# Lecture 15: Transformer-based Encoders



Claire Cardie, Tanya Goyal

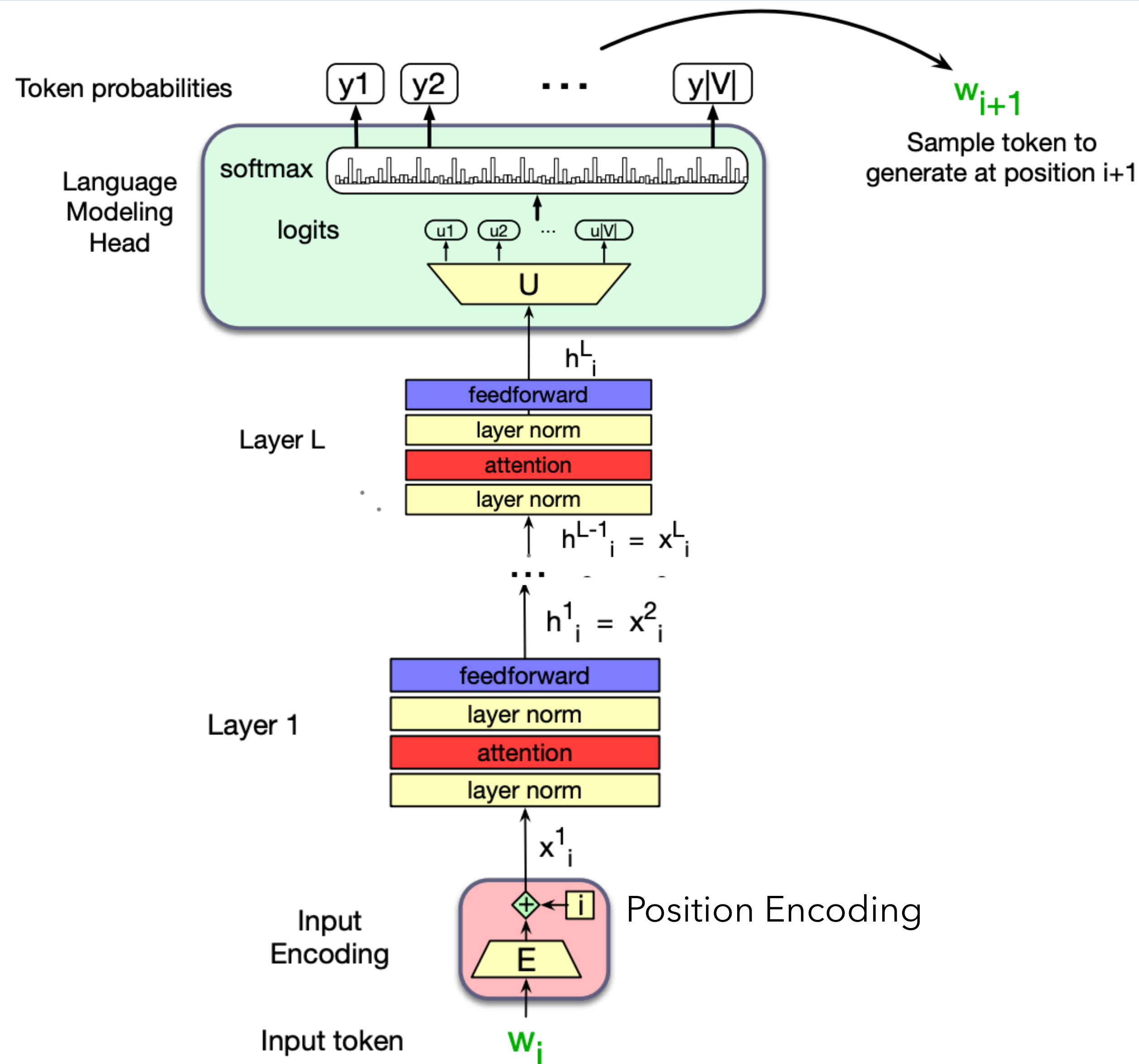CS 4740 (and crosslists): Introduction to Natural Language Processing

# Announcements

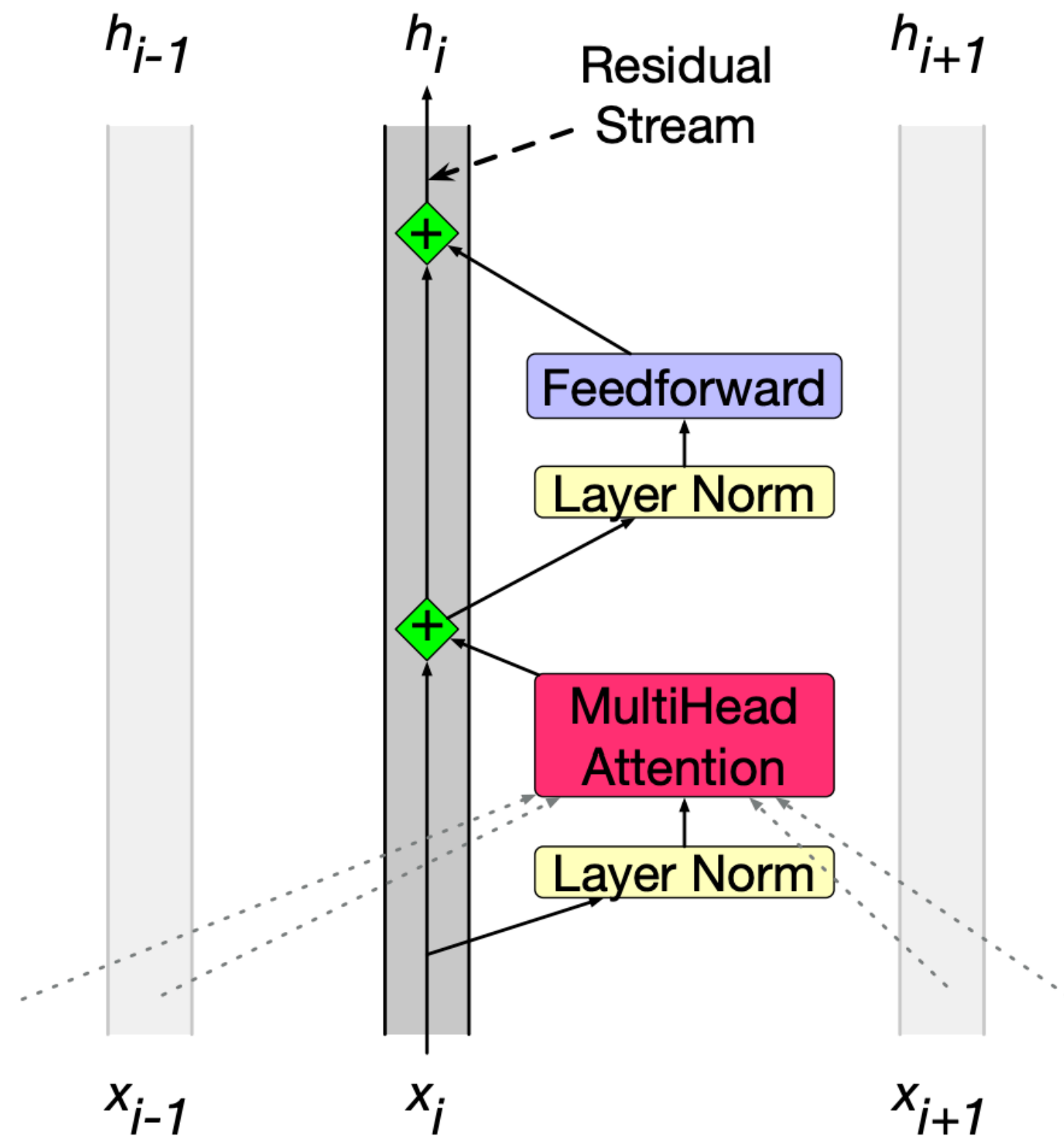- HW3 will be released on Wednesday.

# Today

- Recap: Decoder-only transformer models

- Encoder-decoder transformer models

- Encoder-only models (BERT)

# Recap Decoder-only Transformer



- Main components of a transformer model
  - **(Multi-head) Attention**
  - Feed forward
  - Layer Norm

