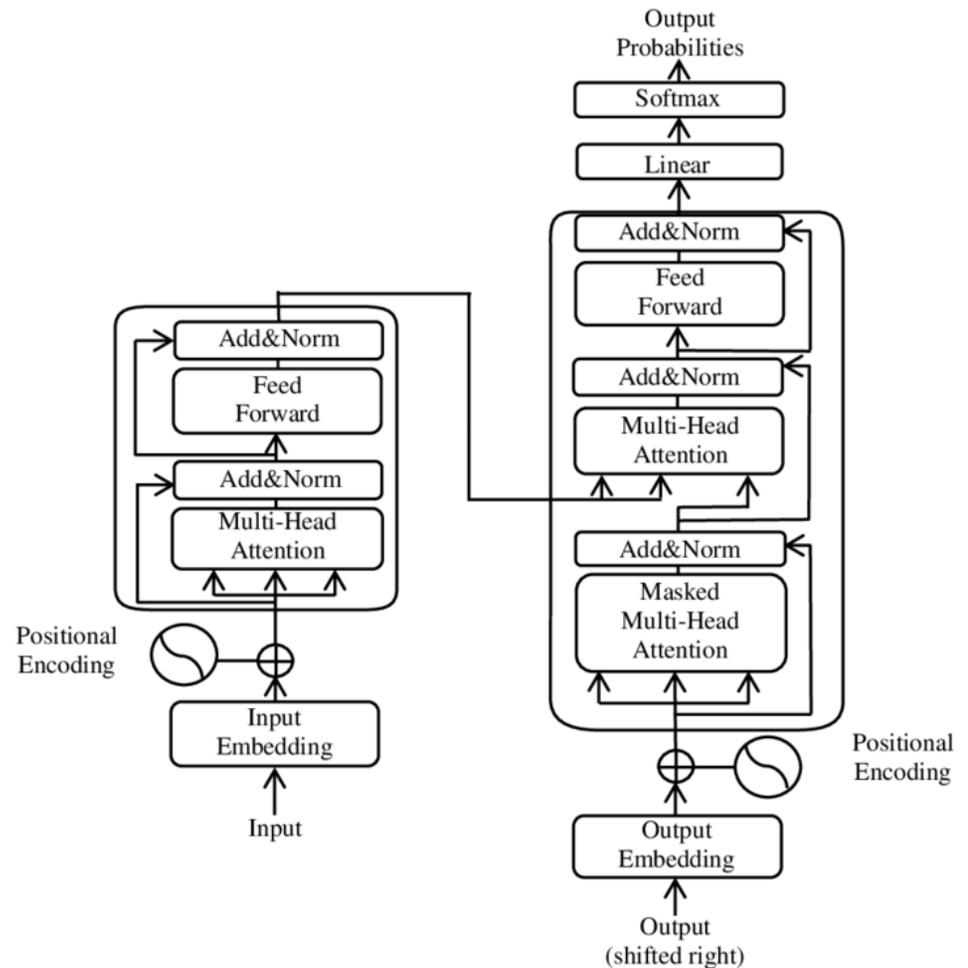


Lecture 14: Transformers

Mark Hasegawa-Johnson

CC0 Public Domain: Re-Use, Re-Mix, Re-distribute at will



Outline

- Text generation
- Attention
- Multi-headed attention
- Positional embedding
- Layer norm
- Residual connections

