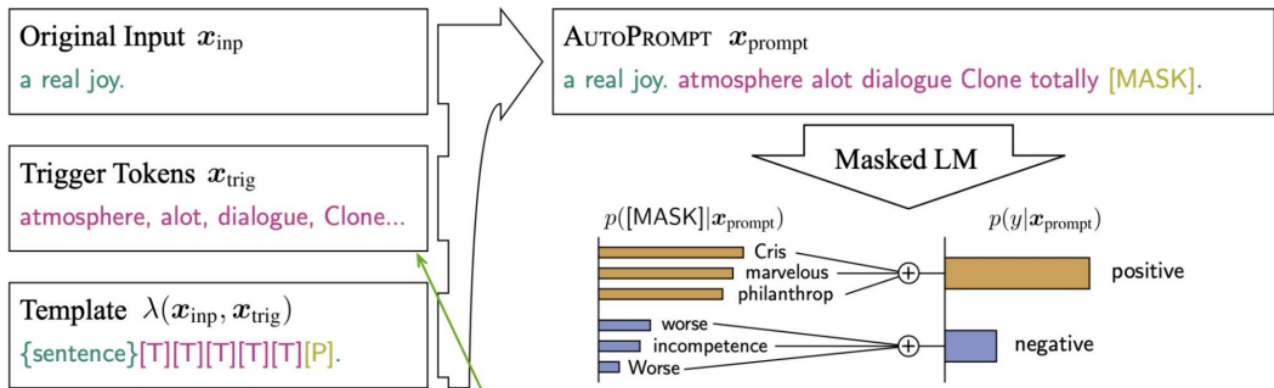


# FOLT Lecture 7, Text Generation

## Intro : AutoPrompt recap



In iterations, change trigger tokens to improve things

## AutoPrompt trigger search

Template  $\lambda(x_{\text{inp}}, x_{\text{trig}})$   
 $\{\text{sentence}\}[\text{T}][\text{T}][\text{T}][\text{T}][\text{T}][\text{P}]$ .

Original Input  $x_{\text{inp}}$   
 a real joy.

[P] or y { good  
bad

Template  $\lambda(x_{\text{inp}}, x_{\text{trig}})$   
 $\{\text{sentence}\}[\text{T}][\text{T}][\text{T}][\text{T}][\text{T}][\text{P}]$ .

Trigger Tokens  $x_{\text{trig}}$   
 atmosphere, alot, dialogue, Clone...

## AutoPrompt – Trigger Search

- Initialize all trigger tokens with [MASK] tokens
- Use a batch of example prompts to compute the likelihood  $p(y|x)$
- Use the gradients to switch one of the tokens
  - Select top-k tokens that improve  $p(y|x)$
- Evaluate on the development set
  - Select the best likelihood on dev
- Continue until the set of tokens is not changed anymore

Candidate Trigger	Batch Of Examples	$p(y x)$
[MASK] [MASK] [MASK]	An artful, intelligent...	0.01
	The uninspiring...	0.05
	It's a beautiful story...	0.03
[MASK] [MASK] [MASK]		
⋮ ⋮ ⋮		
man ##s cameo		
Rating fiennes go		
Rating [MASK] [MASK]	An artful, intelligent...	0.18
	The uninspiring...	0.11
	It's a beautiful story...	0.08
⋮ ⋮ ⋮		
Rating ##omi #!	An artful, intelligent...	0.95
	The inspirational...	0.89
	It's a beautiful story...	0.77

## AutoPrompt – Label Token Search

Template  $\lambda(x_{\text{inp}}, x_{\text{trig}})$   
 {sentence}[T][T][T][T][T][P].

[MASK].

Predict token

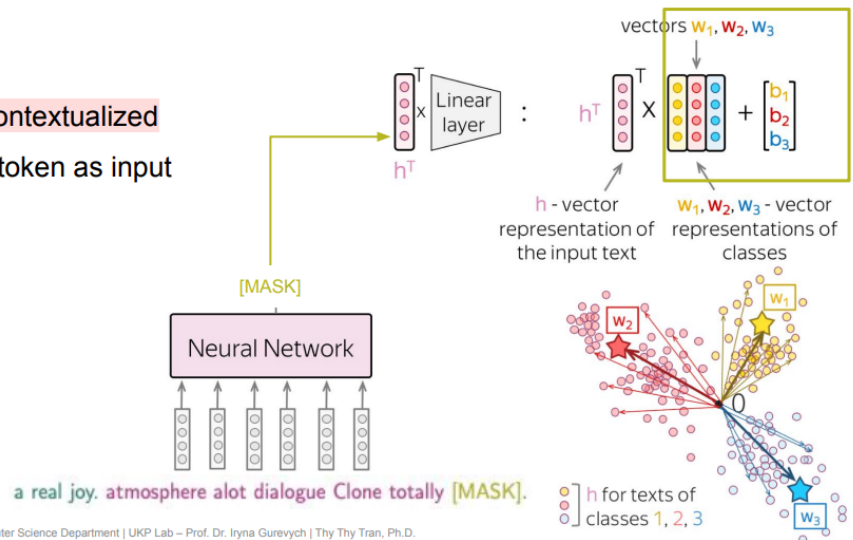
Template  $\lambda(x_{\text{inp}}, x_{\text{trig}})$   
 {sentence}[T][T][T][T][T][P].

Predict token [MASK]

## AutoPrompt – Label Token Search

### Two-step approach

1. Train a classifier using the contextualized embedding  $h$  of the [MASK] token as input



Computer Science Department | UKP Lab – Prof. Dr. Iryna Gurevych | Thy Thy Tran, Ph.D.

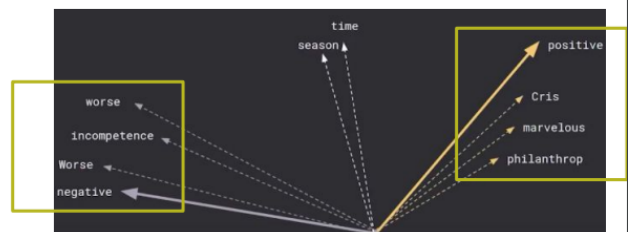
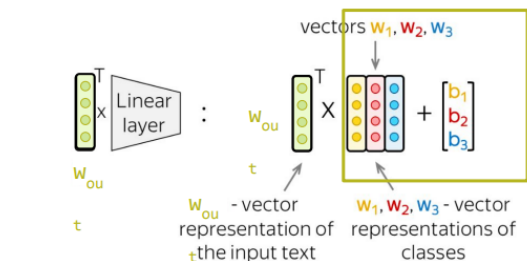
Template  $\lambda(x_{\text{inp}}, x_{\text{trig}})$   
 {sentence}[T][T][T][T][T][P].

Predict token [MASK]

## AutoPrompt – Label Token Search

### Two-step approach

1. Train a classifier using the contextualized embedding  $h$  of the [MASK] token as input
2. Substitute  $h$  with LM output word embeddings  $w_{out}$
3. Select the top-k  $w_{out}$  with highest scores



# Overview

Text gen examples : Machine Translation, Copilot, Text Summarization, Dialogue systems ....

