# Transformers

# Dot Product

- Let's say we have two 3-dimensional vectors, $x_1$ and $x_2$:

$x_1$ = [1, 2, 3] $x_2$ = [4, 5, 6]

$x_1 . x_2$ = (1 x 4) + (2 x 5) + (3 x 6) = 4 + 10 + 18 = 32

- **Activity 1:** Let's say we have two 2-dimensional vectors, $x_1$ and $x_2$, represented as follows: $x_1$ = [3, 5] $x_2$ = [-2, 4], dot product of A and $x_2$ ?

- **Activity2:** Let's say we have two 4-dimensional vectors, $x_1$ and $x_2$, represented as follows: $x_1$ = [1, 2, 3, 4] $x_2$ = [-4, 0, 2, 5], dot product of $x_1$ and $x_2$?

# Self Attention

- Attention is comparison of input to other input elements.

  e.g value of $y_3$ depends on the score of $(x_3, x_1), (x_3, x_2)$ and $(x_3, x_3)$
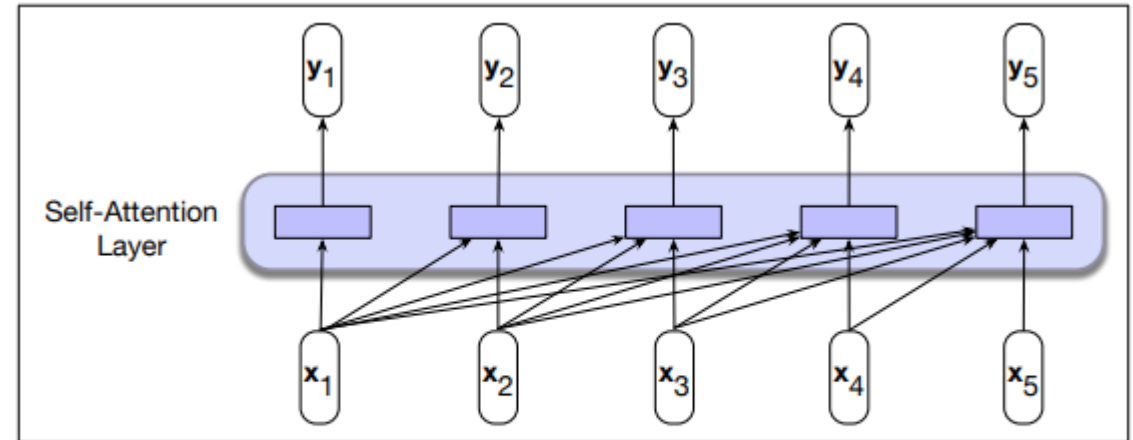


**Figure 10.1** Information flow in a causal (or masked) self-attention model. In processing each element of the sequence, the model attends to all the inputs up to, and including, the current one. Unlike RNNs, the computations at each time step are independent of all the other steps and therefore can be performed in parallel.

$$score(x_i, x_j) = x_i . x_j$$

# Self Attention

$$\alpha_{i,j} = softmax\left(score(\boldsymbol{x}_i, \boldsymbol{x}_j)\right) \quad \forall j \leq i$$

$$\alpha_{i,j} = \frac{\exp\left(score(\boldsymbol{x}_i, \boldsymbol{x}_j)\right)}{\sum_{k=1}^{i} \exp(score(\boldsymbol{x}_i, \boldsymbol{x}_k))} \quad \forall j \leq i$$

$$\boldsymbol{y}_i = \sum_{j \leq i} \alpha_{ij} \boldsymbol{x}_j$$

# Activity

# Query, Key, Value

Considering the output $y_3$

$$\boldsymbol{x}_3 \quad query$$

$$\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{x}_3 \quad key$$

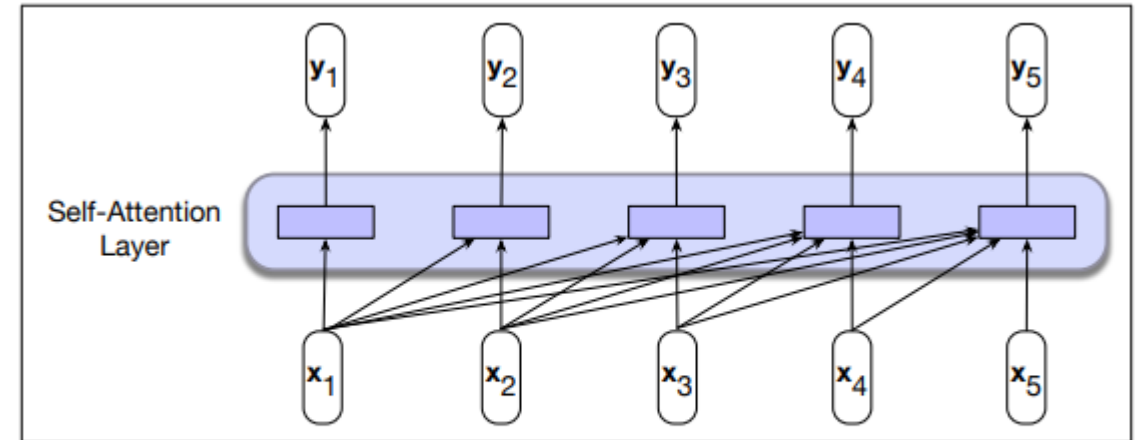$$\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{x}_3 \quad value$$



**Figure 10.1** Information flow in a causal (or masked) self-attention model. In processing each element of the sequence, the model attends to all the inputs up to, and including, the current one. Unlike RNNs, the computations at each time step are independent of all the other steps and therefore can be performed in parallel.

