

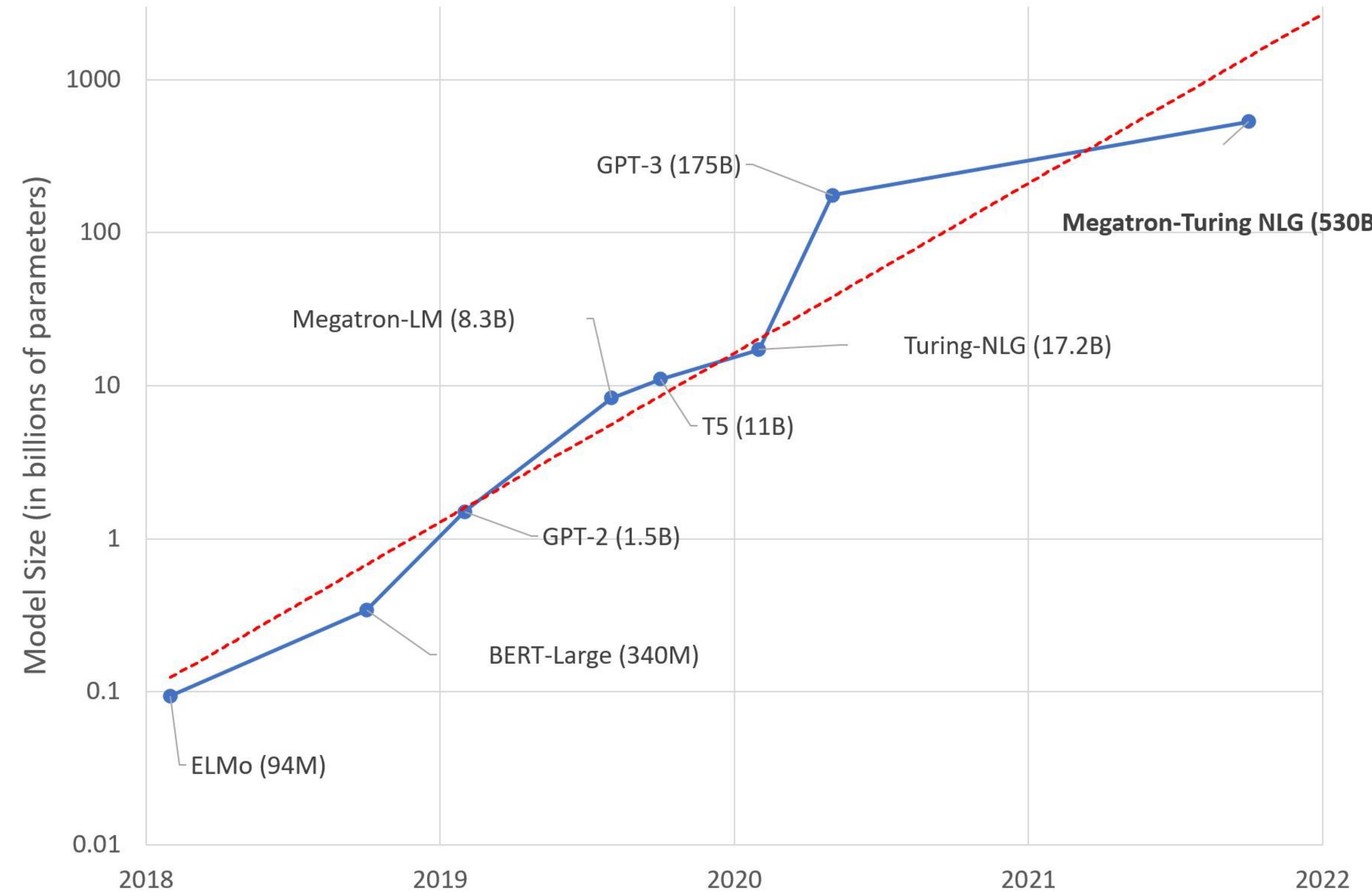
# Large Pretrained Models

COMP3361 – Week 9

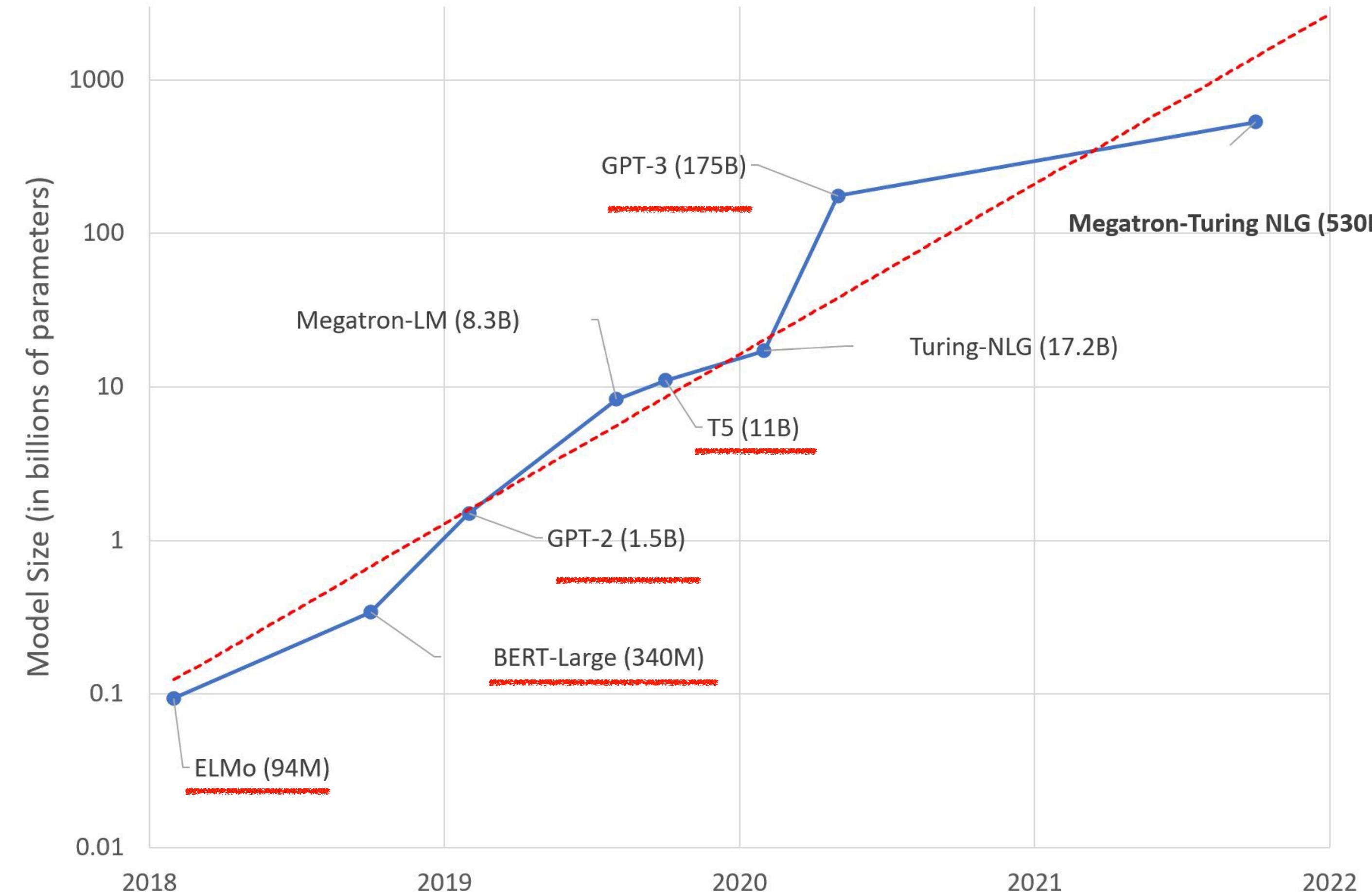
Lingpeng Kong

Department of Computer Science, The University of Hong Kong  
Many materials from Stanford CS224n with special thanks!

# Pretrained Models in the Past Four Years



# Pretrained Models in the Past Four Years



# Pretrained Models are Expensive



One single training run



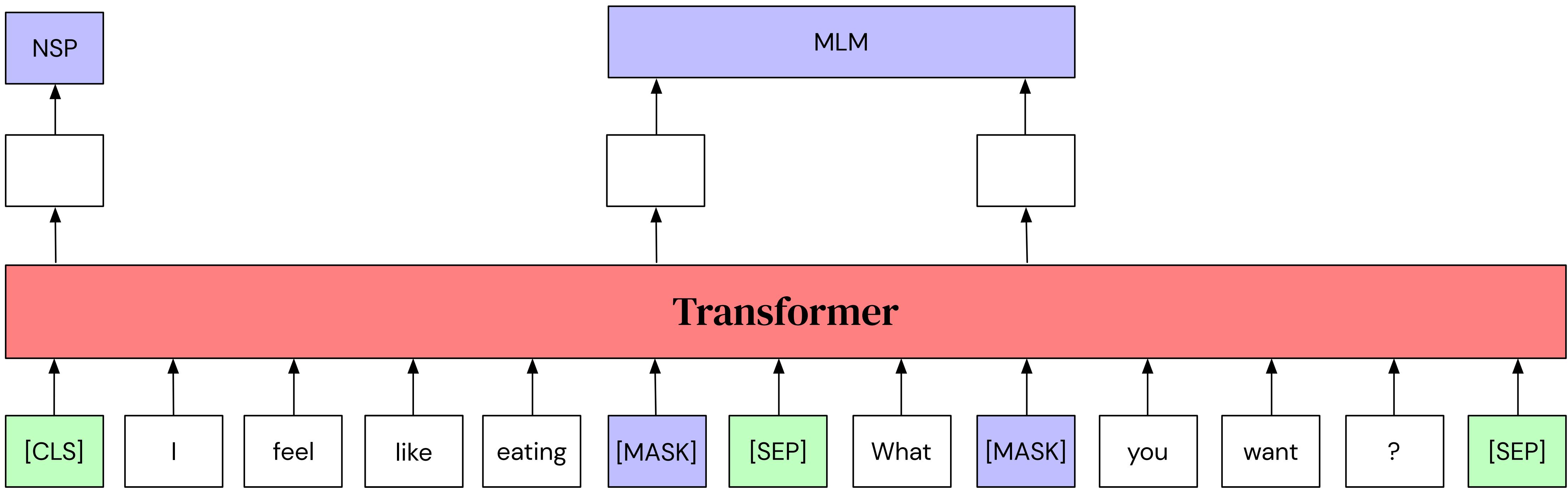
\$12 million



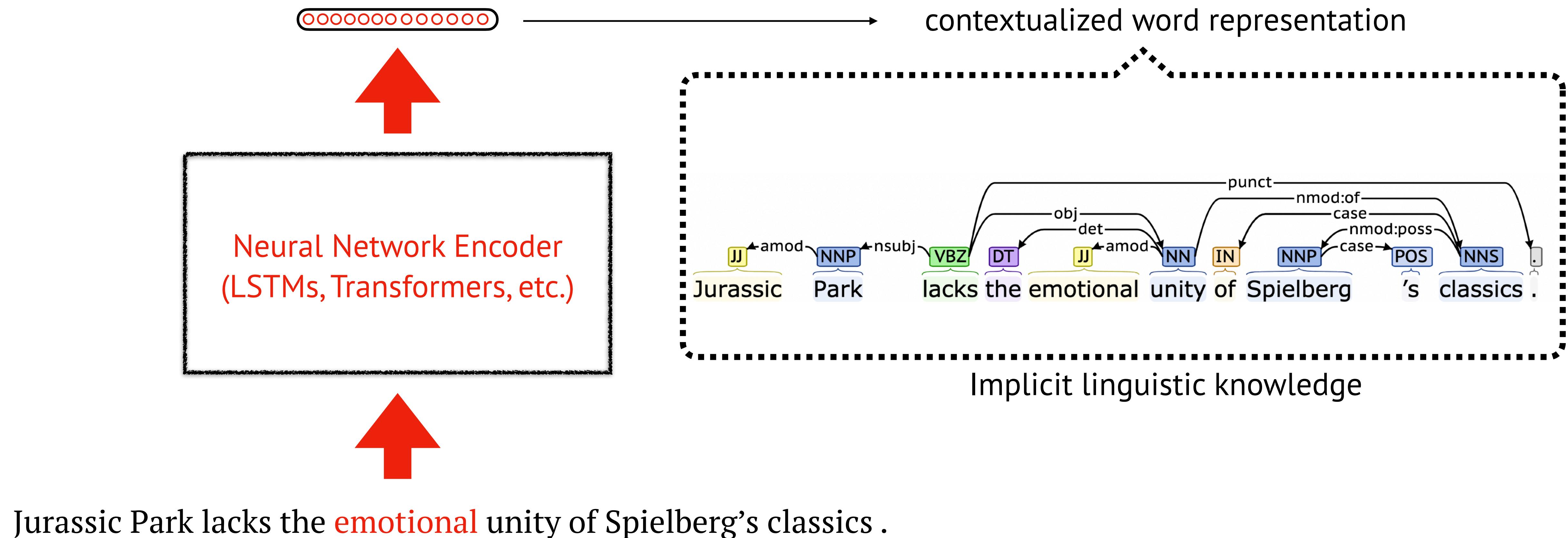
552 metric tons of carbon dioxide  
(120 cars per year)

# Pretraining and Contextualized Word Representations

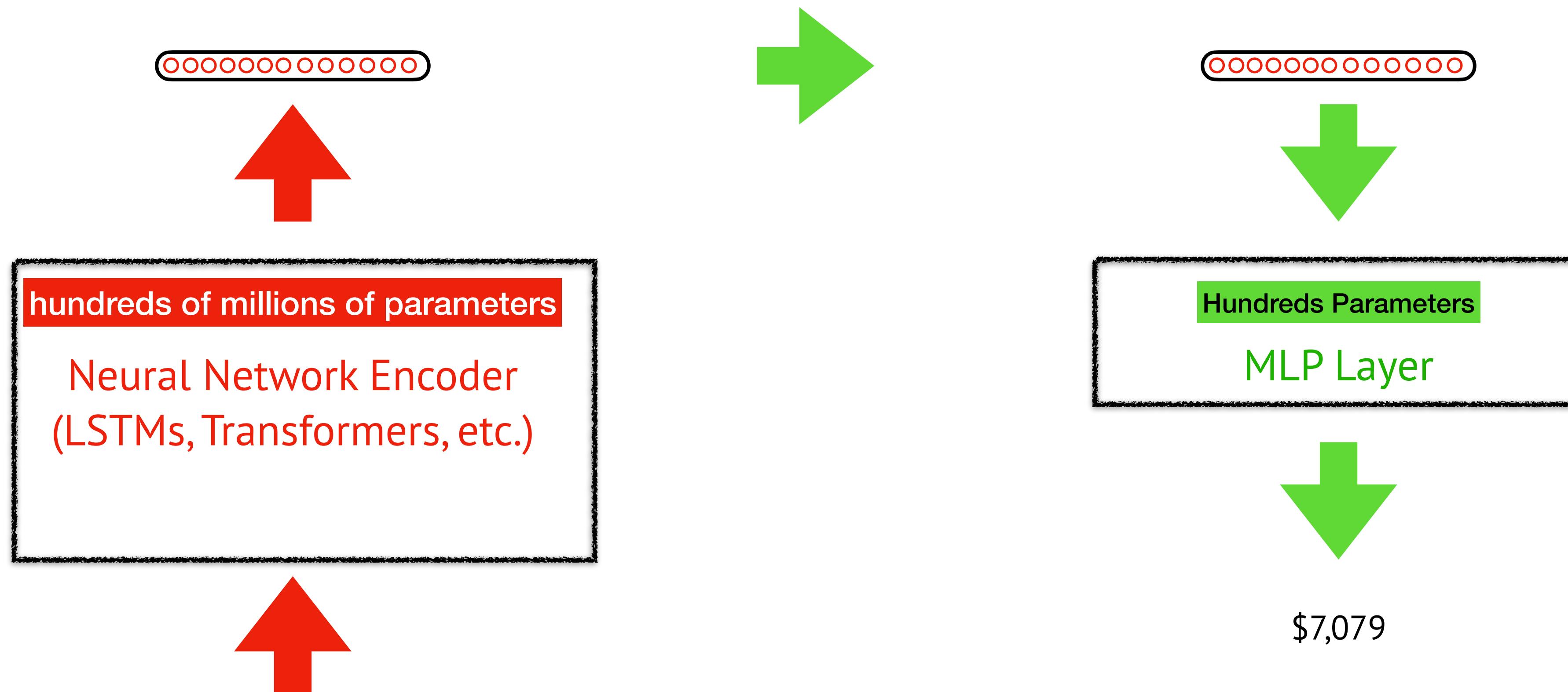
$$\mathbb{E}_{p(x_i, \hat{x}_i)}[p(x_i \mid \hat{x}_i)]$$



# Pretraining and Contextualized Word Representations



# Pretraining and Fine-tuning



Jurassic Park lacks the **emotional** unity of Spielberg's classics .