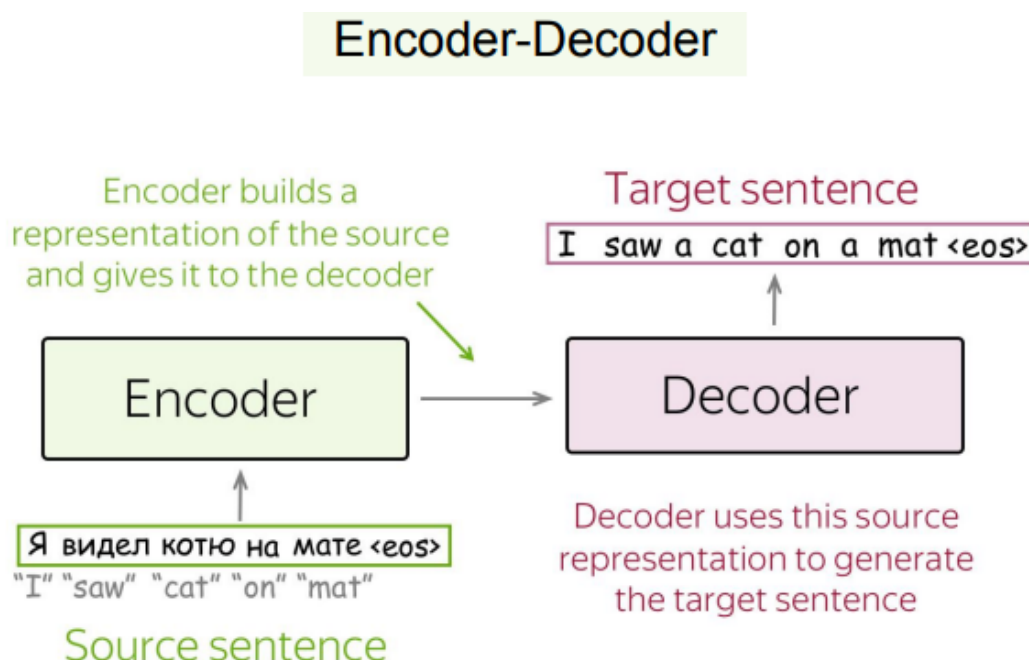
- A text sequence → text sequence in **similar length**
  - E.g., Machine Translation
- A text sequence → text sequence in **much shorter length**
  - E.g., Summarization
- A text sequence → text sequence in **varying lengths**
  - E.g., Conversation/Dialogue

## Key Concepts from previous lectures:

- Language models from lecture 2
- Neural language models lecture 3(Predict next word(token))
- Text generation with neural LM (lecture 2)

## PART1: Generative models

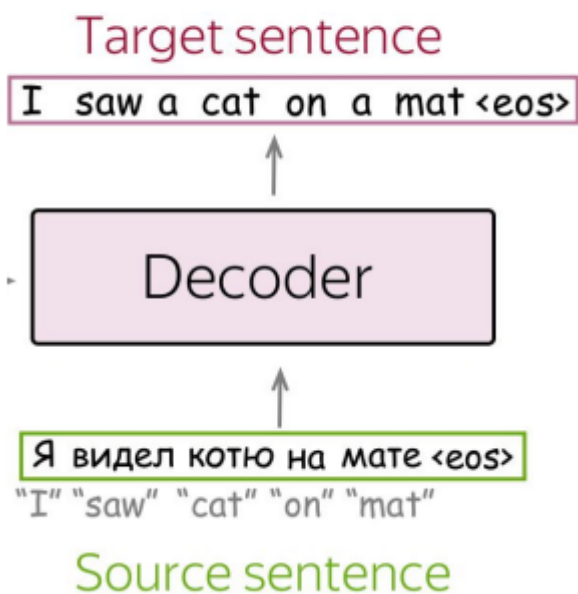
### Examples :



**Encoder** : processes the source sequence and produces its representation(s)

**Decoder** : uses the source representation(s) from the encoder to generate the target sequence

## Decoder-only



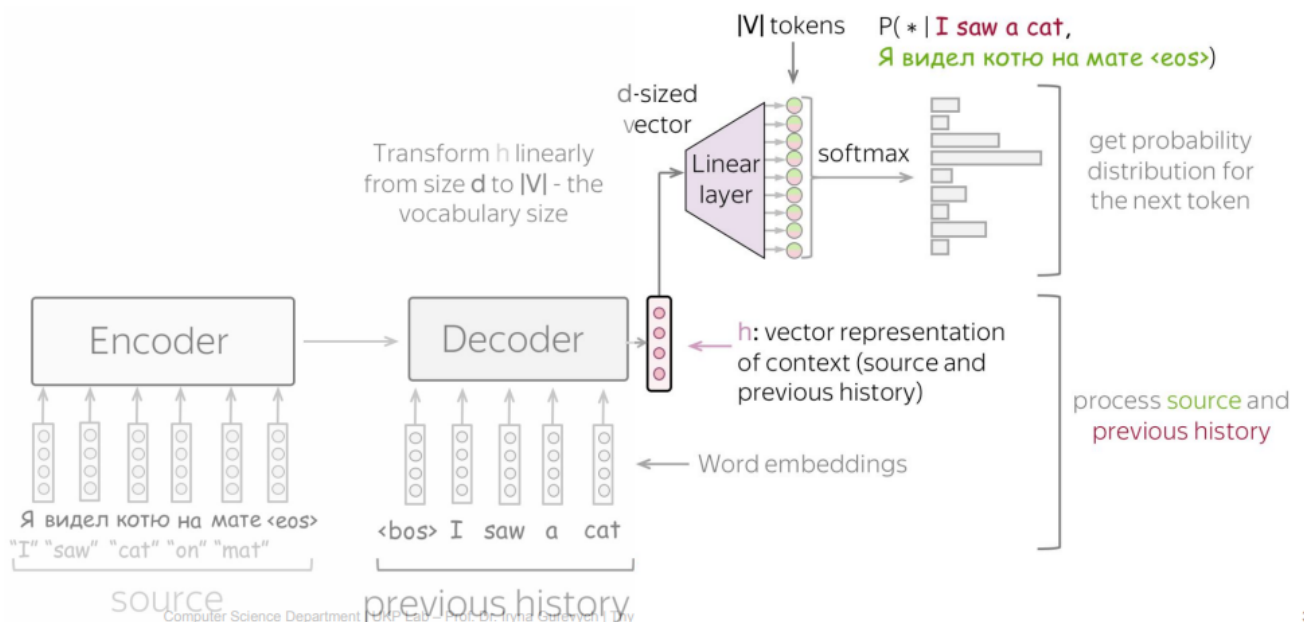
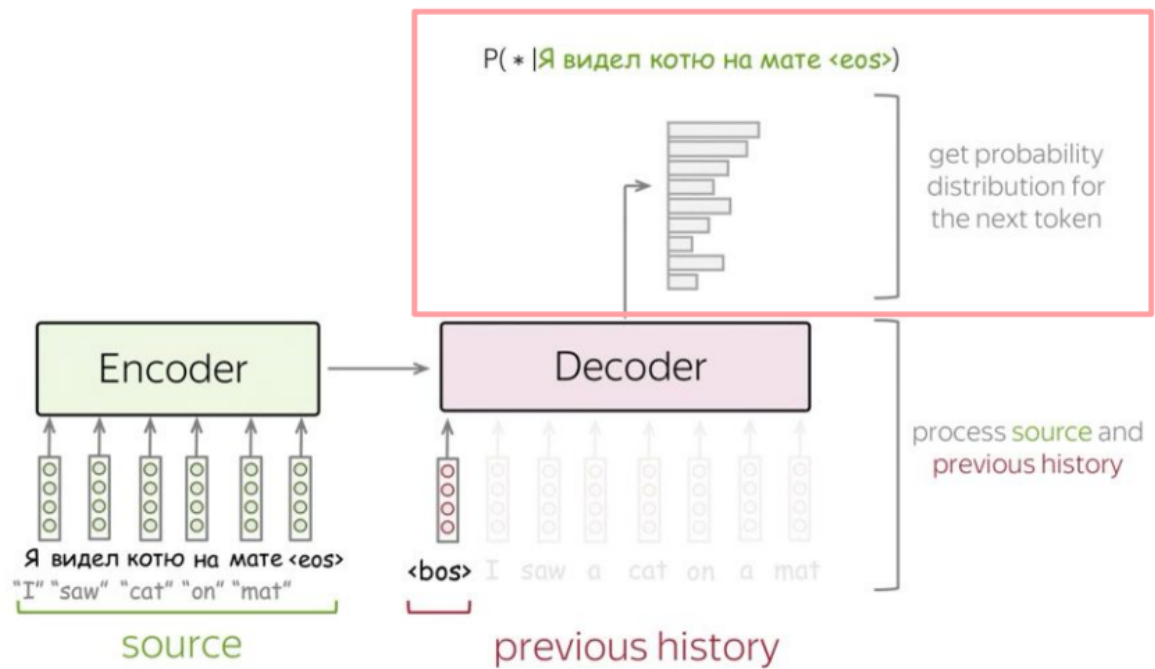
Also there is Encoder-only (e.g., BERT) but not covered in this lecture