Text generation

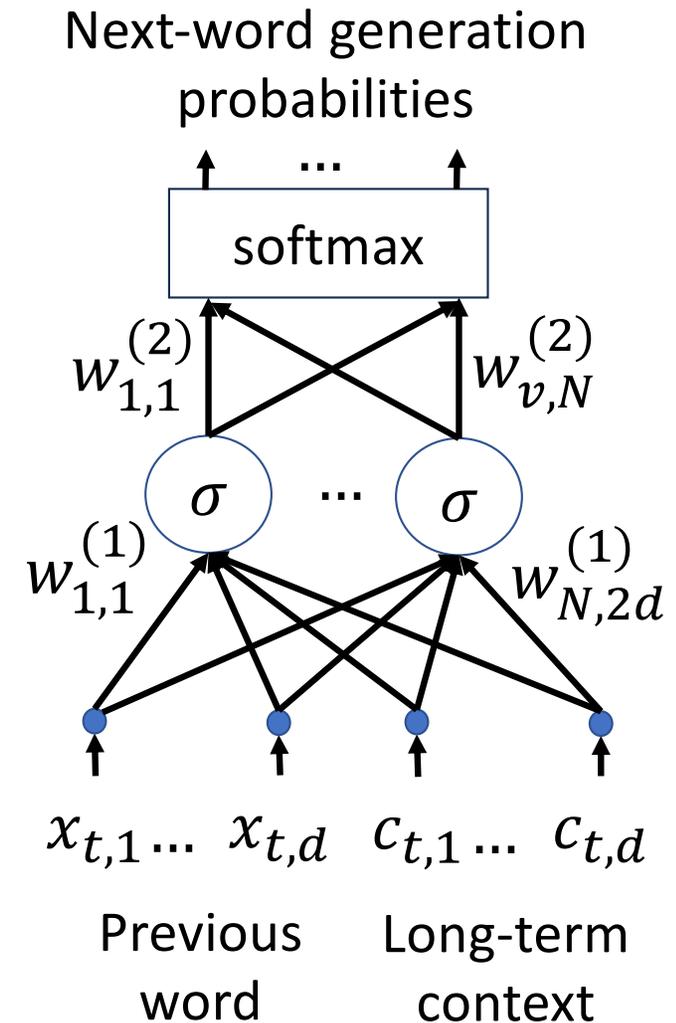
How can a neural net generate text?

- Previous word = \mathbf{x}_t , encode as a vector
- Long-term context = \mathbf{c}_t , encode as a vector
- Next word is $O_{t+1} = y$ with probability

$$P(O_{t+1} = y) = \frac{\exp(z_y)}{\sum_k \exp(z_k)}$$

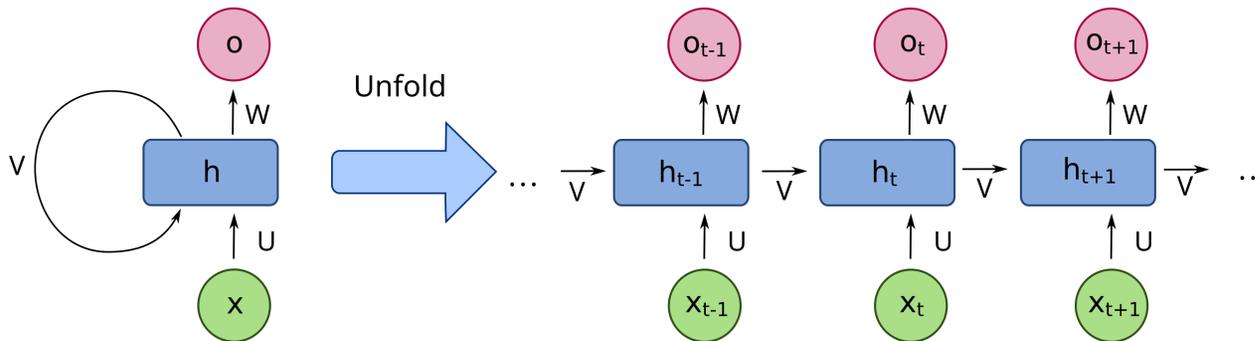
Where z_k is k th output of a feedforward net,

$$\mathbf{z} = \mathbf{W}^{(2)} \sigma(\mathbf{W}^{(1)} [\mathbf{x}_t; \mathbf{c}_t])$$



