# Sequence-to-Sequence Models & Transformers
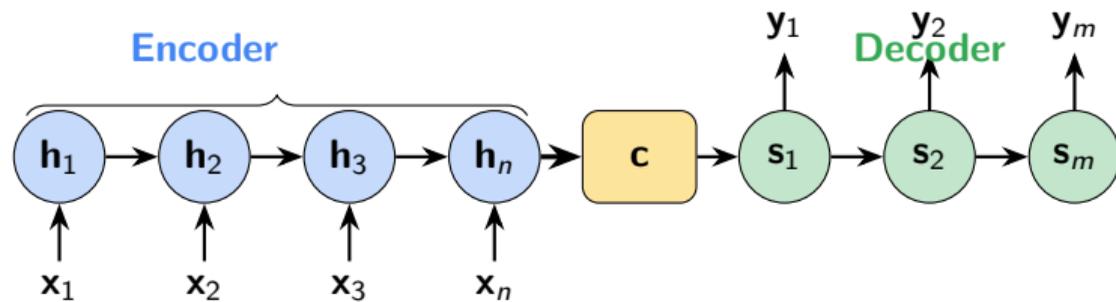
COS 484 Precept

# The Seq2Seq Problem

| Input: "How are you?" | Variable length $n$ |

**Model**

| Output: "Comment allez-vous?" | Variable length $m$ |

**Challenge:** Input and output have **different lengths**

# Encoder-Decoder Architecture



**Encoder:** $\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)$

**Decoder:** $\mathbf{s}_t = g(\mathbf{s}_{t-1}, \mathbf{y}_{t-1}, \mathbf{c})$
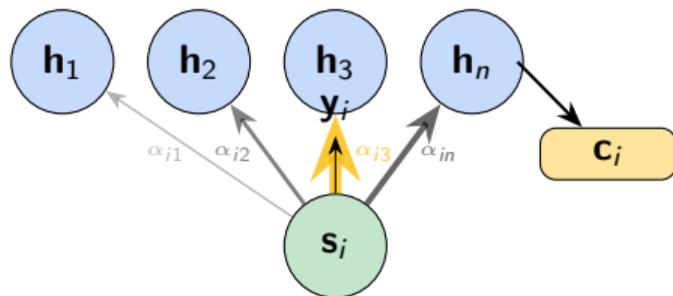
# The Bottleneck Problem



All info compressed here!

## Problem

- Fixed-size vector must encode **entire** input sequence
- Early information gets "forgotten"
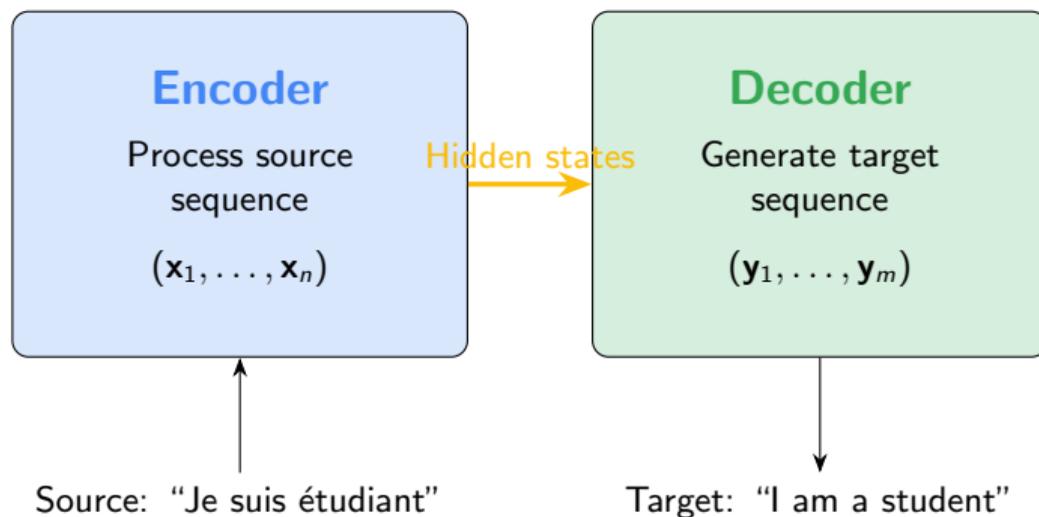- Performance degrades for long sequences

**Intuition:** Each decoder step "looks back" at encoder states

**Dynamic context** $\to$ Different focus for each output
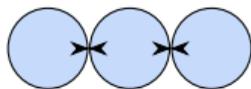
# What Does Each Component Do?