  - Position Encoding

# Recap: Residual Stream view



Input $x_i$ at time step $i$

$$t_i^1 = \text{LayerNorm}(x_i)$$

$$t_i^2 = \text{MultiHead-Attention}(t_i^1, [t_1^1, t_2^1, \ldots t_N^1])$$
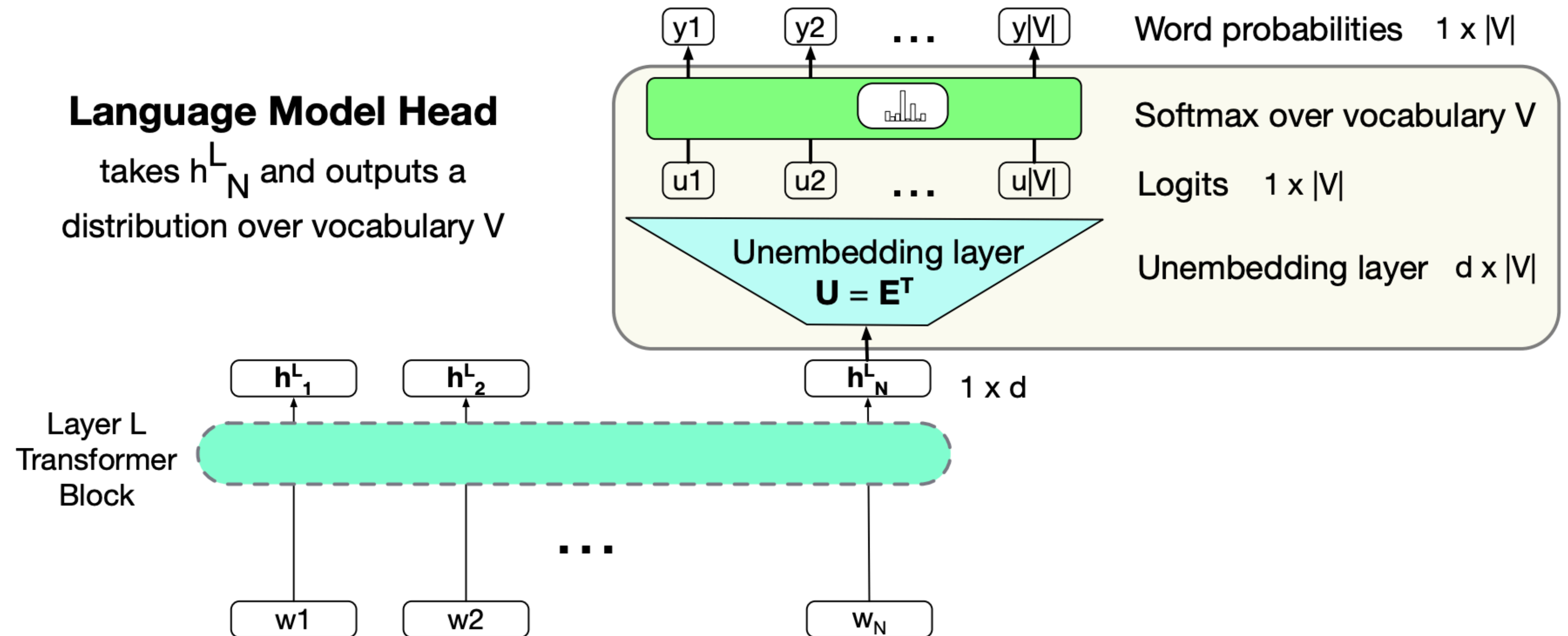
$$t_i^3 = t_i^2 + x_i$$

$$t_i^4 = \text{LayerNorm}(t_i^3)$$

$$t_i^5 = \text{FFNN}(t_i^4)$$

$$h_i = t_i^5 + t_i^3$$

# Recap: Output Layer

- Final output: probability distribution over the vocabulary.

- Training objective: Predict the next token, given preceding tokens.

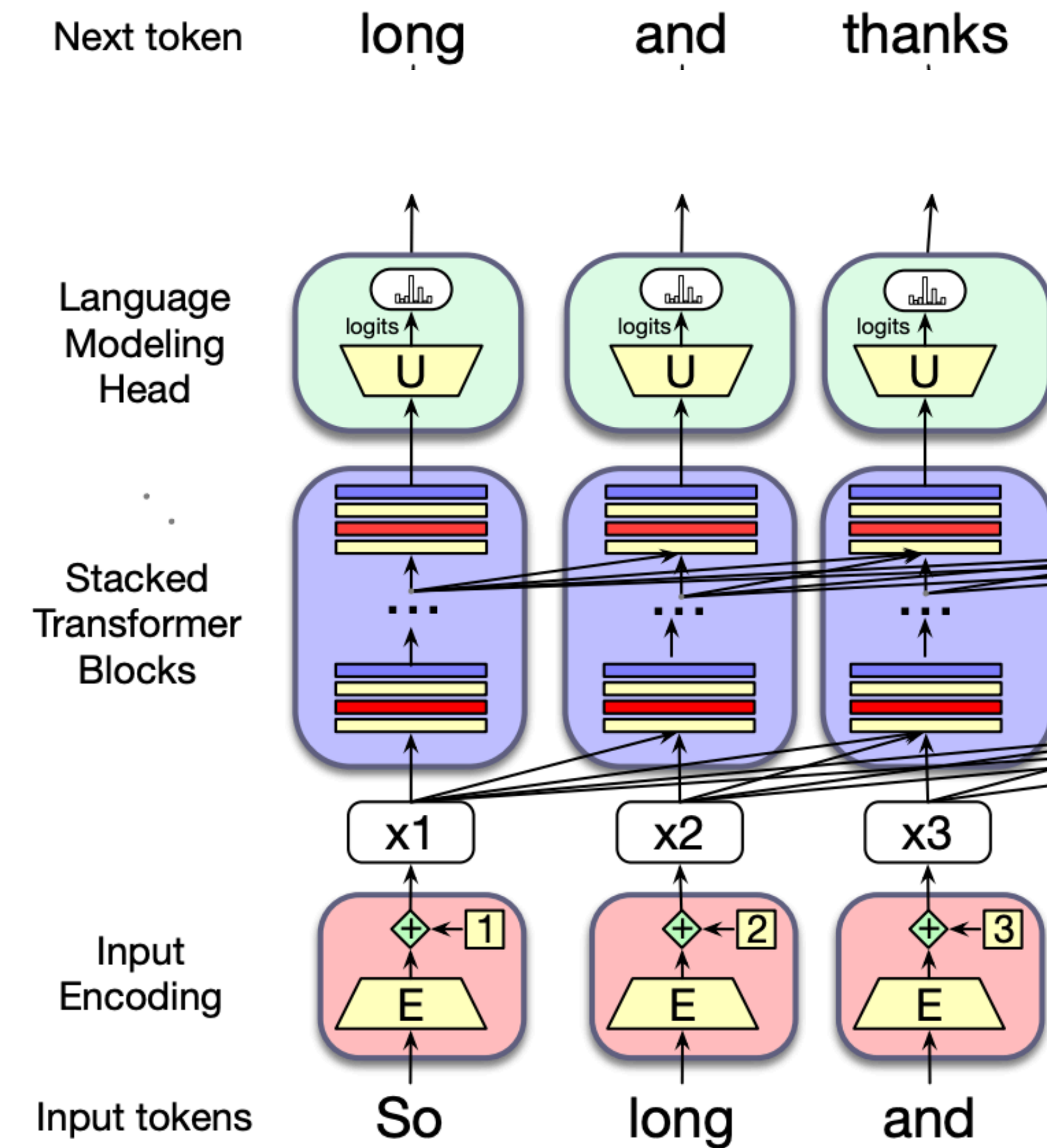# Recap: Language Modeling with Decoder-only Transformers

- Q: How should we model loss?

  - Cross Entropy Loss

  $$L_{CE} = -\sum_{w \in V} \mathbf{y}_t[w] \ \log \hat{\mathbf{y}}_t[w]$$

- Simplifies to:

  $$L_{CE} = -\log \hat{\mathbf{y}}_t[w_{t+1}]$$

# Large Language Models

- Decoder-only transformer models allows us to model the task of language modeling. Why should we care about language modeling?

- Many practical tasks in NLP can be cast as next token prediction.

**Sentiment Analysis:**

The sentiment of the sentence ''I like Jackie Chan" is:

P(positive|The sentiment of the sentence ''I like Jackie Chan" is:)
P(negative|The sentiment of the sentence ''I like Jackie Chan" is:)

# Large Language Models

- Decoder-only transformer models allows us to model the task of language modeling. Why should we care about language modeling?

- Many practical tasks in NLP can be cast as next token prediction.

**Question Answering:**

Q: Who wrote the book ''The Origin of Species"? A:

P(w | Q: Who wrote the book ''The Origin of Species"? A:)

# Large Language Models

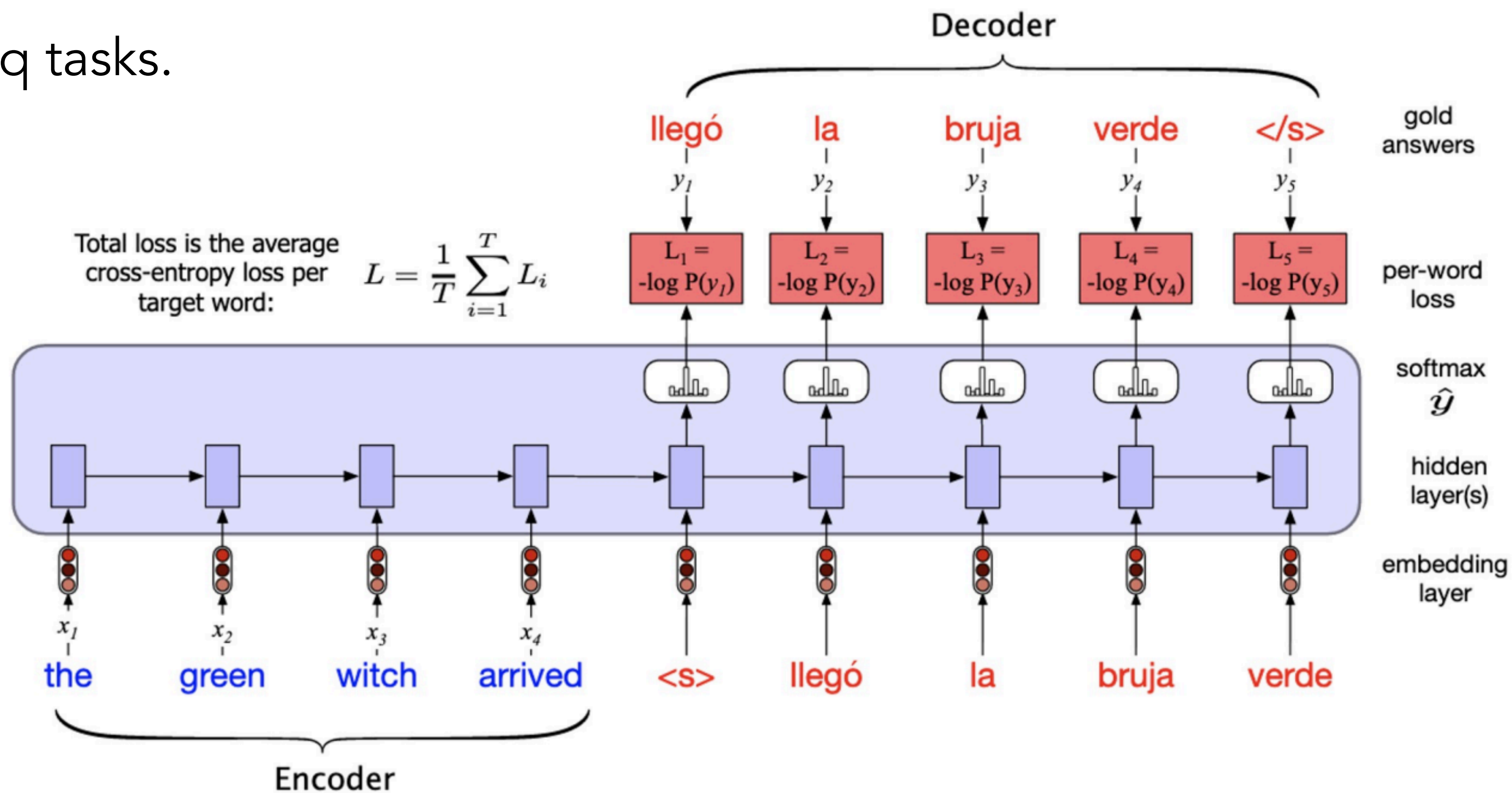- Decoder-only transformer models allows us to model the task of language modeling. Why should we care about language modeling?

