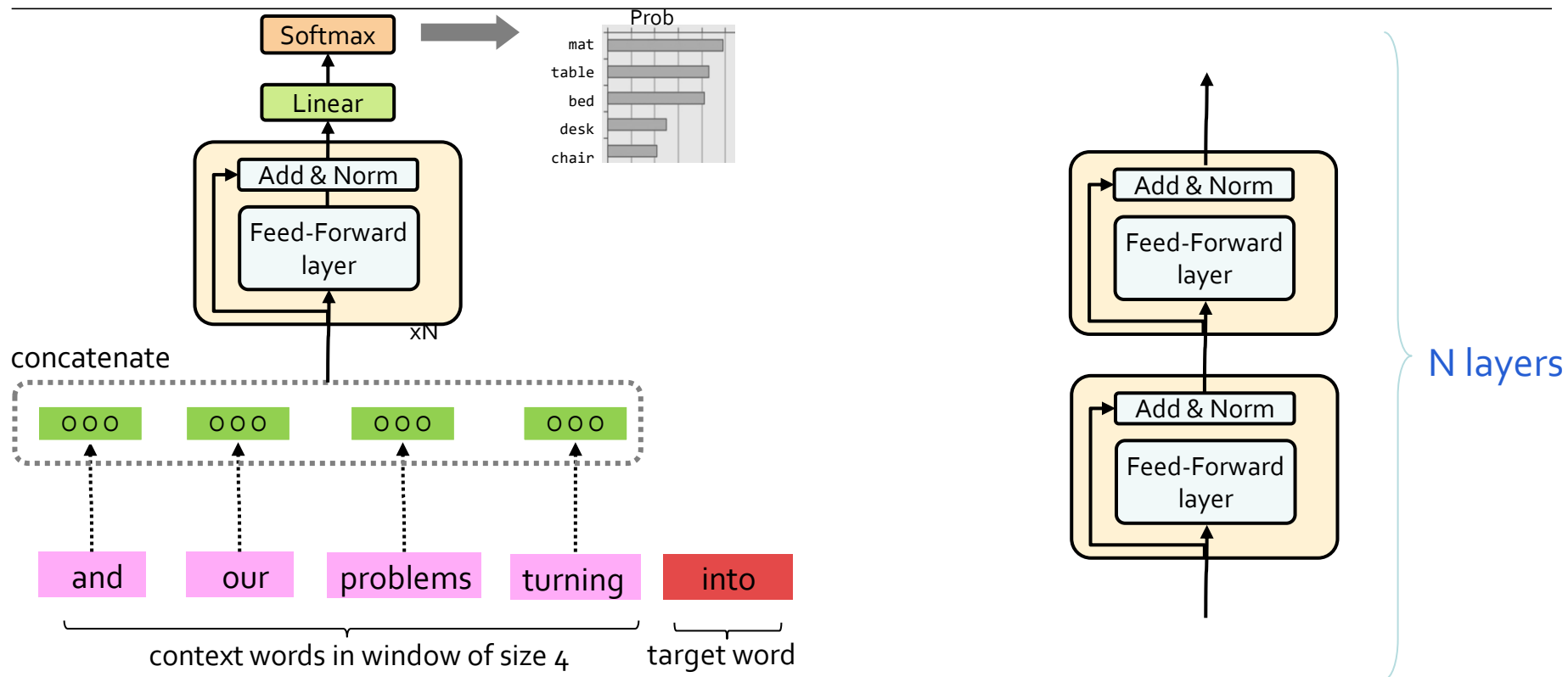
Remaining problems:

- Fixed window is too small
- Enlarging window enlarges  $W$  — Window can never be large enough!
- It's not deep enough to capture nuanced contextual meanings

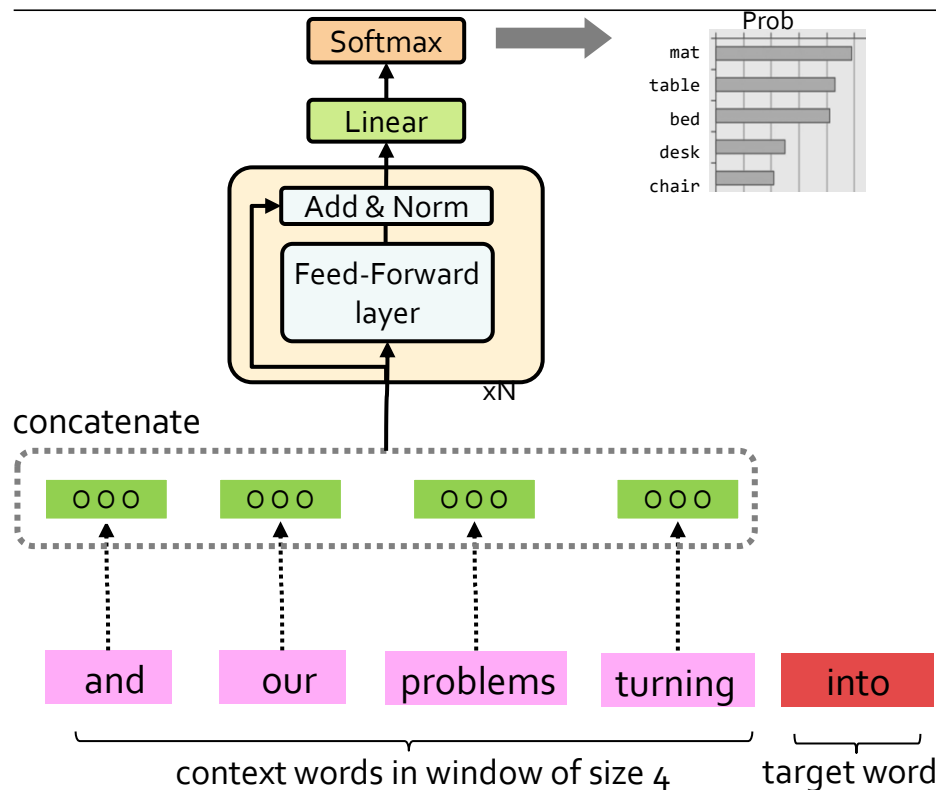




# Sun and Iyer (2021): Revisiting Simple Neural Probabilistic Language Models



# Sun and Iyer (2021): Revisiting Simple Neural Probabilistic Language Models



Model	# Params	Val. perplexity
Transformer	148M	25.0
NPLM-old	32M <sup>2</sup>	216.0
NPLM-old (large)	221M <sup>3</sup>	128.2
NPLM 1L	123M	52.8
NPLM 4L	128M	38.3
NPLM 16L	148M	<b>31.7</b>
- Residual connections	148M	660.0
- Adam, + SGD	148M	418.5
- Layer normalization	148M	33.0

Table 1: NPLM model ablation on WIKITEXT-103.

# What Changed from N-Gram LMs to Neural LMs?

What is the source of Neural LM's strength?

Why sparsity is less of an issue for Neural LMs?

**Answer:** In n-grams, we treat all prefixes independently of each other! (even those that are semantically similar)

students opened their \_\_\_\_  
pupils opened their \_\_\_\_  
scholars opened their \_\_\_\_  
undergraduates opened their \_\_\_\_  
students turned the pages of their \_\_\_\_  
students attentively perused their \_\_\_\_  
...

