Tuesday • October 26, 2021

# The Transformer: Attention-Only Networks
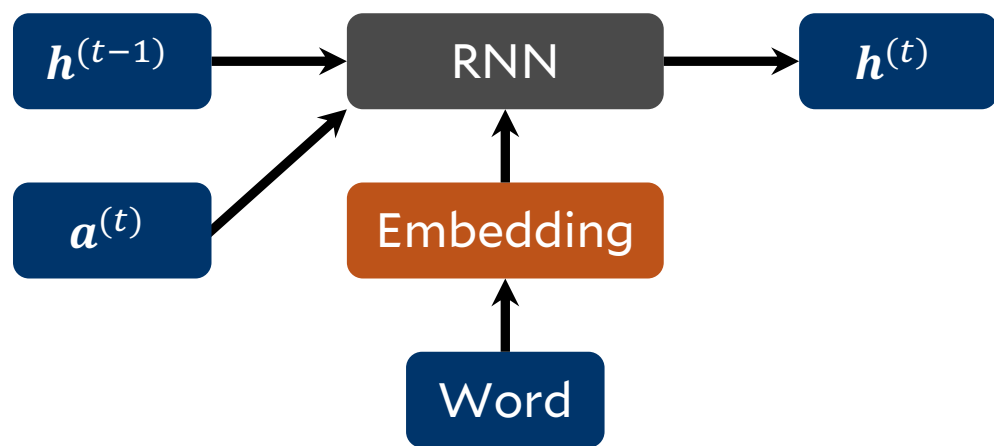
Yale

LING 380/780
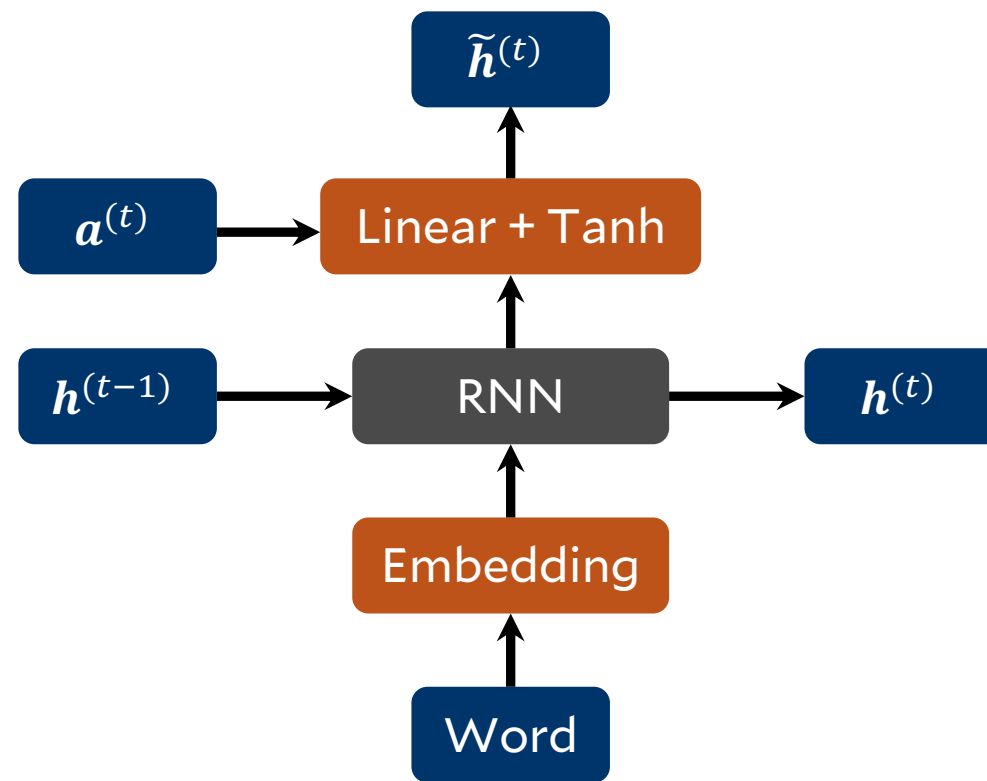
*Neural Network Models of Linguistic Structure*

# RNN Decoder with Attention

Bahdanau-Style:

Luong-Style:

# Computing Attention

First, compute **attention scores**:

Query  Key

$$s_i^{(t)} = \text{score}(\boldsymbol{h}^{(d,T)}, \boldsymbol{h}^{(e,i)})$$

Then, convert the attention scores into **attention weights**:

$$\boldsymbol{\alpha}^{(t)} = \text{softmax}(\boldsymbol{s}^{(t)})$$

Finally, take a weighted average of encoder hidden states:

$$\boldsymbol{a}^{(t)} = \sum_i \alpha_i^{(t)} \, \boldsymbol{h}^{(e,i)} \longleftarrow \text{Value}$$

# Computing Attention

Query  Key

First, compute **attention scores**:

$$s_i^{(t)} = \text{score}(\boldsymbol{q}^{(t)}, \boldsymbol{k}^{(i)})$$

Then, convert the attention scores into **attention weights**:

$$\boldsymbol{\alpha}^{(t)} = \text{softmax}(\boldsymbol{s}^{(t)})$$

Finally, take a weighted average of encoder hidden states:

$$\boldsymbol{a}^{(t)} = \sum_i \alpha_i^{(t)} \boldsymbol{v}^{(i)} \quad \longleftarrow \text{Value}$$

# Intuition behind Attention

- **Attention:** "looking up" previous information

- **Query** ($q^{(t)}$)**:** the thing that we want to look up

- **Keys** ($k^{(i)}$)**:** things that could be looked up

- **Values** ($v^{(i)}$)**:** the result of the lookup

- Score ($s_i^{(t)}$): How well does query $q^{(t)}$ match key $k^{(i)}$?