- Many practical tasks in NLP can be cast as next token prediction.

Language models need to be very powerful perform well at all these tasks!

- **Very** deep network

- Train on a lot of data

    - E.g. GPT-3 model (released in 2020) trained on 300B tokens, LLaMA-3 model trained on 15T tokens.

# Encoder-Decoder Architecture

- Recall RNNs.

- Useful for seq2seq tasks.



Decoder

Total loss is the average cross-entropy loss per target word:

$$L = \frac{1}{T} \sum_{i=1}^{T} L_i$$

gold answers: llegó, la, bruja, verde, </s>

$y_1$, $y_2$, $y_3$, $y_4$, $y_5$

$L_1 = -\log P(y_1)$, $L_2 = -\log P(y_2)$, $L_3 = -\log P(y_3)$, $L_4 = -\log P(y_4)$, $L_5 = -\log P(y_5)$ — per-word loss

softmax $\hat{y}$

hidden layer(s)

embedding layer

$x_1$, $x_2$, $x_3$, $x_4$

the, green, witch, arrived, <s>, llegó, la, bruja, verde

Encoder

# Encoder-Decoder Architecture

- Recall RNNs.

- Useful for seq2seq tasks.

- What about transformers?



**Passed to the decoder**

concat

$h_{1\_back}$

RNN 2 (Right to Left)