# Key Elements in BERT

Transformer

— neural representation learner

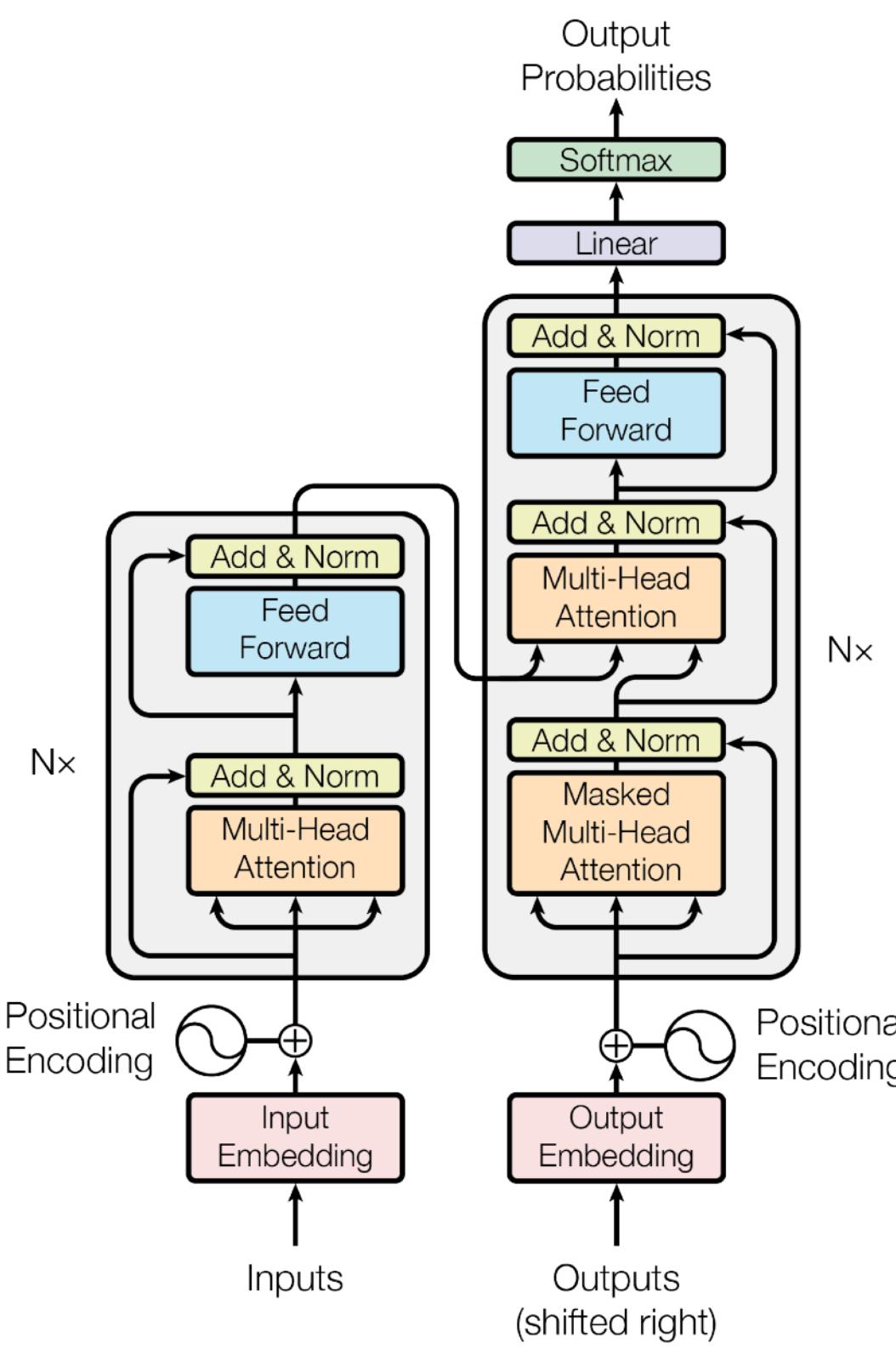
Masked Language Modeling (MLM),  
Next Sentence Prediction (NSP)

— pretraining objective

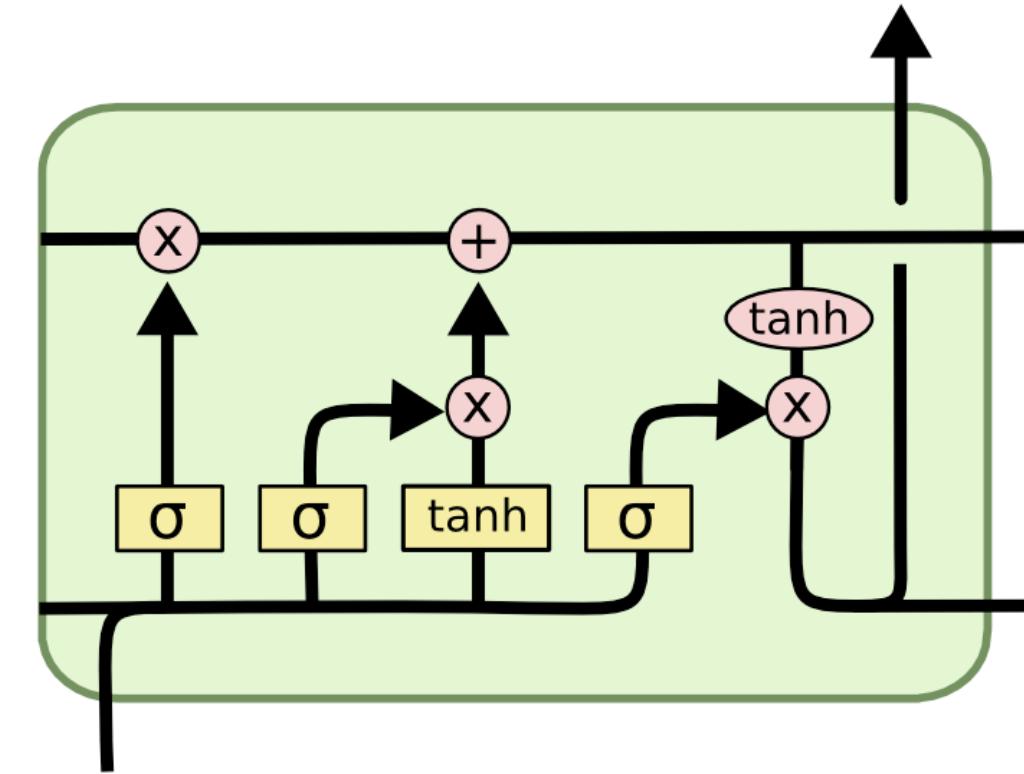
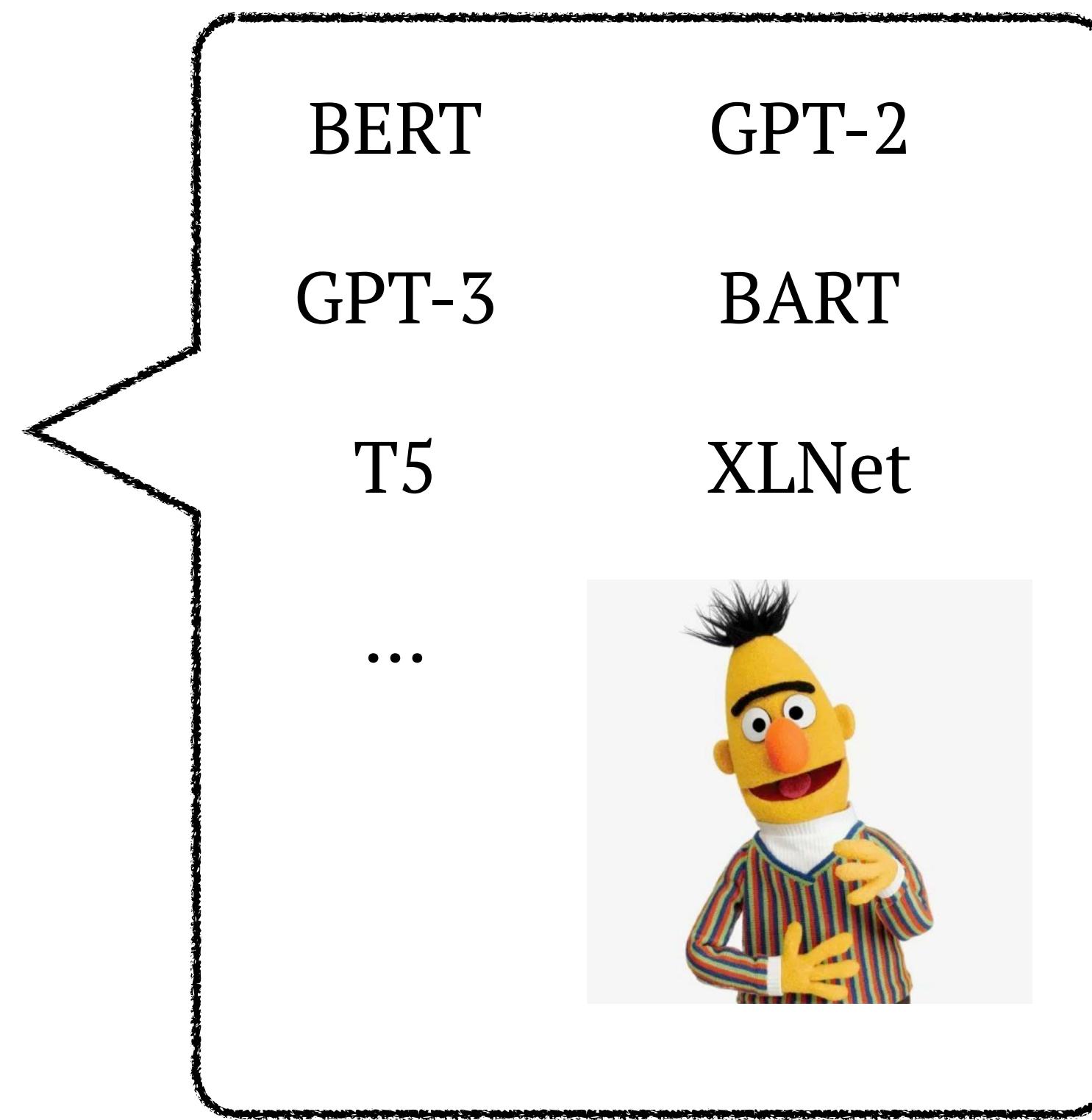
Bidirectional Encoder