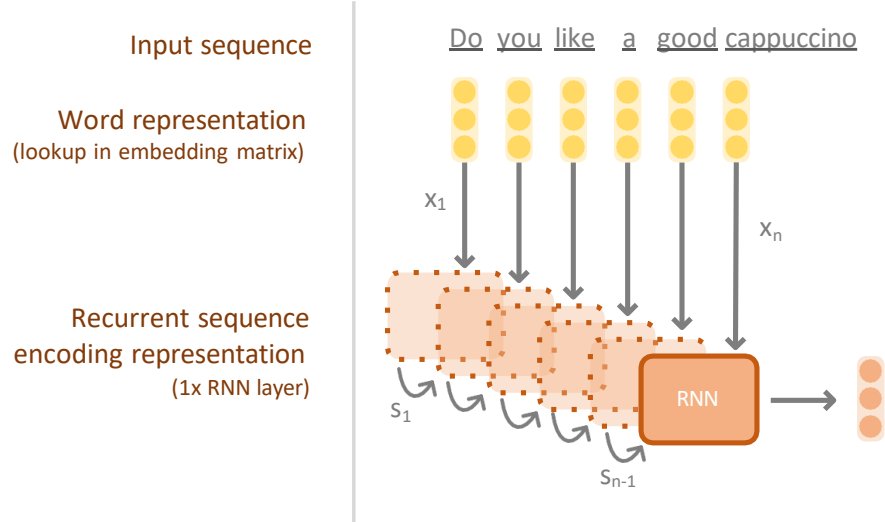
Neural LMs are able to share information across these semantically-similar prefixes and overcome the sparsity issue.

**Recap: LM and Sub-word tokenization**

**Basic Neural Language Models**

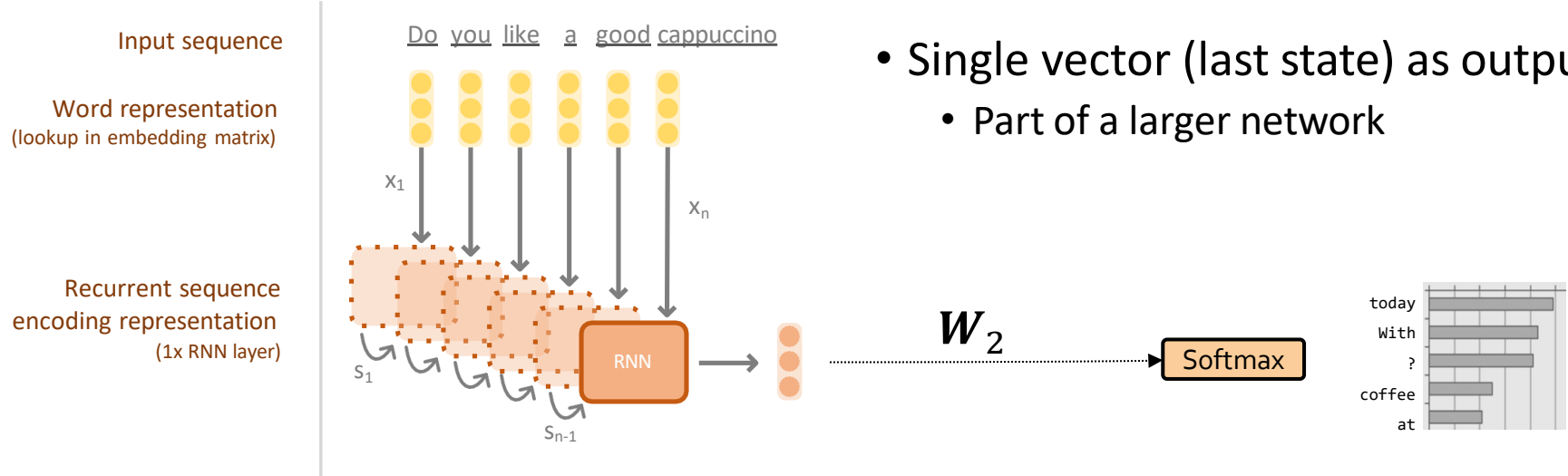
**RNN and Transformer LM**

# RNN Language Model



- Sequence as the input
- Single vector (last state) as output
  - Part of a larger network

# RNN Language Model



- Sequence as the input
- Single vector (last state) as output
  - Part of a larger network

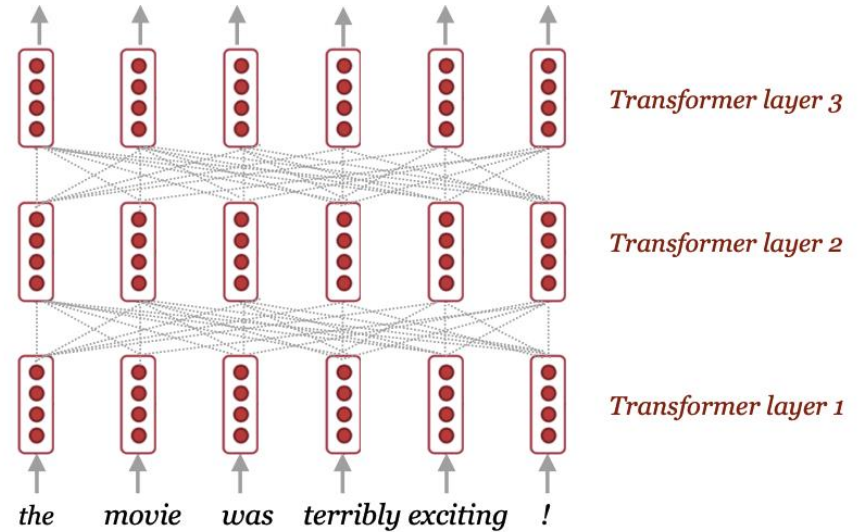
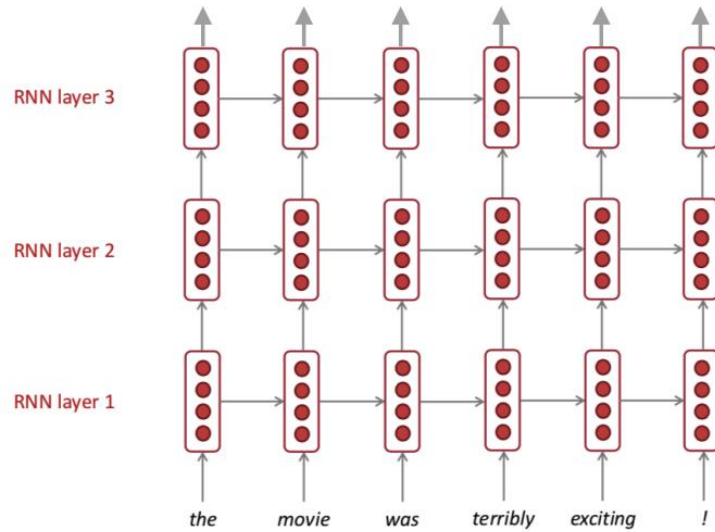
What are the cons?

- While RNNs in theory can represent long sequences, they quickly forget portions of the input.
- Vanishing/exploding gradients
- Difficult to parallelize



What can we do?