# Computing Attention Scores

Computing attention scores

- Dot Product: $\text{score}(\boldsymbol{q}, \boldsymbol{k}) = \boldsymbol{q}^\top \boldsymbol{k}$

- Bilinear: $\text{score}(\boldsymbol{q}, \boldsymbol{k}) = \boldsymbol{q}^\top \boldsymbol{W} \boldsymbol{k}$

- MLP: $\text{score}(\boldsymbol{q}, \boldsymbol{k}) = \text{MLP}\left(\begin{bmatrix} \boldsymbol{q} \\ \boldsymbol{k} \end{bmatrix}\right)$

# Evaluating Machine Translation

- How do RNNs with attention compare to RNNs?

- What does it mean to translate well?

# Evaluating Machine Translation

- **Example Sentence:** 猫はマットにいる。

- **Reference Translation:** The cat is on the mat.

- **Possible Machine Translations:**

  - The cat is on the mat.

  - The cats are on the mat.

  - The mat is where the cat is.

  - The the the the the the.

# Evaluating Machine Translation

n-Gram Precision:

$$\frac{\text{\# of n−grams in the machine and reference translation}}{\text{\# of n−grams in the machine translation}}$$

# Evaluating Machine Translation

- **Reference Translation:** The cat is on the mat.

- **Machine Translation:** The the mat.

- 1-gram precision: 100% (✔ The, ✔ the, ✔ mat, ✔ .)

- 2-gram precision: 67% (✘ The the, ✔ the mat, ✔ mat .)

- 3-gram precision: 50% (✘ The the mat, ✔ the mat .)

- 4-gram precision: 0% (✘ The the mat .)

# Evaluating Machine Translation

Bilingual Evaluation Understudy (BLEU) Score

$$\text{BLEU} = \text{length penalty} \cdot \left( \prod_{n=1}^{4} n\text{-gram precision} \right)^{1/4}$$
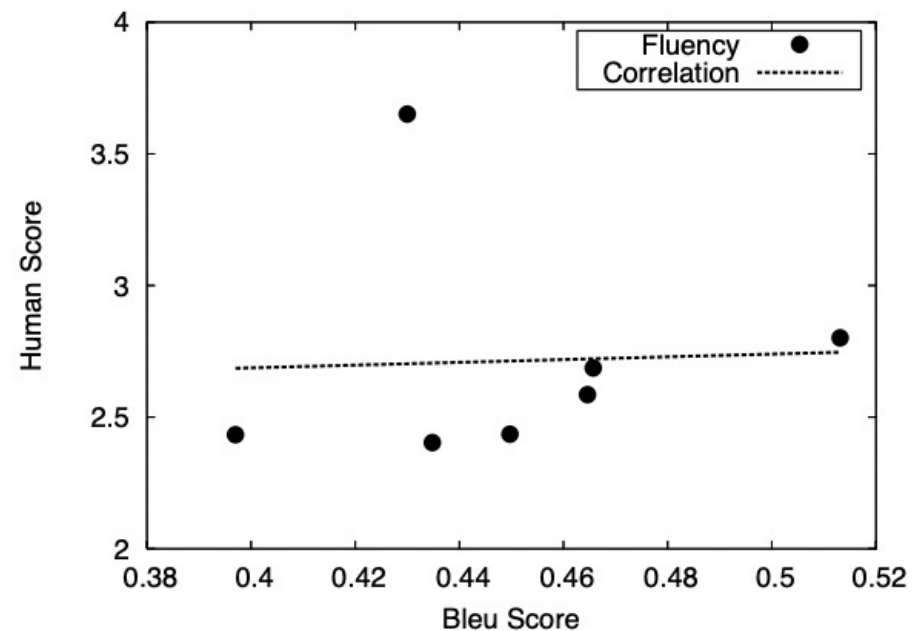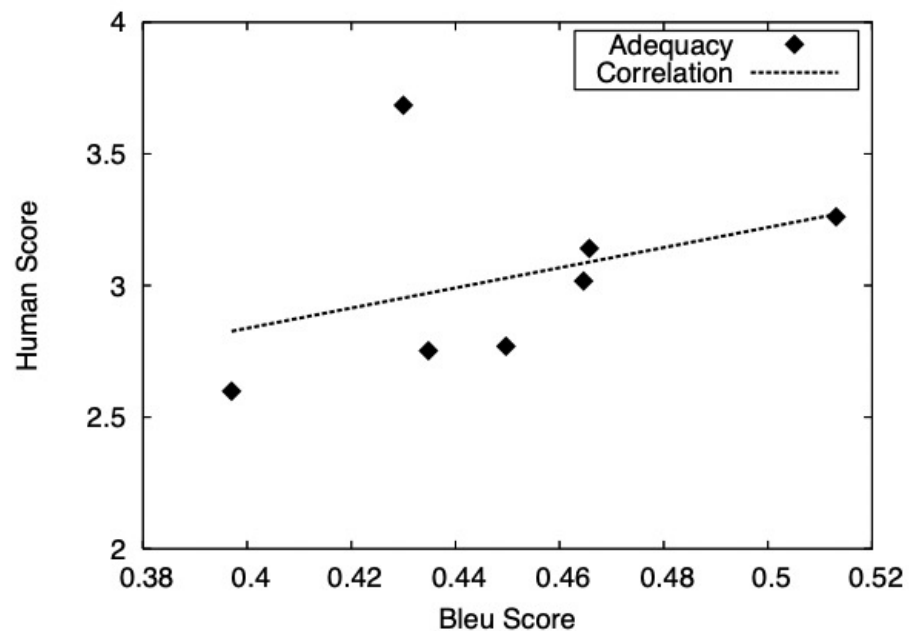
where

$$\text{length penalty} = \min(1, e^{1 - \frac{\text{length of reference translation}}{\text{length of machine translation}}})$$

# Evaluating Machine Translation

- **Example Sentence:** 猫はマットにある。

- **Reference Translation:** The cat is on the mat.