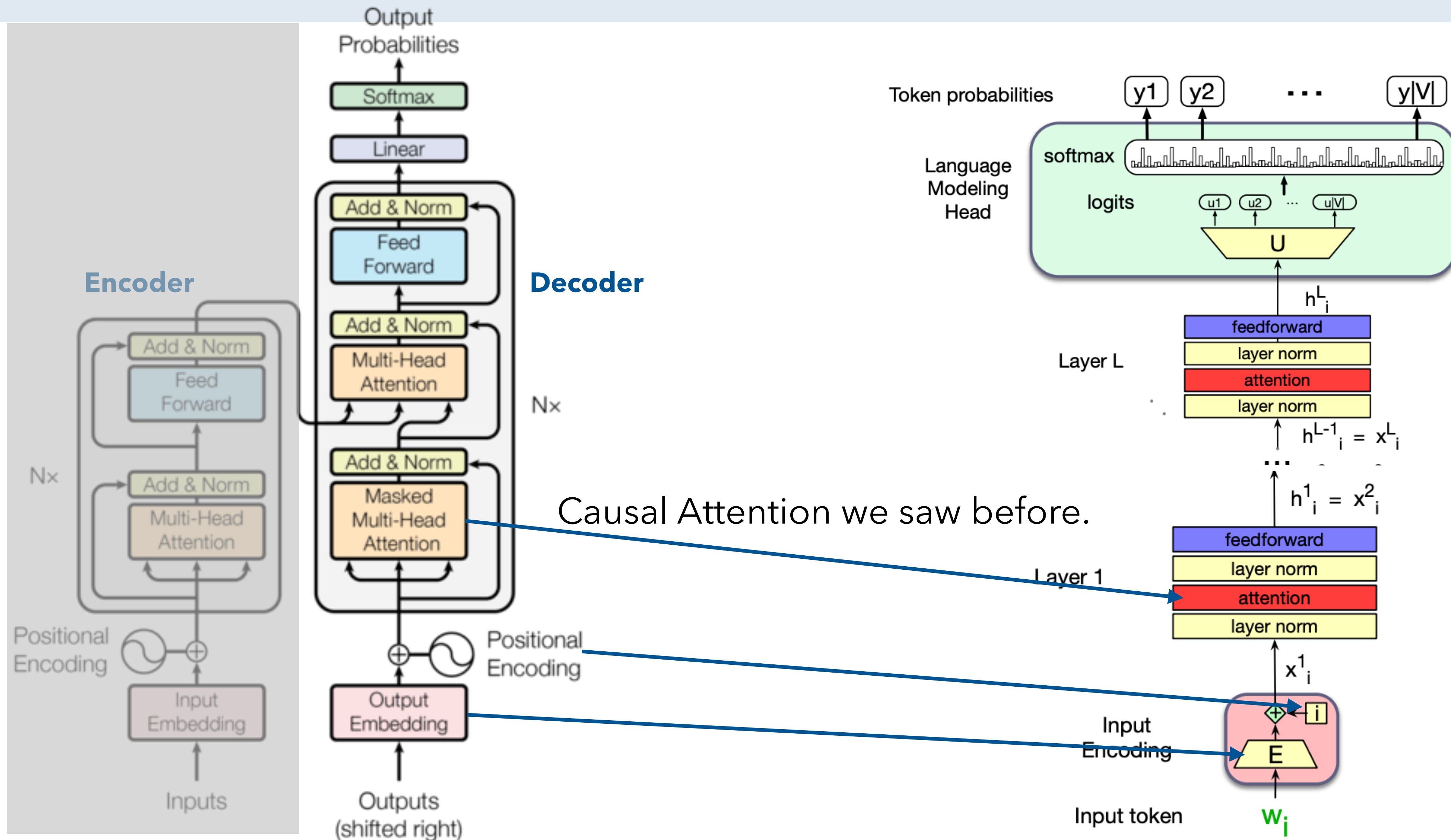RNN 1 (Left to Right)  $h_{n\_forw}$

$x_1$  $x_2$  $x_3$  $x_n$

# Encoder-Decoder Architecture



- The actual figure from "Attention is all you need", Vaswani et al, 2017 paper.
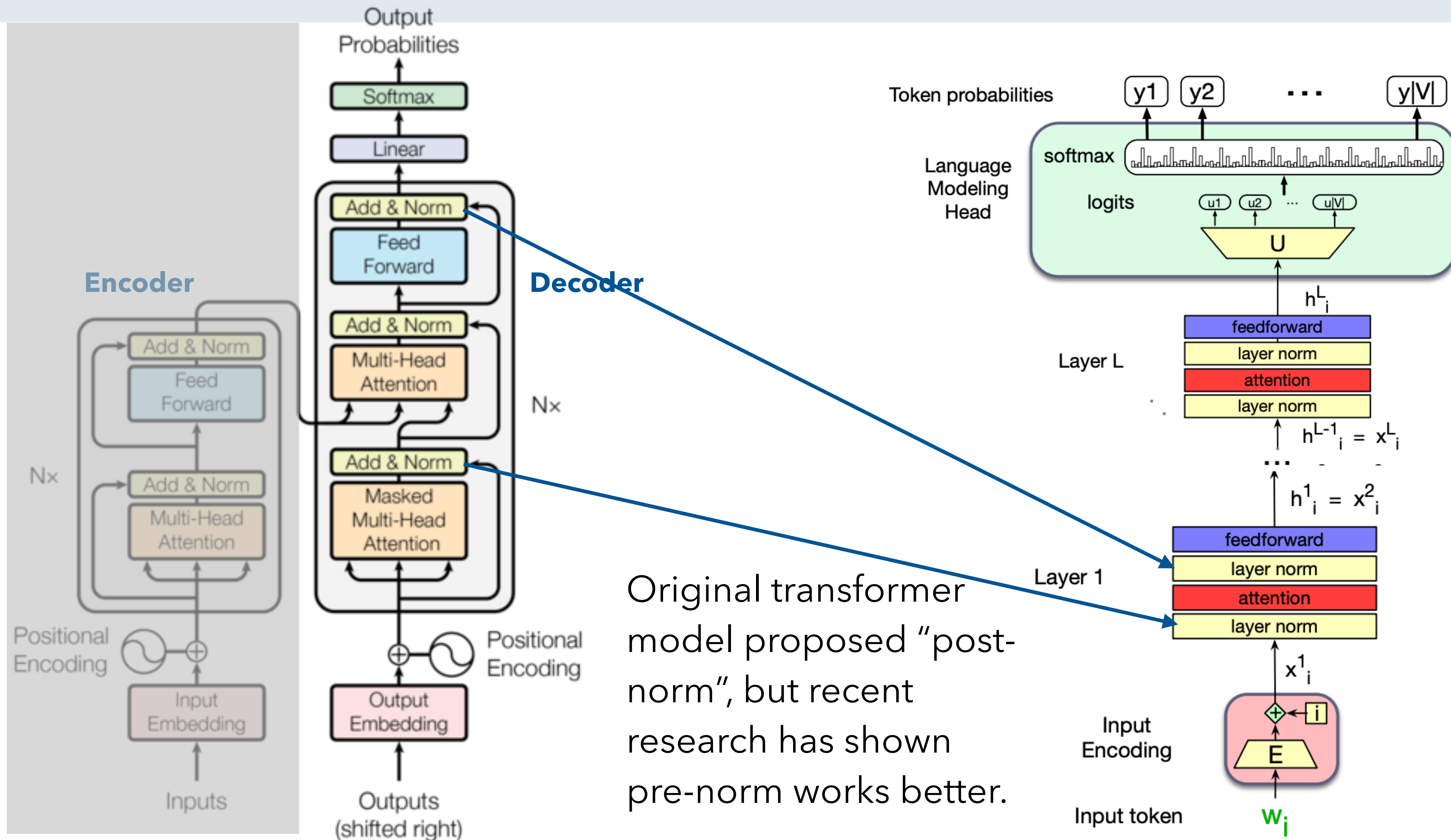
- … and the figure you will see everywhere on the internet
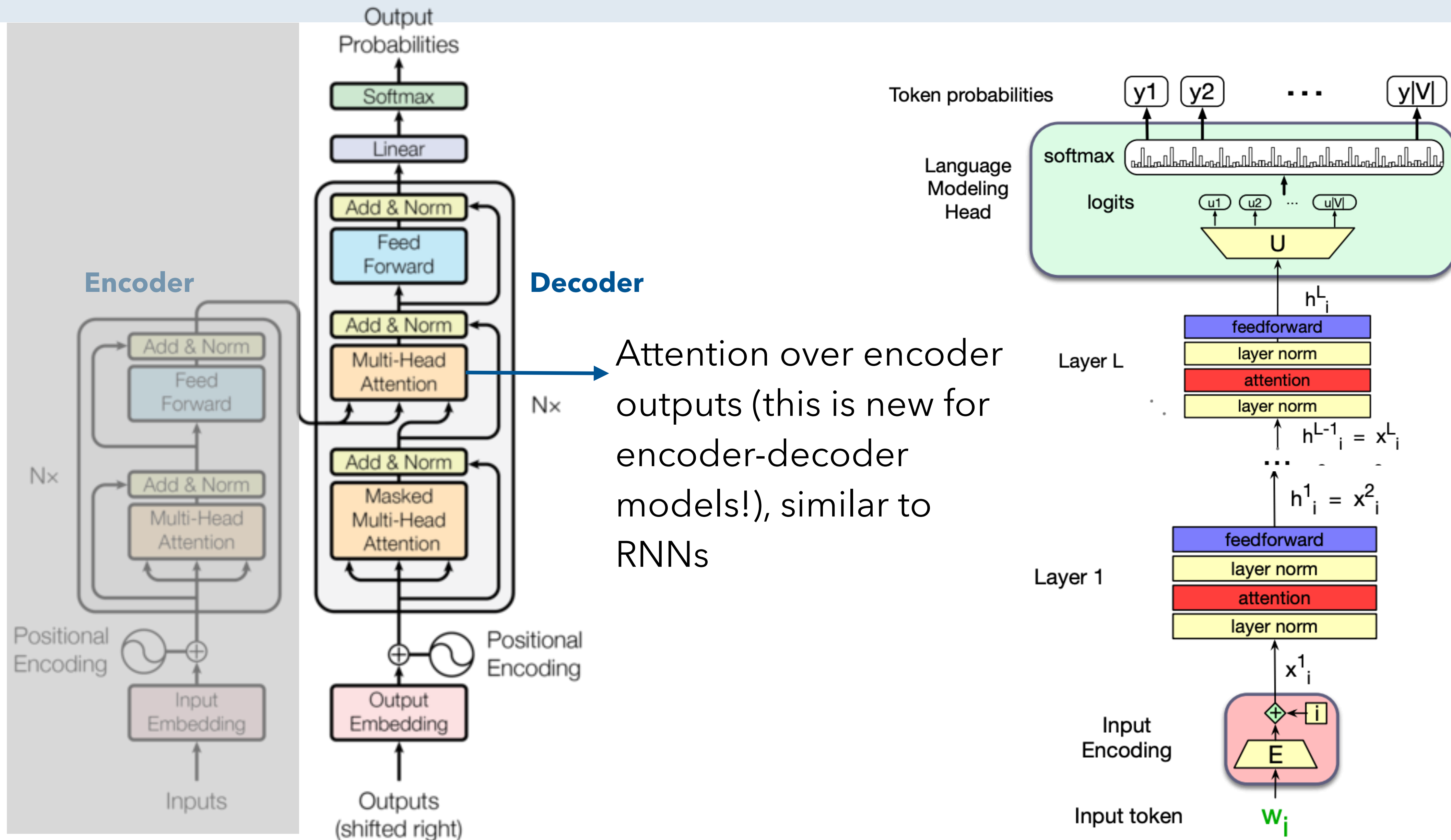
# Encoder-Decoder Architecture



Encoder

Decoder

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Output Embedding

Outputs (shifted right)

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

Causal Attention we saw before.

Token probabilities: $y1$, $y2$, $\cdots$, $y|V|$

Language Modeling Head

softmax

logits: $u1$, $u2$, $\cdots$, $u|V|$

$U$

$h^L_i$

feedforward

layer norm

attention

layer norm

Layer L

$h^{L-1}_i = x^L_i$

$\cdots$

$h^1_i = x^2_i$

feedforward

layer norm

attention

layer norm

Layer 1

$x^1_i$

Input Encoding

$i$

$E$

Input token

$w_i$

# Encoder-Decoder Architecture

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

**Encoder**

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

**Decoder**

N×

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Token probabilities

y1  y2  · · ·  y|V|

Language Modeling Head

softmax

logits   u1  u2  ...  u|V|

U

$h^L_i$

feedforward

layer norm

attention

layer norm

Layer L

$h^{L-1}_i = x^L_i$

· · ·

$h^1_i = x^2_i$

feedforward

layer norm

attention

layer norm

Layer 1

$x^1_i$

Input Encoding

E

Input token   $w_i$

Original transformer model proposed "post-norm", but recent research has shown pre-norm works better.

# Encoder-Decoder Architecture



**Decoder**

Attention over encoder outputs (this is new for encoder-decoder models!), similar to RNNs

# Encoder-Decoder Architecture



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Output Embedding

Outputs (shifted right)

**Encoder**

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

**Decoder**

Residual Connections

Token probabilities   y1   y2   · · ·   y|V|

Language Modeling Head

softmax

logits   u1   u2   ...   u|V|

U

$h^L_i$

feedforward

layer norm

attention

layer norm

Layer L

$h^{L-1}_i = x^L_i$

...

$h^1_i = x^2_i$

feedforward

layer norm

attention

layer norm

Layer 1

$x^1_i$

Input Encoding

$\oplus$   i

E

Input token   $w_i$

# Encoder-Decoder Architecture



- Multi-head attention from decoder states at each layer to **output of the last layer** of the encoder.

# Encoder-Decoder Architecture



Sources: (middle) The Illustrated Transformer
(left and right): Cropped from Attention Is All You Need (AIAYN)