Generating text without context: Bigram RNN

- Without the context vector, $P(O_{t+1})$ only depends on x_t
- As a statistical model, it's called a "Bigram language model"
- As a neural net, that's called a "Recurrent neural network"



By fdeloche - Own work, CC BY-SA 4.0,
<https://commons.wikimedia.org/w/index.php?curid=60109157>

Generating text with context

There are many ways to include long-term context:

- We can expand the hidden state vector so it's larger than the input
- We can use special architectures: LSTM, MAMBA, other special architectures that remember long-term context

Today we focus on the Transformer

- Attention = a special architecture that creates a weighted summary of long-term context
- Layer norm = a type of normalization to make sure attention works
- Residual connections = each layer is a small change

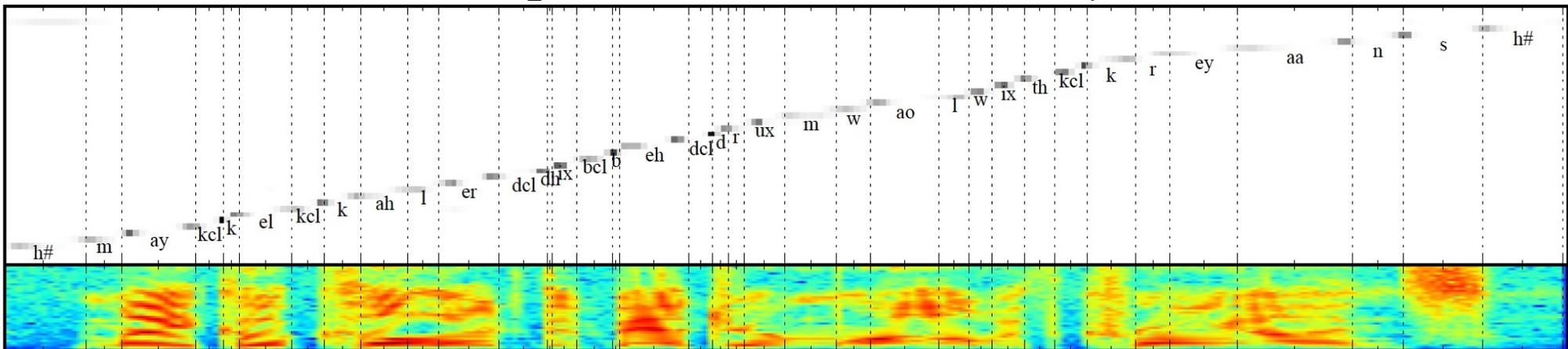
Outline

- Text generation
- Attention
- Multi-headed attention
- Positional embedding
- Layer norm
- Residual connections

Summarizing short-term memory

- Suppose the most recent 1024 words in the conversation are stored in a “short-term memory buffer”
- To compute the i^{th} output, we need context vector $\mathbf{c}_i = \text{summary}(\mathbf{x}_1, \dots, \mathbf{x}_{1024})$
- How can we make that summary?

FDHC0_SX209: Michael colored the bedroom wall with crayons.



Chorowski, Bahdanau, Serdyk, Cho & Bengio, [Attention-Based Models for Speech Recognition](#), Fig. 1

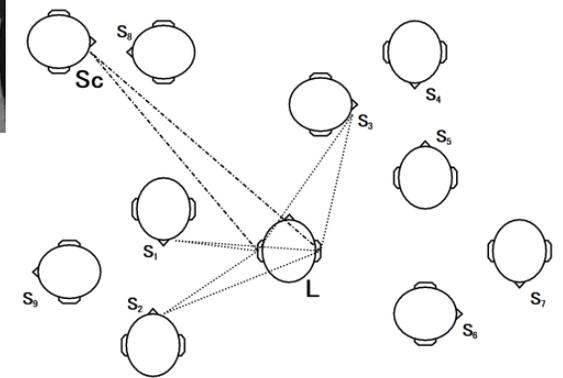
The Cocktail-Party Effect

- If you are focusing on one person's voice, but hear your name spoken by another person, your attention immediately shifts to the second voice.
- This “cocktail-party effect” suggests a model of hearing in which all sounds are processed preconsciously. Trigger sounds in an unattended source will cause attention to re-orient to that source.