- **Possible Machine Translations:**
  - The cat is on the mat. **BLEU:** 100.0%
  - The cats are on the mat. **BLEU:** 43.5%
  - The mat is where the cat is. **BLEU:** 0.0%
  - The the the the the the. **BLEU:** 0.0%

# Is BLEU Reasonable?

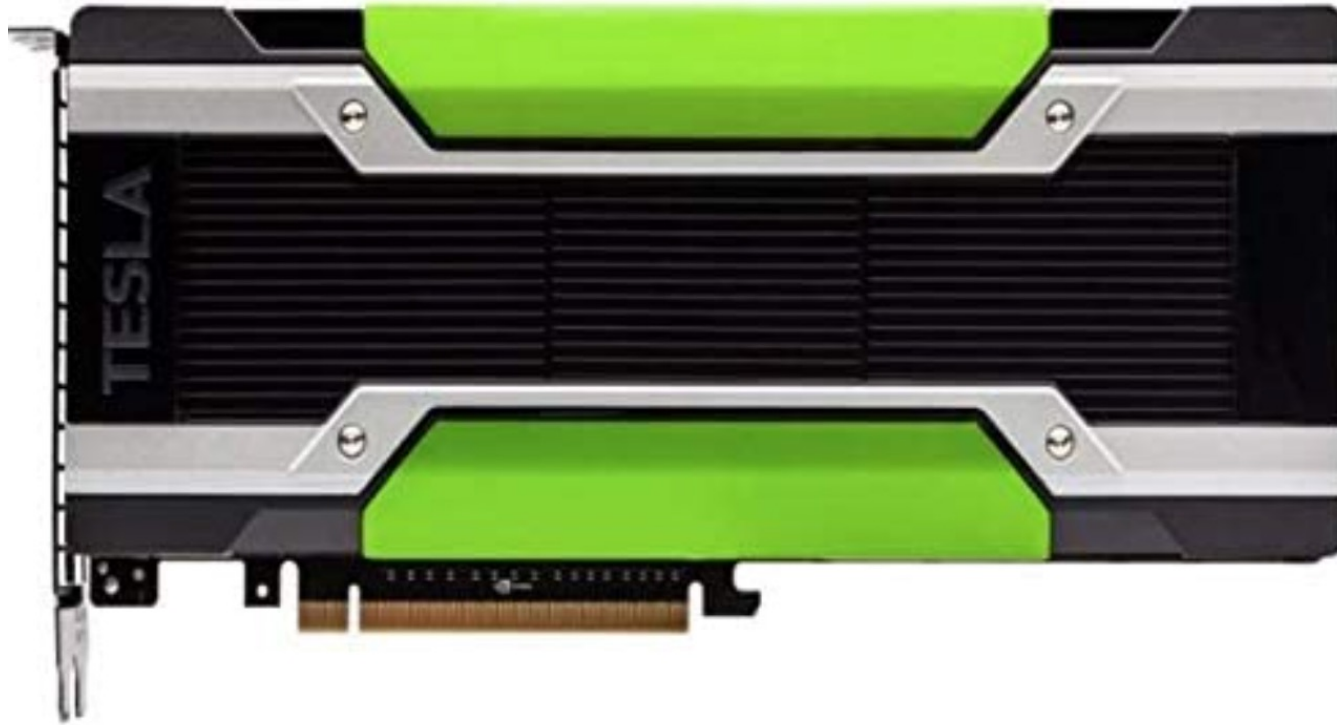BLEU vs. Human Judgements (Callison-Burch et al., 2006)

# EN − FR Machine Translation Performance

| Model | BLEU (%) |
|---|---|
| GRU, 50 Hidden (Bahdanau et al., 2015) | 17.82 |
| GRU, 50 Hidden + Bahdanau Attn. (ibid) | 28.45 |
| GRU, 1000 Hidden (Cho et al., 2014) | 33.87 |
| Google NMT: LSTM, 8 Layers, 1024 Hidden + Bahdanau Attn. (Wu et al., 2016) | 39.92 |

# Cost of Model Training

- The Google NMT model requires **9 days** of training on **96 NVIDIA K80 GPUs**.

- Yale's Grace high performance cluster only has 12 NVIDIA K80 GPUs.

# Cost of Model Training



## Nvidia Tesla K80 24GB GDDR5 CUDA Cores Graphic Cards

Visit the NVIDIA Store

★★★½☆ ∨    48 ratings | 83 answered questions

🏷 8 Price Changes

Price: **$397.99** & FREE Returns

S   Best price   S +

**Pay $66.33/month for 6 months** (plus S&H, tax) with 0% interest equal monthly payments when you're approved for an Amazon Prime Store Card.

May be available at a lower price from other sellers, potentially without free Prime shipping.

| | |
|---|---|
| **Brand** | NVIDIA |
| **Graphics Coprocessor** | Nvidia Tesla |
| **Video Output Interface** | VGA |
| **Graphics Processor** | NVIDIA |

# Parallel Computation

- Graphics processing units (GPUs) are devices that perform **parallel computation**.

- Parallel algorithms **perform multiple computations at the same time**.

# Example: Sequential Matrix Multiplication

Suppose $A, B \in \mathbb{R}^{4 \times 4}$. How do you compute $AB$?

- Initialize $C \leftarrow 0$.