— type of architecture

# Neural Representation Learners



Transformer

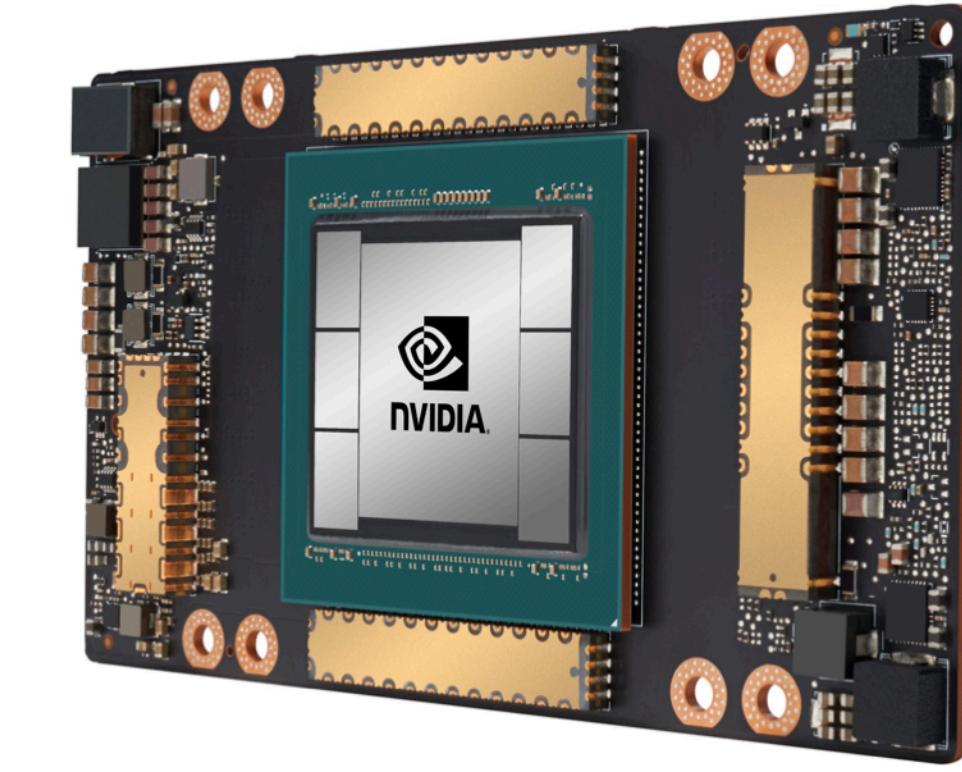
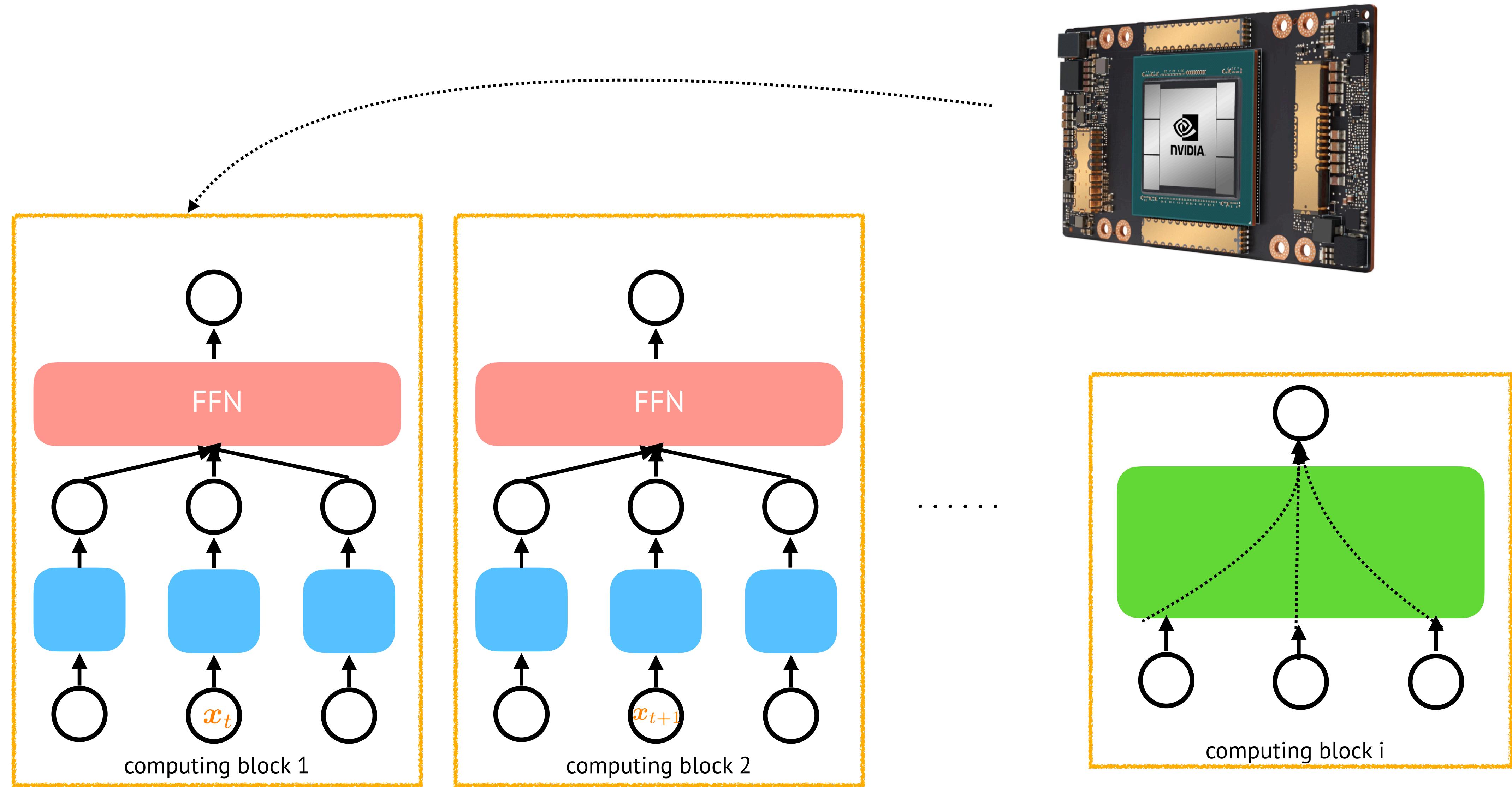


LSTM

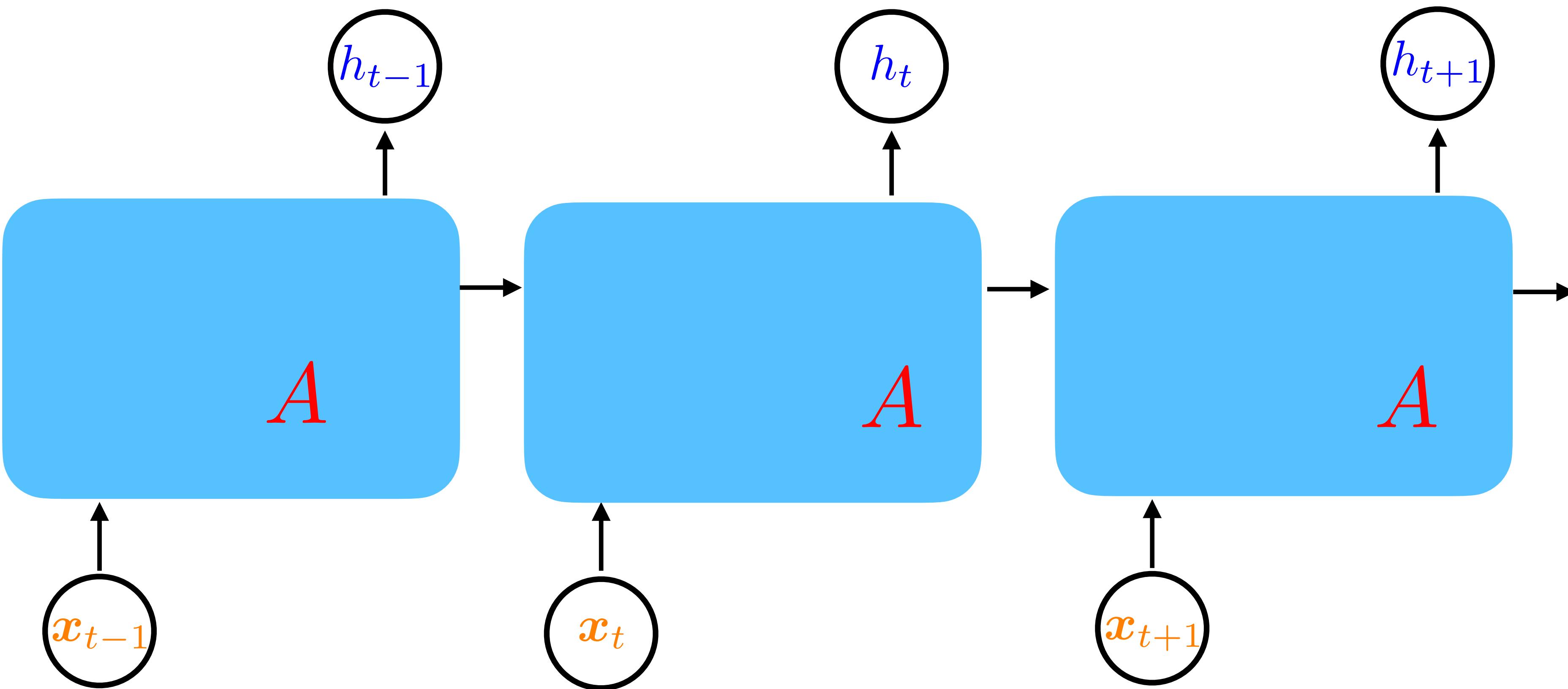


ELMo

# Why Transformers?

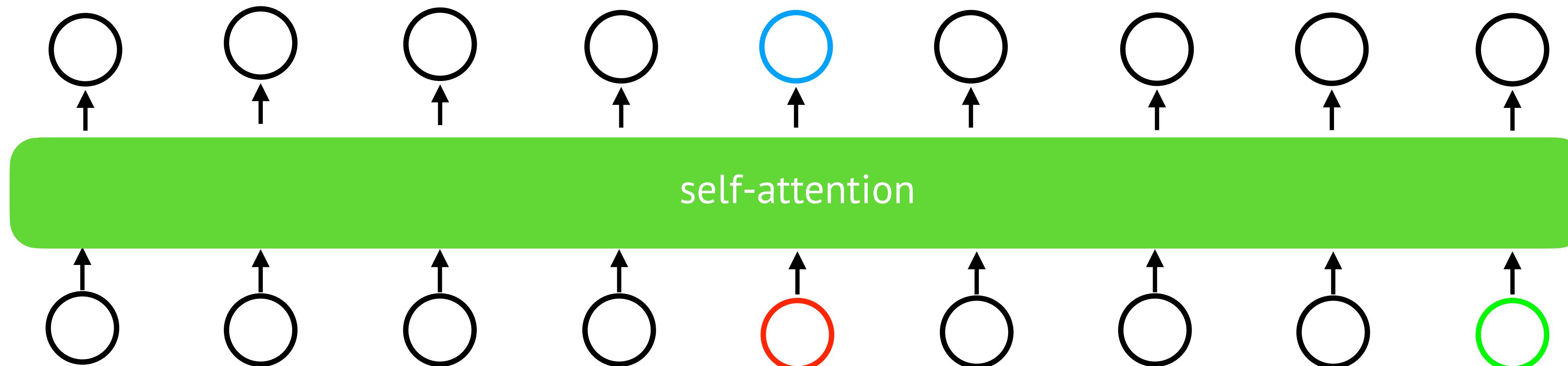


# Why Transformers?



# Why Transformers?

Direct pair-wise interaction between any tokens in the sequence



# Pretraining Objective

training instance (MLM):

x: I feel like eating <MASK> today. What <MASK> you want to eat?

y: noodles, do

training instance (NSP):

x: I feel like eating <MASK> today. ||| What <MASK> you want to eat?

y: True

# Pretraining Objective

What makes a good pretraining objective?

1. No human labeling should be involved.
2. Leads to good representations. (How and why?)

# Mutual Information

$$\begin{aligned} I(A, B) &= H(A) - H(A \mid B) \\ &= H(B) - H(B \mid A). \end{aligned}$$

$$f_{\theta}(a, b) = g_{\psi}(b)^{\top} g_{\omega}(a)$$

$$\theta = \{\omega, \psi\}$$

Goal of Training:

$$I(f(A, B)) \geq \mathbb{E}_{p(a,b)} \left[ f_{\theta}(a, b) - \mathbb{E}_{q(\tilde{\mathcal{B}})} \left[ \log \sum_{\tilde{b} \in \tilde{\mathcal{B}}} \exp f_{\theta}(a, \tilde{b}) \right] \right] + \log |\tilde{\mathcal{B}}|,$$

InfoNCE (Logeswaran & Lee, 2018; van den Oord et al., 2019)

$$\mathbb{E}_{p(a,b)} \left[ f_{\theta}(a, b) - \log \sum_{\tilde{b} \in \mathcal{B}} \exp f_{\theta}(a, \tilde{b}) \right].$$

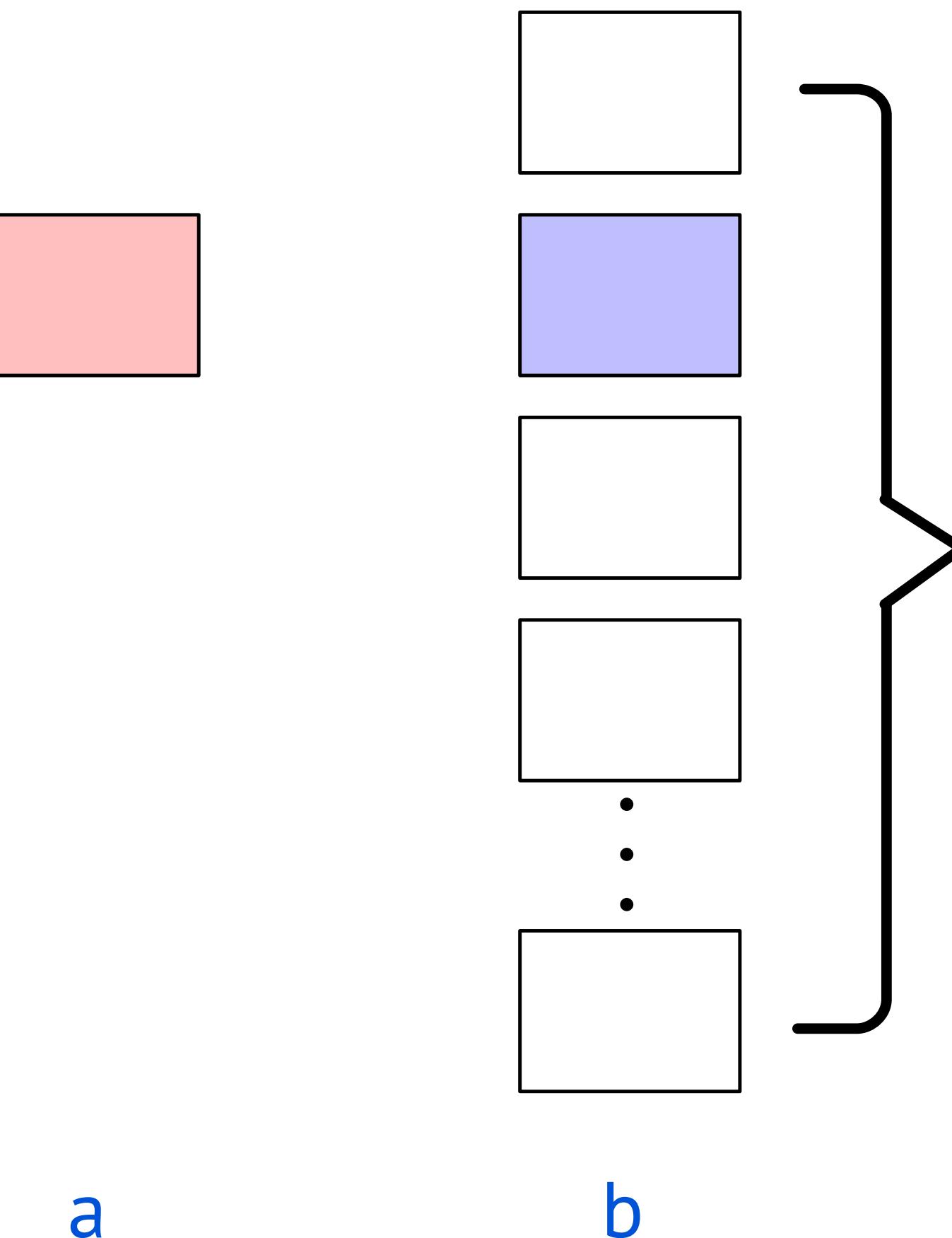
Cross Entropy (Softmax)