### 1) Conditional Language Models

Language Models:  $P(y_1, y_2, \dots, y_n) = \prod_{t=1}^n p(y_t | y_{<t})$   
(left-to-right)

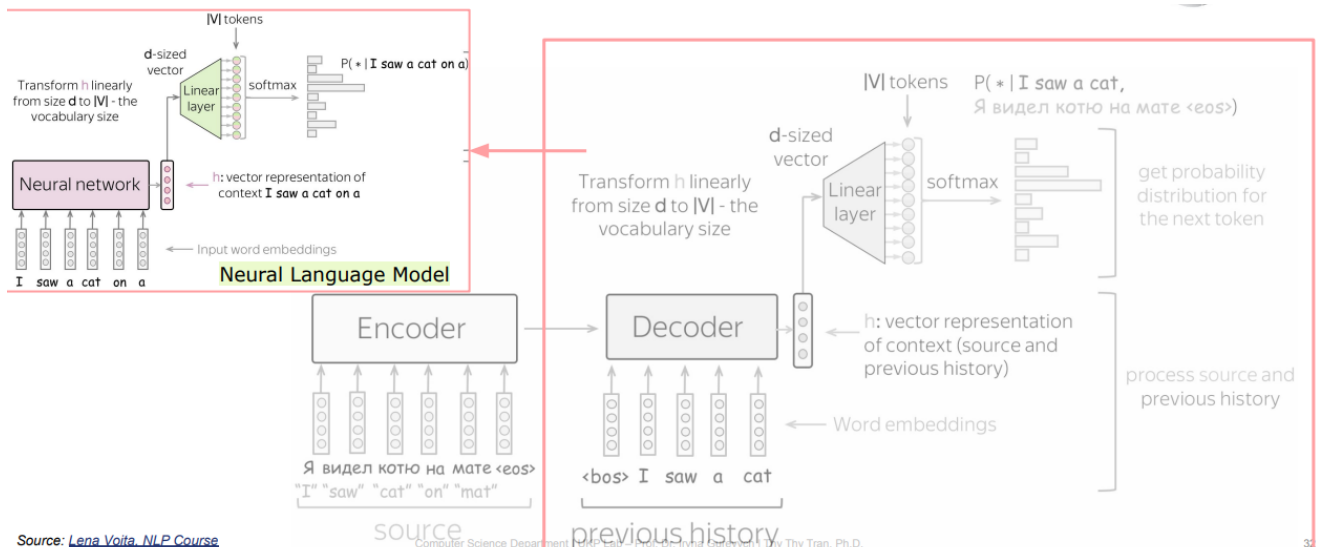
Conditional  
Language Models:  $P(y_1, y_2, \dots, y_n, | \mathbf{x}) = \prod_{t=1}^n p(y_t | y_{<t}, \mathbf{x})$   
condition on source  $\mathbf{x}$