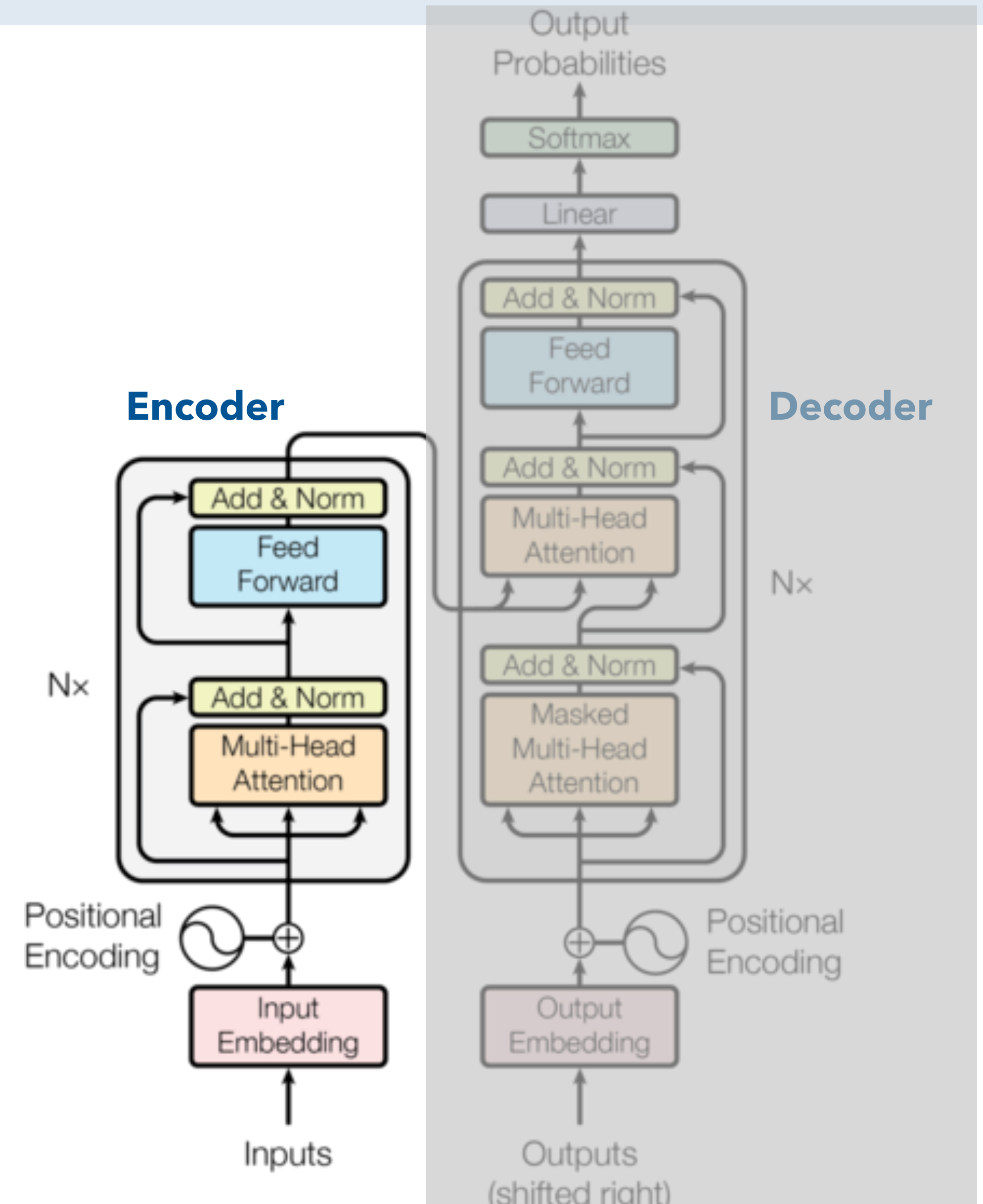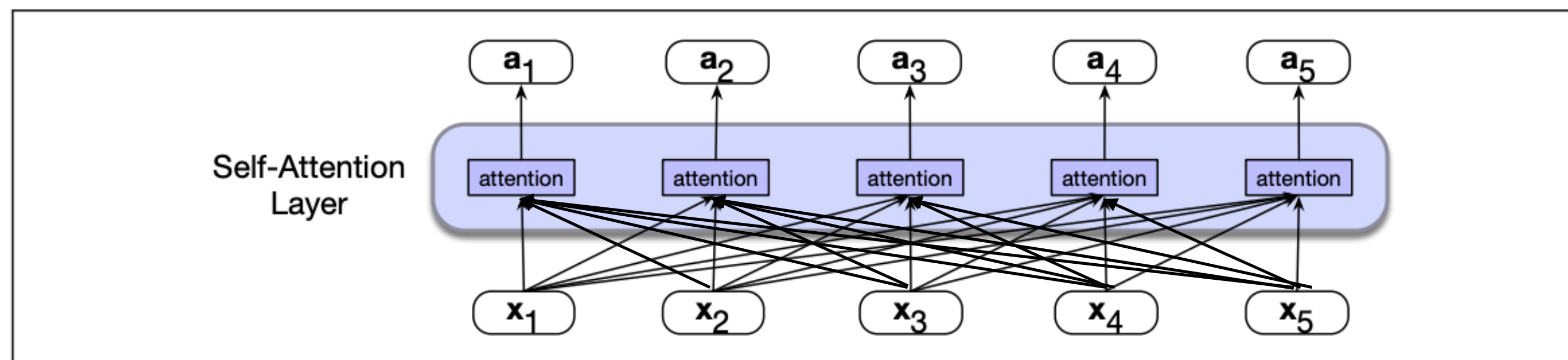# Encoder-Decoder Architecture



- Recall:

$$\mathbf{k}_1 \quad \mathbf{k}_2 \quad \mathbf{k}_3 \qquad \mathbf{q}_3$$

| 0.02 | 0.02 | 0.96 |
|------|------|------|

$$\alpha_{31}\mathbf{v}_1 \quad \alpha_{32}\mathbf{v}_2 \quad \alpha_{33}\mathbf{v}_3$$

$$+ \quad + \quad = \quad W^O = \mathbf{a}_3$$

- What is the source of keys, queries, values in attention from decoder to encoder?

# Encoder-Decoder Architecture

- Encoder architecture is similar to decoder.

- Only difference: This attention is **not** causal.

  - All tokens attend to all other tokens.

# Encoder-Decoder Architecture

- Encoder architecture is similar to decoder.

- Only difference: This attention is **not** causal.

  - All tokens attend to all other tokens.

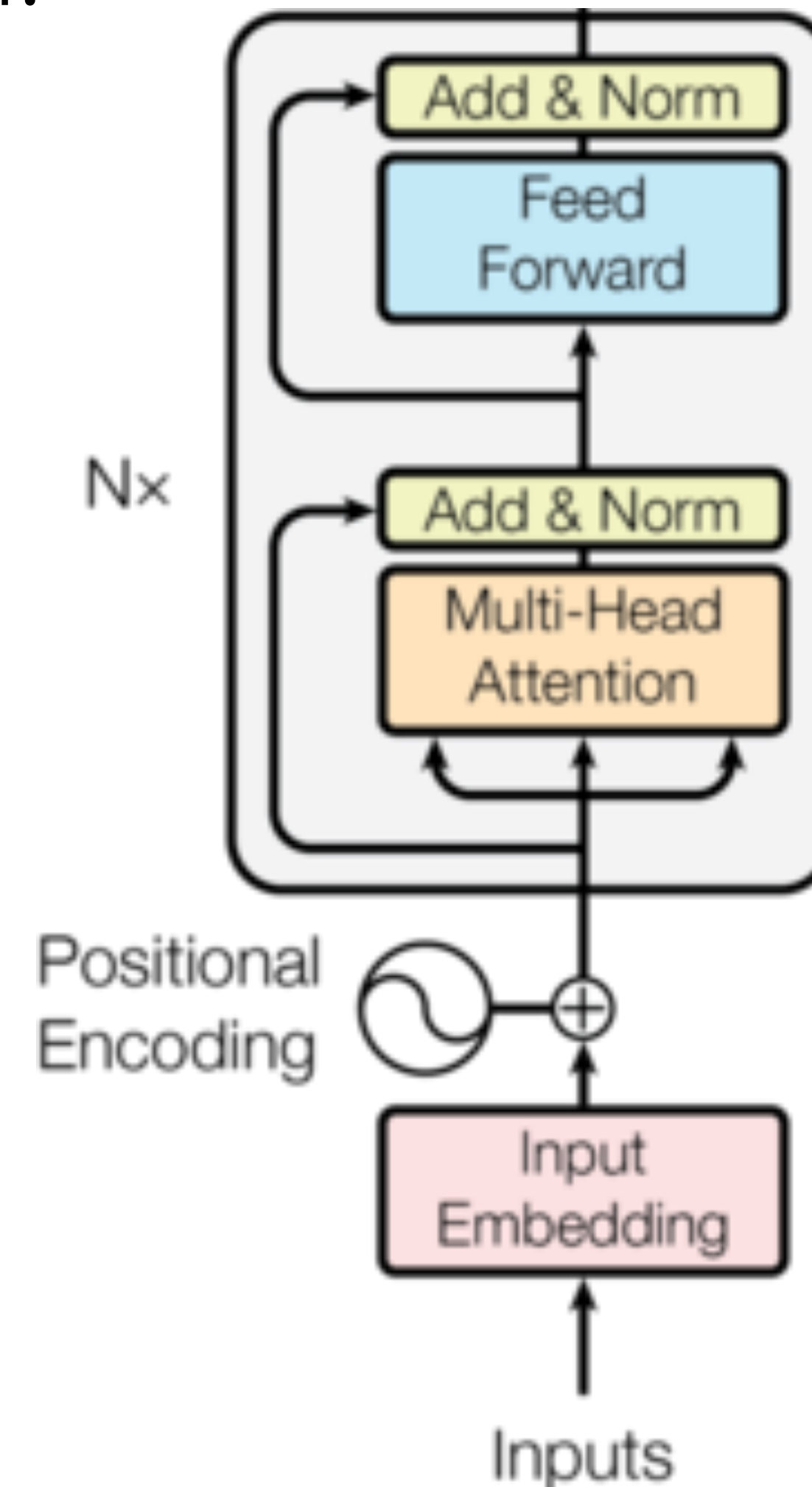- At each time step, the encoder output $h_t$ can be viewed as a "contextual" representation of the input word $w_t$.
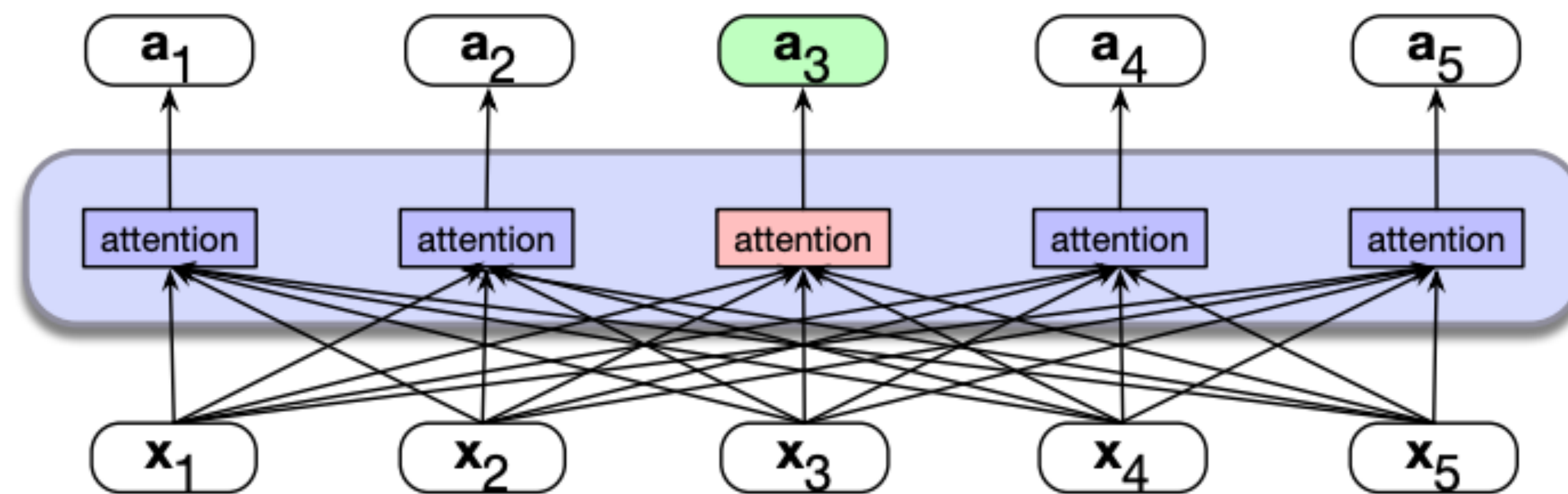
# Encoder-Decoder Architecture