# Mutual Information

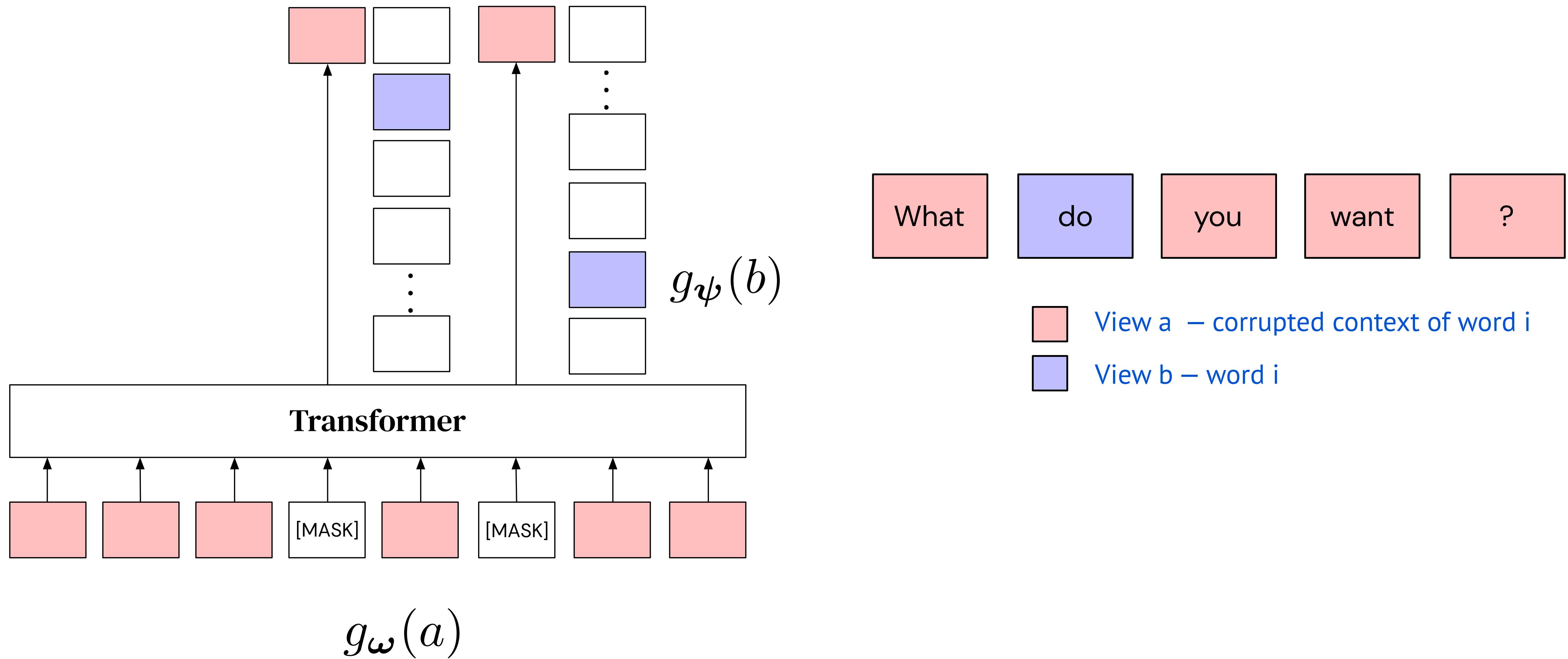
Cross Entropy (Softmax)

$$\mathbb{E}_{p(a,b)} \left[ f_{\theta}(a,b) - \log \sum_{\tilde{b} \in \mathcal{B}} \exp f_{\theta}(a, \tilde{b}) \right].$$

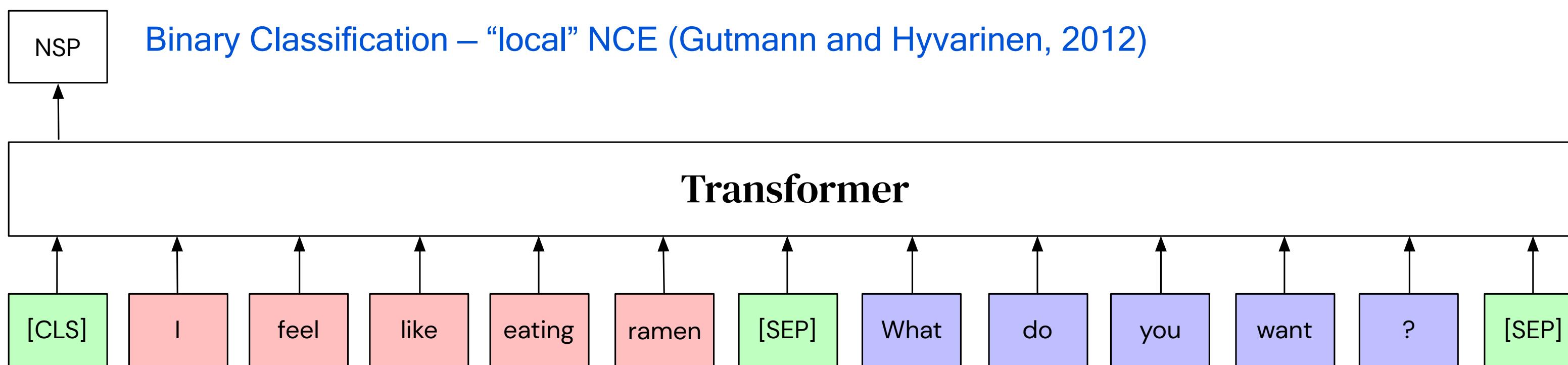
“Hope”                    “Fear”



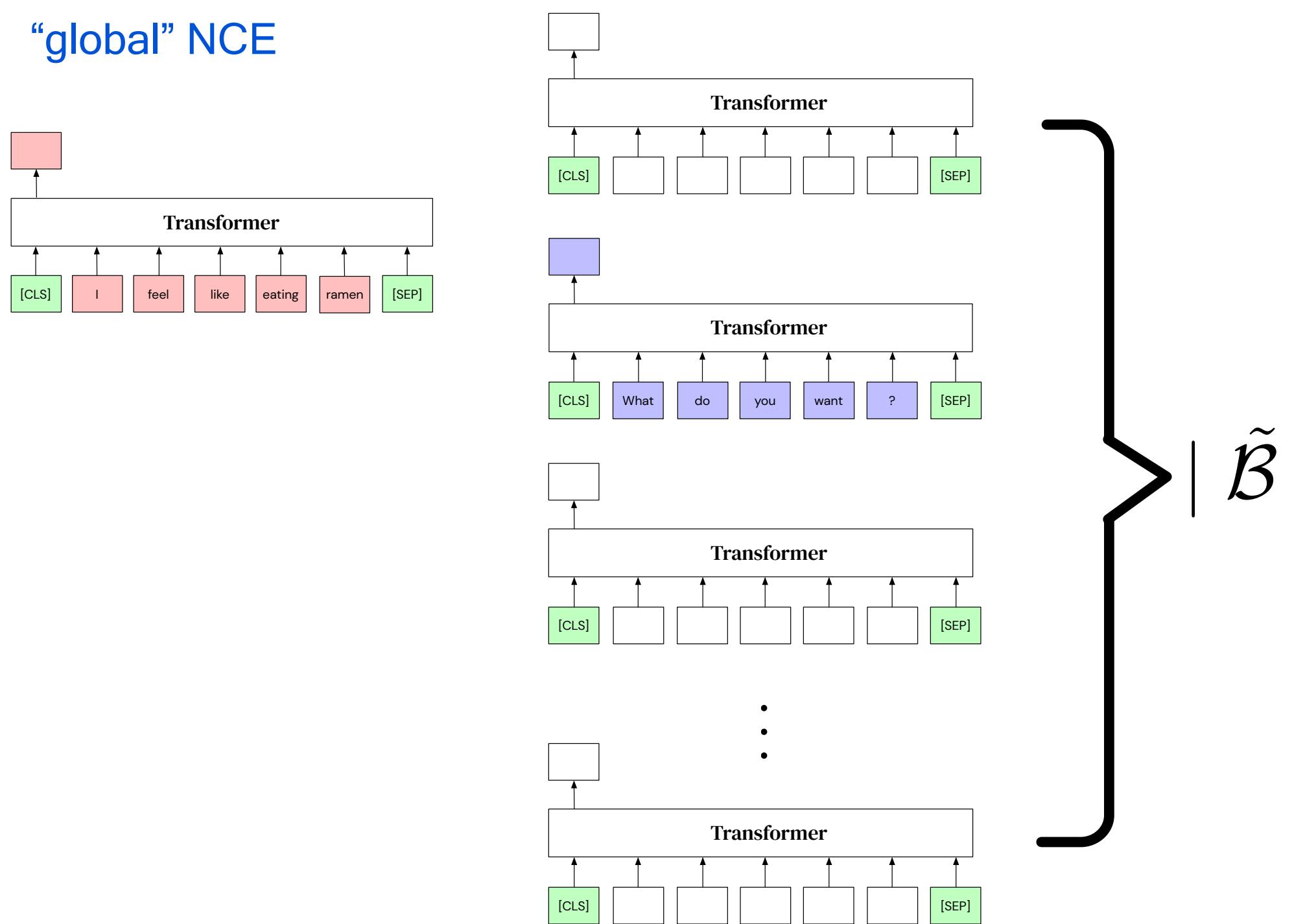
# Masked Language Modeling



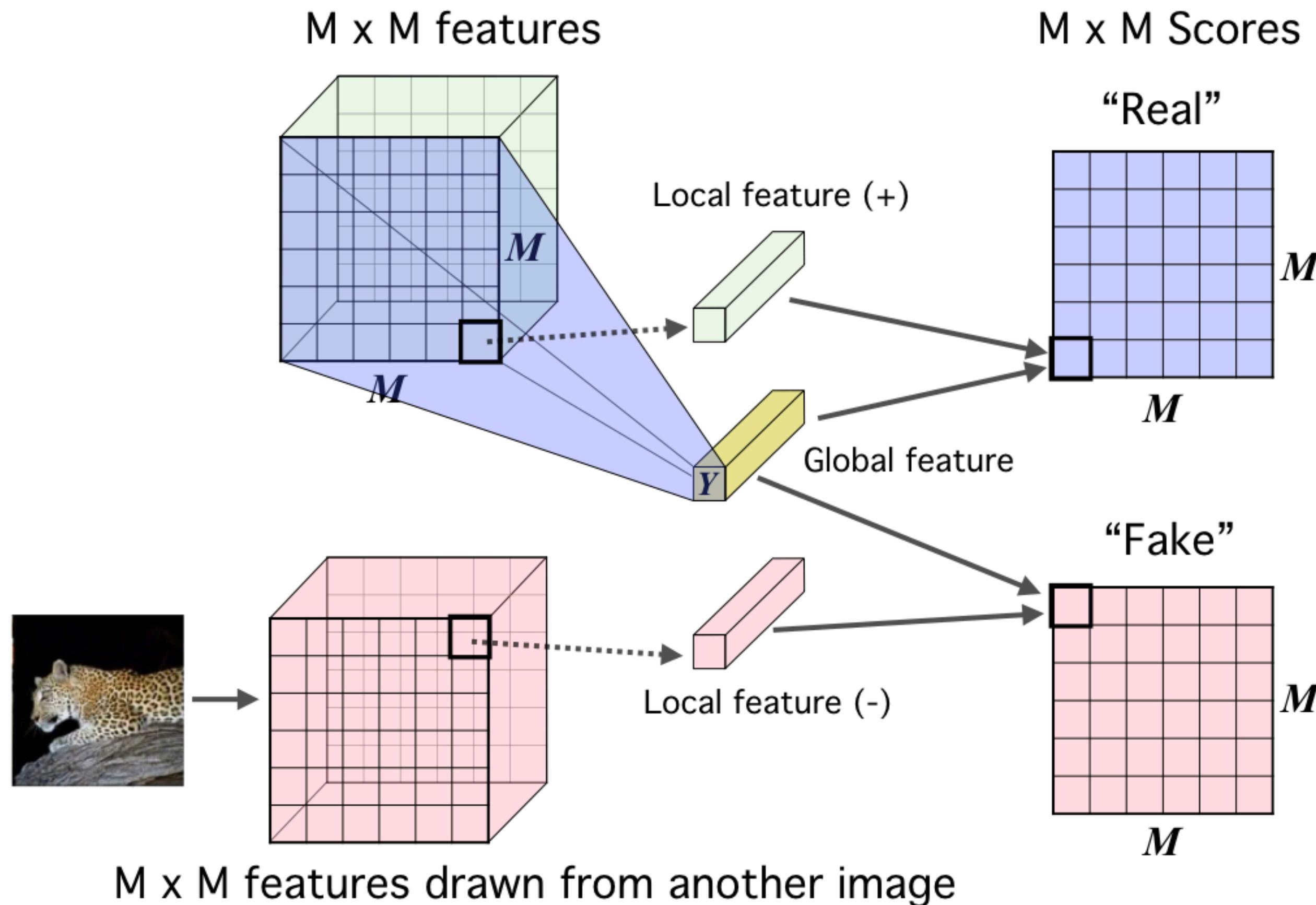
# Next Sentence Prediction



“global” NCE

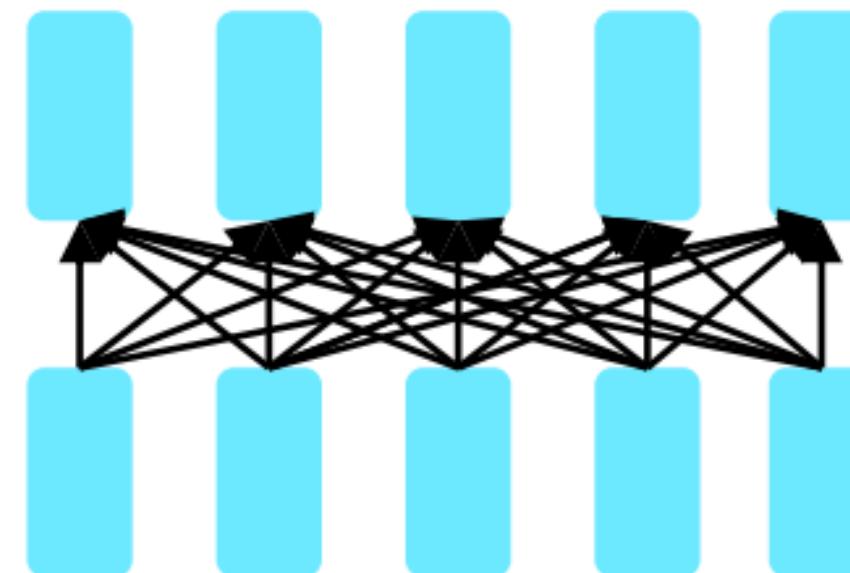


# Connections with Computer Vision



Deep InfoMax (DIM; Hjelm et al., 2019)

# Type of Architecture



Encoders

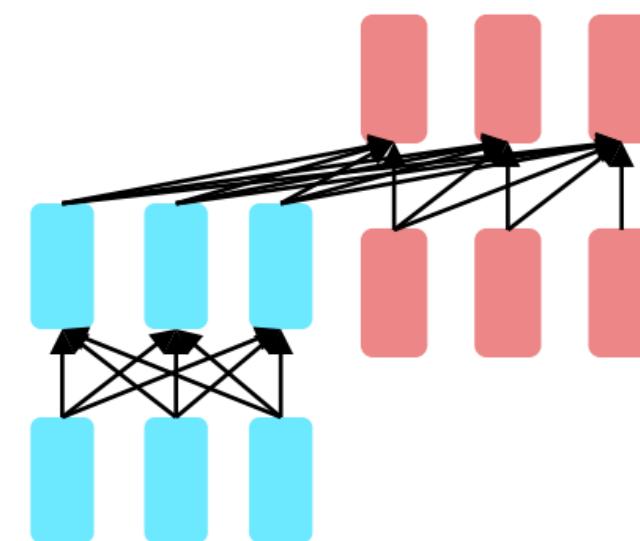
Parameters are what we get from the pretraining process.