*Attention Values*

$$score(\boldsymbol{x}_3, \boldsymbol{x}_1) = \alpha_{31} \qquad score(\boldsymbol{x}_3, \boldsymbol{x}_2) = \alpha_{32} \qquad score(\boldsymbol{x}_3, \boldsymbol{x}_3) = \alpha_{33}$$

$$\alpha_{31} \boldsymbol{x}_1 + \alpha_{32} \boldsymbol{x}_2 + \alpha_{33} \boldsymbol{x}_3 = \boldsymbol{y}_3$$

# Key, Query and Value Matrices

- To capture the roles, transformer has 3 matrices

$$W^Q, W^K, W^V$$

- Transformer project each $x_i$ into **query, key** and **value** vectors

$$q_i = W^Q x_i \; ; k_i = W^K x_i \; : v_i = W^V x_i$$

$$W^Q, W^K, W^V \in R^{d \times d}, \qquad x_i, y_i \in R^{1 \times d}$$

# Self-attention (improved version)

$$score(\boldsymbol{x}_i, \boldsymbol{x}_j) = \boldsymbol{q}_i.\boldsymbol{k}_j$$

$$score(\boldsymbol{x}_i, \boldsymbol{x}_j) = \frac{\boldsymbol{q}_i.\boldsymbol{k}_j}{\sqrt{d}}$$

$$\alpha_{i,j} = \frac{\exp\left(score(\boldsymbol{x}_i,\boldsymbol{x}_j)\right)}{\sum_{k=1}^{i} \exp(score(\boldsymbol{x}_i,\boldsymbol{x}_k))} \quad \forall j \leq i$$

$$\boldsymbol{y}_i = \sum_{j \leq i} \alpha_{ij} \boldsymbol{v}_j$$

# Visual look of calculating $y_3$

- Single output at single time step $i$

- Each output $y_i$ is calculated independently



Output Vector

Weight and Sum value vectors

Softmax $\alpha_{i,j}$

Key/Query Comparisons

Generate key, query, value vectors

$x_1$ $x_2$ $x_3$

**Figure 10.2** Calculating the value of $\mathbf{y}_3$, the third element of a sequence using causal (left-to-right) self-attention.
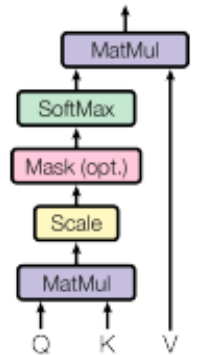
# Vectorization

- Input $X$: A sequence of $N$ tokens with embedding dimension $d$

$$Q = XW^Q \; ; K = XW^K \; : V = XW^V$$

$$W^Q, W^K, W^V \in R^{d \times d}, \qquad X \in R^{n \times d}$$

$$Q, K, V \in R^{n \times d},$$

Scaled Dot-Product Attention



$$SelfAttention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$$

# Attention Values



**Figure 10.3** The $N \times N$ $\mathbf{QK}^\mathsf{T}$ matrix showing the $q_i \cdot k_j$ values, with the upper-triangle portion of the comparisons matrix zeroed out (set to $-\infty$, which the softmax will turn to zero).

Attention is quadratic
In the length of input

# Residual Connection

- In deep networks, residual connections are connections that pass information from a lower layer to a higher layer without going through the intermediate layer.

- Residual connections in transformers are implemented by adding a layer's input vector to its output vector before passing it forward.
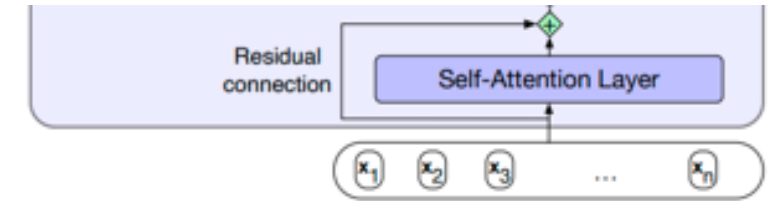
$$x + Self\,Attention(x),$$

# Residual Connection



Residual connection — Self-Attention Layer — $x_1$ $x_2$ $x_3$ ... $x_n$

- Suppose we have an input matrix X representing word embeddings of a sentence:

- Assume we have already calculated the SelfAttention(Q, K, V) matrix using the self-attention mechanism:

- Residual_Connection =

X + SelfAttention(Q, K, V) =

X = [   [1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]   ]

SelfAttention(Q, K, V) =
[   [0.5, 1.0, 1.5],
    [2.0, 2.5, 3.0],
    [3.5, 4.0, 4.5]   ]

[   [1.5, 3.0, 4.5],
    [6.0, 7.5, 9.0],
    [10.5, 12.0, 13.5]   ]