https://commons.wikimedia.org/wiki/File:Cocktail_party_attendees_at_Fuller_Lodge,_1946.jpg



[https://commons.wikimedia.org/wiki/File:Cocktail_Party_At_The_Imperial_Hotel_March_13,_1961_\(Tokyo,_Japan\)_496610682.jpg](https://commons.wikimedia.org/wiki/File:Cocktail_Party_At_The_Imperial_Hotel_March_13,_1961_(Tokyo,_Japan)_496610682.jpg)



https://commons.wikimedia.org/wiki/File:Cocktail-party_effect.svg

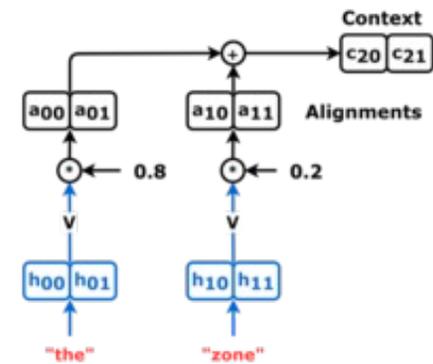
Attention is like probability

- Attention is like probability: You only have a fixed amount of attention, so you need to decide how to distribute it.
- The context vector (\mathbf{c}_i) summarizes input vectors ($\mathbf{x}_1, \dots, \mathbf{x}_{1024}$) using:

$$\mathbf{c}_i = \sum_t \alpha_{i,t} \mathbf{x}_t$$

- where $\alpha_{i,t}$ = importance of \mathbf{x}_t for \mathbf{c}_i :

$$\sum_t \alpha_{i,t} = 1$$

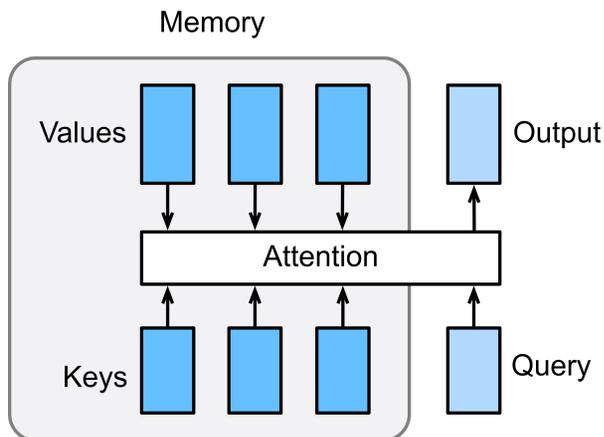


https://commons.wikimedia.org/wiki/File:Encoder_cross-attention,_computing_the_context_vector_by_a_linear_sum.png

What is the “importance of \mathbf{x}_t for \mathbf{c}_i ”?

Let’s break the problem down into smaller pieces:

- “Value vector” \mathbf{v}_t = information from \mathbf{x}_t that will be used in \mathbf{c}_i
- “Key vector” \mathbf{k}_t = describes why \mathbf{v}_t might be useful
- “Query vector” \mathbf{q}_i = describes what we’re looking for