Pros for the “encoders” architecture:

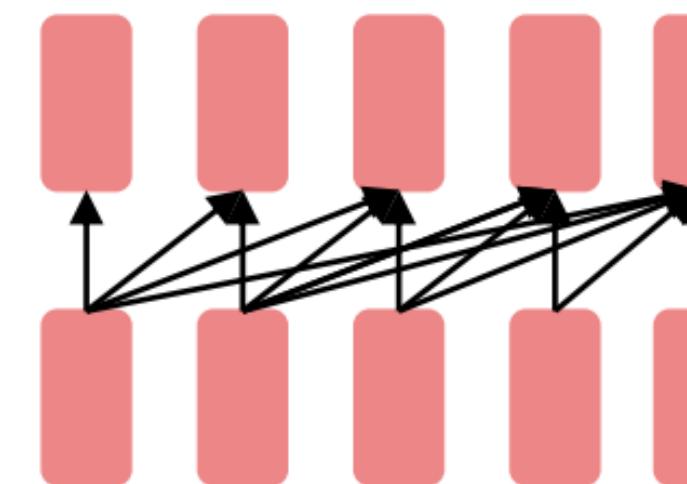
Gets bidirectional context.

Easy to use in language understanding tasks!

Other members in the family:



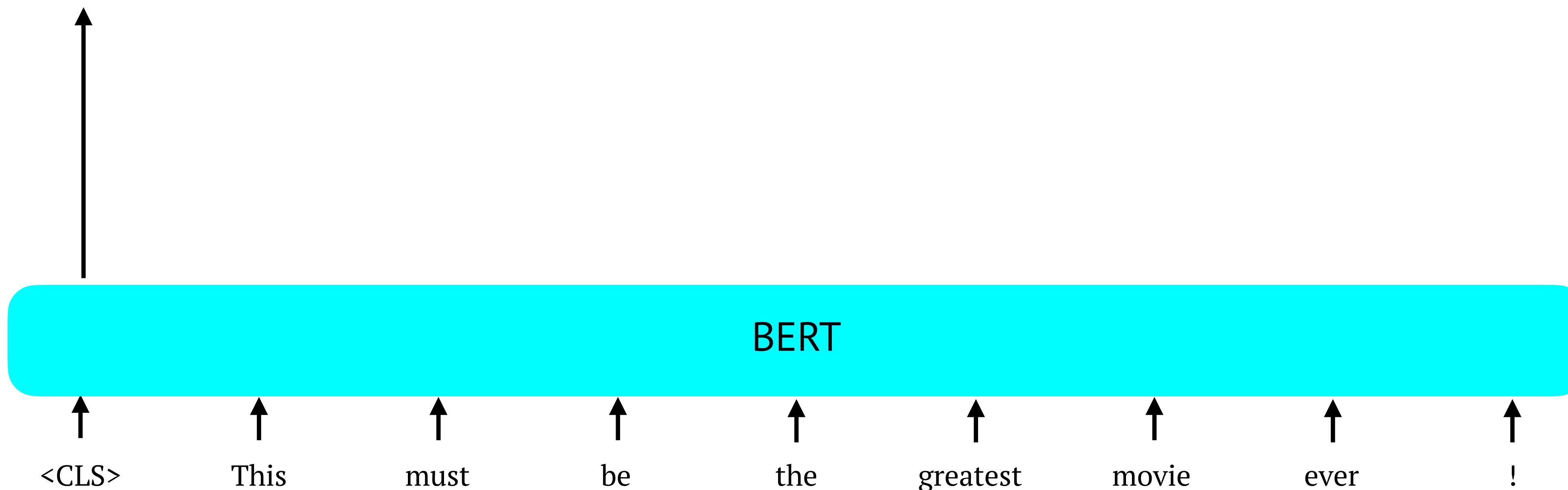
Encoder-Decoders



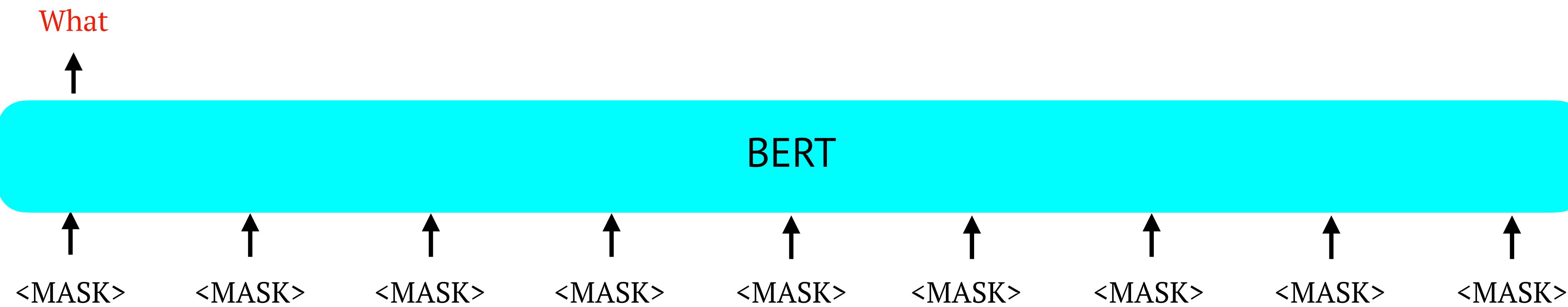
Decoders

# BERT for Understanding

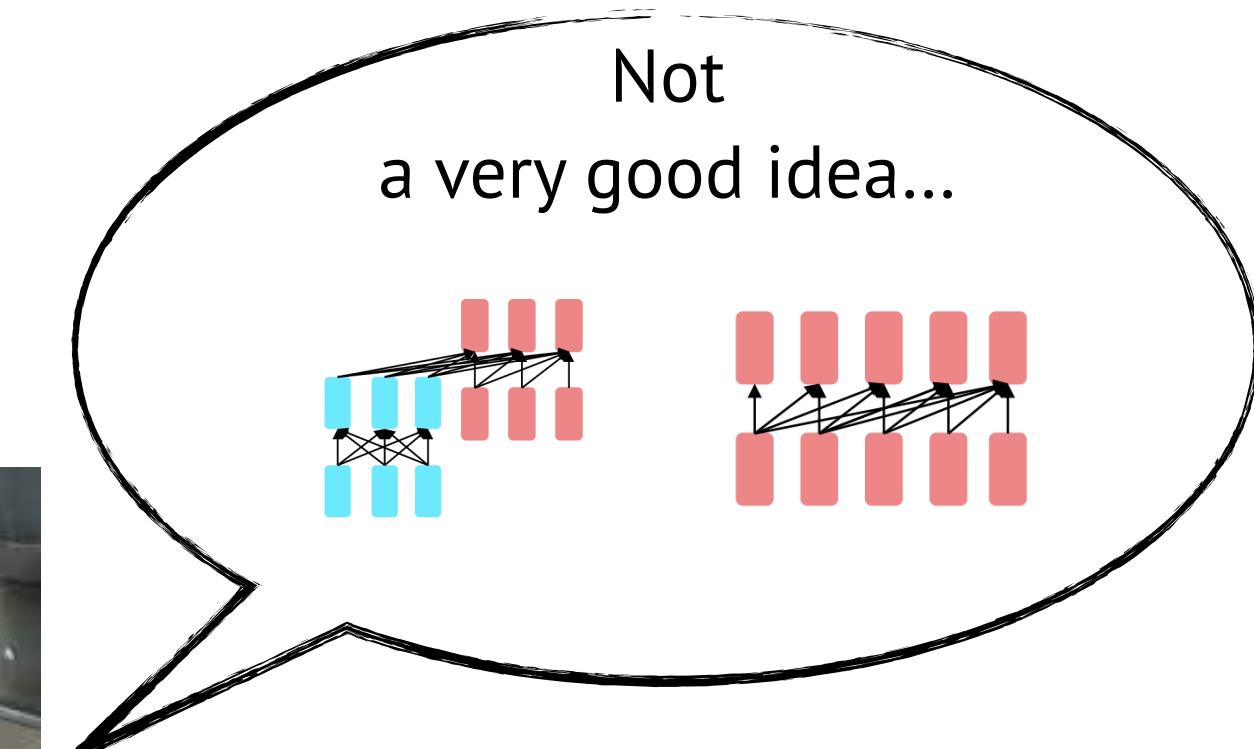
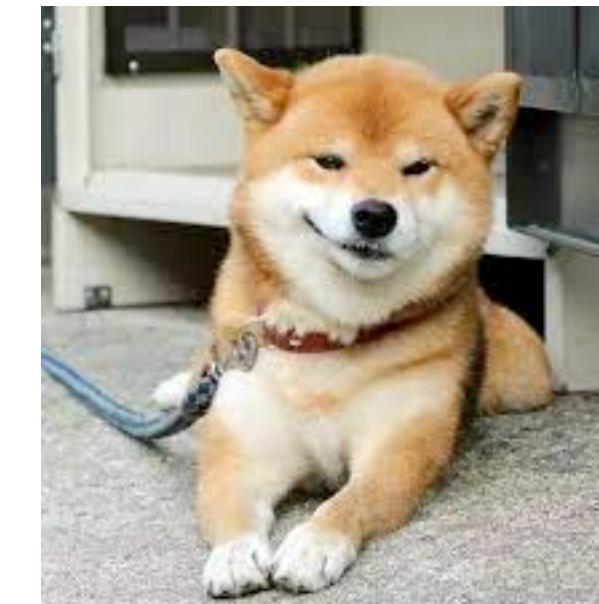
Positive / Negative



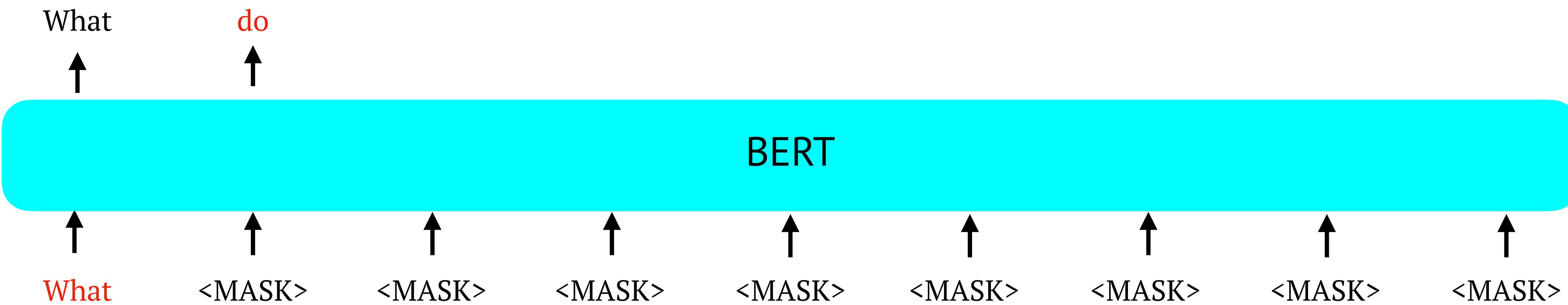
# BERT for Generation



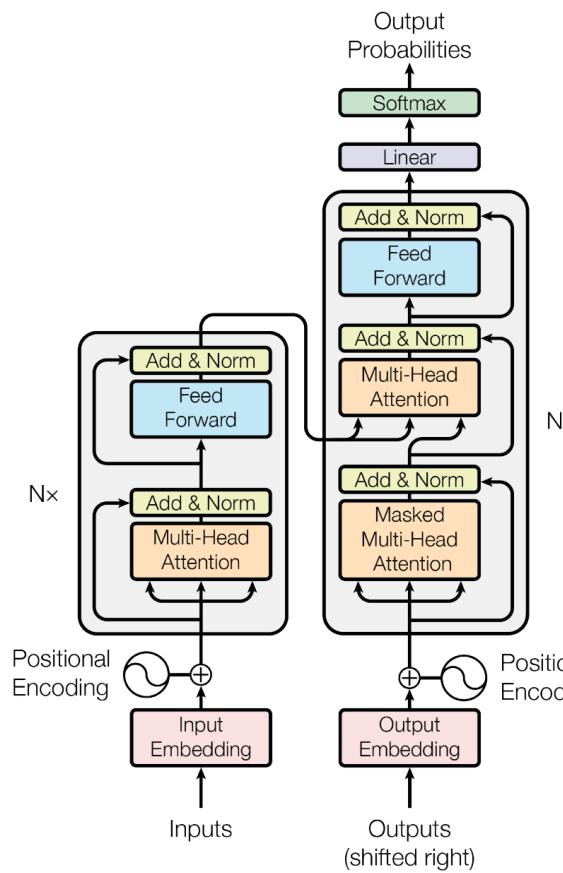
# BERT for Generation



Input has been changed. The representations will need to be recomputed!