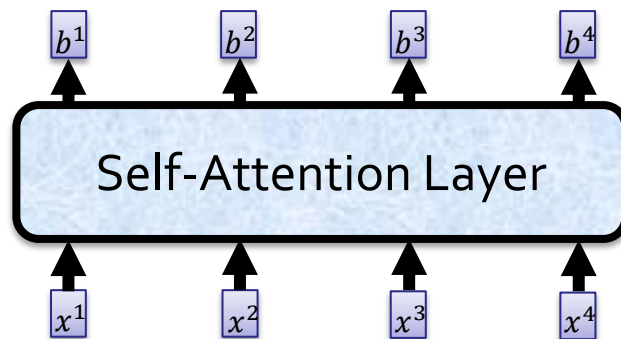
# RNN vs Transformer





# Self-Attention: Back to Big Picture

- **Attention** is a powerful mechanism to create context-aware representations
- A way to focus on select parts of the input



- Better at maintaining **long-distance dependencies** in the context.

[[Attention Is All You Need, Vaswani et al. 2017](#)]

# Properties of Self-Attention

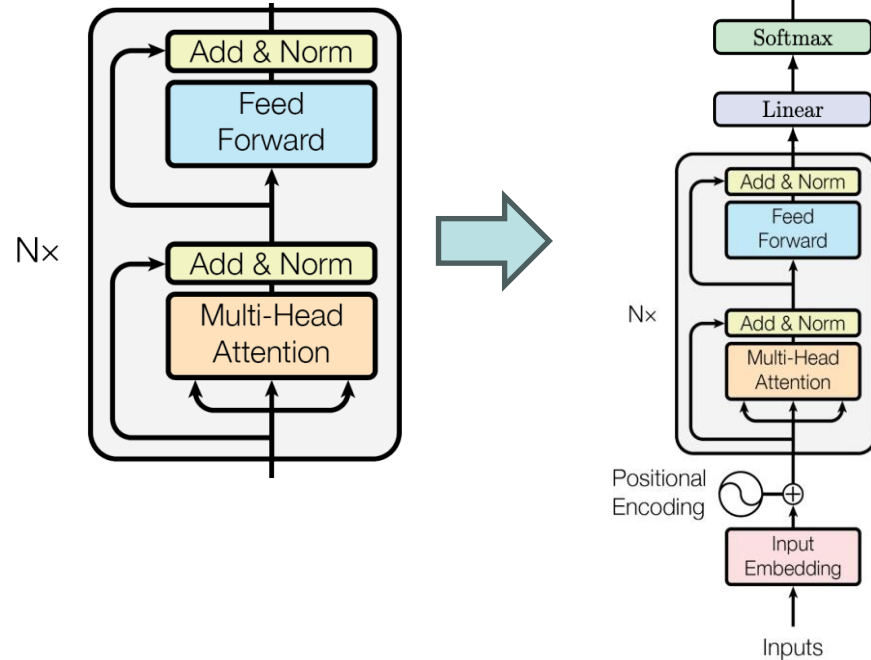
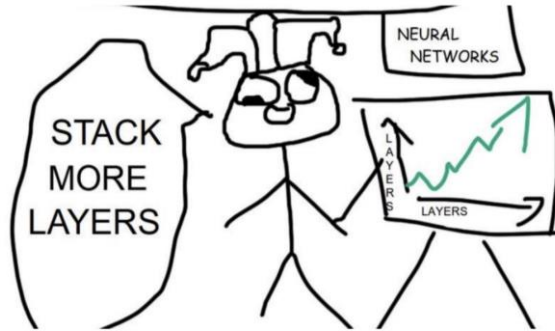
Layer Type	Complexity per Layer	Sequential Operations
Self-Attention	$O(n^2 \cdot d)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$

- $n$  = sequence length,  $d$  = hidden dimension
- Quadratic complexity, but:
  - $O(1)$  sequential operations (not linear like in RNN)
- Efficient implementations

[[Attention Is All You Need, Vaswani et al. 2017](#)]

# How Do We Make it **Deep**?

- Step 1: Stack more layers!
- Step 2: ...
- Step 3: Profit!



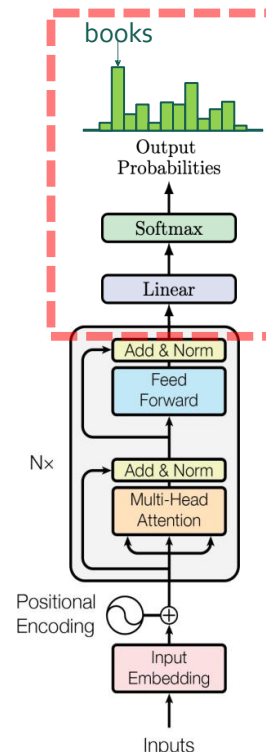
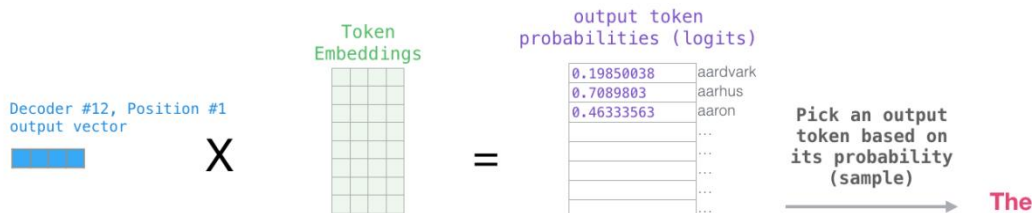
# From Representations to Prediction

- To perform prediction, add a classification head on top of the final layer of the transformer.
- This can be per token (Language modeling)
- Or can be for the entire sequence (only one token)

$\text{out} \in \mathbb{R}^{S \times d}$  (S: Sequence length)

$\text{logits} = \text{Linear}_{(d, V)}(\text{out}) = f(\text{out} \cdot W_V) \in \mathbb{R}^{S \times V}$

$\text{probabilities} = \text{softmax}(\text{logits}) \in \mathbb{R}^{S \times V}$