# Layer Normalization

- Given a batch of input vectors X with dimensions $(batch\_size, feature\_dim)$, a scale parameter vector γ and a bias parameter vector β, the layer normalization of X is computed as:

$$\mu = \frac{1}{batch\_size} \sum_{i=1}^{batch\_size} X_i$$

$$\sigma^2 = \frac{1}{batch\_size} \sum_{i=1}^{batch\_size} X_i - \mu$$

$$\hat{X} = \frac{X - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$LayerNorm = \alpha\hat{X} + \beta$$

Where $\alpha$ and $\beta$ are learnable parameters and $\epsilon$ is small constant

# Layer normalization

Layer Normalize

- Step 1:  Calculate the mean and standard deviation of each feature.
- Step 2: Normalize each feature using the mean and standard deviation

-

Residual_Connection = [
[1.5, 3.0, 4.5],
[6.0, 7.5, 9.0],
[10.5, 12.0, 13.5]
]

Mean = [
(1.5+6.0+10.5)/3, (3.0+7.5+12.0)/3, (4.5+9.0+13.5)/3
] = [6.0, 7.5, 9.0]

Standard Deviation = [
sqrt((($1.5-6.0)^2+(6.0-6.0)^2+(10.5-6.0)^2)/3),
sqrt((($3.0-7.5)^2+(7.5-7.5)^2+(12.0-7.5)^2)/3),
sqrt((($4.5-9.0)^2+(9.0-9.0)^2+(13.5-9.0)^2)/3)
] ≈ [3.674, 3.674, 3.674]

Normalized = [
[(1.5-6.0)/3.674, (3.0-7.5)/3.674, (4.5-9.0)/3.674],
[(6.0-6.0)/3.674, (7.5-7.5)/3.674, (9.0-9.0)/3.674],
[(10.5-6.0)/3.674, (12.0-7.5)/3.674, (13.5-9.0)/3.674]
] ≈ [
[-1.22, -1.22, -1.22],
[0.00, 0.00, 0.00],
[1.22, 1.22, 1.22]
]

# Layer normalization

Layer Normalize

• Step 3: Apply learnable scale and shift parameters

Layer normalization also includes learnable parameters (gamma) for scaling and (beta) for shifting. For simplicity, assume gamma = 1 and beta = 0.

Layer_Normalized = gamma * Normalized + beta ≈

[ [-1.22, -1.22, -1.22],

[0.00, 0.00, 0.00],

[1.22, 1.22, 1.22]  ]

# Transformer Block



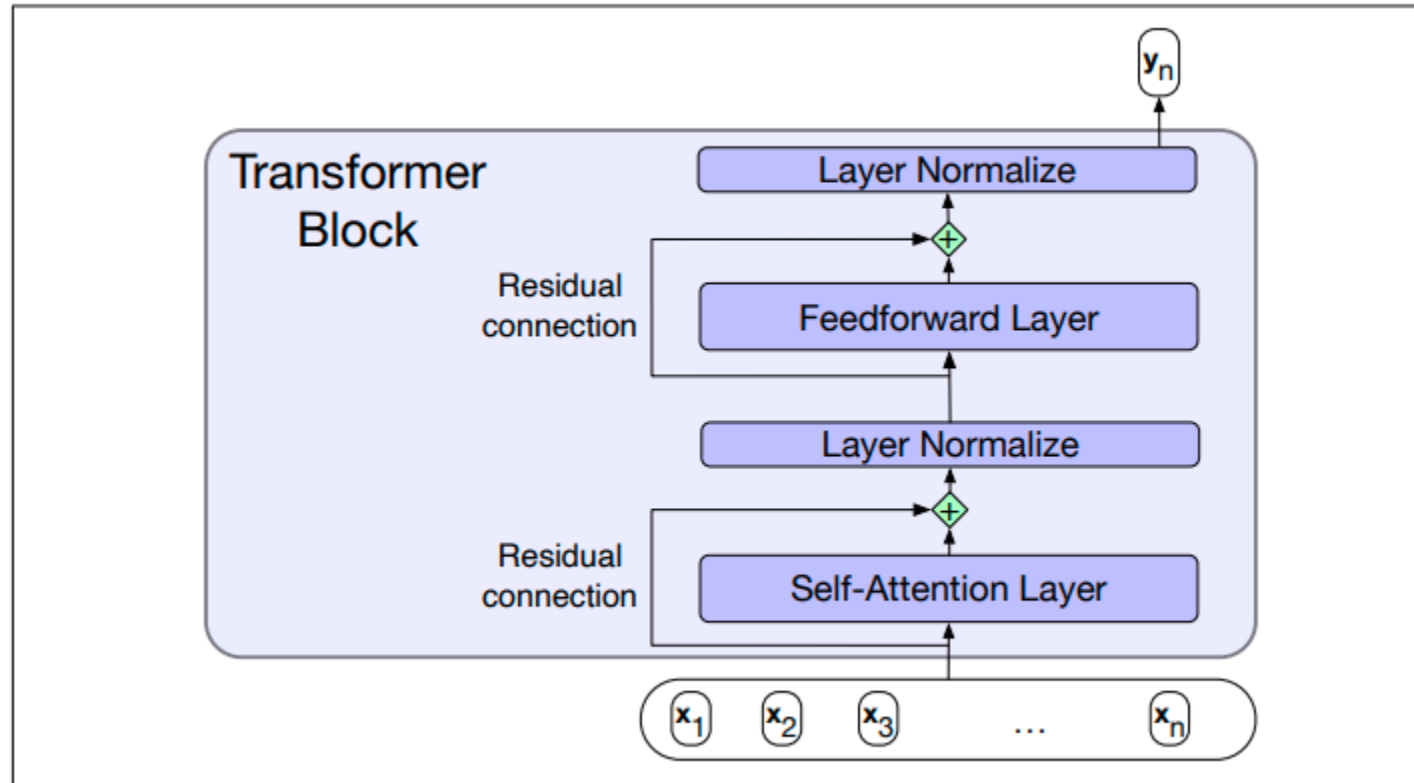**Figure 10.4** A transformer block showing all the layers.