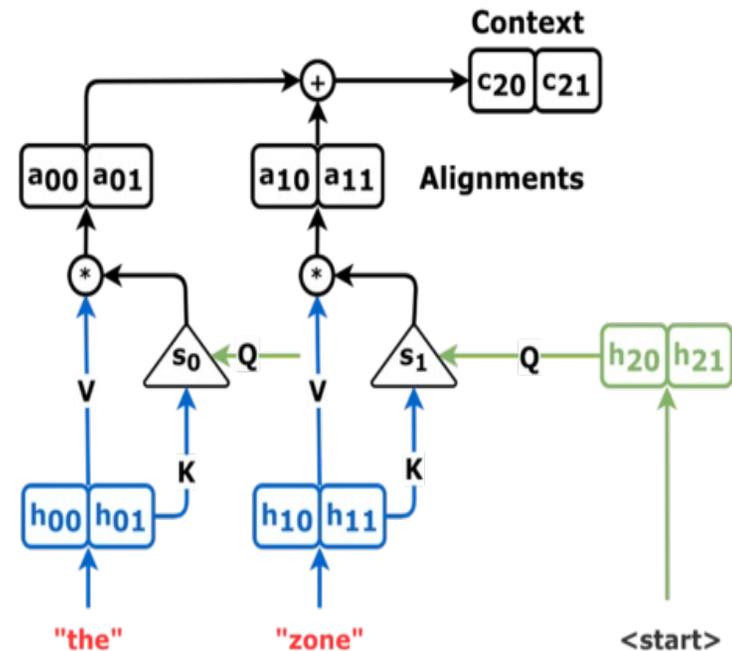


By Zhang, Aston and Lipton, Zachary C. and Li, Mu and Smola, Alexander J. - <https://github.com/d2l-ai/d2l-en>, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=15226571>

What is the “importance of \mathbf{x}_t for \mathbf{c}_i ”?

Let’s break the problem down into smaller pieces:

- “Value vector” \mathbf{v}_t = information from \mathbf{x}_t that will be used in \mathbf{c}_i
- “Key vector” \mathbf{k}_t = describes why \mathbf{v}_t might be useful
- “Query vector” \mathbf{q}_i = describes what we’re looking for

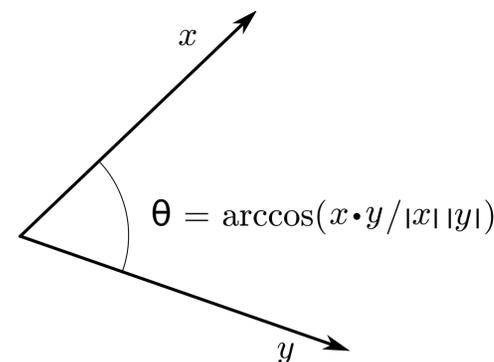


Dot-product attention

How can you decide which value vectors, v_t are most relevant to a particular query? Answer:

1. Create a key vector, \mathbf{k}_t , such that $\mathbf{q}_i^T \mathbf{k}_t > 0$ if v_t is relevant to \mathbf{q}_i , otherwise $\mathbf{q}_i^T \mathbf{k}_t < 0$.
2. Assume that each the key is roughly zero-mean, unit-variance, so $|\mathbf{k}_t|^2 \approx d$, the length of the vector. Assume that the query is unit-length, $|\mathbf{q}_i|^2 \approx 1$ More later about why these are reasonable assumptions.
3. On that assumption,

$$\cos(\mathbf{k}_t, \mathbf{q}_i) \approx \mathbf{q}_i^T \mathbf{k}_t / \sqrt{d}$$



CC BY-SA 3.0
<https://commons.wikimedia.org/w/index.php?curid=49972362>

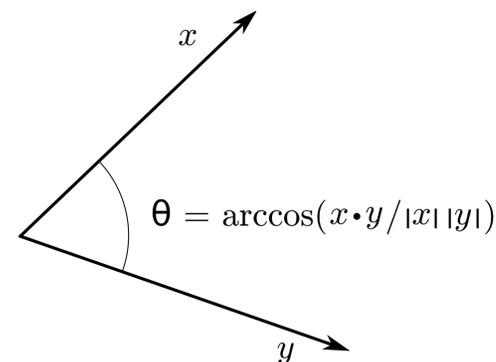
Dot-product attention

Similarity between the query and the key is:

$$\cos(\mathbf{k}_t, \mathbf{q}_i) \approx \mathbf{q}_i^T \mathbf{k}_t / \sqrt{d}$$