# Transformer-based Language Modeling

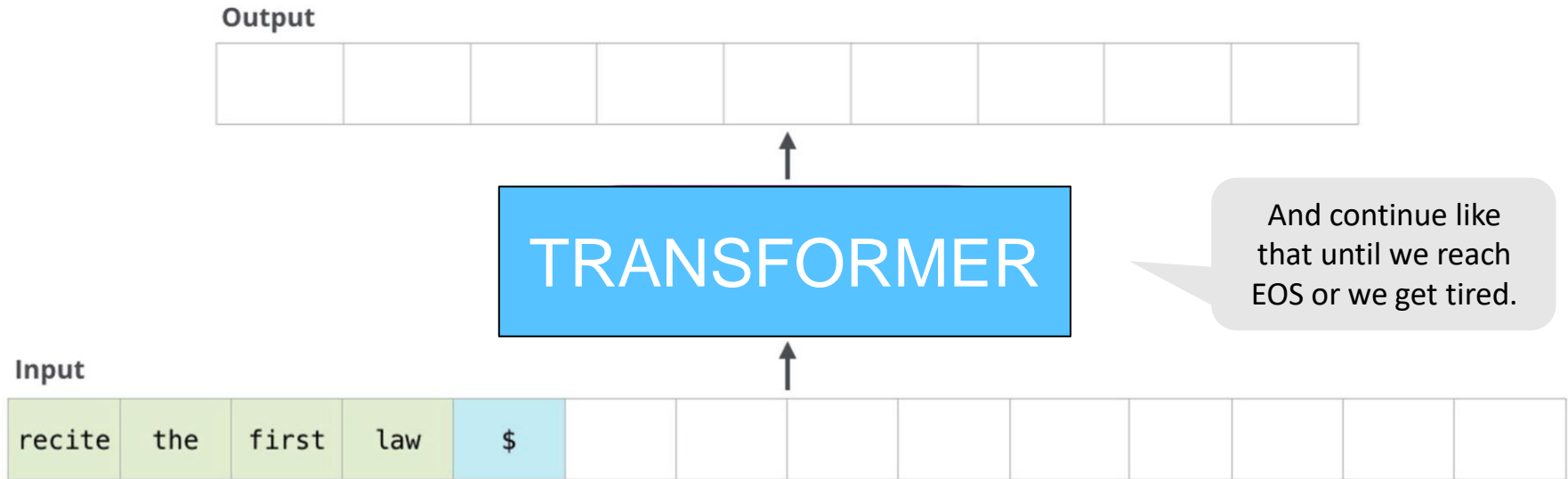


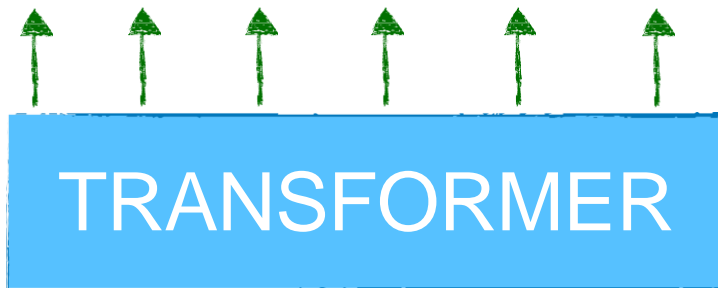
Image by <http://jalammar.github.io/illustrated-gpt2/>

# Training a Transformer Language Model

- **Goal:** Train a Transformer for language modeling (i.e., predicting the next word).
- **Approach:** Train it so that each position is predictor of the next (right) token.
  - We just shift the input to right by one, and use as labels

(gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad </s>$

EOS special token



$X = \text{the} \quad \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat}$

```
X = text[:, :-1]
Y = text[:, 1:]
```

[Slide credit: Arman Cohan]

# Training a Transformer Language Model

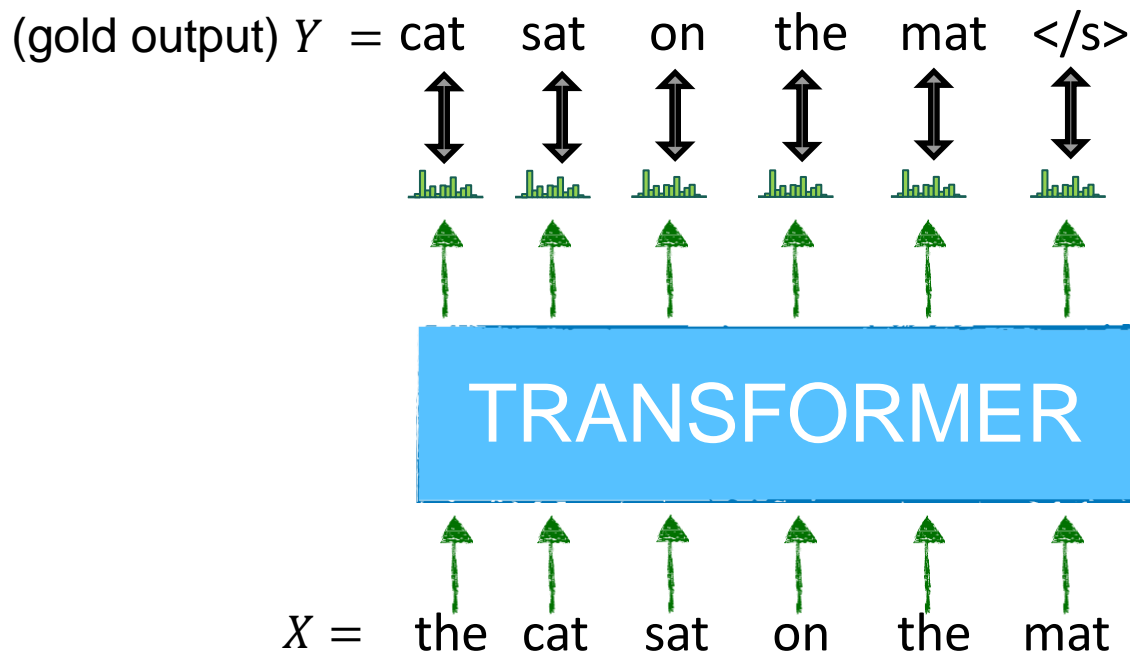
- For each position, compute their corresponding **distribution** over the whole vocab.

(gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad </s>$



# Training a Transformer Language Model

- For each position, compute the **loss** between the distribution and the gold output label.

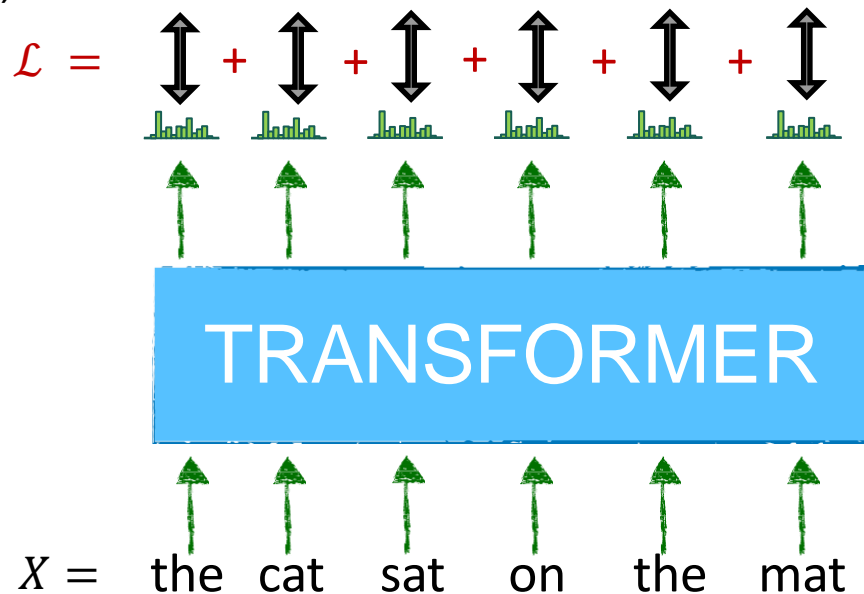




# Training a Transformer Language Model

- Sum the position-wise loss values to obtain a **global loss**.