### 2) Generation with Encoder-Decoder



Computer Science Department | NLP Lab | Prof. Dr. Irina Serebrenik | Duy

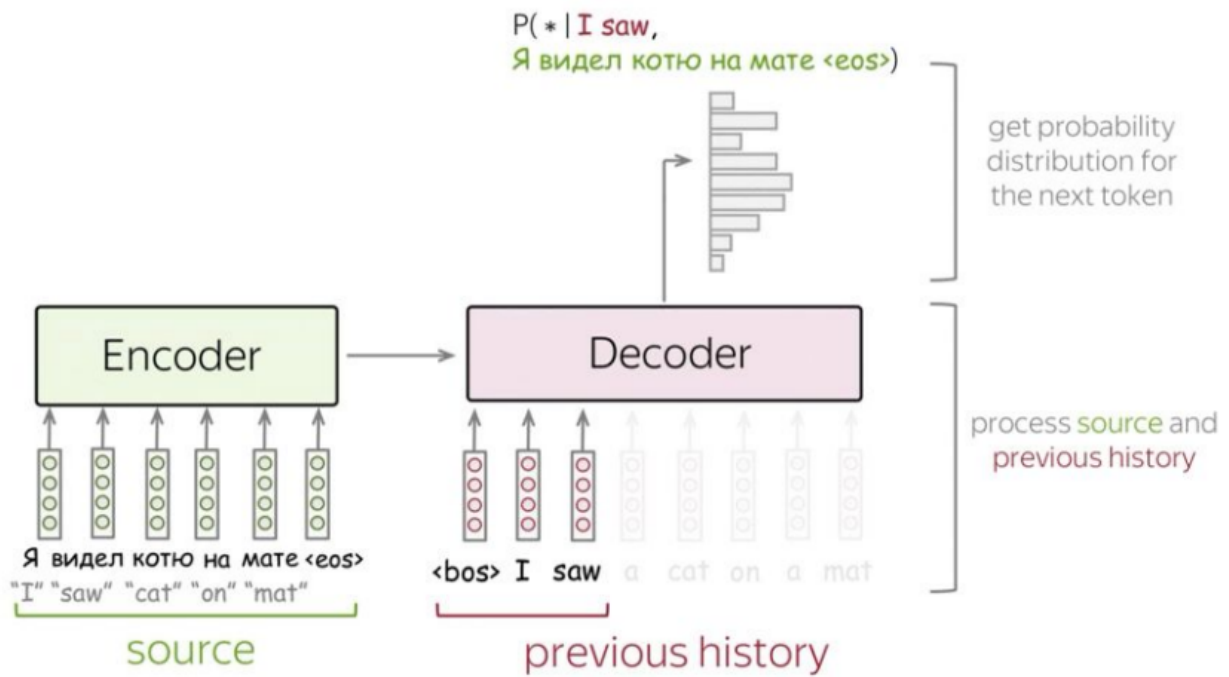
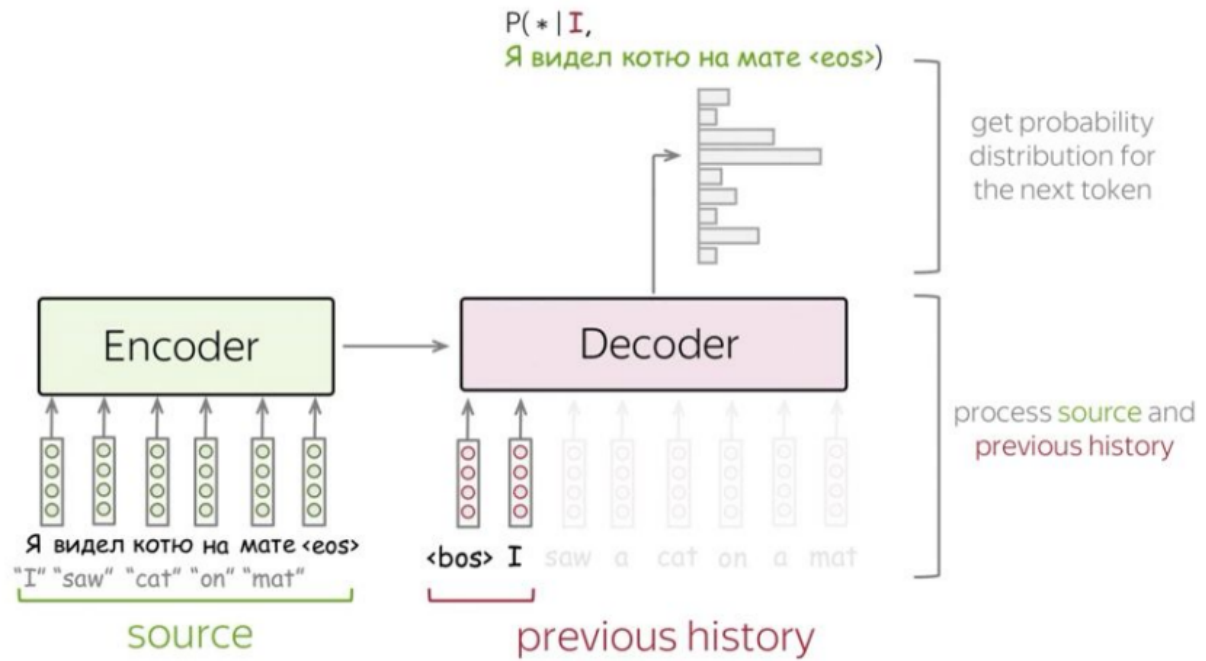
31



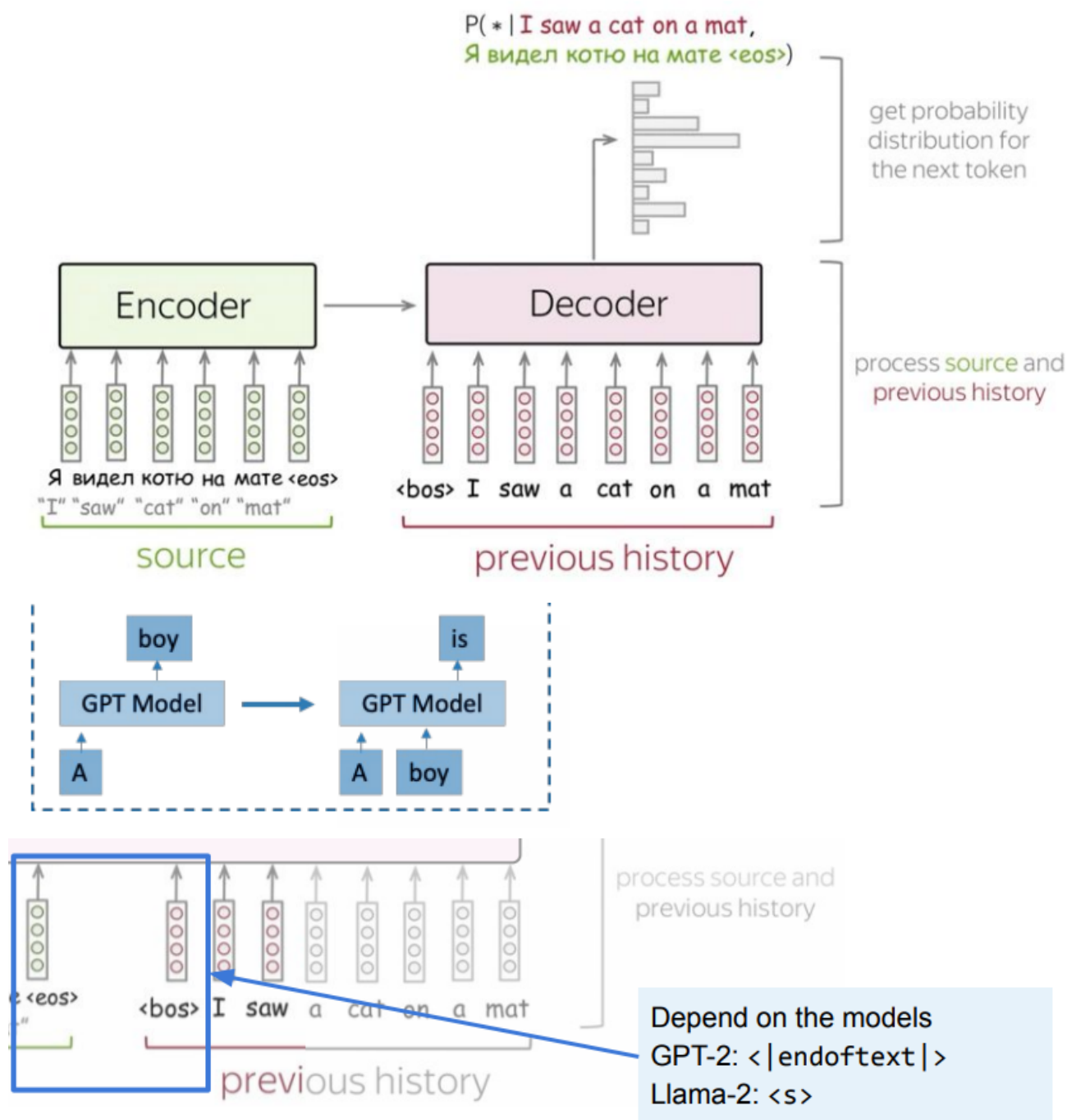
Source: [Lena Voita, NLP Course](#)

Computer Science Department | NLP Lab | Prof. Dr. Irina Serebrenik | Duy Thy Tran, Ph.D.