## Encoder

- Reads the **entire** source sequence
- Creates **contextualized representations**
- Output: Hidden states $\mathbf{h}_1, \ldots, \mathbf{h}_n \in \mathbb{R}^{n \times d}$
- Bidirectional: can see past & future

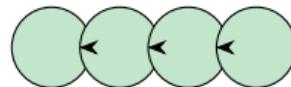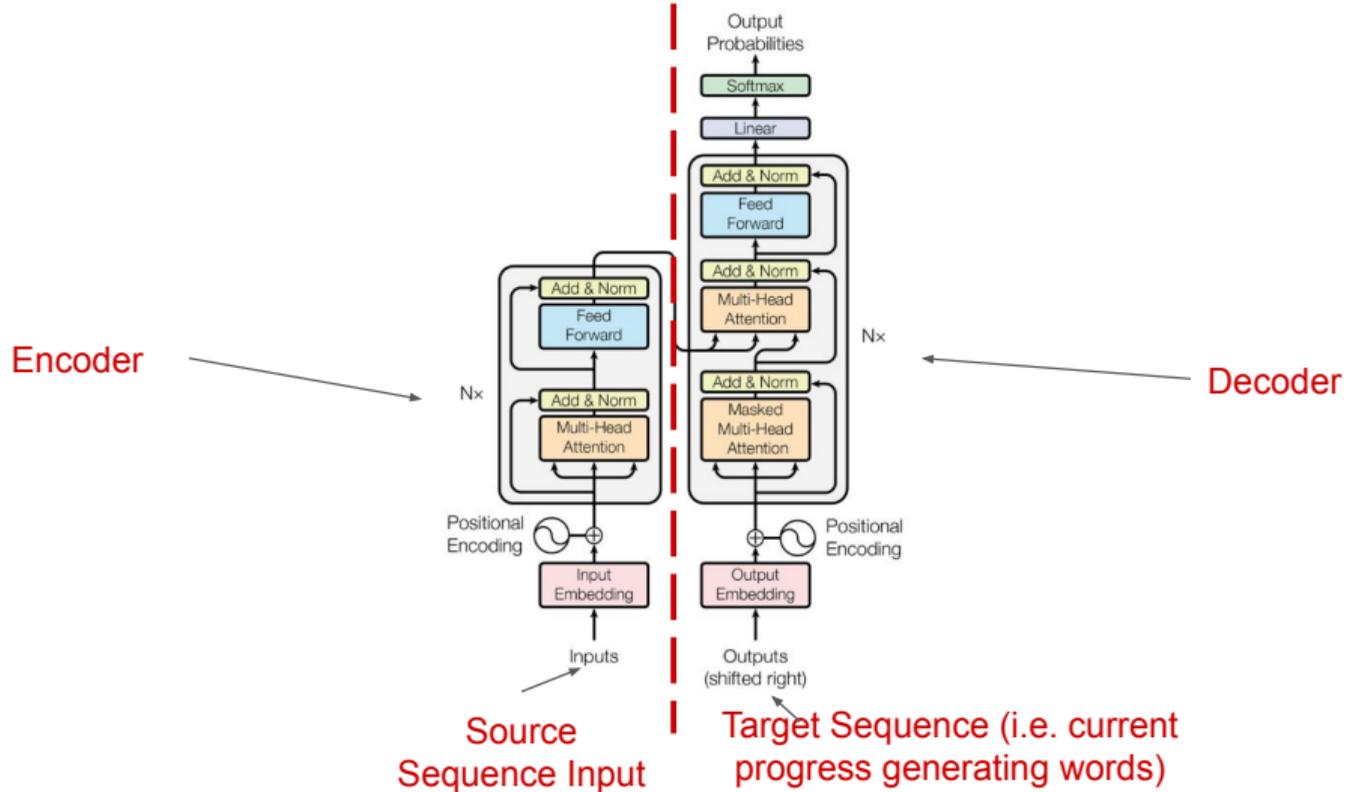Je   suis étudiant