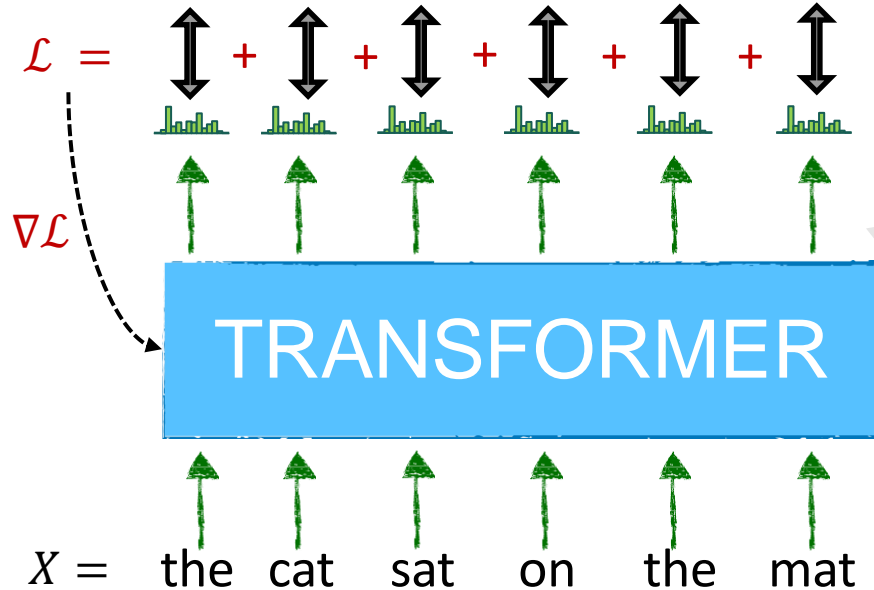
(gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad </s>$



# Training a Transformer Language Model

- Using this loss, do **Backprop** and **update** the Transformer parameters.

(gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad </s>$

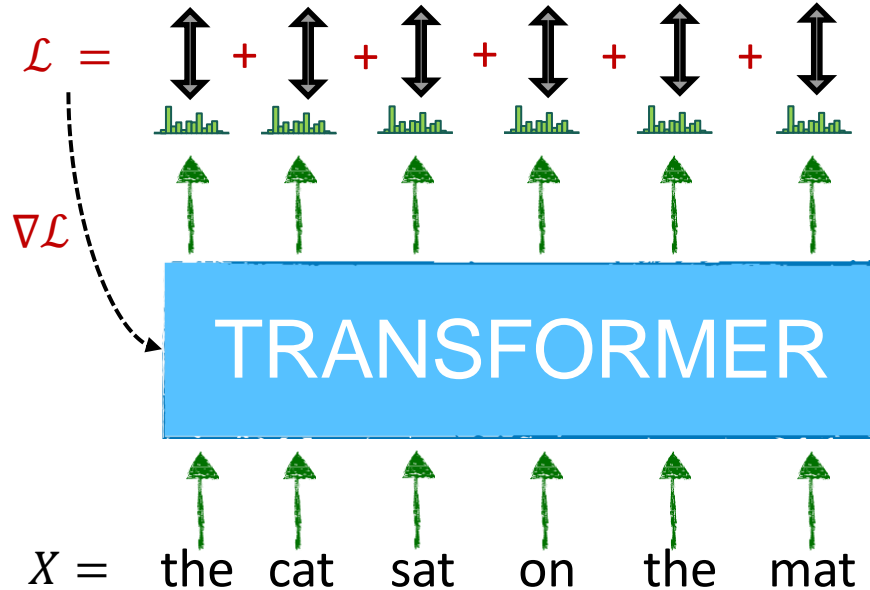


Well, this is not quite right ? ...  
what is the problem with this?

# Training a Transformer Language Model

- The model would solve the task by **copying** the next token to output (data leakage).

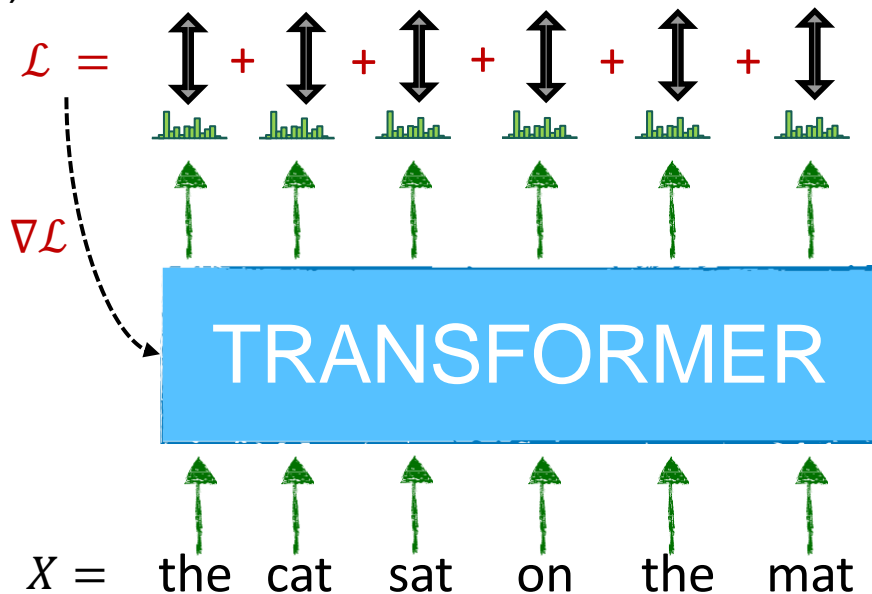
(gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad \text{</s>}$



# Training a Transformer Language Model

- We need to **prevent information leakage** from future tokens! How?

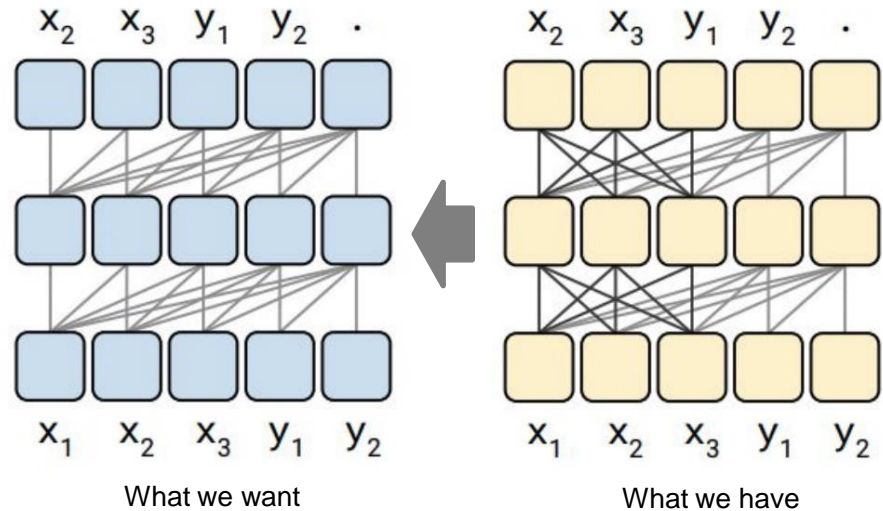
(gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad </s>$