- For $i, j \in \{1, 2, 3, 4\}$: *16 repetitions*

  - Set $C_{i,j} \leftarrow A_{i,1}B_{1,j} + A_{i,2}B_{2,j} + A_{i,3}B_{3,j} + A_{i,4}B_{4,j}$. *7 operations*

- Return $C$.

Total time: 112 FLOPs (floating-point operations)

# Example: Parallel Matrix Multiplication

- We can speed up matrix multiplication by doing it in parallel.

- Divide the matrices into 4 blocks:

$$AB = \begin{bmatrix} A_{:2,:2}B_{:2,:2} + A_{:2,3:}B_{3:,:2} & A_{:2,:2}B_{:2,3:} + A_{:2,3:}B_{3:,3:} \\ A_{3:,:2}B_{:2,:2} + A_{3:,3:}B_{3:,:2} & A_{3:,:2}B_{:2,3:} + A_{3:,3:}B_{3:,3:} \end{bmatrix}$$

# Example: Parallel Matrix Multiplication

- Initialize $C \leftarrow 0, D \leftarrow 0$.

- Fork the process into 8 threads: *12 operations each*

  - $C_{:2,:2} \leftarrow A_{:2,:2}B_{:2,:2}$     $C_{:2,3:} \leftarrow A_{:2,:2}B_{:2,3:}$     $C_{3:,2:} \leftarrow A_{3:,:2}B_{:2,:2}$,
  - $C_{3:,3:} \leftarrow A_{3:,:2}B_{:2,3:}$     $D_{:2,:2} \leftarrow A_{:2,3:}B_{3:,:2}$     $D_{:2,3:} \leftarrow A_{:2,3:}B_{3:,3:}$,
  - $D_{3:,2:} \leftarrow A_{3:,3:}B_{3:,:2}$     $D_{3:,3:} \leftarrow A_{3:,3:}B_{3:,3:}$

- Return $C + D$. *16 operations*

Total time: 28 FLOPs

# Neural Networks and Parallel Computation

- Feedforward (i.e., non-recurrent) networks like MLPs can be **trained very quickly** using parallel computation with GPUs.

- RNNs cannot, since inputs must be processed one at a time.

- **Solution:** Create a **feedforward** architecture with **attention**!

# Self Attention

The **self-attention head** is an attention layer that pays attention **to its own input**.

- Input: $x^{(1)}, x^{(2)}, \ldots, x^{(n)}$

- **Queries:** $q^{(t)} = \text{Linear}(x^{(t)})$

- **Keys:** $k^{(i)} = \text{Linear}(x^{(i)})$

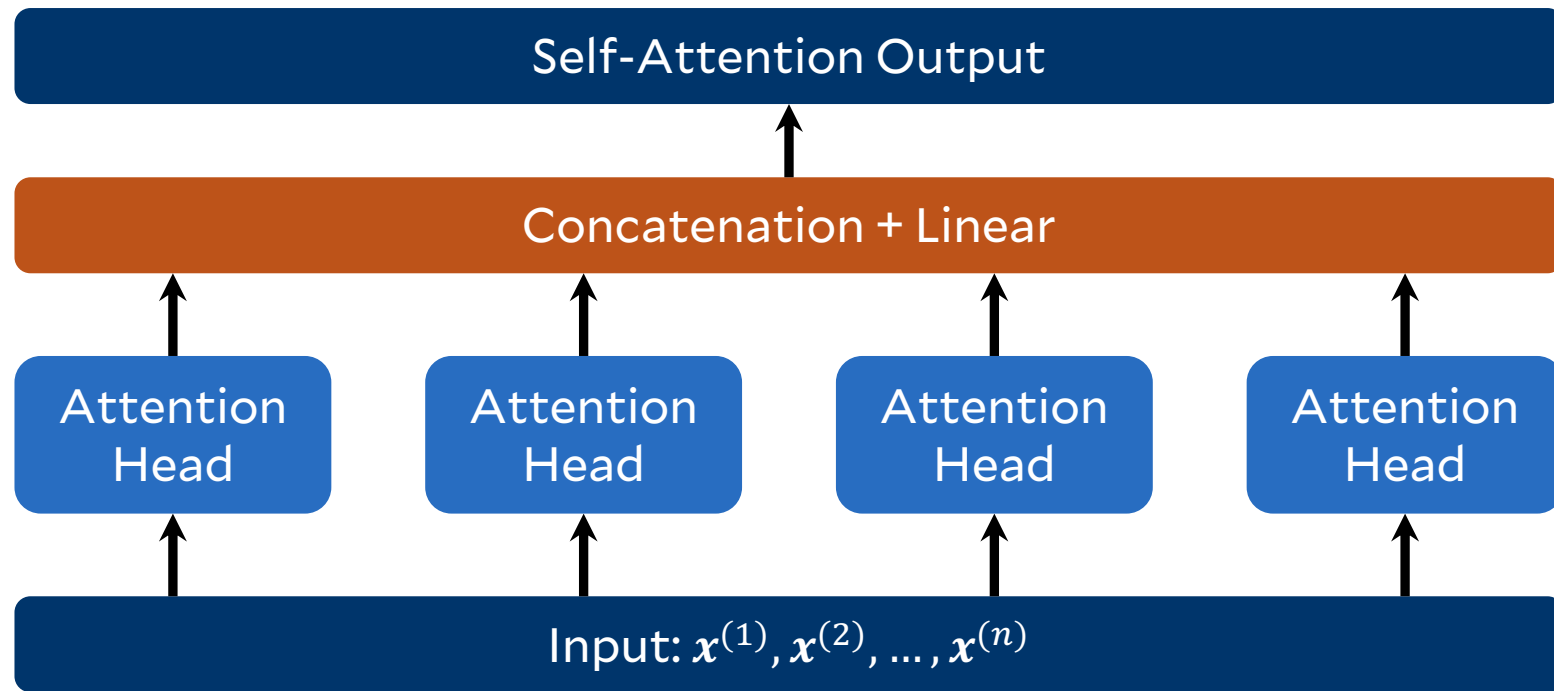- **Values:** $v^{(i)} = \text{Linear}(x^{(i)})$

# Self Attention

Attention scores are computed using **scaled dot-product attention**:

$$\text{score}(\boldsymbol{q}, \boldsymbol{k}) = \frac{\boldsymbol{q}^\top \boldsymbol{k}}{\sqrt{d}}$$

where $\boldsymbol{q}, \boldsymbol{k} \in \mathbb{R}^d$.