## Decoder

- Generates output **one token at a time**
- Conditions on encoder output
- Output: $P(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})$
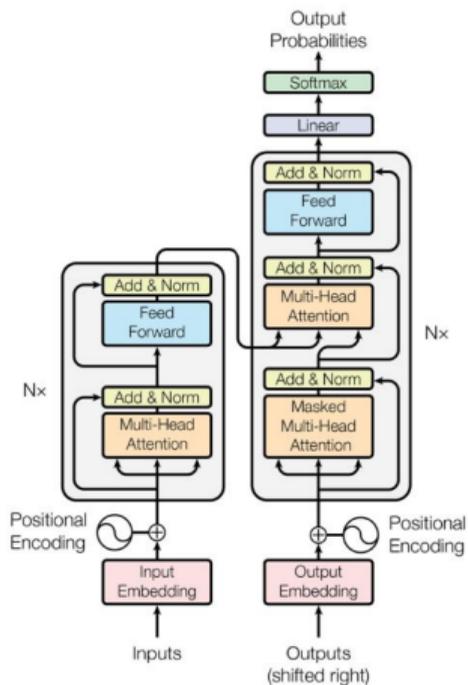- Autoregressive: only sees past

I   am   a student

# Transformer Architecture

# Transformer Encoder



Source
sequence
$(x_1, \ldots, x_n)$

# Transformer Encoder: Positional + Word Embedding

**Input and Positional Embedding**

| | | |
|---|---|---|
| EMBEDDING WITH TIME SIGNAL | $x_1$ ☐☐☐☐ | $x_2$ ☐☐☐☐ | $x_3$ ☐☐☐☐ |
| | = | = | = |
| POSITIONAL ENCODING | $t_1$ ☐☐☐☐ | $t_2$ ☐☐☐☐ | $t_3$ ☐☐☐☐ |
| | + | + | + |
| EMBEDDINGS | $x_1$ ☐☐☐☐ | $x_2$ ☐☐☐☐ | $x_3$ ☐☐☐☐ |
| INPUT | Je | suis | étudiant |

Embedded source sequence
$\mathbb{R}^{n \times d_1}$

Positional Encoding

Input Embedding

Inputs

Source sequence
$(x_1, \ldots, x_n)$

Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention

Nx

# Transformer Encoder: Multi-Head Self Attention

**Self-Attention:** $W_i^Q \in \mathbb{R}^{d_1 \times d_q}, W_i^K \in \mathbb{R}^{d_1 \times d_k}, W_i^V \in \mathbb{R}^{d_1 \times d_v}$

**Step 1:**

**Step 2:**

After Multi-Head Attention
$\mathbb{R}^{n \times d_2}$

Embedded source sequence
$\mathbb{R}^{n \times d_1}$

Positional Encoding

Input Embedding

Inputs

Source sequence
$(x_1, ..., x_n)$



**MultiHead Attention:** $W^O \in \mathbb{R}^{d \times d_2}$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{head}_i = \text{Attention}(XW_i^Q, XW_i^K, XW_i^V)$$

# Transformer Encoder: Multi-Head Self Attention

**Self-Attention:** $W_i^Q \in \mathbb{R}^{d_1 \times d_q}, W_i^K \in \mathbb{R}^{d_1 \times d_k}, W_i^V \in \mathbb{R}^{d_1 \times d_v}$

**Step 1:**

**Step 2:**



In practice, $d_1 = d_2$

After Multi-Head Attention
$\mathbb{R}^{n \times d_1}$

Embedded source sequence
$\mathbb{R}^{n \times d_1}$

"Self" attention means Q, K, V are all computed from a single sequence

Source sequence
$(x_1, ..., x_n)$

**MultiHead Attention:** $W^O \in \mathbb{R}^{d \times d_2}$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
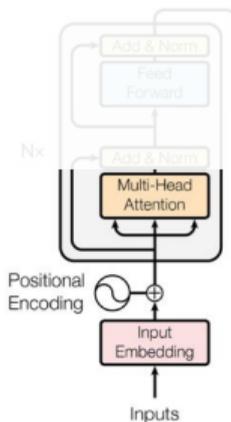$$\text{head}_i = \text{Attention}(XW_i^Q, XW_i^K, XW_i^V)$$