32



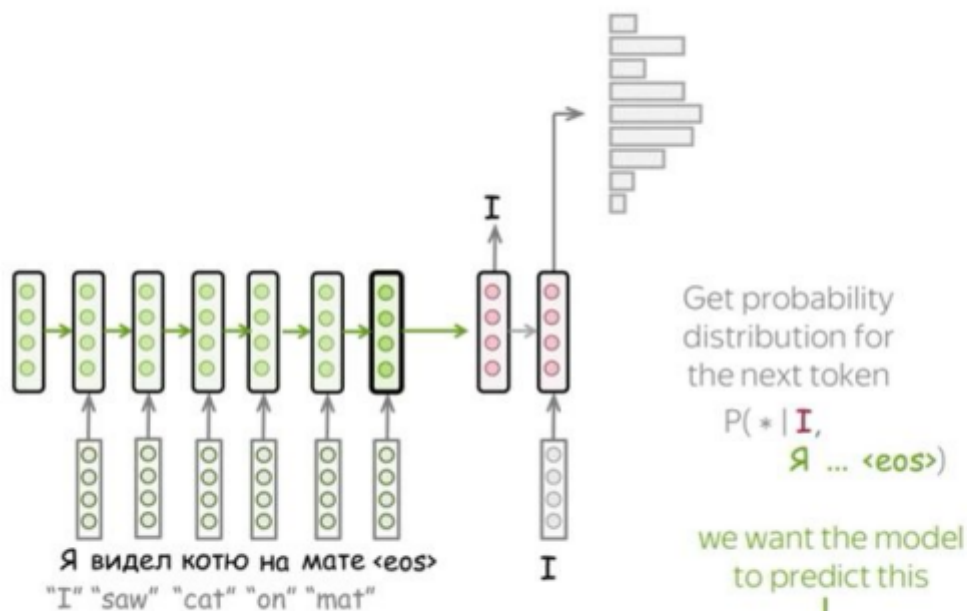




### 3) Training text generation

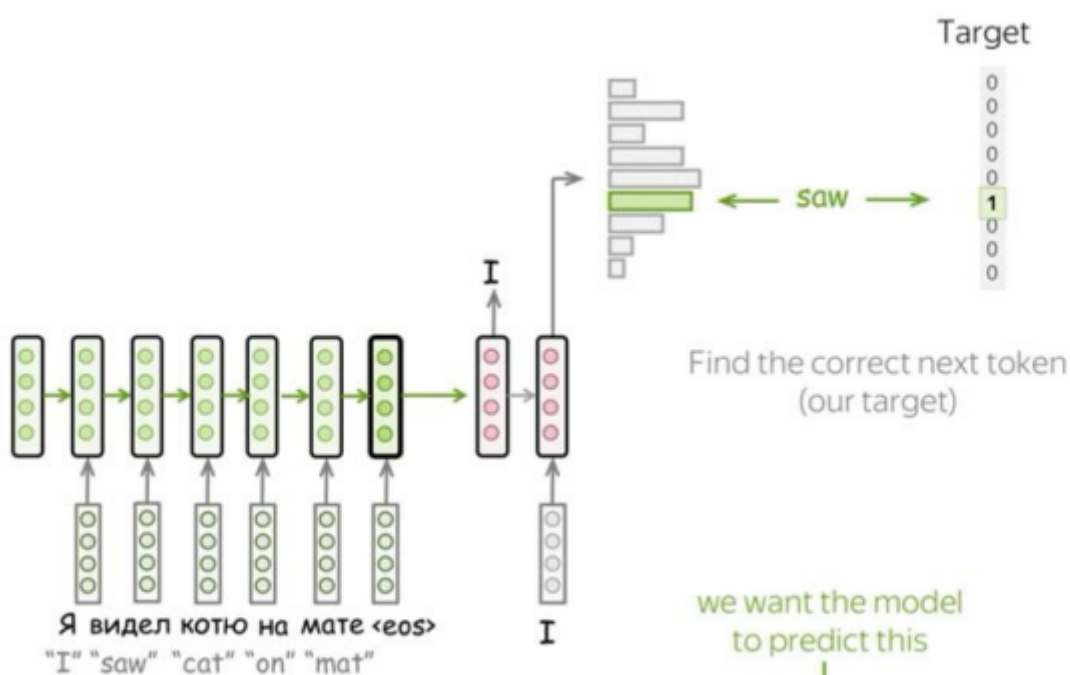
Similar to neural LMs, neural text generation models are trained to maximize the probability distribution of the next token given the previous context





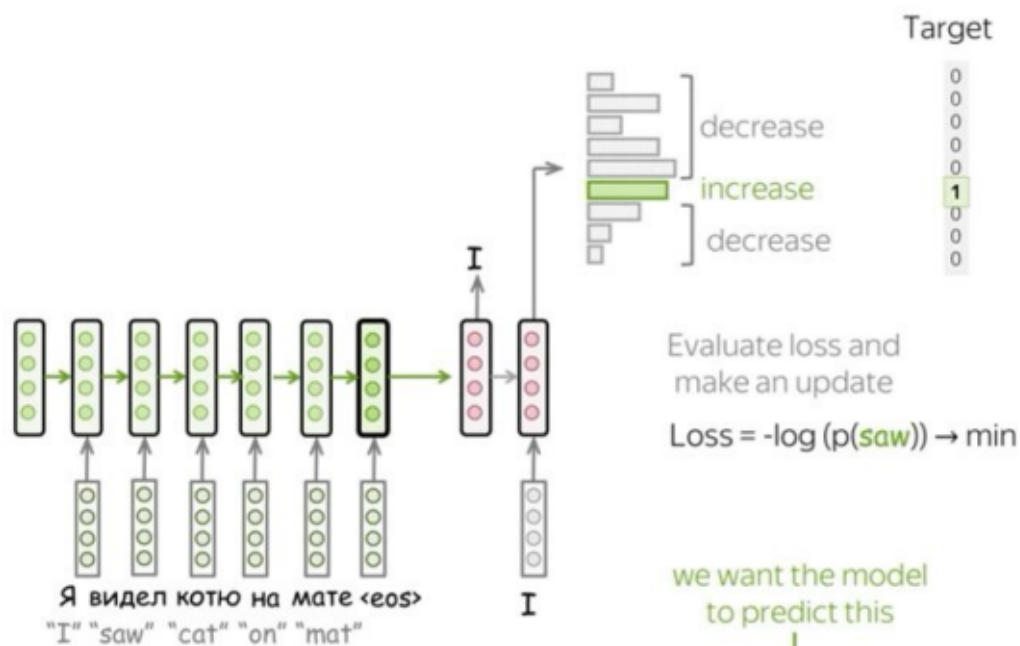
Source: Я видел котю на мате <eos>  
 "I" "saw" "cat" "on" "mat"