# Attention Heads

**Multi-head attention** combines several self-attention heads.
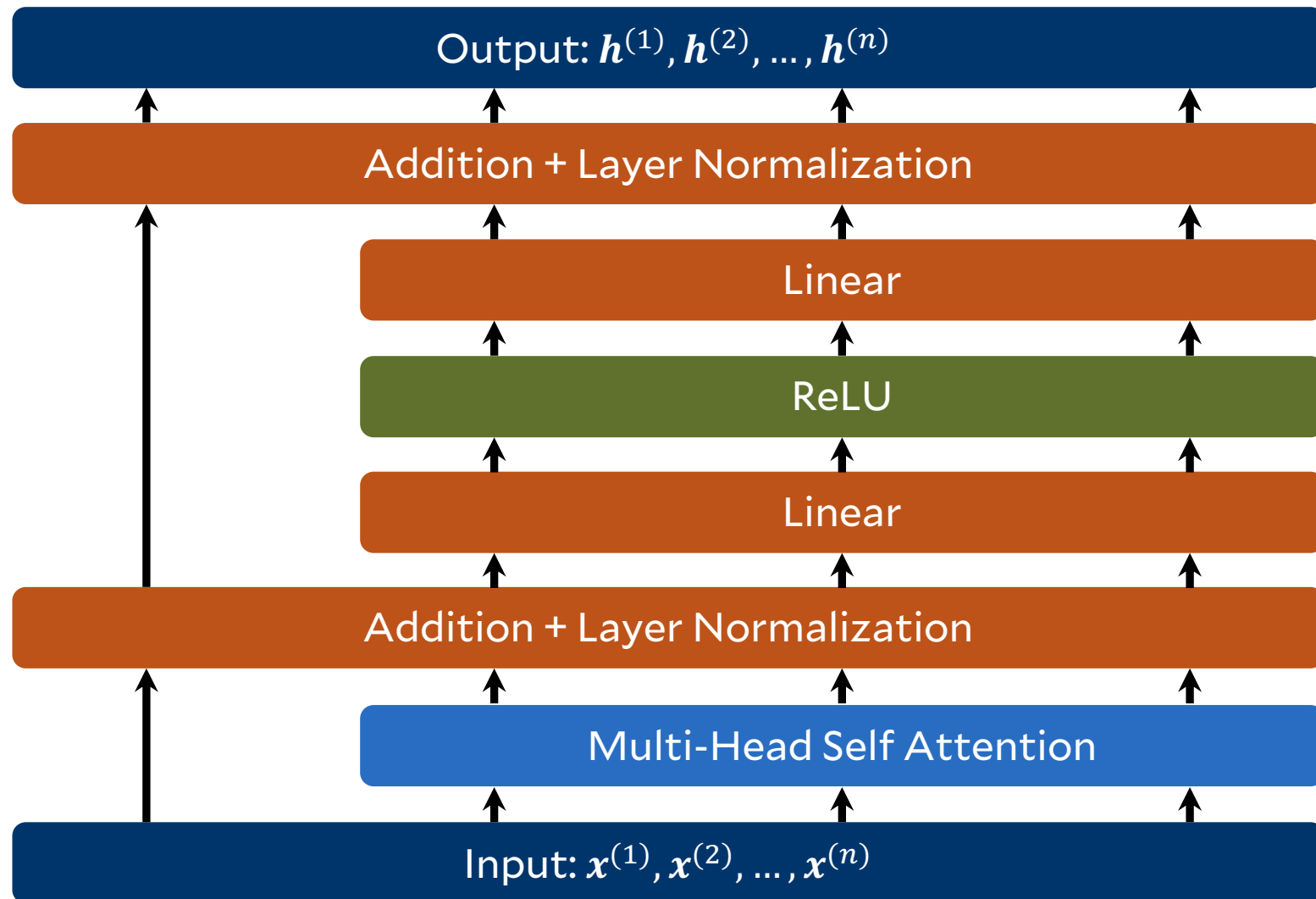
# The Transformer Architecture

The **Transformer** is a self-attention-based architecture made up of **encoder blocks**.

Each encoder contains:

- Several attention heads

- **Layer normalization** layers

$$\text{LayerNorm}(\boldsymbol{x}) = \gamma \frac{\boldsymbol{x} - \text{mean}(\boldsymbol{x})}{\sqrt{\text{var}(\boldsymbol{x}) + \varepsilon}} + \beta$$

# Encoder Block



Output: $\boldsymbol{h}^{(1)}, \boldsymbol{h}^{(2)}, \ldots, \boldsymbol{h}^{(n)}$

Addition + Layer Normalization

Linear

ReLU

Linear

Addition + Layer Normalization

Multi-Head Self Attention

Input: $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(n)}$

# Transformer Encoder

The **Transformer encoder** combines word embeddings with a **positional encoding** before passing them through several encoder blocks.