# Transformer Encoder: Add & Norm

**Add & Norm:**

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

After Add & Norm
$\mathbb{R}^{n \times d_1}$

After Multi-Head Attention
$\mathbb{R}^{n \times d_1}$

Embedded source sequence
$\mathbb{R}^{n \times d_1}$

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Source sequence
$(x_1, \ldots, x_n)$

**LayerNorm**

$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

# Transformer Encoder: Feed Forward

**Feed Forward**

$$\text{FFN}(\mathbf{x}_i) = \text{ReLU}(\mathbf{x}_i\mathbf{W_1} + \mathbf{b_1})\mathbf{W_2} + \mathbf{b_2}$$

$$\mathbf{W_1} \in \mathbb{R}^{d \times d_{ff}}, \mathbf{b_1} \in \mathbb{R}^{d_{ff}}$$

$$\mathbf{W_2} \in \mathbb{R}^{d_{ff} \times d}, \mathbf{b_2} \in \mathbb{R}^{d}$$

Compute transformation over each value in the sequence **independently**
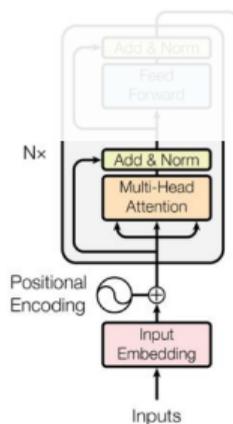
After Feed Forward
$\mathbb{R}^{n \times d_1}$

After Add & Norm
$\mathbb{R}^{n \times d_1}$

After Multi-Head Attention
$\mathbb{R}^{n \times d_1}$

Embedded source sequence
$\mathbb{R}^{n \times d_1}$

N×

Feed Forward

Add & Norm
Multi-Head Attention

Positional Encoding ⊕

Input Embedding

Inputs

Source sequence
$(x_1, \ldots, x_n)$

# Transformer Encoder: Final Add & Norm

After Final Add & Norm
$$\mathbb{R}^{n \times d_1}$$

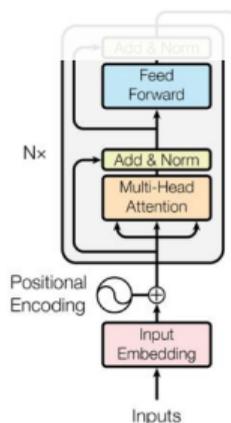After Feed Forward
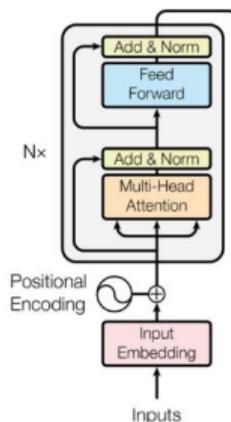$$\mathbb{R}^{n \times d_1}$$

After Add & Norm
$$\mathbb{R}^{n \times d_1}$$

After Multi-Head Attention
$$\mathbb{R}^{n \times d_1}$$

Embedded source sequence
$$\mathbb{R}^{n \times d_1}$$



N×

Add & Norm
Feed Forward

Add & Norm
Multi-Head Attention

Positional Encoding

Input Embedding
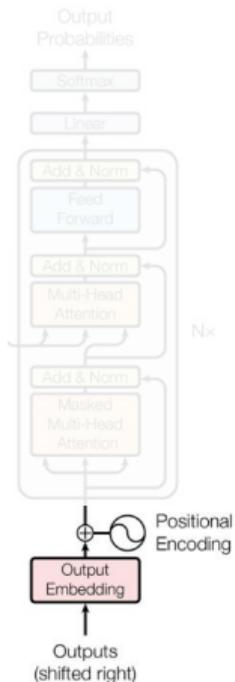
Inputs

Source sequence
$(x_1, \ldots, x_n)$

**Add & Norm:**

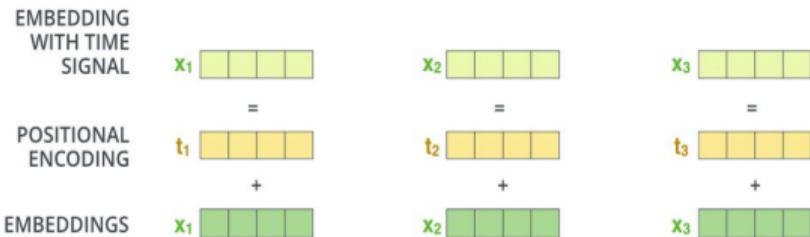$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