Target: I saw a cat on a mat <eos>



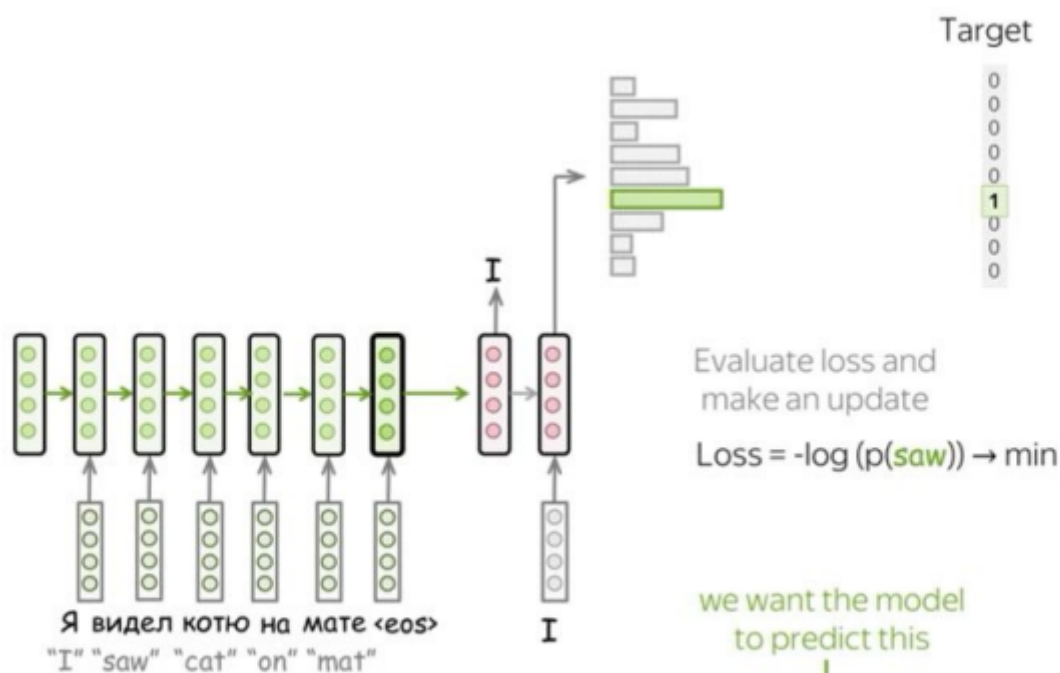
Source: Я видел котю на мате <eos>  
 "I" "saw" "cat" "on" "mat"

Target: I saw a cat on a mat <eos>



Source: Я видел котю на мате <eos>  
"I" "saw" "cat" "on" "mat"

Target: I saw a cat on a mat <eos>



Source: Я видел котю на мате <eos>  
"I" "saw" "cat" "on" "mat"

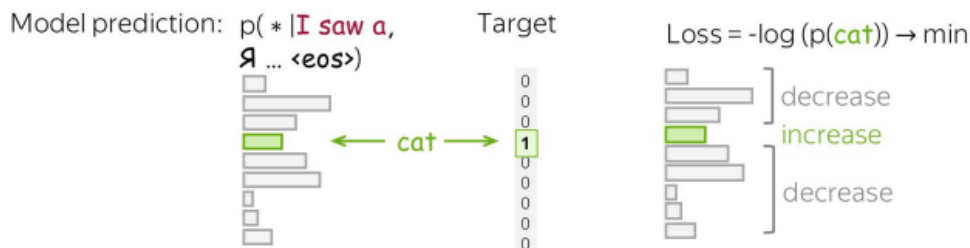
Target: I saw a cat on a mat <eos>

Source sequence: Я видел котю на мате <eos>  
 "I" "saw" "cat" "on" "mat"

Target sequence: I saw a cat on a mat <eos>  
 "I" "saw" "a" "cat" "on" "a" "mat"

← one training example  
 ← one step for this example

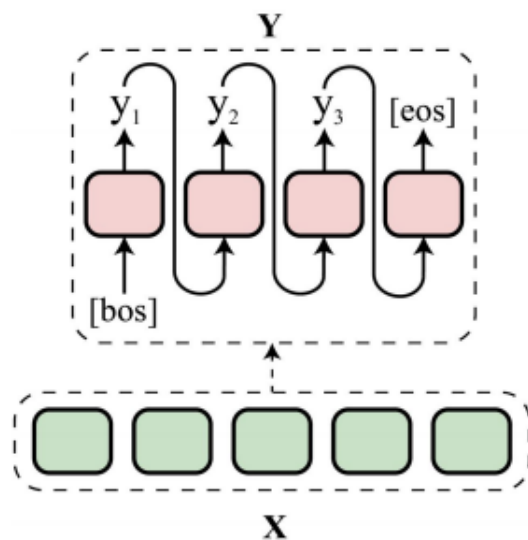
previous tokens we want the model to predict this



## PART 2: Decoding Strategies(Generating text)

### 1)Autoregressive Generation

*Similar to text generation with neural LMs*



Autoregressive Generation [AG]

- Starts with [bos] (begin-of-sequence)
- At each step
  - process previous generated tokens
  - get probability distribution for the next token
- Stops by [eos] (end-of-sequence)
  - Terminate when [eos] is predicted
  - Or stop generating text when max target sequence length is reached
- max target sequence length or max new tokens: expected maximum length of Y
- max token length: expected maximum length of X + Y

### 2) Probability for Next Token

We have:

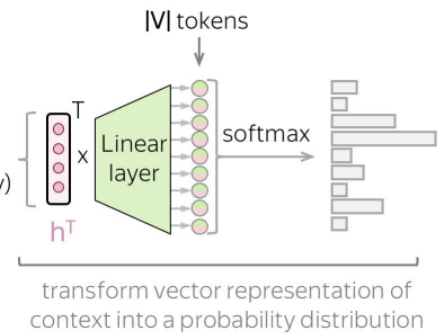
- $\mathbf{h}$  - vector of size  $d$

We need:

- vector of size  $|V|$  – probabilities for all tokens in the vocabulary

Transform linearly  
from size  $d$  to size  $|V|$

$d$  features  
(vector  
dimensionality)



Conditional

Language Models:  $P(y_1, y_2, \dots, y_n | \mathbf{x}) = \prod_{t=1}^n p(y_t | y_{<t}, \mathbf{x})$   
 condition on source  $\mathbf{x}$

$$\mathbf{y} = \text{softmax}(\mathbf{h}\mathbf{W} + \mathbf{b})$$

$$\mathbf{h} \in \mathbb{R}^d, \mathbf{W} \in \mathbb{R}^{d \times |V|}, \mathbf{b} \in \mathbb{R}^{|V|}$$

$$\frac{\exp(\mathbf{z}_c)}{\sum_{c' \in C} \exp(\mathbf{z}_{c'})}$$

$$p(y_t = w | y_{<t}, \mathbf{x}) \propto \exp(\text{score}(w))$$

### 3) Temperature – Tempered Sampling

$$p(y_t = w | y_{<t}, \mathbf{x}) \propto \exp(\text{score}(w))$$

$$q(y_t = w | y_{<t}, \mathbf{x}) \propto \exp(\text{score}(w)/T) \quad \text{where } T \in (0, +\infty)$$

- Typically we choose  $T \in (0, 1)$ , which makes the distribution more peaky.