# Transformer Encoder

Positional Encoding: Given sequence $w_1, w_2, \ldots, w_n$, the representation of $w_i$ is

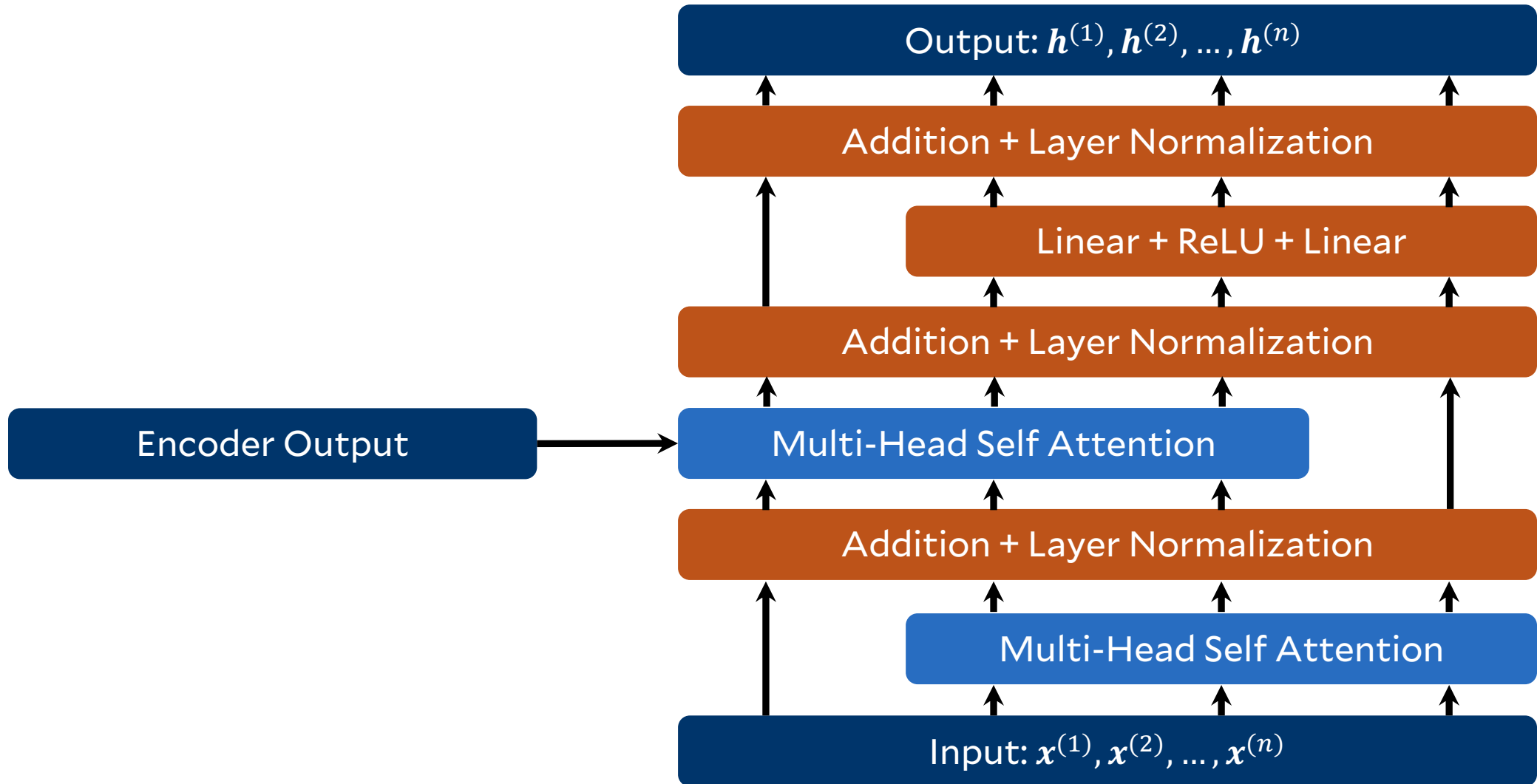$$\text{Embedding}(w_i) + \boldsymbol{p}^{(i)}$$

where the embedding and hidden size is $h$ and

$$p_j^{(i)} = \begin{cases} \sin\left(i \cdot 10000^{-2j/h}\right), j \text{ is even} \\ \cos\left(i \cdot 10000^{-2j/h}\right), j \text{ is odd} \end{cases}$$