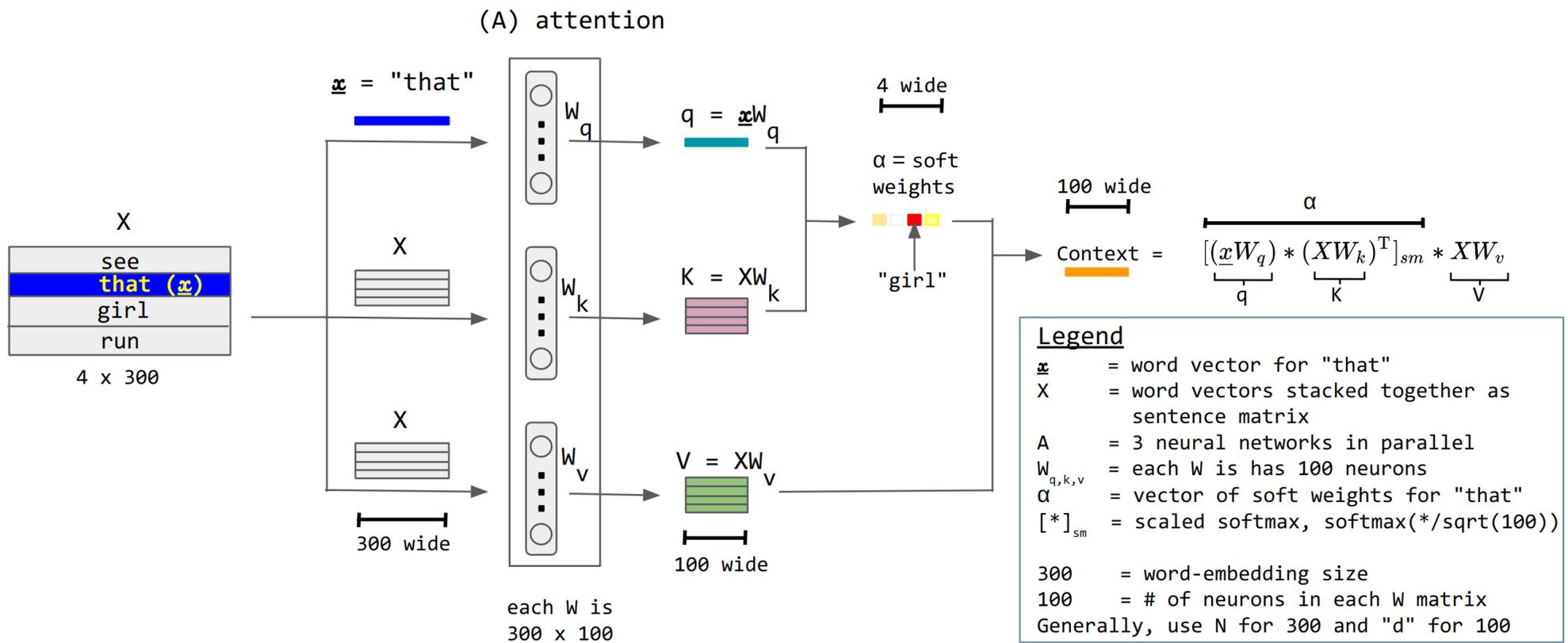
Convert the similarity measures into a probability-like quantity (called attention) using softmax:

$$\alpha_{i,t} = \frac{\exp(\mathbf{q}_i^T \mathbf{k}_t / \sqrt{d})}{\sum_{\tau} \exp(\mathbf{q}_i^T \mathbf{k}_{\tau} / \sqrt{d})}$$



CC BY-SA 3.0

<https://commons.wikimedia.org/w/index.php?curid=49972362>



By Numiri - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=131663062>

Attention

- Stack the input vectors in a matrix, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$
- Multiply by \mathbf{W}_Q , \mathbf{W}_K , and \mathbf{W}_V to create query, key, and value vectors
$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_V$$