What happens when  $T \rightarrow 0$  and  $T \rightarrow +\infty$ ?

### 4) Repetition Penalty

[Keskar et al., 2019] Intuition: to avoid generating duplicate substrings

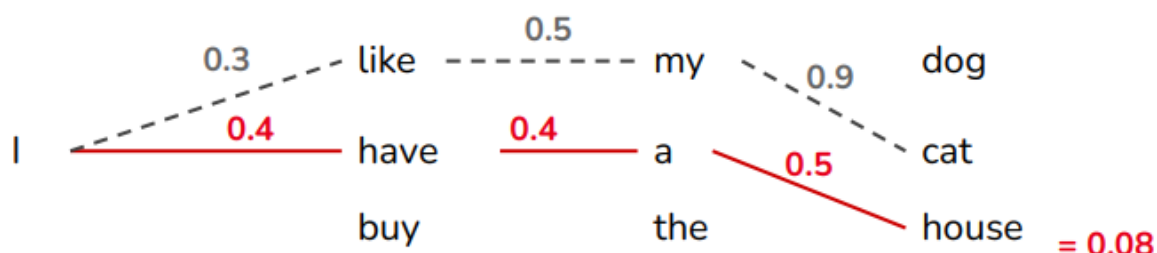
$$p_i = \frac{\exp(x_i / (T \cdot I(i \in g)))}{\sum_j \exp(x_j / (T \cdot I(j \in g)))} \quad I(c) = \theta \text{ if } c \text{ is True else } 1$$

- $g$  contains a set of previously generated tokens
- $\mathbb{1}$  is an identity function
- $\theta = 1.2$  is found to yield a good balance between less repetition and truthful generation.

## 5) Decoding Strategies

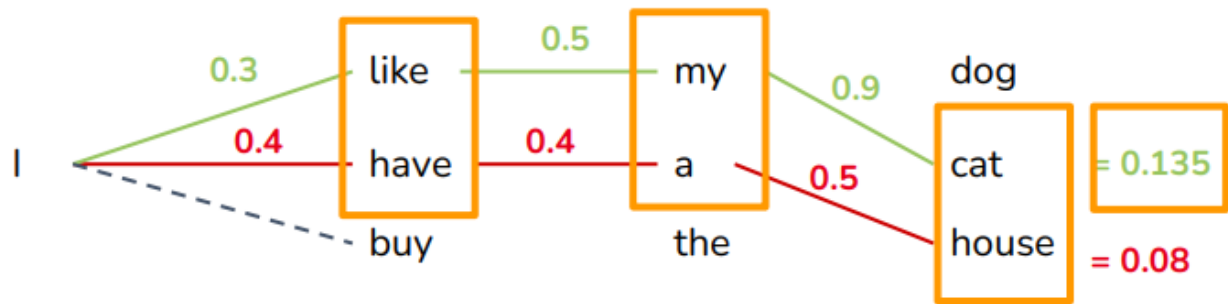
	Greedy	Beam Search	Top-k Sampling	Top-p (Nucleus) Sampling
At each step	Pick the best word	Try a few best words	Random sample from top-k	smallest set with cumulative probability > p
Output	One sequence	Several partial sequences	One sequence	One sequence

### a) Greedy Decoding



Weakness : Repetition as it always selects the most frequent token

### b) Beam Search



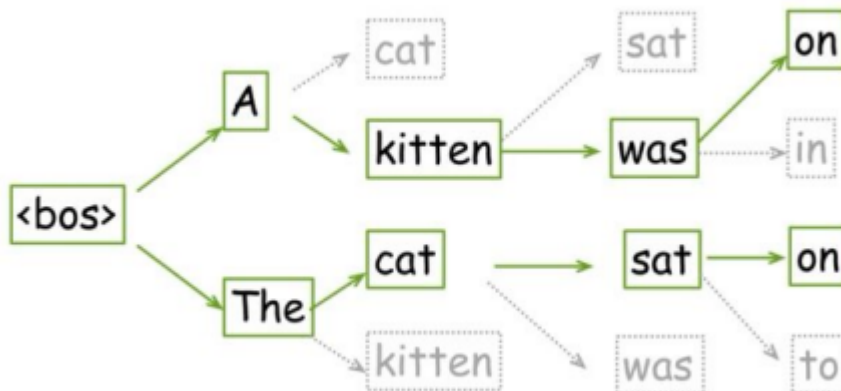
- Top-N is called beam size, usually 4-10.
- Increasing beam size is computationally inefficient and may lead to worse quality.

