

# Language Models

Task: generating text

Ex: Input: "Once upon a time, in a land far away,"

Output: "there lived a wise old dragon who loved to read books."

Ex: Input: "Why did the chicken cross the road?"

Output: "To get to the other side."

Ex: How Chat GPT works

Input: System prompt

"You are an AI assistant trained to provide accurate, concise, and helpful responses to user inquiries."

+

"Why did the chicken cross the road?"

Output: "To get to the other side."

We will only study auto-regressive models

Ex: ChatGPT, Gemini, Claude

Auto-regressive: generates output sequentially by using its previous outputs and inputs

Ex: Input: "Why did the chicken cross the road?"

Output: "To"

Input: "Why did the chicken cross the road?  
To"

Output: "get"

Input: "Why did the chicken cross the road?  
To get"

Output: "to"

⋮

Input: "Why did the chicken cross the road?  
To get to the other side."

Output: "<EOS>" End Of sequence token

We want a model  $f: \mathcal{X} \rightarrow \mathcal{Y}$

$x \in \mathcal{X}$  is a sequence of ~~words~~ tokens

$f(x) \in \mathcal{Y}$  is the next ~~word~~ token

$\mathcal{Y} = \{ \text{all words} + \text{punctuation} + \langle \text{EOS} \rangle + \langle \text{PAD} \rangle \}$

$= \{1, \dots, K\}$

Vocabulary

$|\mathcal{Y}| = K$

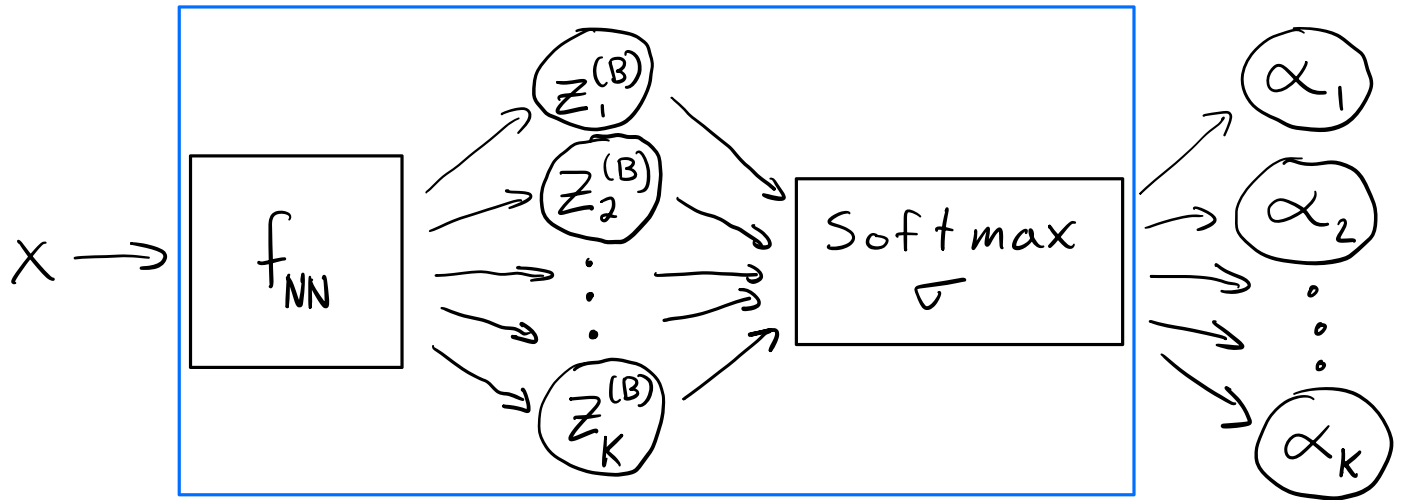
Discrete labels make optimization hard

We want a model  $f_{\text{prob}}: \mathcal{X} \rightarrow \mathcal{Y}_{\text{prob}}^{\leftarrow} = [0, 1]^K$