# Attention mask

Attention raw scores

0	-0.08	1.24	0.69	-0.98	1.43	-0.6	0.7	0.16	0.93	1.28	-1.61	-1.1
1	-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	-1.01	1.01
2	0.86	1.19	1.59	0.86	-0.13	-0.15	-2.13	-0.98	-0.87	-1.72	1.87	-0.72
3	0.12	-0.03	-0.02	0.88	-0.46	-0.7	0.54	-0.42	-1.89	-0.38	0.04	-0.84
4	0.51	0.17	0.13	-1.64	0.24	-0.02	1.68	-0.36	0.64	0.36	0.27	0.66
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-0.31	1.54	1.66	1.14	0.58	-1.44
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.96	-0.17	-0.9	-1.57	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	0.5	-0.3
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-0.43
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11
	1	2	3	4	5	6	7	8	9	10	11	12



# Attention mask

Attention raw scores

0	-0.08	1.24	0.69	-0.98	1.43	-0.6	0.7	0.16	0.93	1.28	-1.61	-1.1
1	-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	-1.01	1.01
2	0.86	1.19	1.59	0.86	-0.13	-0.15	-2.13	-0.98	-0.87	-1.72	1.87	-0.72
3	0.12	-0.03	-0.02	0.88	-0.46	-0.7	0.54	-0.42	-1.89	-0.38	0.04	-0.84
4	0.51	0.17	0.13	-1.64	0.24	-0.02	1.68	-0.36	0.64	0.36	0.27	0.66
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-0.31	1.54	1.66	1.14	0.58	-1.44
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.96	-0.17	-0.9	-1.57	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	0.5	-0.3
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-0.43
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11
	1	2	3	4	5	6	7	8	9	10	11	12

Attention mask



# Attention mask

Attention raw scores

0	-0.08	1.24	0.69	-0.98	1.43	-0.6	0.7	0.16	0.93	1.28	-1.61	-1.1
1	-0.09	-0.0	-0.7	0.06	0.25	0.23	0.26	0.18	0.78	-0.21	-1.01	1.01
2	0.86	1.19	1.59	0.86	-0.13	-0.15	-2.13	-0.98	-0.87	-1.72	1.87	-0.72
3	0.12	-0.03	-0.02	0.88	-0.46	-0.7	0.54	-0.42	-1.89	-0.38	0.04	-0.84
4	0.51	0.17	0.13	-1.64	0.24	-0.02	1.68	-0.36	0.64	0.36	0.27	0.66
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-0.31	1.54	1.66	1.14	0.58	-1.44
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.96	-0.17	-0.9	-1.57	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	0.5	-0.3
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-0.43
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11
	1	2	3	4	5	6	7	8	9	10		

X

Attention mask



Note matrix multiplication is quite fast in GPUs.

# Attention mask

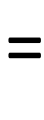
Attention raw scores

0	0.09	1.24	0.89	-0.96	1.43	-0.4	0.7	0.16	0.89	1.28	-1.01	-1.1
1	-0.09	-0.1	-0.7	0.06	0.25	0.23	0.26	0.16	0.78	-0.21	-1.01	1.01
2	0.98	1.19	1.06	0.88	-0.13	-0.19	-2.79	0.98	-0.87	-1.22	1.07	-0.72
3	0.12	0.03	0.02	0.96	-0.46	-0.27	0.56	-0.42	0.16	0.26	0.24	0.96
4	0.51	0.17	0.13	-1.84	0.26	-0.12	1.08	-0.36	0.04	0.26	0.27	0.06
5	0.24	0.16	0.43	0.74	0.86	1.79	-0.31	1.34	1.06	1.14	0.56	-1.49
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-0.36	0.17	0.9	-1.07	0.22
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	-inf	-inf	-inf
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	-inf	-inf
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-inf
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11
	1	2	3	4	5	6	7	8	9	10	11	12

X

Raw attention scores

0	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
1	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
2	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
3	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
4	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
5	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
6	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
7	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
8	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
9	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
10	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
11	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
	1	2	3	4	5	6	7	8	9	10	11	12



=

Masked attention raw scores

0	-0.08	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
1	-0.09	-0.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
2	0.86	1.19	1.59	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
3	0.12	-0.03	-0.02	0.88	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
4	0.51	0.17	0.13	-1.64	0.24	-inf	-inf	-inf	-inf	-inf	-inf	-inf
5	0.24	-1.44	0.43	0.74	0.96	-1.21	-inf	-inf	-inf	-inf	-inf	-inf
6	0.26	-0.1	0.93	0.72	-0.38	1.65	0.47	-inf	-inf	-inf	-inf	-inf
7	-0.55	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-inf	-inf	-inf	-inf
8	0.74	-0.76	-0.44	-0.08	-1.38	-0.13	1.25	-1.37	1.84	-inf	-inf	-inf
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.52	1.06	1.87	-inf	-inf
10	1.56	0.9	0.39	1.46	1.44	-1.05	0.9	-0.73	0.36	-0.67	-0.62	-inf
11	0.32	0.74	0.44	-0.1	1.19	0.83	0.29	2.06	0.51	-0.26	1.51	0.11
	1	2	3	4	5	6	7	8	9	10	11	12



# Attention mask

Attention raw scores

0	0.09	1.24	0.89	-0.96	1.43	-0.4	0.7	0.16	0.89	1.28	1.61	-1.1
1	-0.09	-0.0	-0.7	0.66	0.25	0.23	0.26	0.16	0.78	-0.21	1.01	1.01
2	0.88	1.19	1.06	0.88	-0.13	-0.19	2.79	0.89	-0.87	1.22	1.87	-0.72
3	0.12	0.03	0.02	0.96	-0.46	-0.27	0.56	-0.42	0.8	0.26	0.24	0.96
4	0.35	0.17	0.13	-1.84	0.24	-0.02	1.68	-0.36	0.84	0.26	0.27	0.86
5	0.24	1.31	0.43	0.74	0.86	1.79	-0.31	1.34	1.66	1.14	0.56	1.44
6	0.26	-0.1	0.89	0.72	-0.38	1.81	0.87	-0.86	0.17	0.9	-1.07	0.22
7	-0.85	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-0.7	-0.04	1.54	0.81
8	0.74	-0.76	0.44	-0.06	-1.06	-0.13	1.28	-1.1	1.84	0.3	0.57	0.74
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.1	1.06	1.87	0.5	-0.3
10	1.96	0.9	0.36	1.66	1.64	1.08	0.19	-0.73	0.36	0.87	0.62	-0.43
11	0.02	0.74	0.44	-0.1	1.19	0.88	0.26	0.26	0.51	-0.26	1.01	0.11
	1	2	3	4	5	6	7	8	9	10	11	12