# MultiHead Attention

$$\text{MultiheadAttention}(X) = (head_1 \oplus head_2 \oplus \cdots head_h )W^O$$

$$Q = XW_i^Q : \; K = XW_i^K : \;\; V = XW_i^V$$

$$head_i = Self\,Attention(Q, K, V)$$

$$W^Q \in R^{d \times d_k}, W^K \in R^{d \times d_k}, W^V \in R^{d \times d_v}, \qquad X \in R^{n \times d}$$

$$Q \in R^{N \times d_k}, K \in R^{N \times d_k}, V \in R^{N \times d_v},$$

$$W^O \in R^{hd_v \times h}$$

# Multihead Attention Transformer Block



**Figure 10.4** A transformer block showing all the layers.

# Word Order

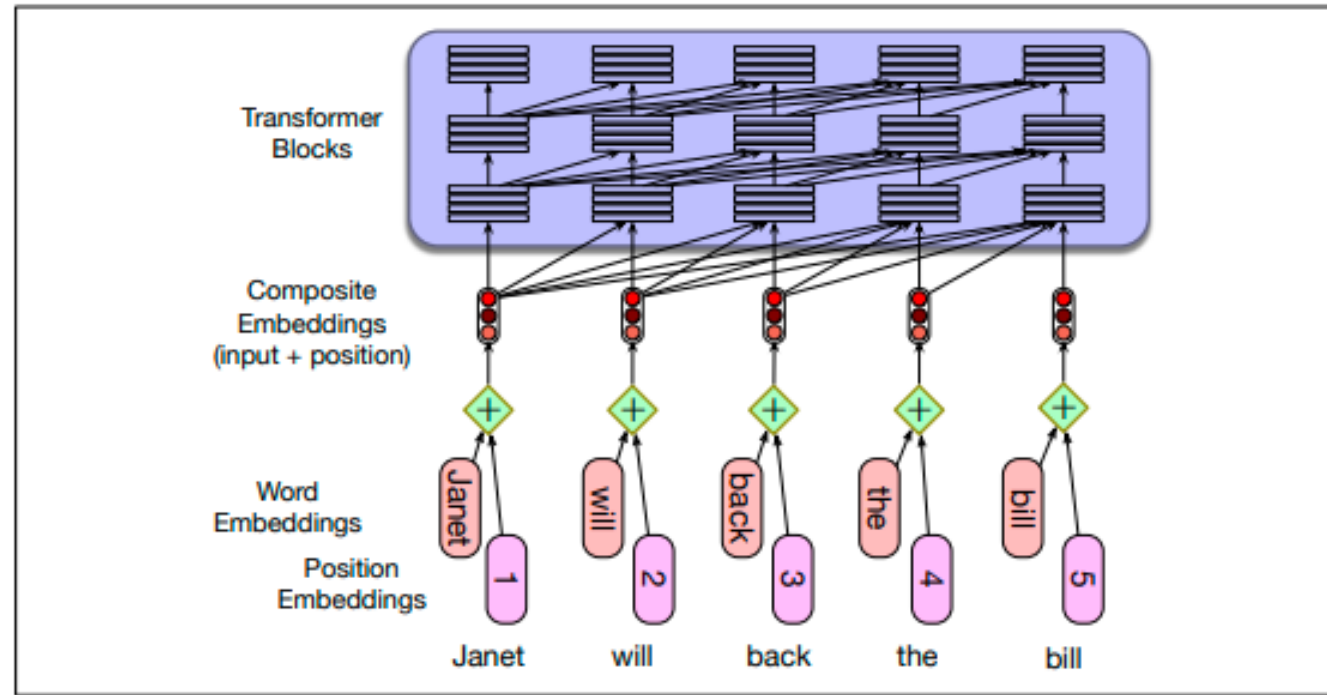- Does transformer model the position of each token in the input sequence ?