$x \in \mathcal{X}$  is a sequence of tokens

$f_{\text{prob}}(x) \in \mathcal{Y}_{\text{prob}}$  is a vector of  $K$  probabilities

with  $d^{(B)} = K$   
 We can use a NN  $f_{NN}(\vec{x})$  with  $h^{(B)}(z) = z$  and then  
 apply softmax to get a probability



$$f_{\text{prob}}(\vec{x}) = \sigma(f_{NN}(x))$$

$$\sum_{q=1}^K \alpha_q = 1, \alpha_q \in [0, 1]$$

How do we represent a sequence of tokens  $s$  as a vector  $\vec{x} \in \mathbb{R}^{d+1}$ ?

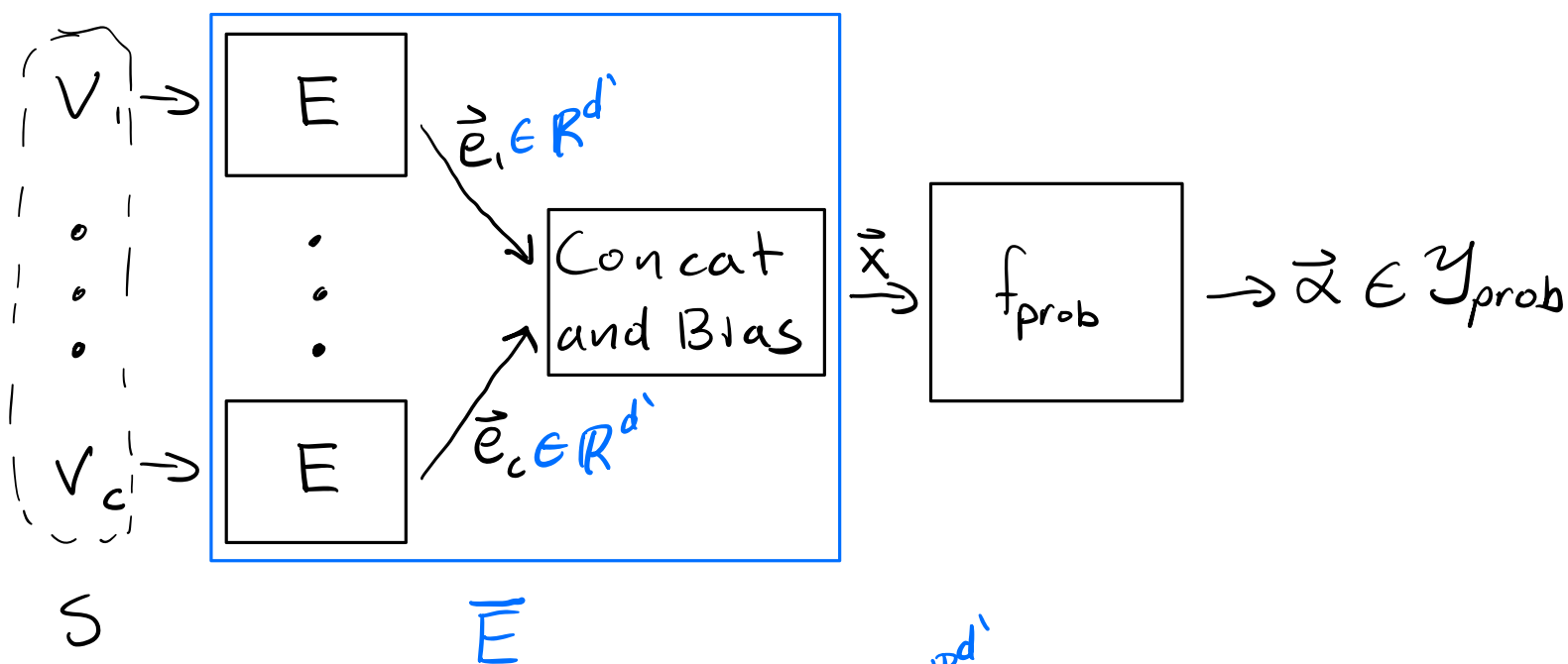
$$s = (v_1, \dots, v_c) \in \mathcal{Y}^c$$