X

Raw attention scores

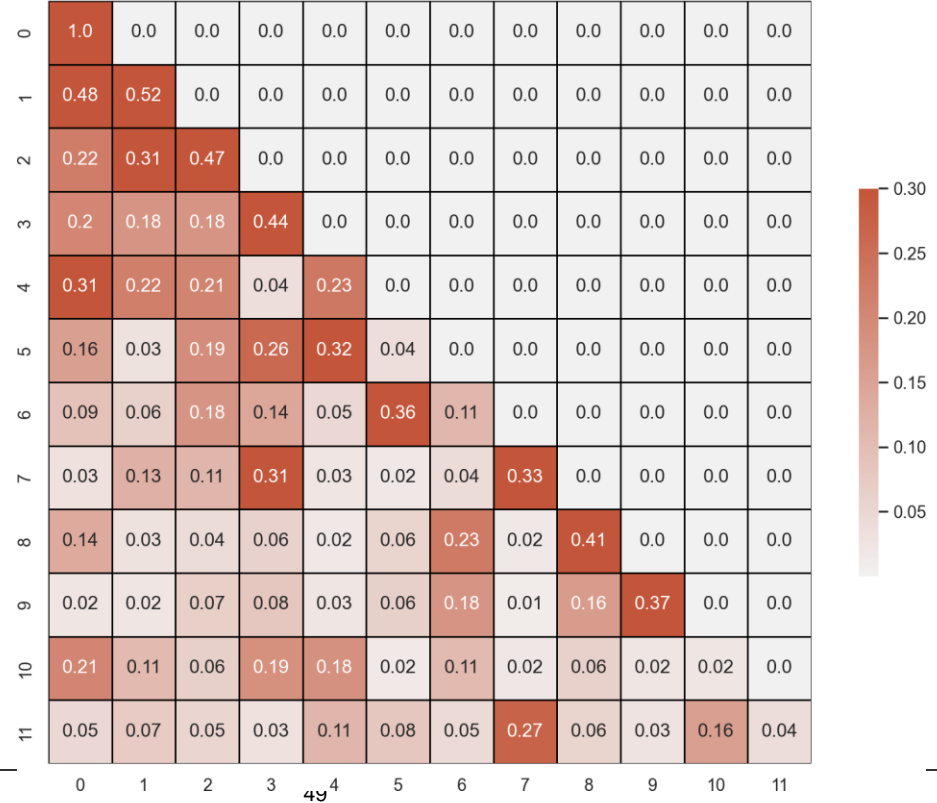
0	1.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
1	1.0	1.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
2	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
3	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
4	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf	-inf	-inf
5	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf	-inf
6	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf	-inf
7	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf	-inf
8	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf	-inf
9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf	-inf
10	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-inf
11	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
	0	1	2	3	4	5	6	7	8	9	10	11

softmax

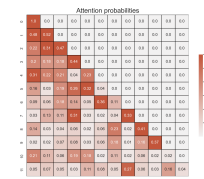
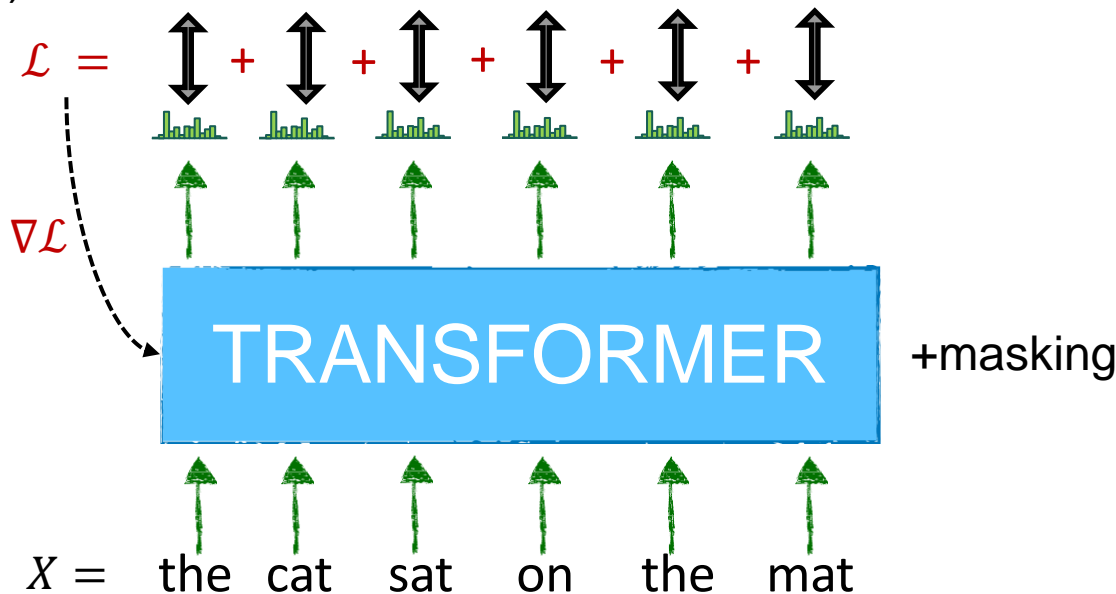
Masked attention raw scores

0	-0.08	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
1	-0.09	-0.3	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
2	0.86	1.19	1.06	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
3	0.12	-0.03	-0.02	0.88	-inf	-inf	-inf	-inf	-inf	-inf	-inf	-inf
4	0.31	0.17	0.13	1.44	0.24	-inf	-inf	-inf	-inf	-inf	-inf	-inf
5	0.24	-1.44	0.43	0.74	0.86	-1.21	-inf	-inf	-inf	-inf	-inf	-inf
6	0.26	-0.1	0.83	0.72	0.38	1.65	0.47	-inf	-inf	-inf	-inf	-inf
7	-0.86	0.81	0.71	1.7	-0.8	-1.14	-0.32	1.78	-inf	-inf	-inf	-inf
8	0.74	-0.76	-0.46	-0.06	-1.06	-0.13	1.28	-1.1	1.84	-inf	-inf	-inf
9	-0.97	-0.91	0.15	0.35	-0.81	0.11	1.14	-1.1	1.06	1.87	-inf	-inf
10	1.96	0.9	0.36	1.66	1.64	1.08	0.19	-0.73	0.36	0.87	-0.42	-inf
11	0.02	0.74	0.44	-0.1	1.19	0.88	0.26	0.26	0.51	-0.26	1.01	0.11
	1	2	3	4	5	6	7	8	9	10	11	12

Attention probabilities

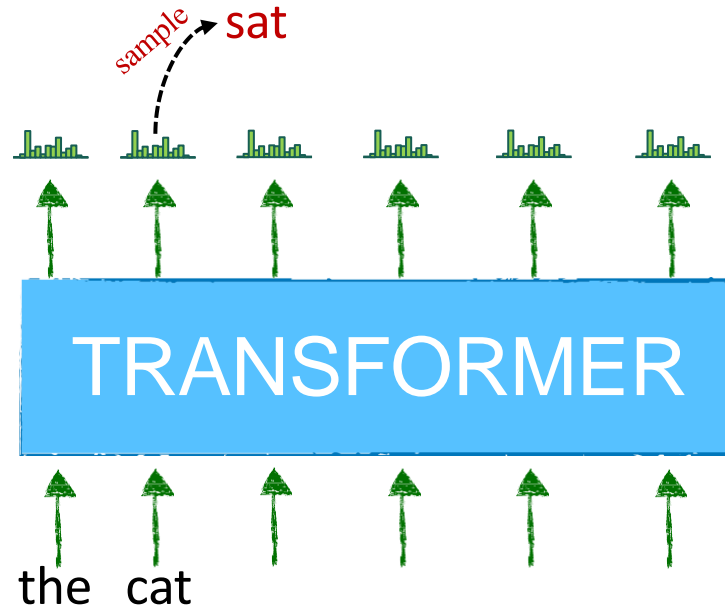


- (gold output)  $Y = \text{cat} \quad \text{sat} \quad \text{on} \quad \text{the} \quad \text{mat} \quad \text{</s>}$