- Training

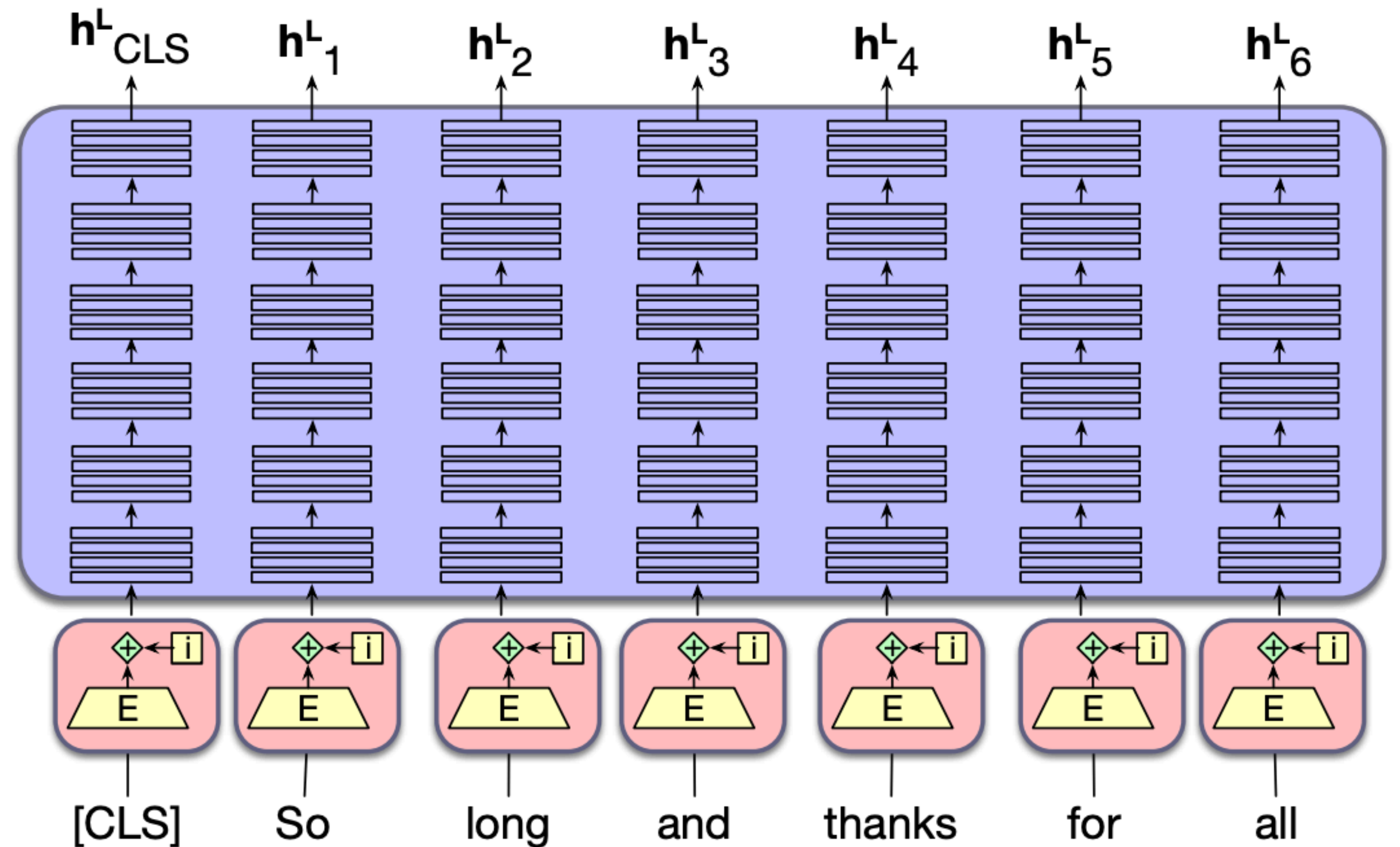  - Next-token prediction at the decoder output.

# Encoder-only architectures

- Same architecture components as the decoder.

- **Difference:** the attention is bidirectional, not causal

# Encoder-only Architecture

- What is this useful for?

  - Gives contextual representations for each input word.
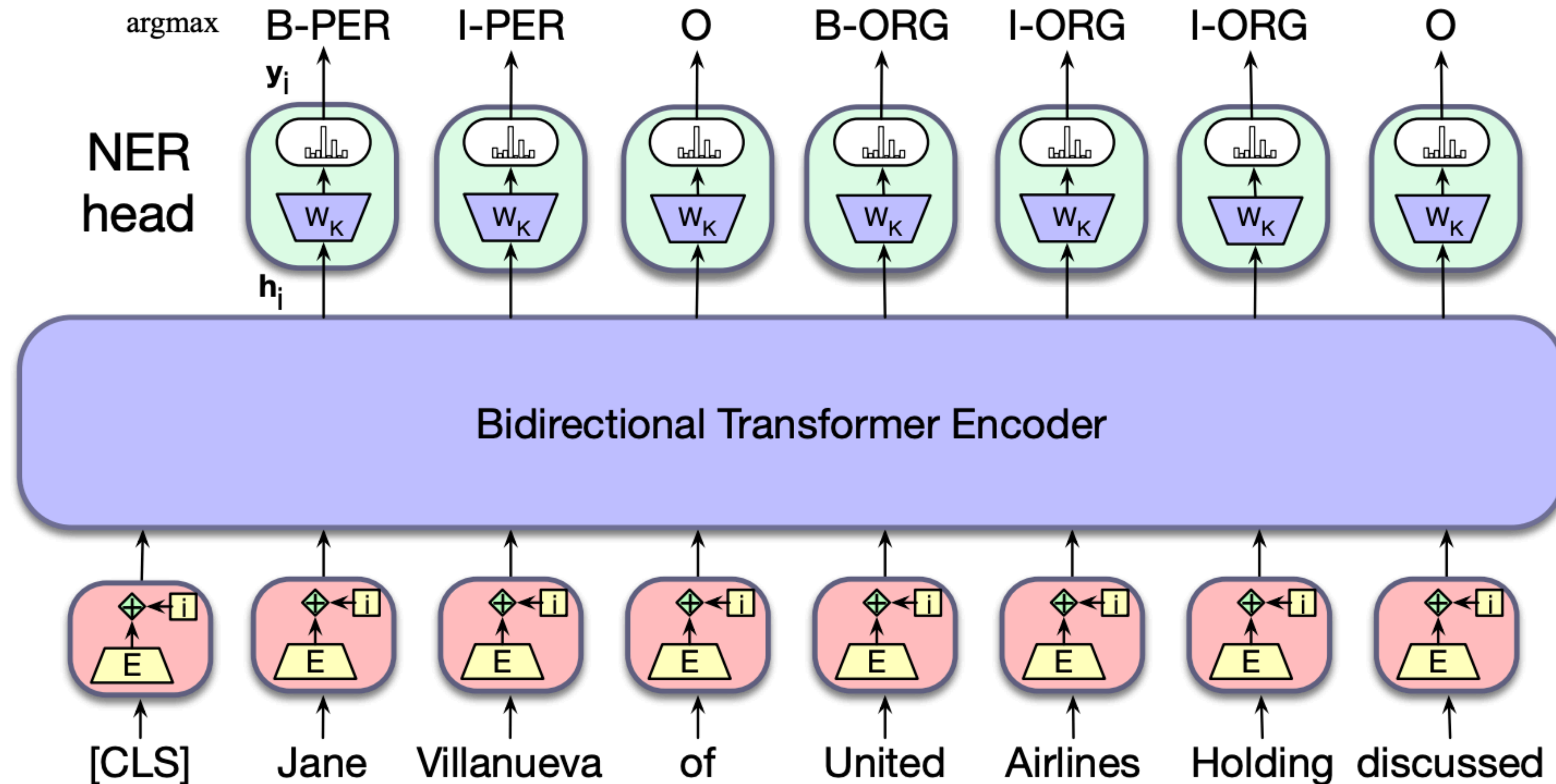
# Encoder-only Architecture

- Word sense disambiguation

  - A sense (or word sense) is a discrete representation of one aspect of the meaning of a word



German article "die"

Was der Fall ist, **die** Tatsache, ist das Bestehen von Sachverhalten.

über **die** Verhandlungen der Königl.

single person dies ←→ multiple people die

a playing die

Chernenko became the first Soviet leader to **die** in less than three years

Over 60 people **die** and over 100 are unaccounted for.