**Figure 10.6** A simple way to model position: simply adding an embedding representation of the absolute position to the input word embedding to produce a new embedding of the same dimenionality.
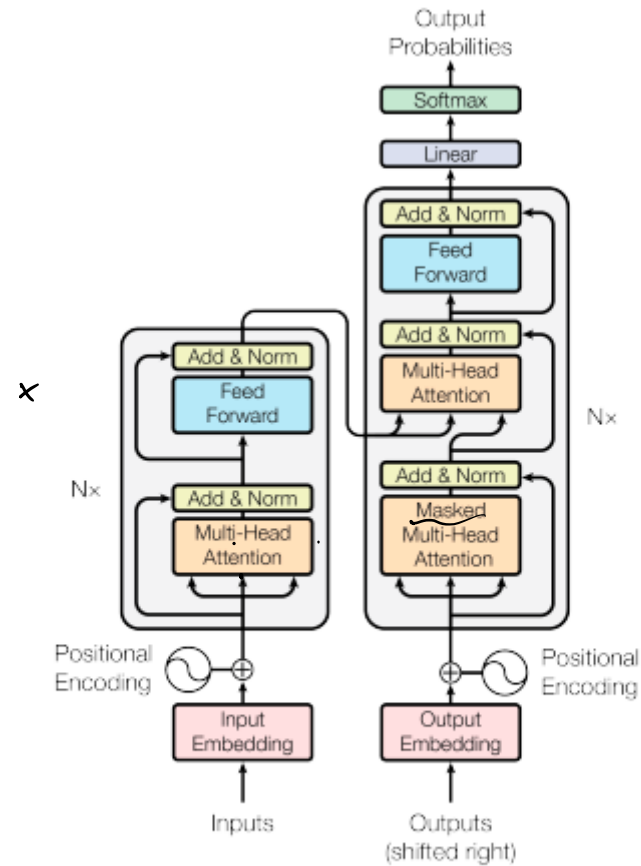
# Transformer architecture



Figure 1: The Transformer - model architecture.