**LayerNorm**

$$y = \frac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + \beta$$

# Transformer Decoder:

**Output and Positional Embedding**

EMBEDDING WITH TIME SIGNAL $x_1$ $x_2$ $x_3$

=

POSITIONAL ENCODING $t_1$ $t_2$ $t_3$

+

EMBEDDINGS $x_1$ $x_2$ $x_3$

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

⊕ → Positional Encoding

Output Embedding

Outputs (shifted right)

Embedded target sequence
$\mathbb{R}^{m \times d_1}$

Target sequence
(<bos>, $x_1$, ..., $x_m$)

# Transformer Decoder: Masked Multi-Head Attention

**Masked Self-Attention:** $W_i^Q \in \mathbb{R}^{d_1 \times d_q}, W_i^K \in \mathbb{R}^{d_1 \times d_k}, W_i^V \in \mathbb{R}^{d_1 \times d_v}$



**Step 1:**

X × W<sup>Q</sup> = Q

X × W<sup>K</sup> = K

X × W<sup>V</sup> = V

**Step 2:**

$$\frac{Q \times K^T}{\sqrt{d_k}}$$

Elementwise Multiply by Mask
(equivalent to setting masked indices to -∞)

$\odot$ =

**Step 3:**

softmax( ) V

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

Embedded target sequence
$\mathbb{R}^{m \times d_1}$

Output
Probabilities
Softmax
Linear
Add & Norm
Feed
Forward
Add & Norm
Multi-Head
Attention
N×
Add & Norm
Masked
Multi-Head
Attention
Positional
Encoding
Output
Embedding
Outputs
(shifted right)

Target sequence
(<bos>, $x_1$, ..., $x_m$)

**MultiHead Attention:** $W^O \in \mathbb{R}^{d \times d_2}$

$\mathrm{MultiHead}(Q, K, V) = \mathrm{Concat}(\mathrm{head}_1, ..., \mathrm{head}_h)W^O$

$\mathrm{head}_i = \mathrm{Attention}(XW_i^Q, XW_i^K, XW_i^V)$

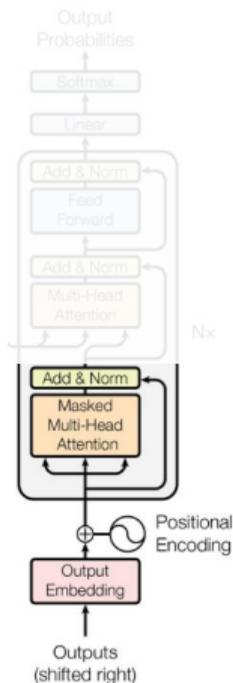# Transformer Decoder:



After Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

**Embedded target sequence**
$\mathbb{R}^{m \times d_1}$

Target sequence
($<bos>$, $x_1$, ..., $x_m$)

**Add & Norm:**

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

**LayerNorm**

$$y = \frac{x - \text{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

# Transformer Decoder: Multi-Head (Cross) Attention



Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

After Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

Embedded target sequence
$\mathbb{R}^{m \times d_1}$

Target sequence
($<bos>$, $x_1$, ..., $x_m$)

**Cross-Attention:** $W_i^Q \in \mathbb{R}^{d_1 \times d_q}, W_i^K \in \mathbb{R}^{d_1 \times d_k}, W_i^V \in \mathbb{R}^{d_1 \times d_v}$