Ex: "why did the chicken"

$$s = (\text{"why"}, \text{"did"}, \text{"the"}, \text{"chicken"}) \in \mathcal{Y}^4$$

$$E: \mathcal{Y} \rightarrow \mathbb{R}^{d'} \quad \text{embedding function}$$

assume it is given



$$\bar{E}(s) = \bar{E}(v_1, \dots, v_c) = (1, E(v_1), E(v_2), \dots, E(v_c))^T = \vec{x} \in \mathbb{R}^{d+1}$$

$$d = cd'$$

context length

If a sequence is longer than  $c$  tokens  
take the last  $c$  tokens

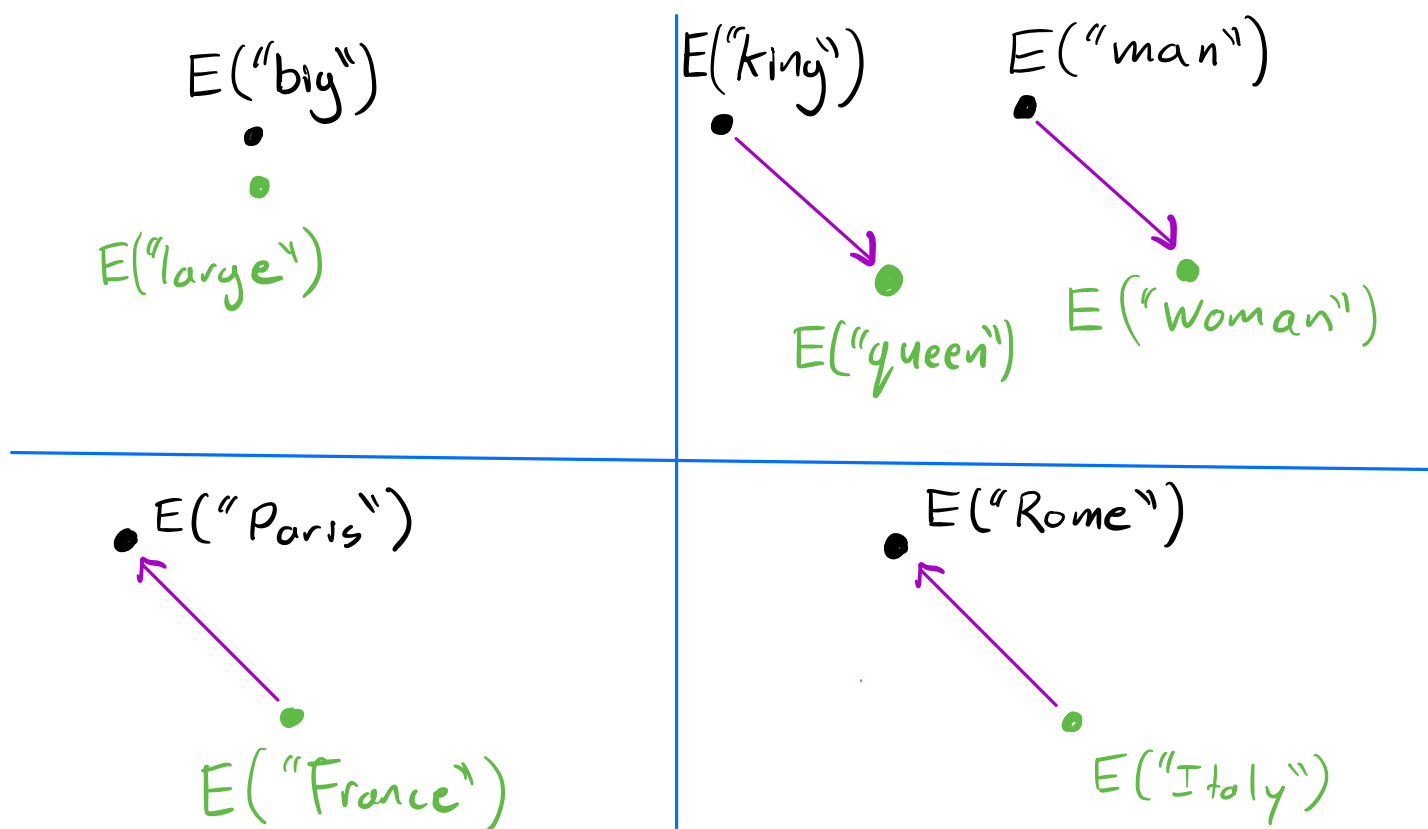
To handle sequences shorter than  $c$  tokens  
a padding token "<PAD>" is added