# Transformers as Language Model

$$L_{CE}(\hat{y}_t, y_t) = -\log \hat{y}_t[w_{t+1}]$$

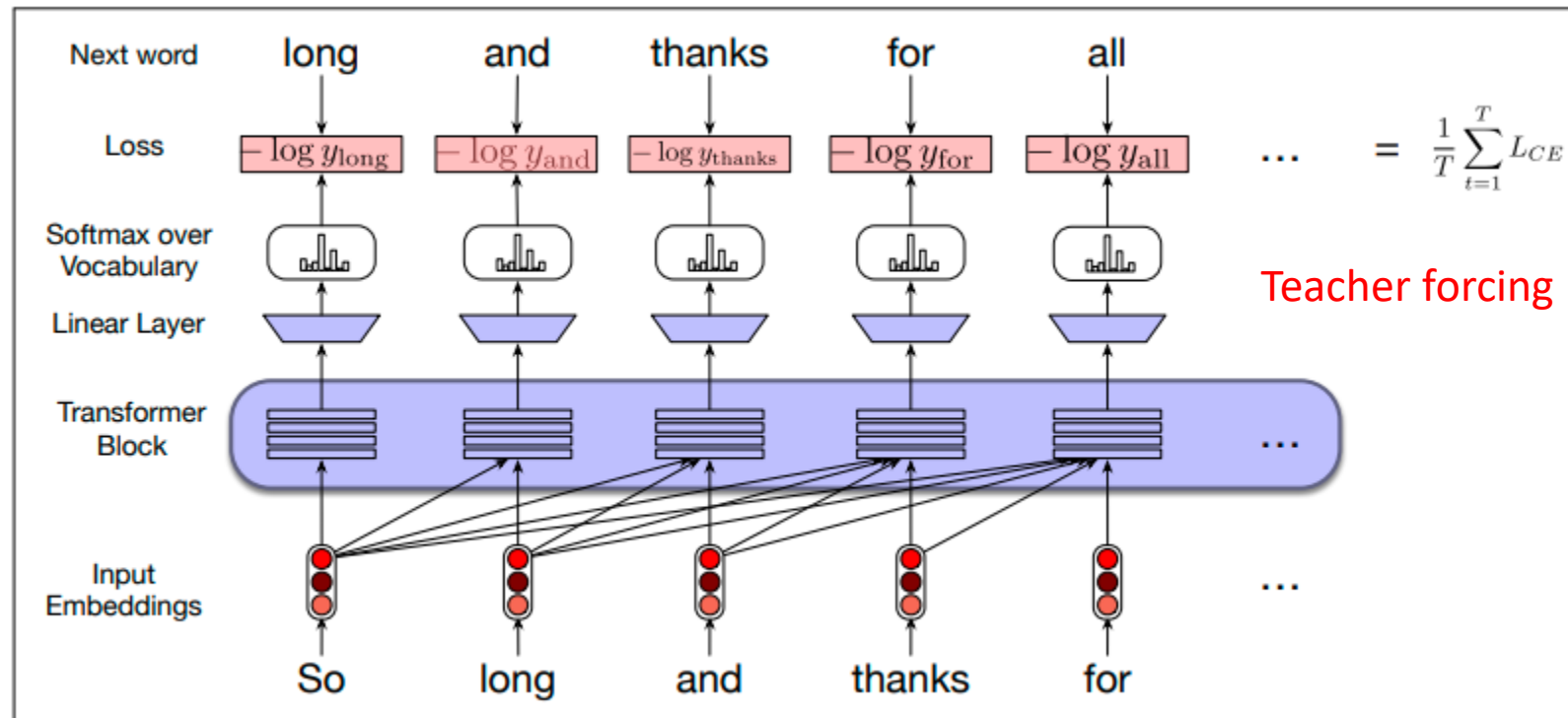Where $w_{t+1}$ is the correct next word in the input sequence
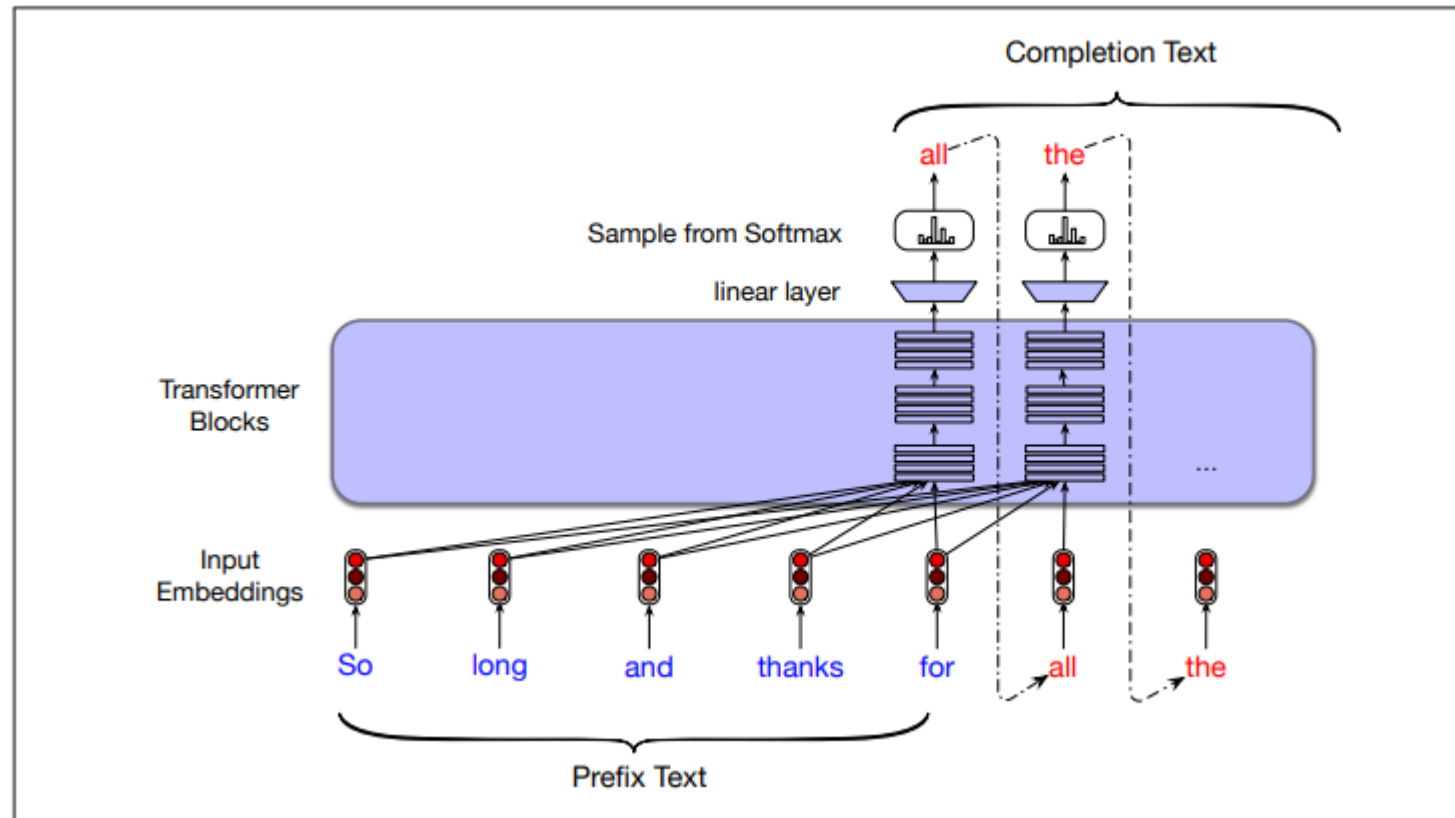


Figure 10.7 Training a transformer as a language model.

# Autoregressive text Completion

# Generating a text from Language model

- Once language model trained we generate the text in the following way
  - ✓ Sample a word in the output from the softmax distribution that results from using the beginning of sentence marker, <s>, as the first input.
  - ✓ Use the word embedding for that first word as the input to the network at the next time step, and then sample the next word in the same fashion.
  - ✓ Continue generating until the end of sentence marker, </s>, is sampled or a fixed length limit is reached.