**Step 1:**

X × $W^Q$ = Q

X × $W^K$ = K

X × $W^V$ = V

**Step 2:**

$$\text{softmax}\left( \frac{Q \times K^T}{\sqrt{d_k}} \right) \quad V$$

"Cross" attention means
Q, K, V are computed
from **separate** sequences

**MultiHead Attention:** $W^O \in \mathbb{R}^{d \times d_2}$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

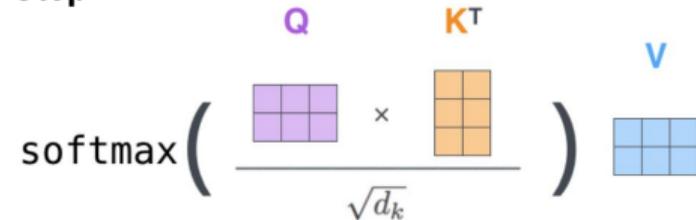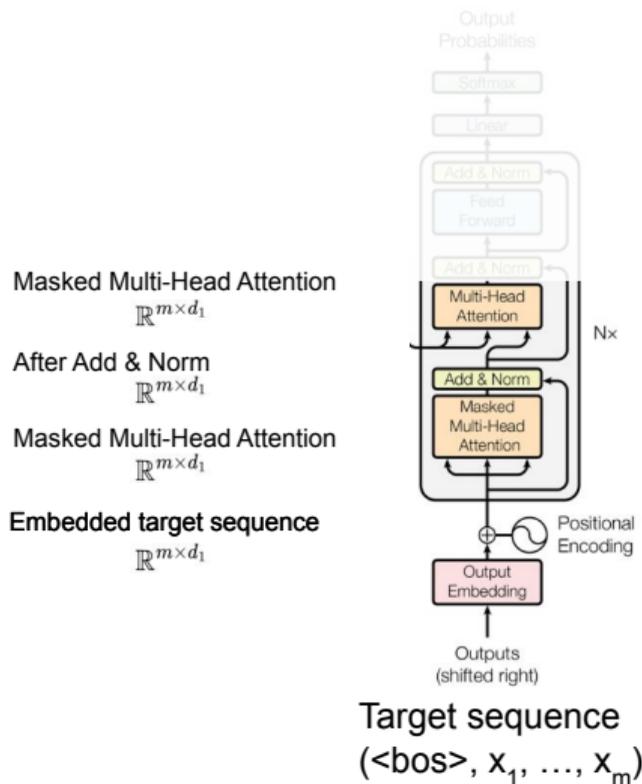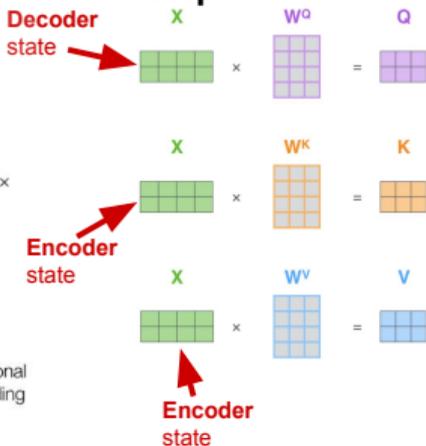$$\text{head}_i = \text{Attention}(XW_i^Q, XW_i^K, XW_i^V)$$

# Transformer Decoder: Multi-Head (Cross) Attention

**Cross-Attention:** $W_i^Q \in \mathbb{R}^{d_1 \times d_q}, W_i^K \in \mathbb{R}^{d_1 \times d_k}, W_i^V \in \mathbb{R}^{d_1 \times d_v}$
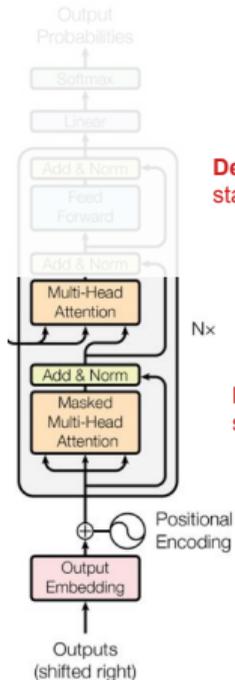
**Step 1:**

**Step 2:**



Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

After Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

Embedded target sequence
$\mathbb{R}^{m \times d_1}$

Target sequence
(<bos>, $x_1$, ..., $x_m$)

"Cross" attention means
Q, K, V are computed
from **separate** sequences

**MultiHead Attention:** $W^O \in \mathbb{R}^{d \times d_2}$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h) W^O$$

$$\text{head}_i = \text{Attention}(XW_i^Q, XW_i^K, XW_i^V)$$

# Transformer Decoder: Add & Norm

Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

After Add & Norm
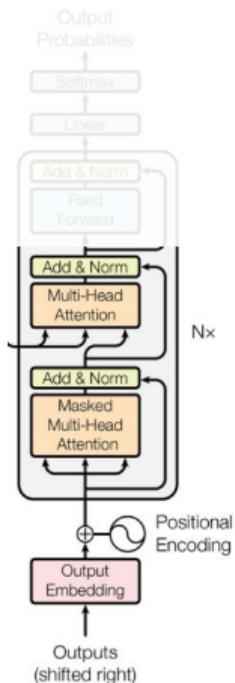$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