# Decoder Block

The Transformer also has a **decoder block** for auto-regressive decoding.
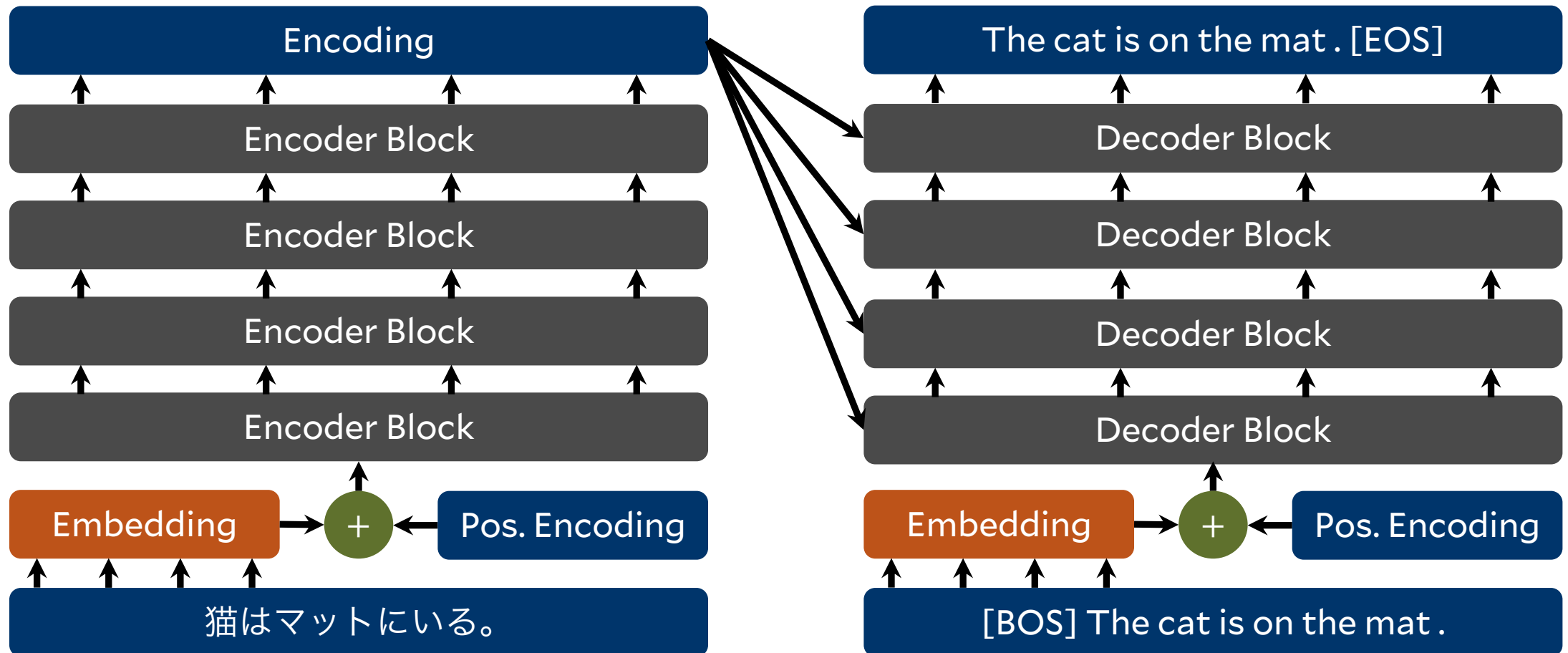
# Decoder Block

# Decoder Block

- The first self-attention layer is **masked**: no position may attend to a future position. If $i \geq t$, then

$$s_i^{(t)} = -\infty$$

- In the second self-attention layer, keys and values are obtained from the encoder outputs $\boldsymbol{o}^{(i)}$ :

  - **Keys:** $\boldsymbol{k}^{(i)} = \text{Linear}(\boldsymbol{o}^{(i)})$

  - **Values:** $\boldsymbol{v}^{(i)} = \text{Linear}(\boldsymbol{o}^{(i)})$

# Transformer Encoder－Decoder

# Teacher Forcing

- During training, the entire reference translation is provided to the decoder. This is called **teacher forcing**.

- Why is this okay?

- The decoder makes all predictions simultaneously, making training more parallel.

- The decoder is never trained to make predictions based on incorrectly predicted previous words.

# EN − FR Machine Translation Performance
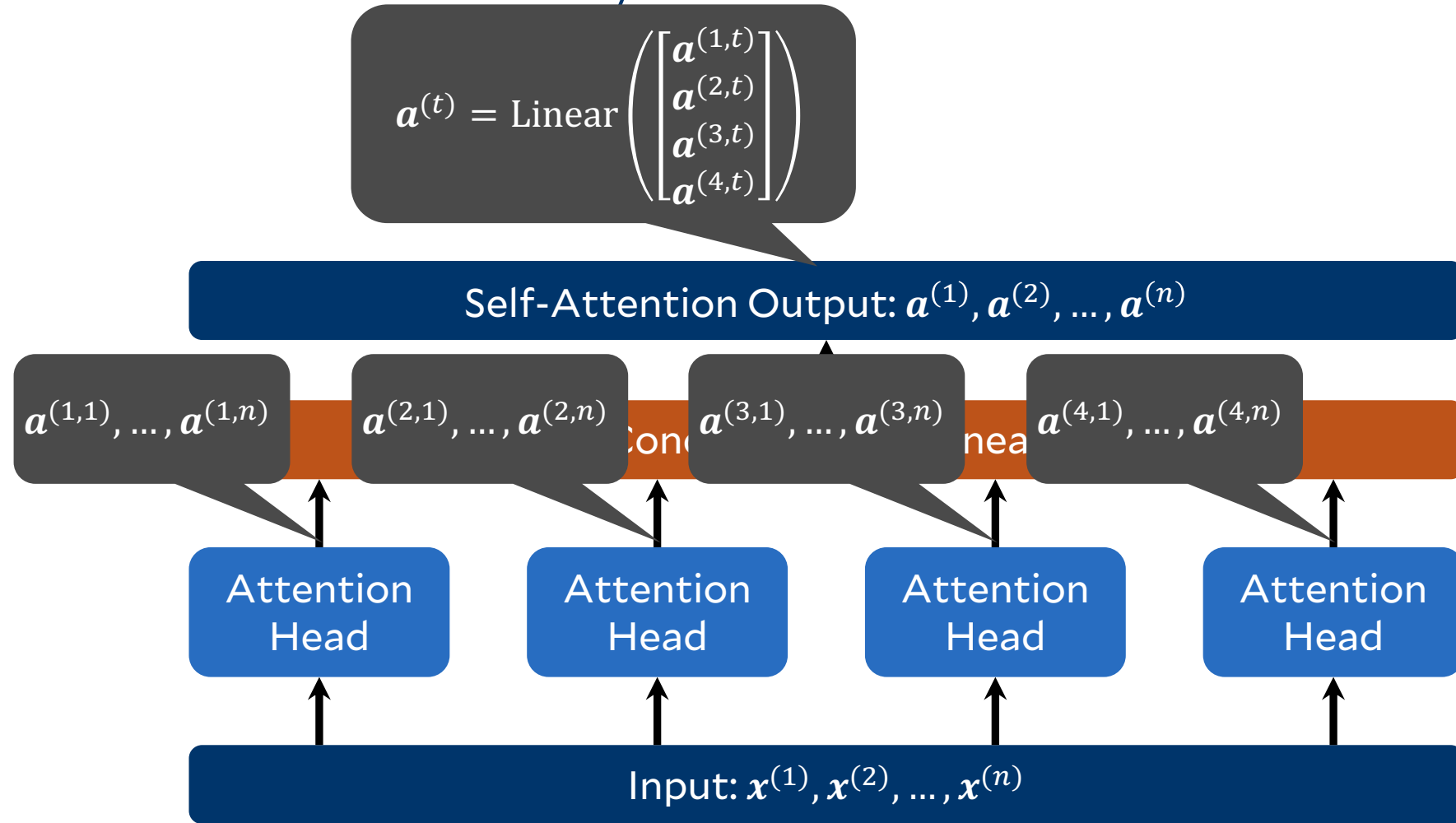
| Model | BLEU (%) |
|---|---|
| GRU, 50 Hidden (Bahdanau et al., 2015) | 17.82 |
| GRU, 50 Hidden + Bahdanau Attn. (ibid) | 28.45 |
| GRU, 1000 Hidden (Cho et al., 2014) | 33.87 |
| Google NMT: LSTM, 8 Layers, 1024 Hidden + Bahdanau Attn. (Wu et al., 2016) | 39.92 |
| Transformer, 512 Hidden, 8 Heads (Vaswani et al., 2017) | 38.1 |
| Transformer, 1024 Hidden, 16 Heads (ibid) | 41.0 |