$$\widehat{y_t} = argmax_{w \in V} P(w|y_1 y_2 \ldots y_{t-1})$$
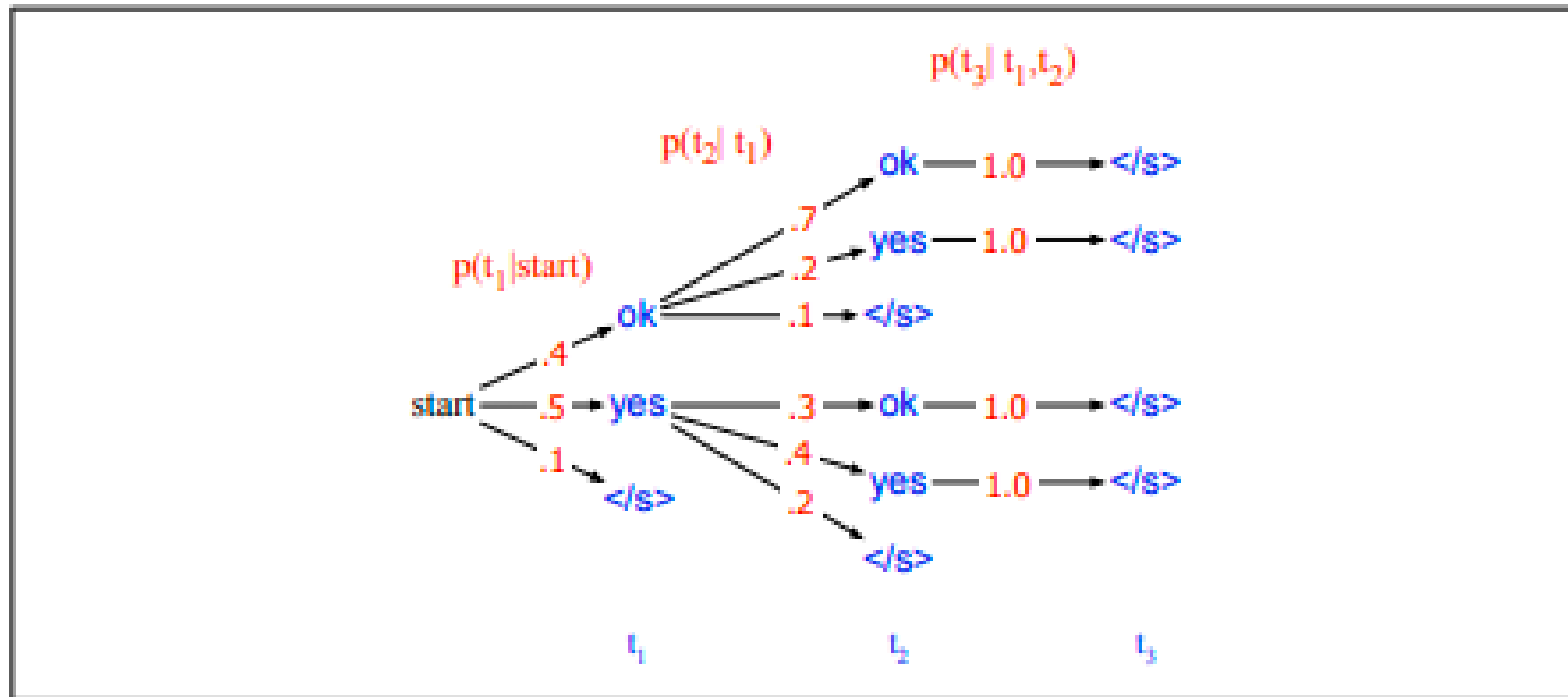
**greedy decoding**