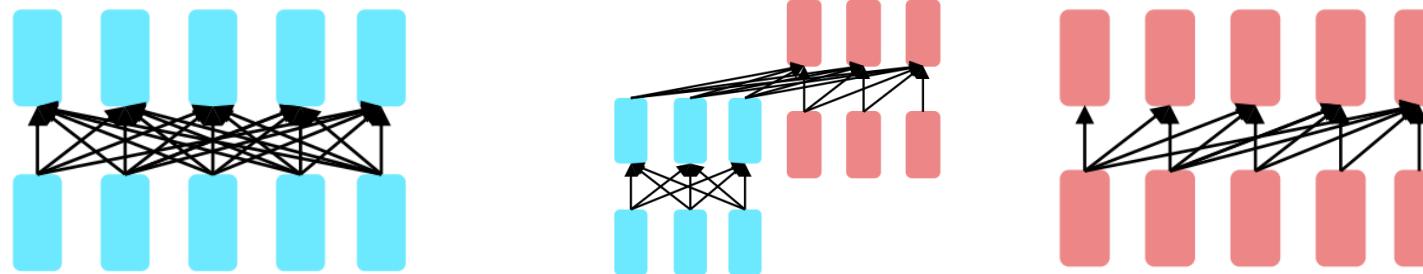
# Pretrained Models



— neural representation learner

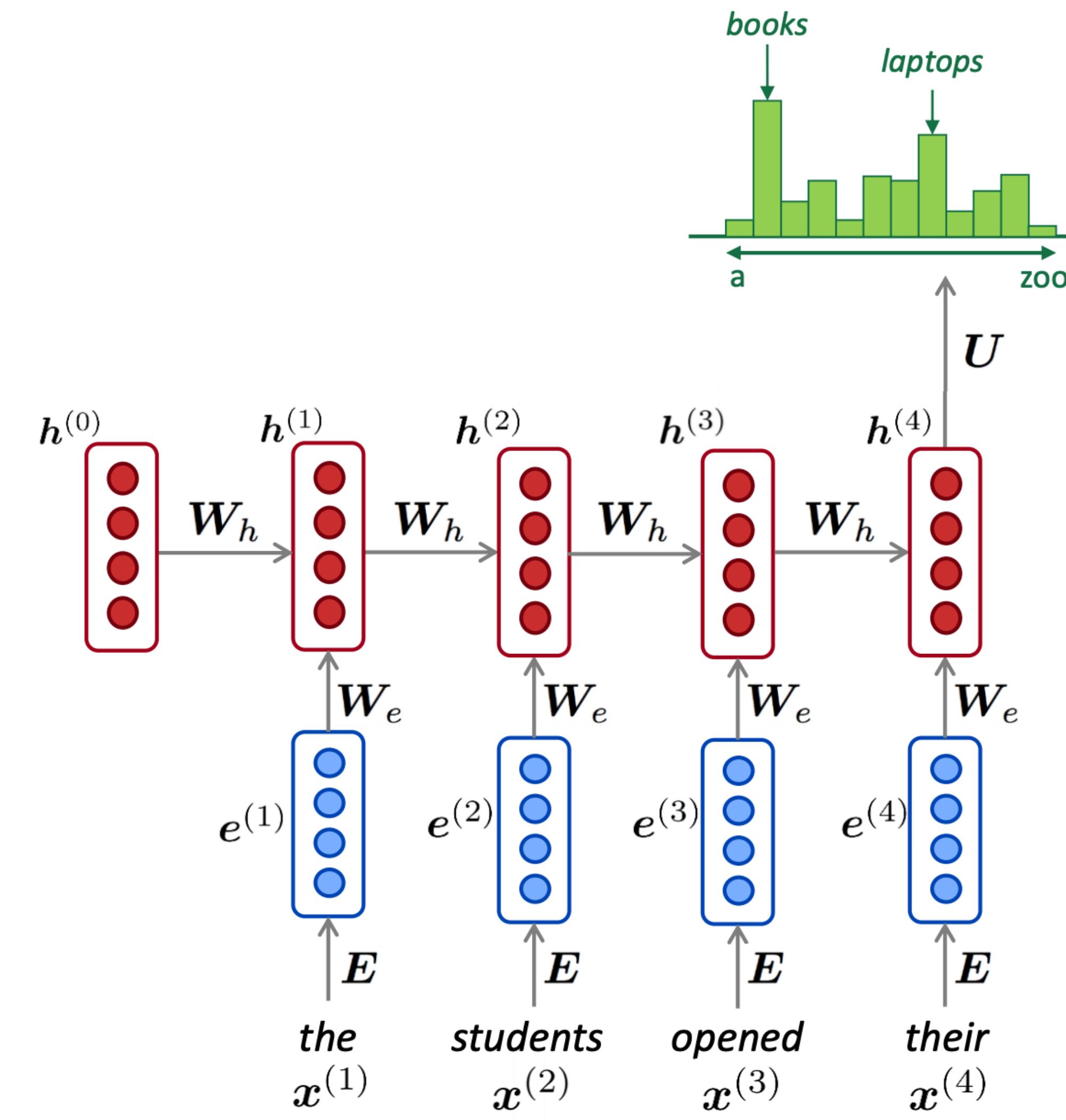
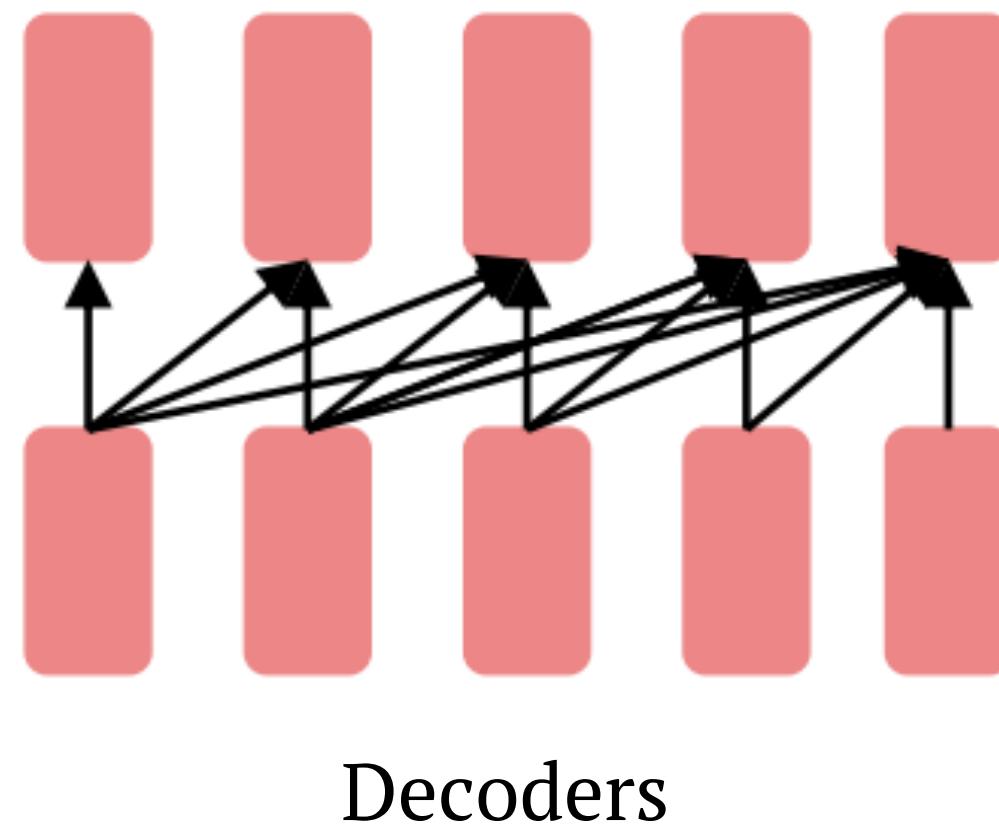
$$\mathbb{E}_{p(x_i, \hat{x}_i)}[p(x_i \mid \hat{x}_i)]$$

— pretraining objective



— type of architecture

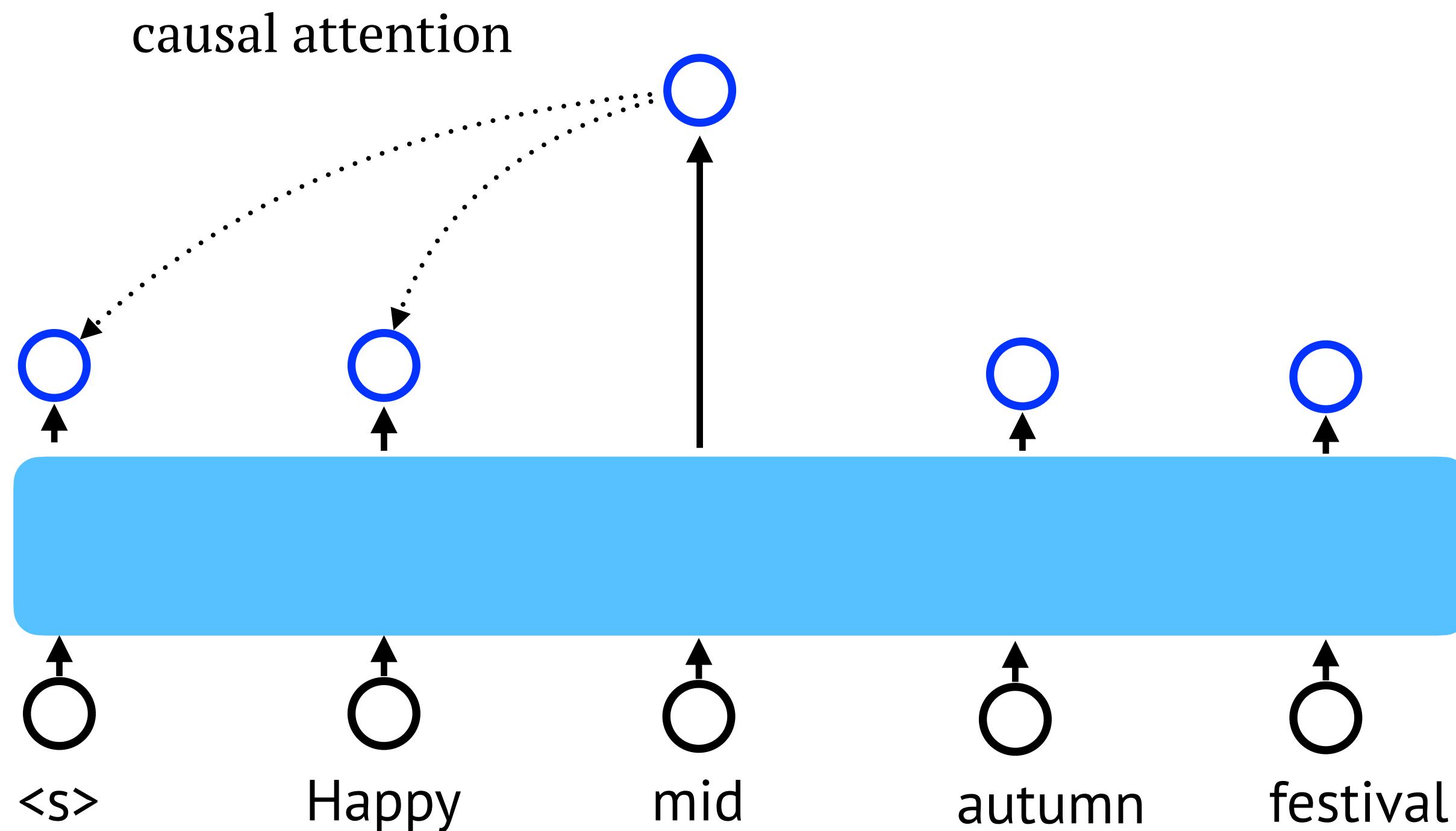
# GPT (Generative Pretrained Transformer)



Radford et al., 2018

# Transformer as Decoder

Need to prevent the attention the future words.



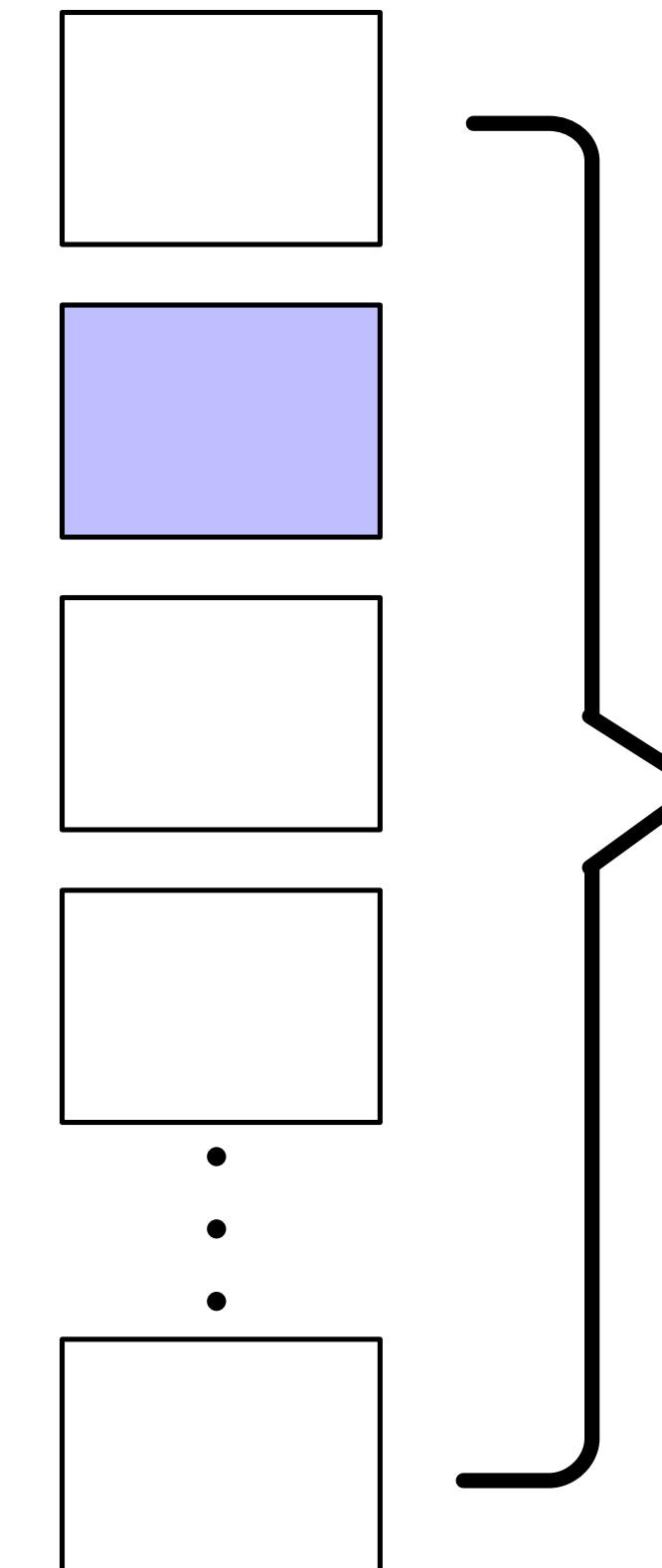
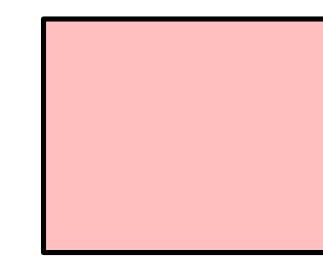
	Happy	mid	autumn	festival
Happy	−∞	−∞	−∞	−∞
mid		−∞	−∞	−∞
autumn			−∞	−∞
festival				−∞

$$e_{ij} = \begin{cases} q_i^\top k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

# GPT (Generative Pretrained Transformer)



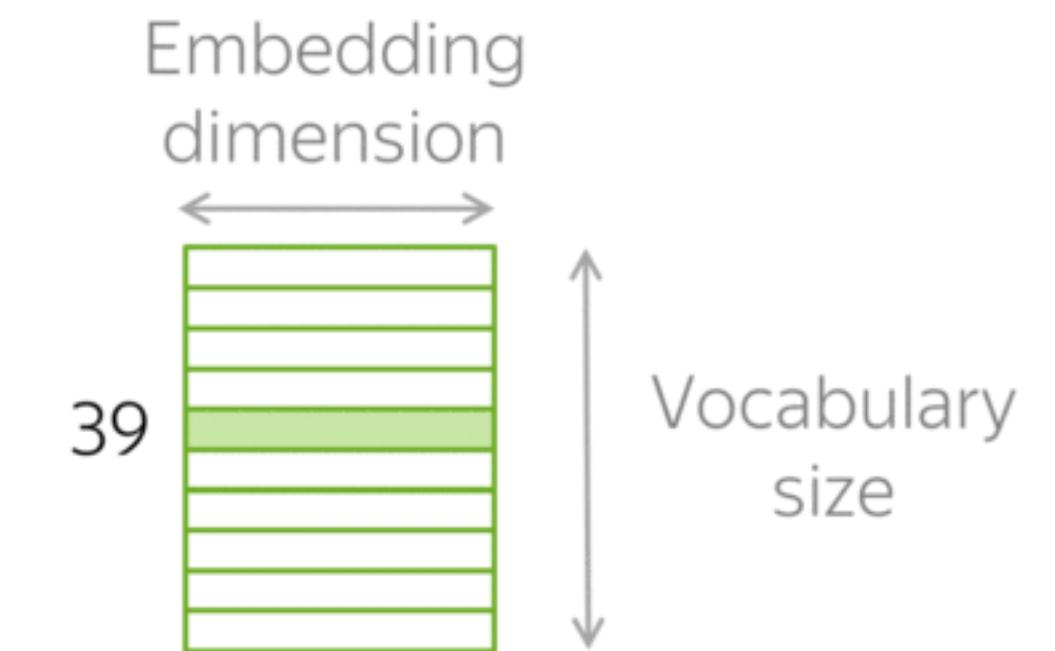
Previous Context



Next Word

Token index in  
the vocabulary

39 1592 10 2548 5  
I saw a cat .