# Self-Attention Head

- Input: $x^{(1)}, x^{(2)}, \ldots, x^{(n)}$
- Output: $a^{(1)}, a^{(2)}, \ldots, a^{(n)}$
- Queries:

$$q^{(t)} = \text{Linear}(x^{(t)}) \in \mathbb{R}^d$$

- Keys:

$$k^{(i)} = \text{Linear}(x^{(i)}) \in \mathbb{R}^d$$

- Values:

$$v^{(i)} = \text{Linear}(x^{(i)})$$

- Attention Scores:

$$s_i^{(t)} = \text{score}(q^{(t)}, k^{(i)})$$
$$= (q^{(t)})^\top k^{(i)} / \sqrt{d}$$

- Attention Weights:

$$\alpha^{(t)} = \text{softmax}(s^{(t)})$$

- Attention Vectors:

$$a^{(t)} = \sum_i \alpha_i^{(t)} v^{(i)}$$

# Self-Attention Layer

# Encoder Block



$\boldsymbol{h}^{(t)} = \text{LayerNorm}(\boldsymbol{b}^{(t)} + \boldsymbol{c}^{(t)})$

Output: $\boldsymbol{h}^{(1)}, \boldsymbol{h}^{(2)}, \dots, \boldsymbol{h}^{(n)}$

Addition + Layer Normalization

$\boldsymbol{c}^{(t)} = \text{MLP}(\boldsymbol{b}^{(t)})$

$\boldsymbol{b}^{(t)} = \text{LayerNorm}(\boldsymbol{x}^{(t)} + \boldsymbol{a}^{(t)})$

MLP (Linear + ReLU + Linear)

Addition + Layer Normalization

$\boldsymbol{a}^{(1)}, \boldsymbol{a}^{(2)}, \dots, \boldsymbol{a}^{(n)}$

Multi-Head Self Attention

Input: $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \dots, \boldsymbol{x}^{(n)}$

# Decoder Block



$h^{(t)} = \text{LayerNorm}(d^{(t)} + e^{(t)})$

Output: $h^{(1)}, h^{(2)}, \ldots, h^{(n)}$

$e^{(t)} = \text{MLP}(d^{(t)})$

Addition + Layer Normalization

$d^{(t)} = \text{LayerNorm}(b^{(t)} + c^{(t)})$

MLP (Linear + ReLU + Linear)

$c^{(1)}, c^{(2)}, \ldots, c^{(n)}$
$k^{(i)} = \text{Linear}(o^{(i)})$
$v^{(i)} = \text{Linear}(o^{(i)})$

Addition + Layer Normalization

Encoder Output:
$o^{(1)}, o^{(2)}, \ldots, o^{(n)}$

Multi-Head Self Attention

$a^{(1)}, a^{(2)}, \ldots, a^{(n)}$

Addition + Layer Normalization

$b^{(t)} = \text{LayerNorm}(x^{(t)} + a^{(t)})$

Multi-Head Self Attention

Input: $x^{(1)}, x^{(2)}, \ldots, x^{(n)}$

# Transformer Encoder – Decoder

# Types of Recurrent Networks

|  | Single Input | Sequence Input |
|---|---|---|
| **Single Output** | MLP | RNN Classifier |
| **Sequence Output** | RNN Generator | RNN Encoder-Decoder |

# Types of Transformer Networks

|  | Single Input | Sequence Input |
|---|---|---|
| **Single Output** | MLP | Transformer Classifier |
| **Sequence Output** | Transformer Generator | Transformer Encoder–Decoder |

# Transformer Generator

# Transformer Classifier

Output

Linear + Softmax

$o^{(0)}$

Encoder Block

Encoder Block

Encoder Block

Embedding  +  Pos. Encoding

Input

# Transformer vs. RNN Training Time

| Model | BLEU (%) | Time (PetaFLOPs) |
|---|---|---:|
| Google NMT | 39.92 | 140,000 |
| Transformer (Base) | 38.1 | 3,300 |
| Transformer (Big) | 41.0 | 23,000 |