Ex:  $c=3$

$S = ("Why", "did") \Rightarrow ("<PAD>", "Why", "did")$

# Embeddings

$E_x$ :  $d' = 2$        $E: \mathcal{Y} \rightarrow \mathbb{R}^{d'}$



# Creating a Dataset

$$S = (V_1, V_2, V_3, \dots, V_c, V_{c+1}, \dots, V_{n+1})$$

$$S_1 = (\overbrace{\langle \text{PAD} \rangle, \dots, \langle \text{PAD} \rangle}^{c-1}, V_1), \quad y_1 = V_2$$

$$S_2 = (\overbrace{\langle \text{PAD} \rangle, \dots, \langle \text{PAD} \rangle}^{c-2}, V_1, V_2), \quad y_2 = V_3$$

⋮

⋮

$$S_c = (V_1, \dots, V_c),$$

$$y_c = V_{c+1}$$

⋮

$$S_n = (V_{n-c+1}, \dots, V_n),$$

$$y_n = V_{n+1}$$

Ex:  $C=2$ ,  $S = (\text{"Why"}, \text{"did"}, \text{"the"})$

$$S_1 = (\langle \text{PAD} \rangle, \text{"Why"}), \quad y_1 = \text{"did"}$$

$$S_2 = (\text{"Why"}, \text{"did"}), \quad y_2 = \text{"the"}$$



$$\vec{x}_1 = \vec{E}(s_1), \vec{x}_2 = \vec{E}(s_2), \dots, \vec{x}_n = \vec{E}(s_n) \in \mathbb{R}^{d+1}$$

$$\vec{y}_1 = \text{onehot}(y_1) \in \{0,1\}^K \subseteq [0,1]^K$$

$$\vdots$$

$$\vec{y}_n = \text{onehot}(y_n) \in \{0,1\}^K$$

$$D = ((\vec{x}_1, \vec{y}_1), \dots, (\vec{x}_n, \vec{y}_n))$$

## ERM Learner:

$$\mathcal{A}(\mathcal{D}) = \hat{f}_{\text{prob}} = \arg \min_{f \in \mathcal{F}} \hat{L}(f)$$

$$\mathcal{F} = \{ f \mid f: \mathcal{X} \rightarrow \mathcal{Y}_{\text{prob}} \text{ where } f = \sigma(f_{\text{NN}}) \text{ and } f_{\text{NN}} \text{ is a NN with a fixed architecture} \}$$

$\ell$  is multiclass cross-entropy loss

For a new sequence  $s = (v_1, \dots, v_c) \in \mathcal{Y}^c$

The Language model is:

$$\hat{f}_{\text{LM}}(s) = \arg \max_{q \in \{1, \dots, k\}} \hat{y}_q$$

$$\hat{\tilde{y}} = \hat{f}_{\text{prob}}(\tilde{E}(s))$$

# Notes

- Transformers are the NN architecture currently used
- Embedding  $E$  can be learned
- Vocabulary can use characters or sub-words instead of words "tokenization"
- Most probable word is not always chosen, instead can sample based on probabilities
- "Large" language models (LLM)  
means a NN with a lot of weights

Ex: GPT-3 has 175 Billion weights

GPT-5 Publically estimated to have  
 $\approx 2$  Trillion weights

The brain has  $\approx 100$  trillion  
connections between neurons