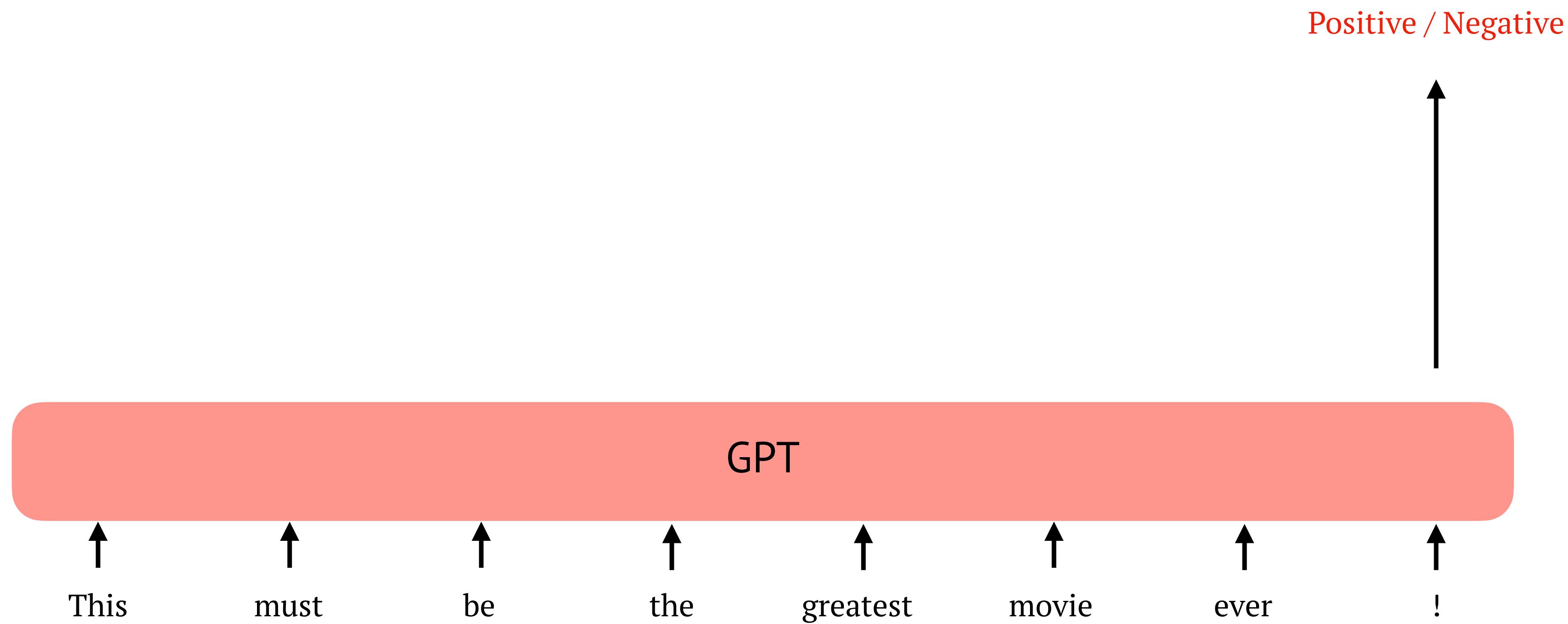
lookup table

Transformer

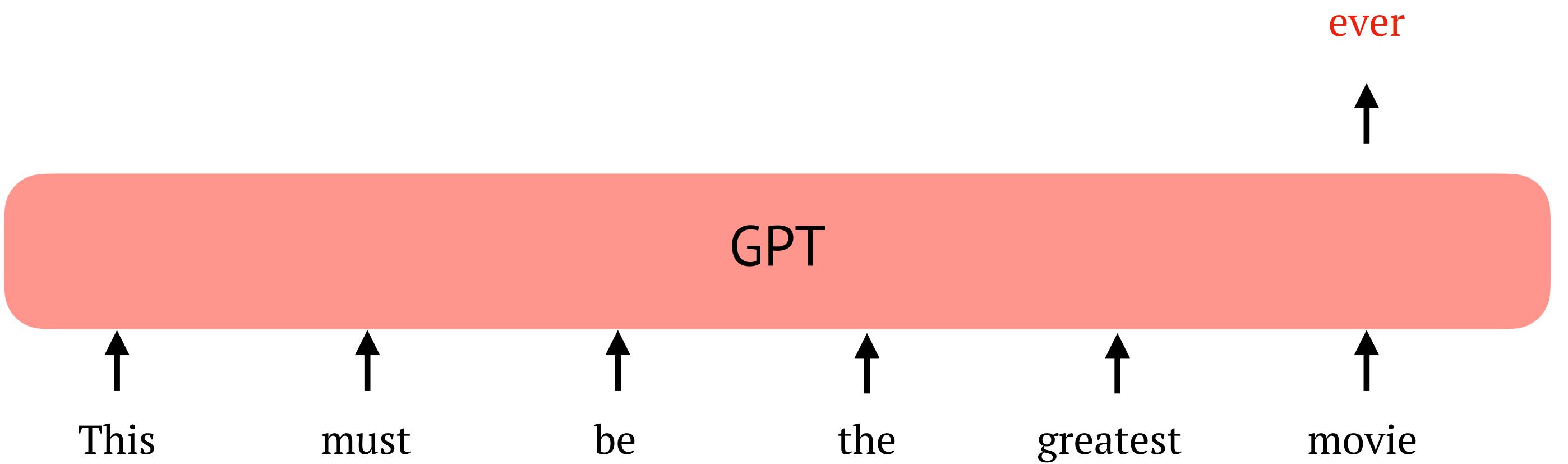
Radford et al., 2018

GIF credit: [Lena Voita](#)

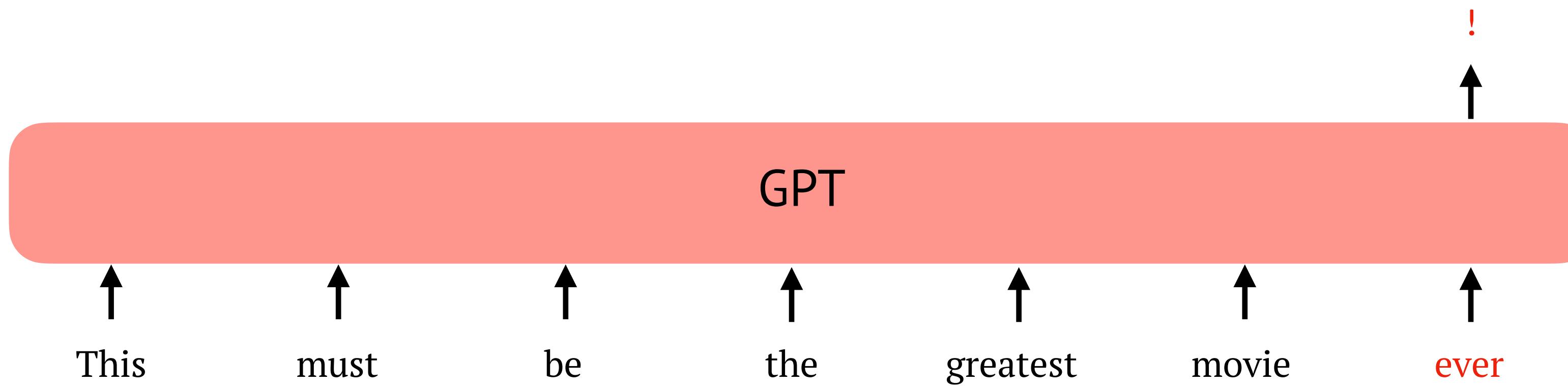
# GPT for Understanding



# GPT for Generation

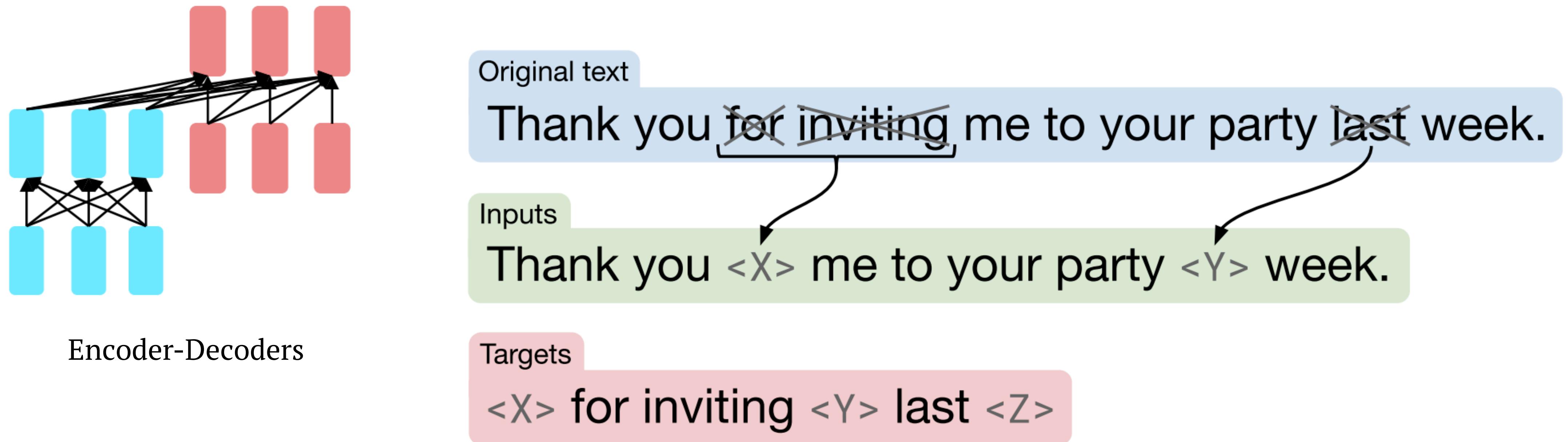


# GPT for Generation



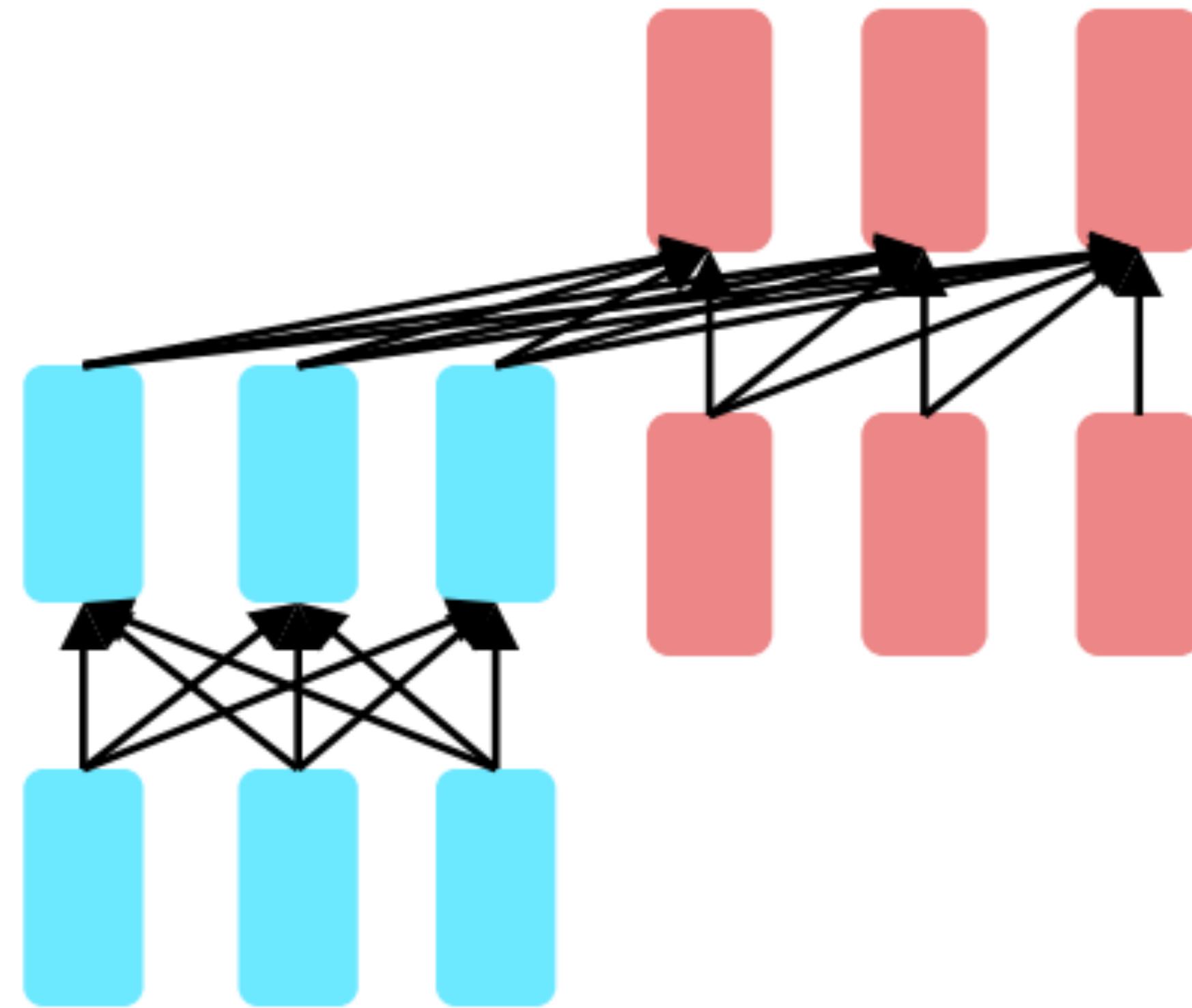
Just “grow” the transformer!