Players must always move a token according to the **die** value

Vaughan's ultimate fantasy was to **die** in a head-on collision with movie star Elizabeth Taylor

Many more **die** from radiation sickness, starvation and cold.

The faces of a **die** may be placed clockwise or counterclockwise

# Encoder-only Architecture
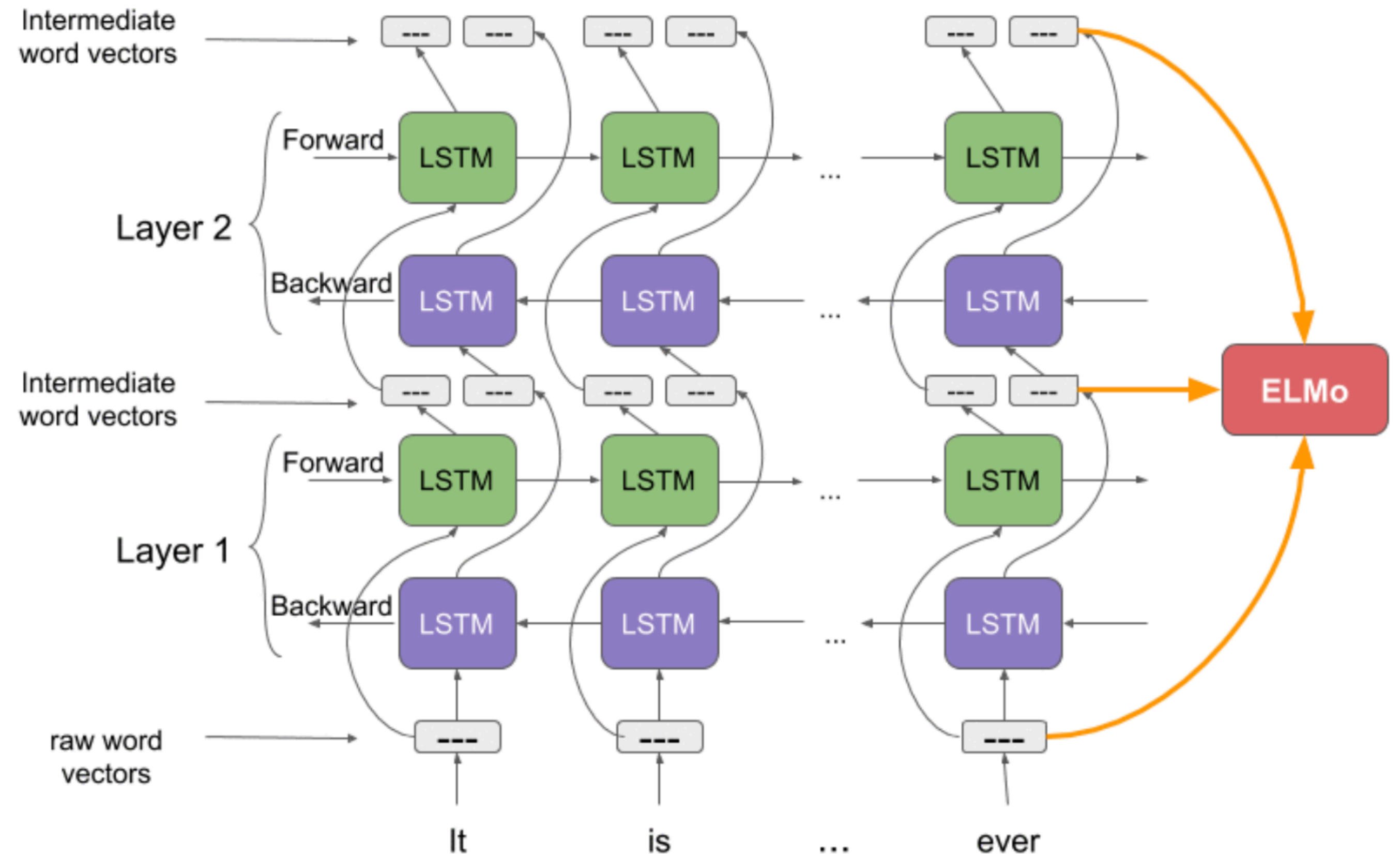
- Classification?

  - Sequence Labeling

# Why are these more powerful than bidirectional RNNs?

- Predecessor: ELMo (Peters et al., 2018)

Uses bi-LSTM-based encoder-decoder.

Combines hidden vectors from different layers.

# Why are these more powerful than bidirectional RNNs?

- BiLSTMs are not quite the same as full self-attention:

Try to predict the italicized word from just the left or just the right context.
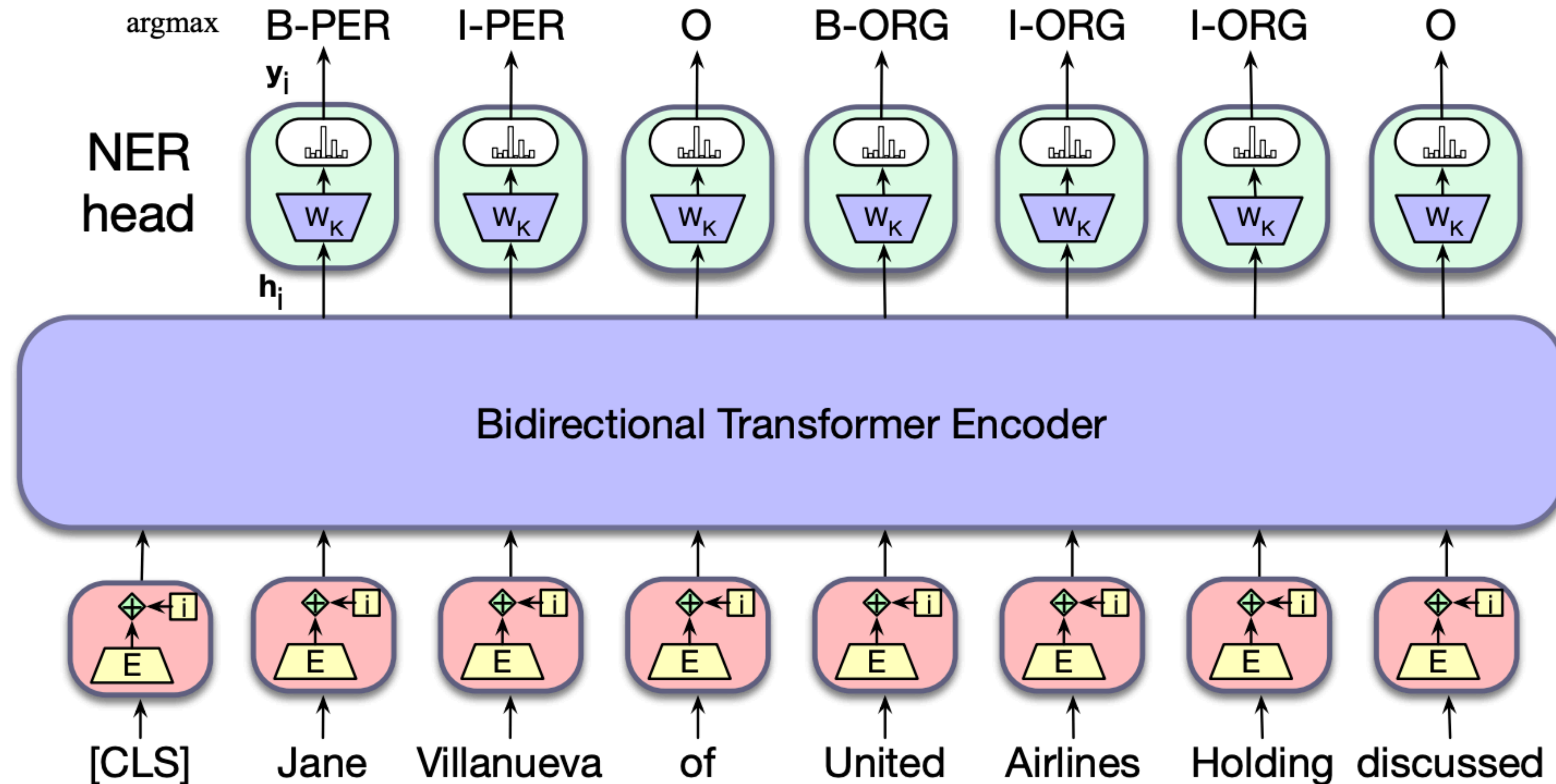The celebrity , Michael ***Jordan*** , was a player in the NBA .