#### Strategy:

Start with the beginning of the sentence token or with an empty sequence

<bos>

Pick top **beam\_size** hypos ,terminate the test

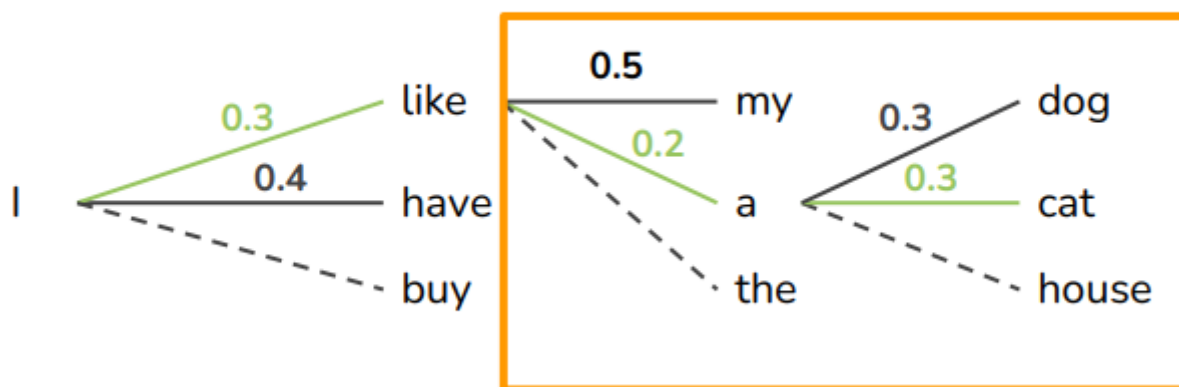


#### Weakness:

- Short sentences
- Less diversity

Beam search Text is not surprising

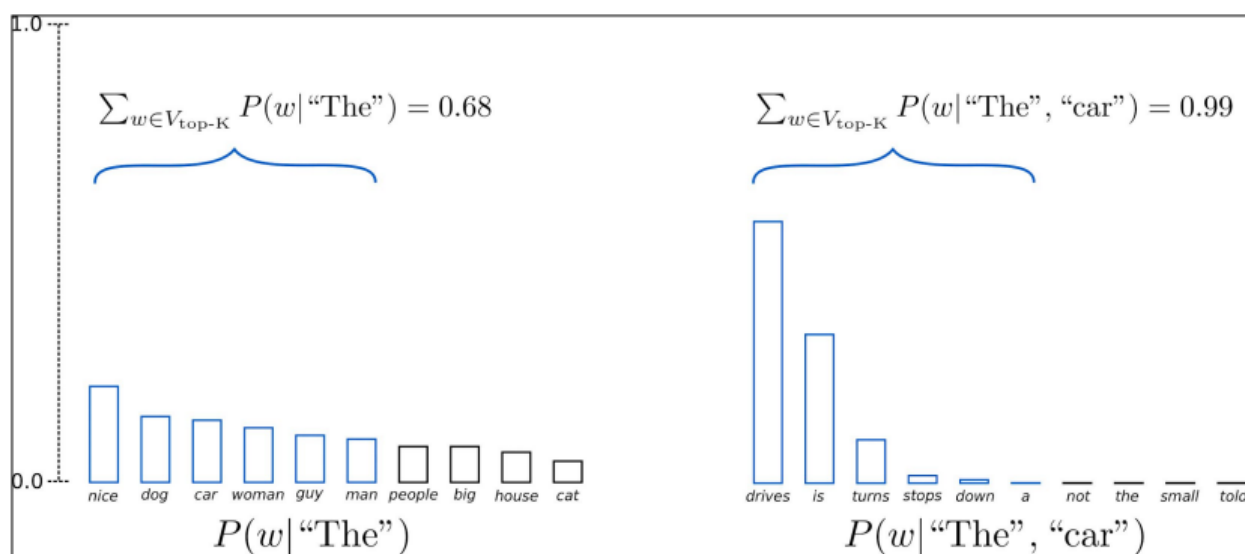
#### c)Top-k Sampling



- Sort the probabilities of the vocabulary at each step
- Select the top-k words, k is often 5 – 20
- $k=1 \Rightarrow$  greedy decoding
- Increase  $k \rightarrow$  have more diverse, also more risky
- Decrease  $k \rightarrow$  have more safe choices but less diverse

**Which k to choose?**

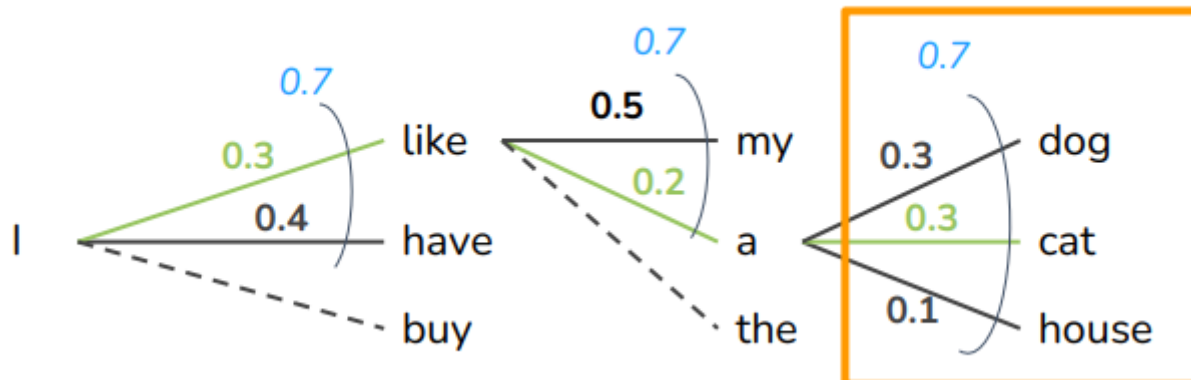
Which k to choose?



**Weakness** • Weird n-grams may occur due to random picking of top-k words  $\Rightarrow$  the output may not be coherent

## d)Top-p (Nucleus) Sampling



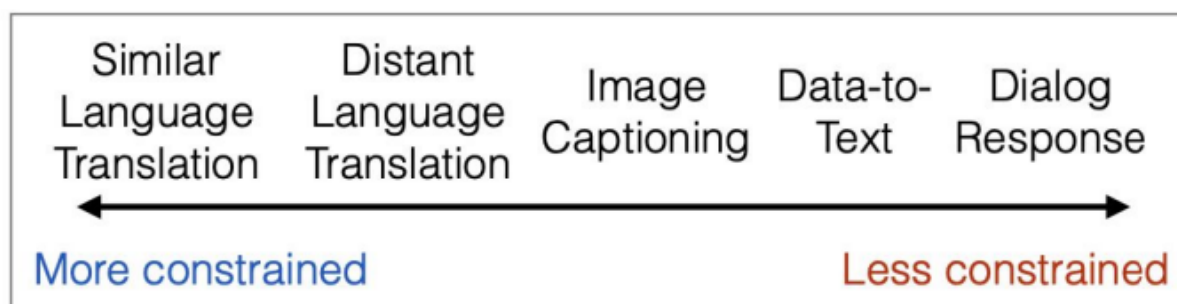


- Commonly-used probability  $p$  is 0.95

**Weakness** • (Surprisingly) may not include surprise words

### e) Decoding in practise:

- Can combine different strategies
  - e.g., temperature + beam search, temperature + top-k
- Use beam search with small beam size for tasks that exists a correct answer (more constrained)
- Use top-k or top-p for open-ended generation (less constrained)
- As models getting better/larger, sampling-based methods tend to work better



More freedom = more flexibility, but often more difficulty in modeling and evaluation

### f) Controlled Generation

- Add a further constraint in addition to content-based ones
- Politeness/Style Control: Take an input  $X$  and a label indicating style, etc.

source	Give me the telephone!
reference	Gib mir das Telefon! [T]
none	Gib mir das Telefon! [T]
polite	Geben Sie mir das Telefon! [V]
informal	Gib mir das Telefon! [T]

- Personalization: Take an input X and a side information about the speaker

**English Sentence:** Accordingly , I consider it essential that both the identification of cattle and the labelling of beef be introduced as quickly as possible on a compulsory basis .

**German Sentence:** Entsprechend halte ich es auch für notwendig , daß die Kennzeichnung möglichst schnell und verpflichtend eingeführt wird , und zwar für Rinder und für Rindfleisch .

**Meta Info:** EUROID="2209" NAME="Schierhuber" LANGUAGE="DE" GENDER="FEMALE" DATE\_OF\_BIRTH="31 May 1946" SESSION\_DATE="97-02-19" AGE="50"

---