**Embedded target sequence**
$\mathbb{R}^{m \times d_1}$

Target sequence
(<bos>, $x_1$, …, $x_m$)

**Add & Norm:**

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

**LayerNorm**

$$y = \frac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + \beta$$

# Transformer Decoder: Feed Forward

Feed Forward
$\mathbb{R}^{m \times d_1}$

Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

After Add & Norm
$\mathbb{R}^{m \times d_1}$
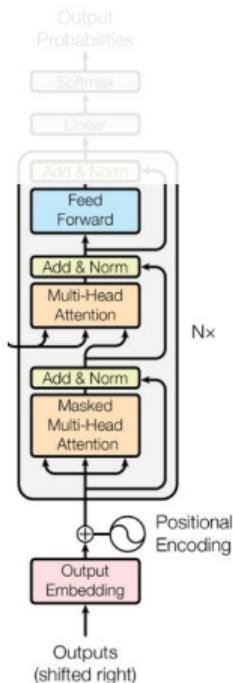
Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

**Embedded target sequence**
$\mathbb{R}^{m \times d_1}$

Target sequence
(<bos>, $x_1$, …, $x_m$)

**Feed Forward**

$$\mathrm{FFN}(\mathbf{x}_i) = \mathrm{ReLU}(\mathbf{x}_i \mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$$

$$\mathbf{W}_1 \in \mathbb{R}^{d \times d_{ff}}, \mathbf{b}_1 \in \mathbb{R}^{d_{ff}}$$

$$\mathbf{W}_2 \in \mathbb{R}^{d_{ff} \times d}, \mathbf{b}_2 \in \mathbb{R}^d$$

# Transformer Decoder: Add & Norm

Add & Norm
$\mathbb{R}^{m \times d_1}$

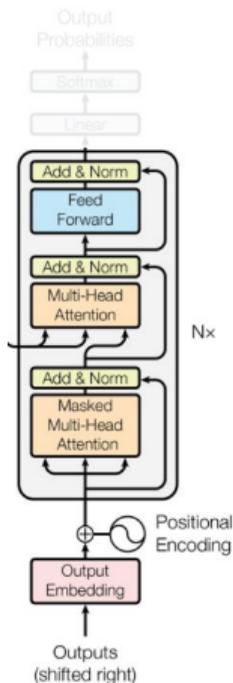Feed Forward
$\mathbb{R}^{m \times d_1}$

Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

After Add & Norm
$\mathbb{R}^{m \times d_1}$

Masked Multi-Head Attention
$\mathbb{R}^{m \times d_1}$

**Embedded target sequence**
$\mathbb{R}^{m \times d_1}$



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Output Embedding

Outputs (shifted right)

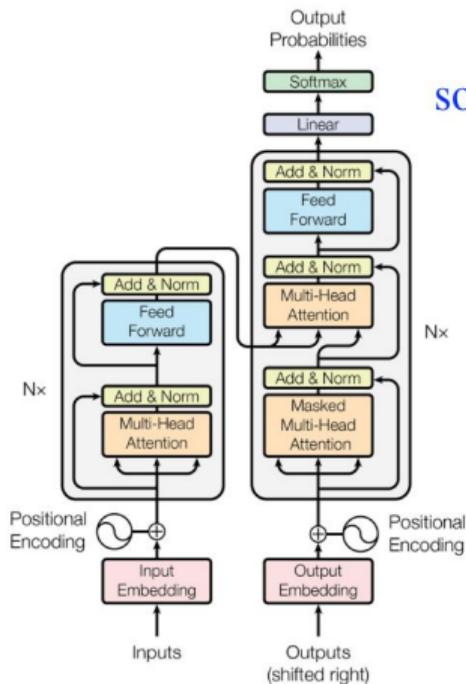Target sequence
(<bos>, $x_1$, …, $x_m$)

**Add & Norm:**

$\text{LayerNorm}(x + \text{Sublayer}(x))$

**LayerNorm**

$$y = \frac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + \beta$$

# Transformer: Final output



$$\text{softmax}(\mathbf{W}_o \mathbf{h}_i)$$

Compute transformation over **concatenated states**