# T5 (Text-to-Text Transfer Transformer)



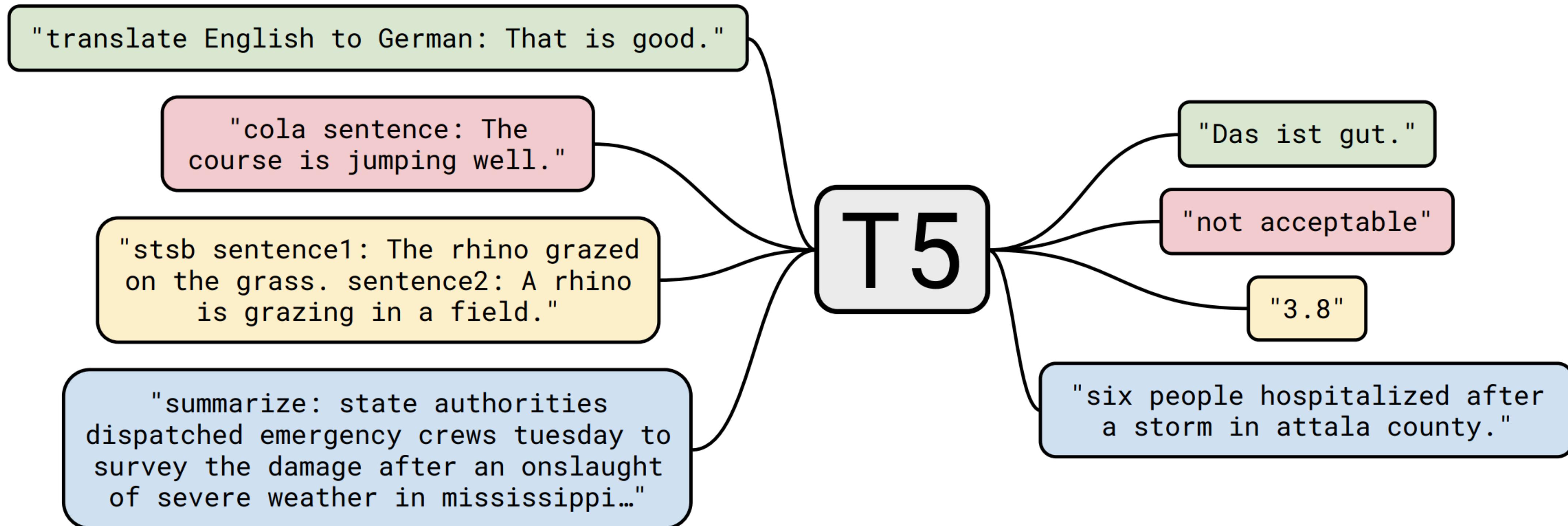
# T5 (Text-to-Text Transfer Transformer)

<X> for inviting <Y> last <Z>



Thank you <X> me to your party <Y> week.

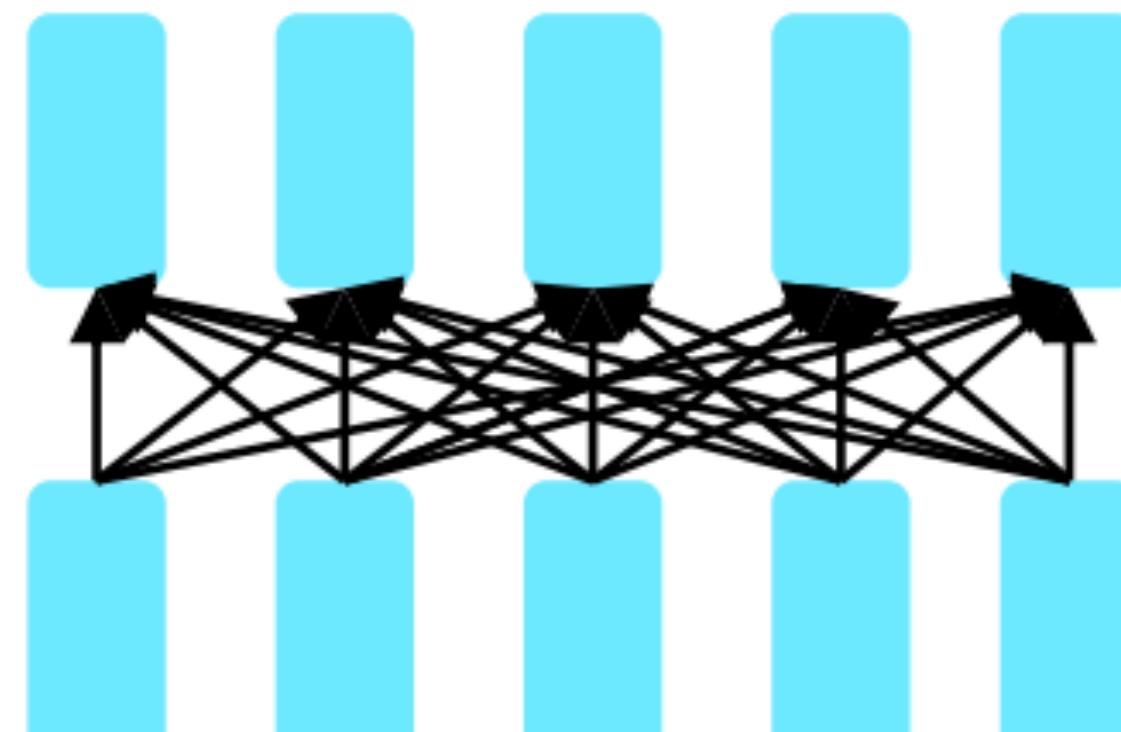
# T5 (Text-to-Text Transfer Transformer)



# T5 (Text-to-Text Transfer Transformer)

Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style <a href="#">Devlin et al. (2018)</a>	Thank you <M> <M> me to your party apple week .	( <i>original text</i> )
Deshuffling	party me for your to . last fun you inviting week Thank	( <i>original text</i> )
MASS-style <a href="#">Song et al. (2019)</a>	Thank you <M> <M> me to your party <M> week .	( <i>original text</i> )
I.i.d. noise, replace spans	Thank you <X> me to your party <Y> week .	<X> for inviting <Y> last <Z>
I.i.d. noise, drop tokens	Thank you me to your party week .	for inviting last
Random spans	Thank you <X> to <Y> week .	<X> for inviting me <Y> your party last <Z>

# ELMo (Embeddings from Language Models)

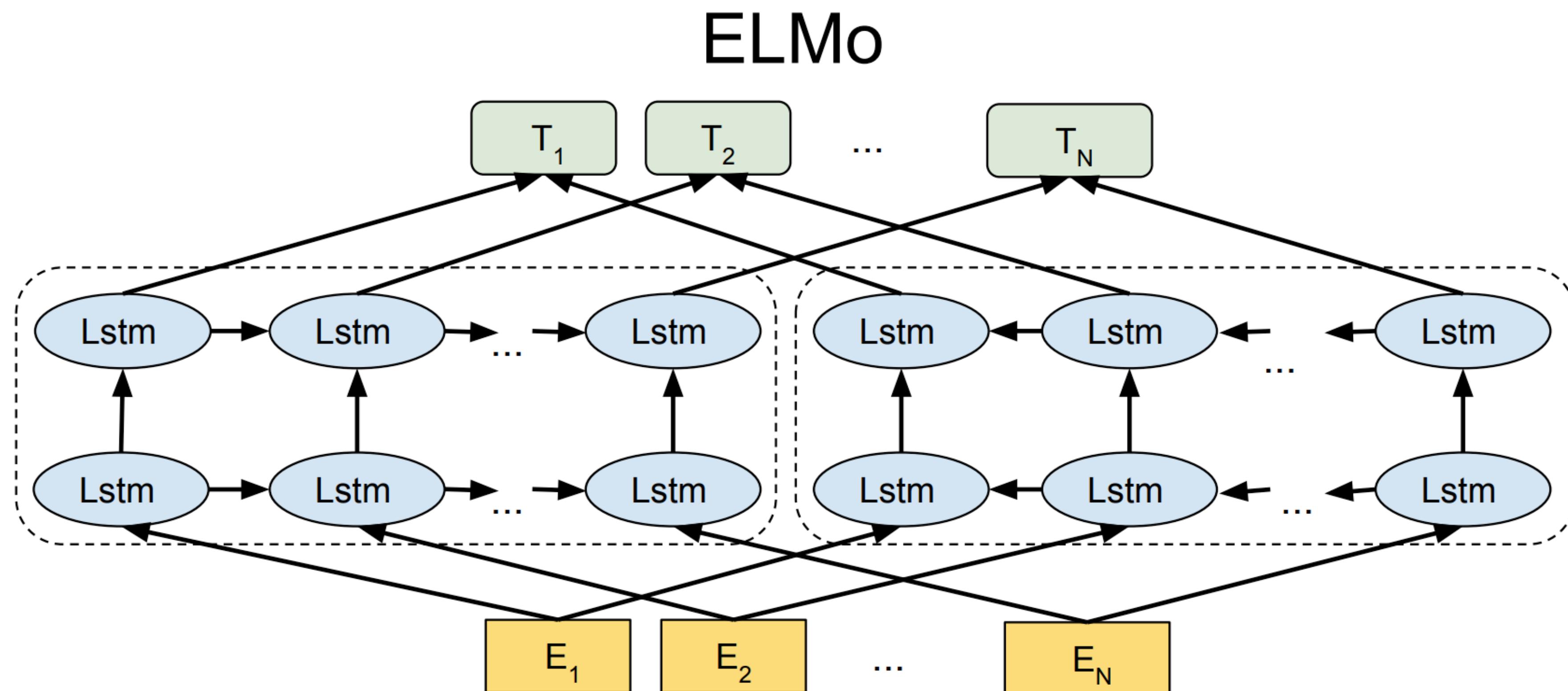


Encoders

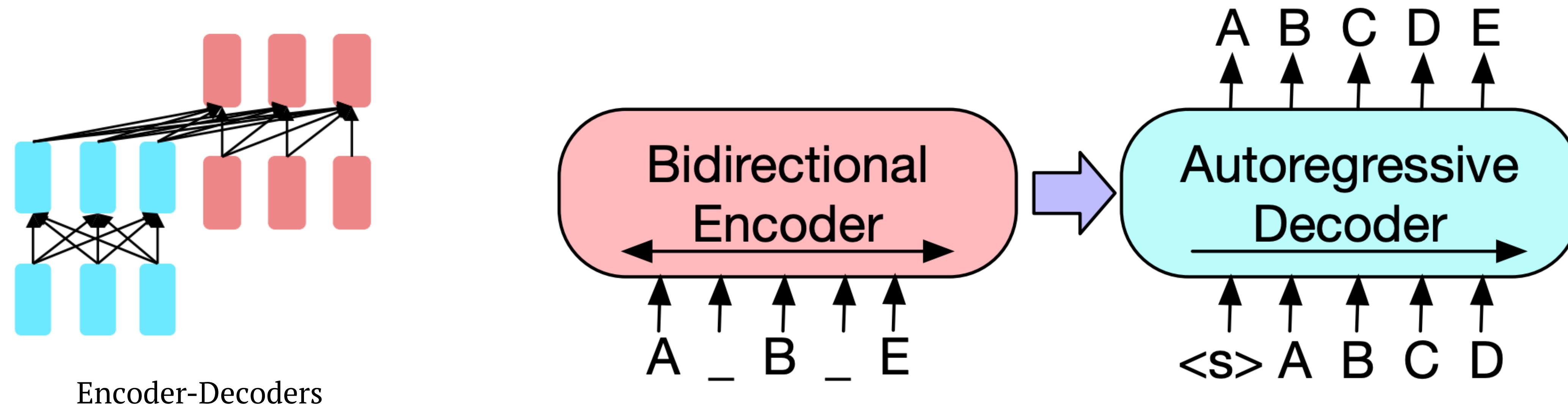
$$\sum_{k=1}^N \left( \log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s) \right)$$

Bidirectional Language Model

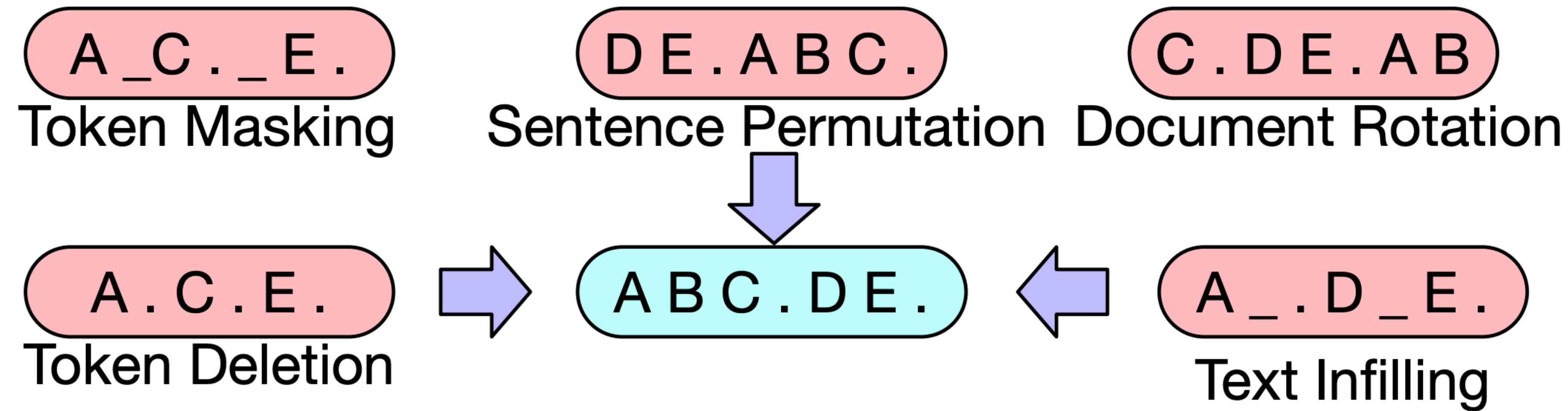
# ELMo (Embeddings from Language Models)



# BART (Denoising Sequence-to-Sequence Pre-training )



# BART (Denoising Sequence-to-Sequence Pre-training )



# InfoWord