# Greedy approach problem



**Figure 10.8**   A search tree for generating the target string $T = t_1, t_2, \ldots$ from the vocabulary $V = \{yes, ok, <s>\}$, showing the probability of generating each token from that state. Greedy search would choose *yes* at the first time step followed by *yes*, instead of the globally most probable sequence *ok ok*.
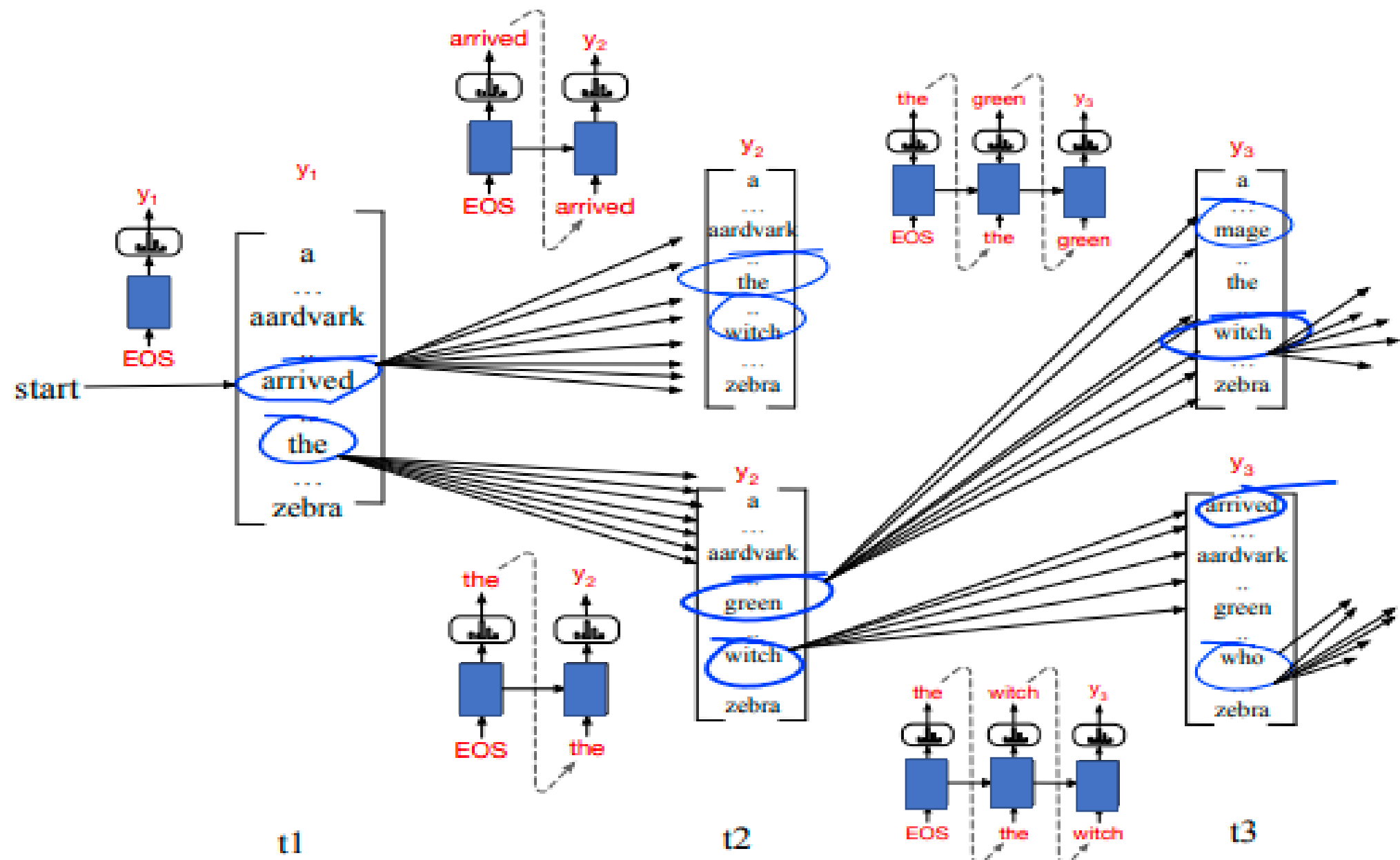
**Figure 10.9** Beam search decoding with a beam width of $k = 2$. At each time step, we choose the $k$ best hypotheses, compute the $V$ possible extensions of each hypothesis, score the resulting $k * V$ possible hypotheses

# Calculating probabilities

$$score(y) = \log P(y|x)$$
$$= \log\left(P(y_1|x)P(y_2|y_1,x)P(y_3|y_1,y_2,x)...P(y_t|y_1,...,y_{t-1},x)\right)$$
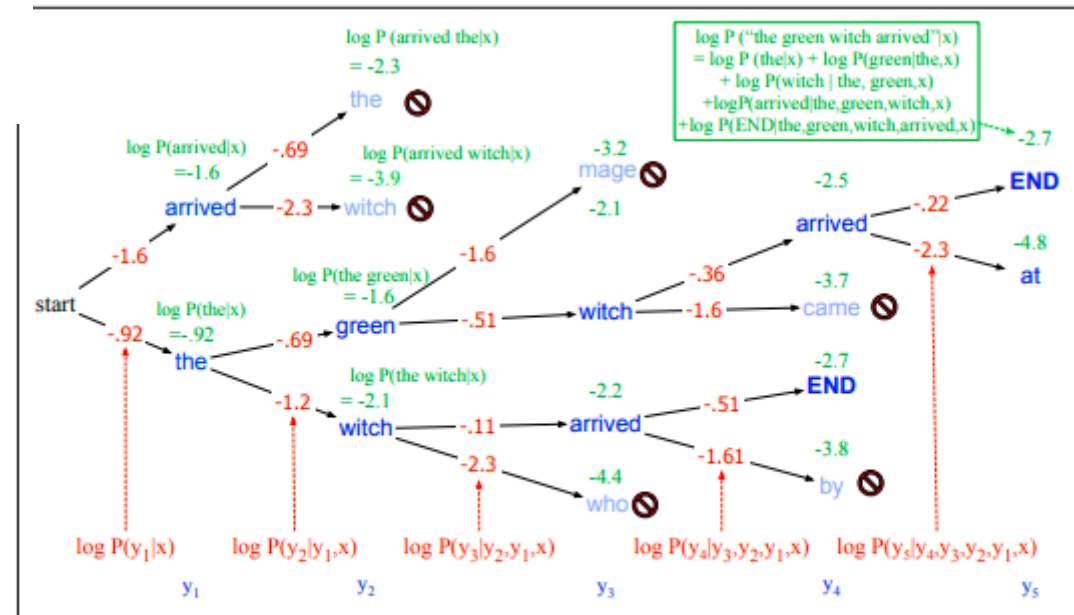$$= \sum_{i=1}^{t} \log P(y_i|y_1,...,y_{i-1},x) \tag{10.22}$$



**Figure 10.10** Scoring for beam search decoding with a beam width of $k = 2$. We maintain the log probability of each hypothesis in the beam by incrementally adding the logprob of generating each next token. Only the top $k$ paths are extended to the next step.