# How to use the model to generate text?

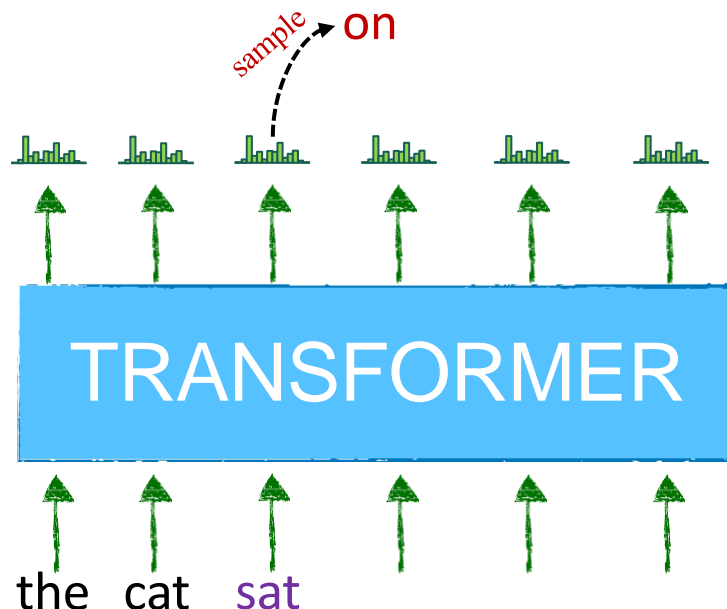
- Use the output of previous step as input to the next step repeatedly



# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

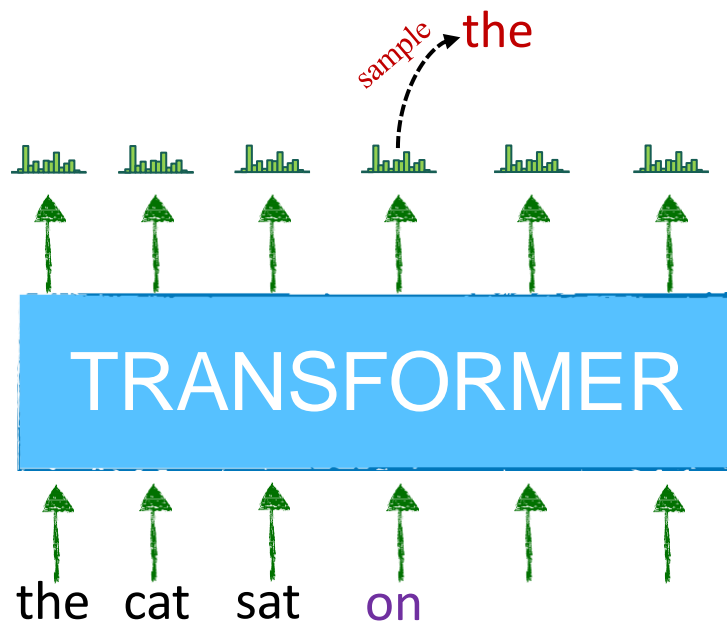
The probabilities get revised upon adding a new token to the input.



# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

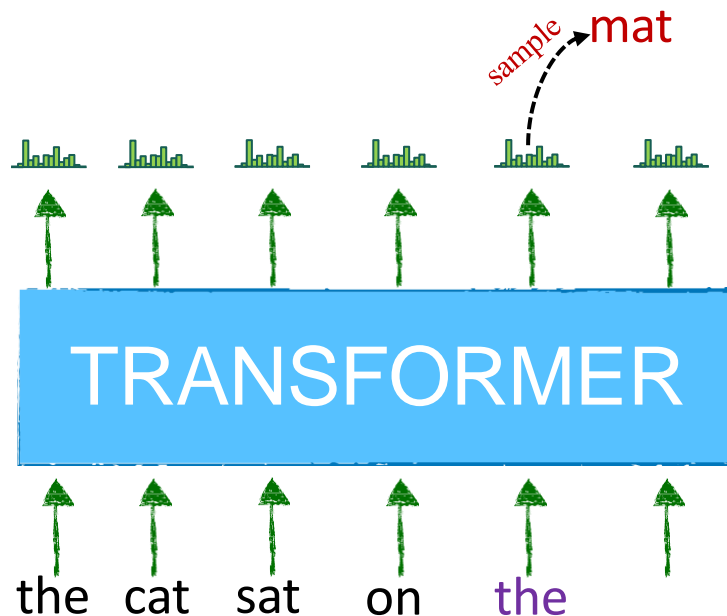
The probabilities get revised upon adding a new token to the input.



# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

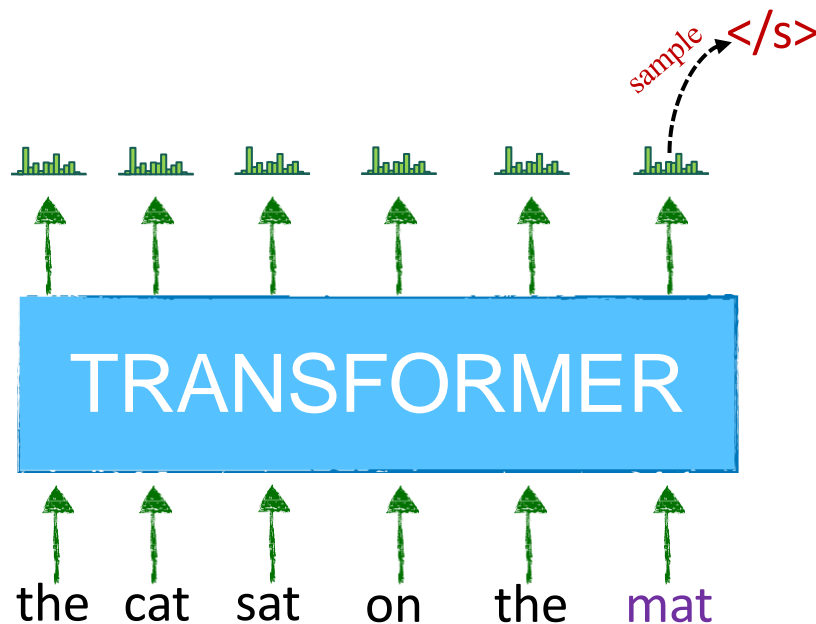
The probabilities get revised upon adding a new token to the input.



# How to use the model to generate text?

- Use the output of previous step as input to the next step repeatedly

The probabilities get revised upon adding a new token to the input.



- Neural models overcome n-gram limitations (sparsity, fixed windows) and model long-range dependencies.
- Fixed-window models reduce sparsity but lack depth and scalability for larger contexts.
- RNNs handle sequences but struggle with vanishing gradients and parallelization.
- Transformers use self-attention for efficient, context-aware representations, excelling in long-distance dependencies.
- Masking prevents future token leakage during training, enabling autoregressive text generation.





# Adaptation & Prompting