A left-to-right model could also answer "Jackson" and be sensible but wrong. (Reference to singer/songwriter Michael Jackson)

A right-to-left model could also answer "Curry" and be sensible but wrong. (Reference to NBA player Steph Curry)
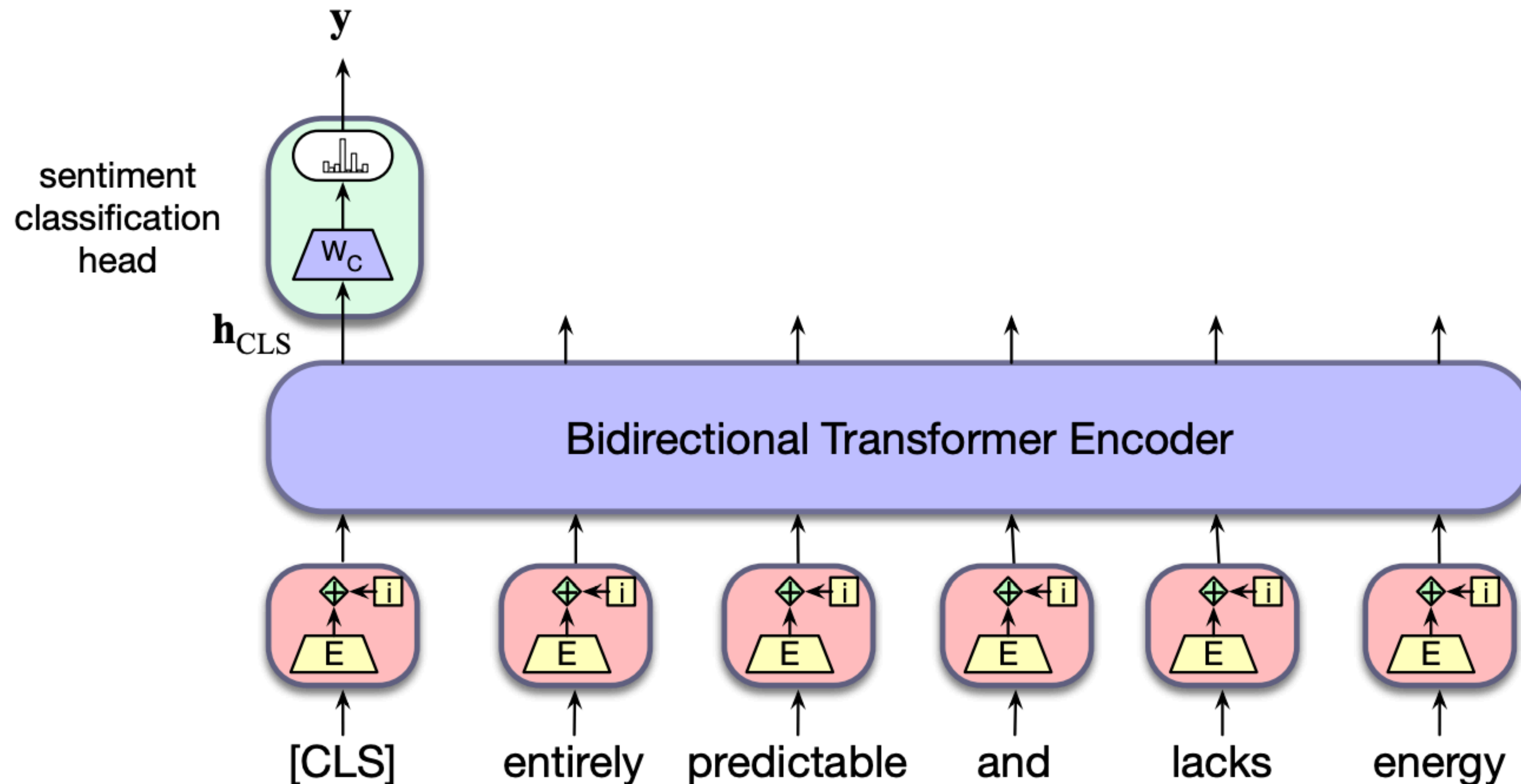
# Encoder-only Architecture

- Classification?

  - Sequence Tagging
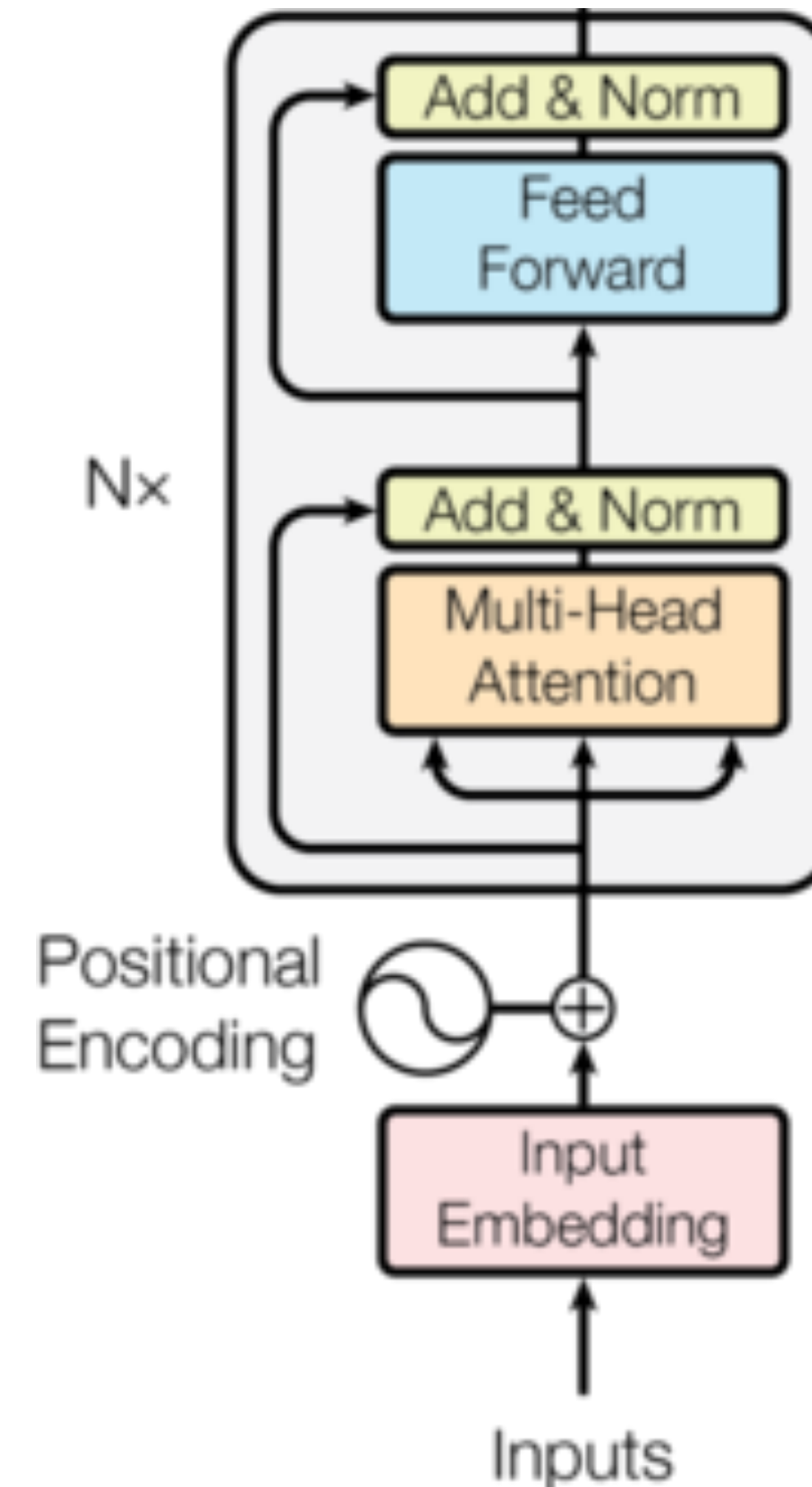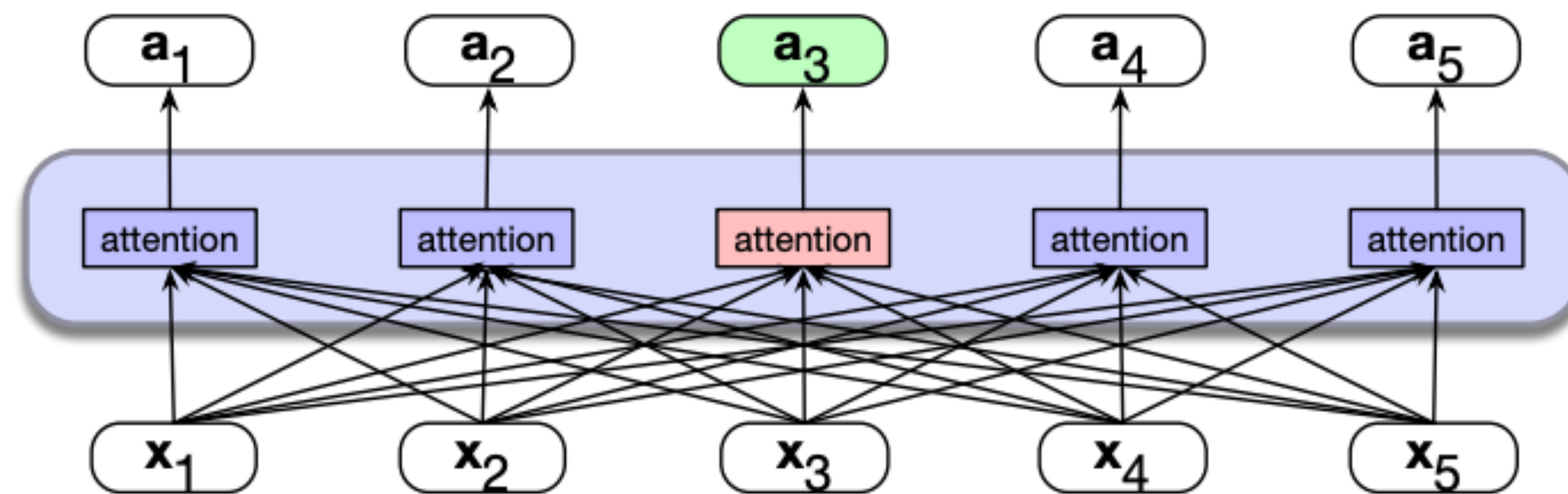
# Encoder-only Architecture

- Classification?

  - Sequence-level classification

# Encoder-only architectures

- Can we use the same next-token prediction task to train encoder models?

  No! The desired output is part of the input!

# Slide Acknowledgements

▸ Earlier versions of this course offerings including materials from Claire Cardie, Marten van Schijndel, Lillian Lee.