- Dot-product similarity:

$$\mathbf{E} = \mathbf{Q}\mathbf{K}^T / \sqrt{d} = \frac{1}{\sqrt{d}} \begin{bmatrix} \mathbf{q}_1^T \mathbf{k}_1 & \cdots & \mathbf{q}_1^T \mathbf{k}_n \\ \vdots & \ddots & \vdots \\ \mathbf{q}_n^T \mathbf{k}_1 & \cdots & \mathbf{q}_n^T \mathbf{k}_n \end{bmatrix}$$

- Attention

$$\mathbf{A} = \begin{bmatrix} \alpha_{1,1} & \cdots & \alpha_{1,n} \\ \vdots & \ddots & \vdots \\ \alpha_{n,1} & \cdots & \alpha_{n,n} \end{bmatrix} = \text{softmax}(\mathbf{E}), \quad \alpha_{i,t} = \frac{\exp(\mathbf{q}_i^T \mathbf{k}_t / \sqrt{d})}{\sum_{\tau} \exp(\mathbf{q}_i^T \mathbf{k}_{\tau} / \sqrt{d})}$$

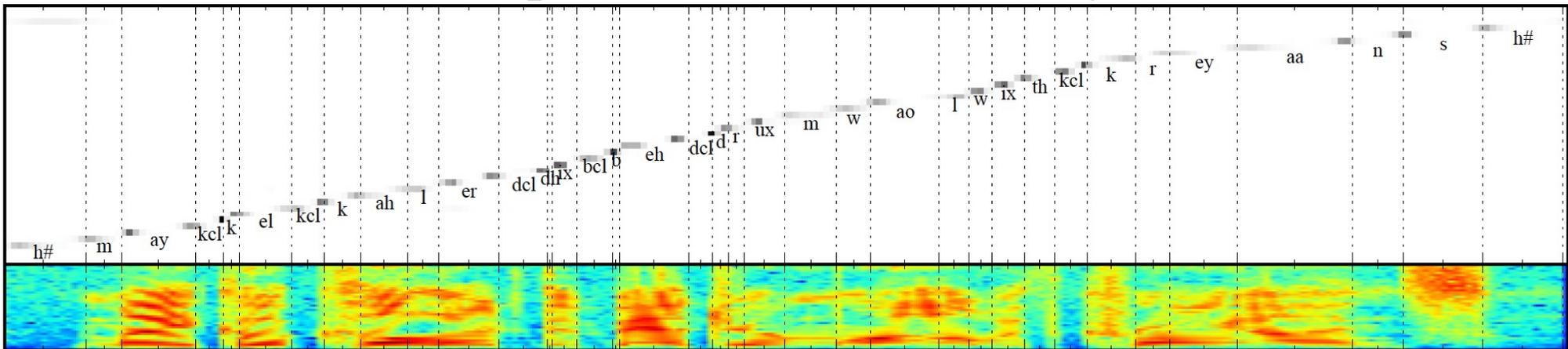
- Context vectors are the dot product of attention times value:

$$\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_n]^T = \mathbf{A}\mathbf{V}, \quad \mathbf{c}_i = \sum_t \alpha_{i,t} \mathbf{v}_t$$

Cross-attention visualization

This plot shows $\alpha_{i,t}$ where i = output character, and t = input spectrum

FDHC0_SX209: Michael colored the bedroom wall with crayons.



Chorowski, Bahdanau, Serdyk, Cho & Bengio, [Attention-Based Models for Speech Recognition](#), Fig. 1

Quiz!

Go to [prairielearn](#) and try the quiz!

Outline

- Text generation
- Attention
- Multi-headed attention
- Positional embedding
- Layer norm
- Residual connections

What is the “importance of \mathbf{x}_t for \mathbf{c}_i ”?

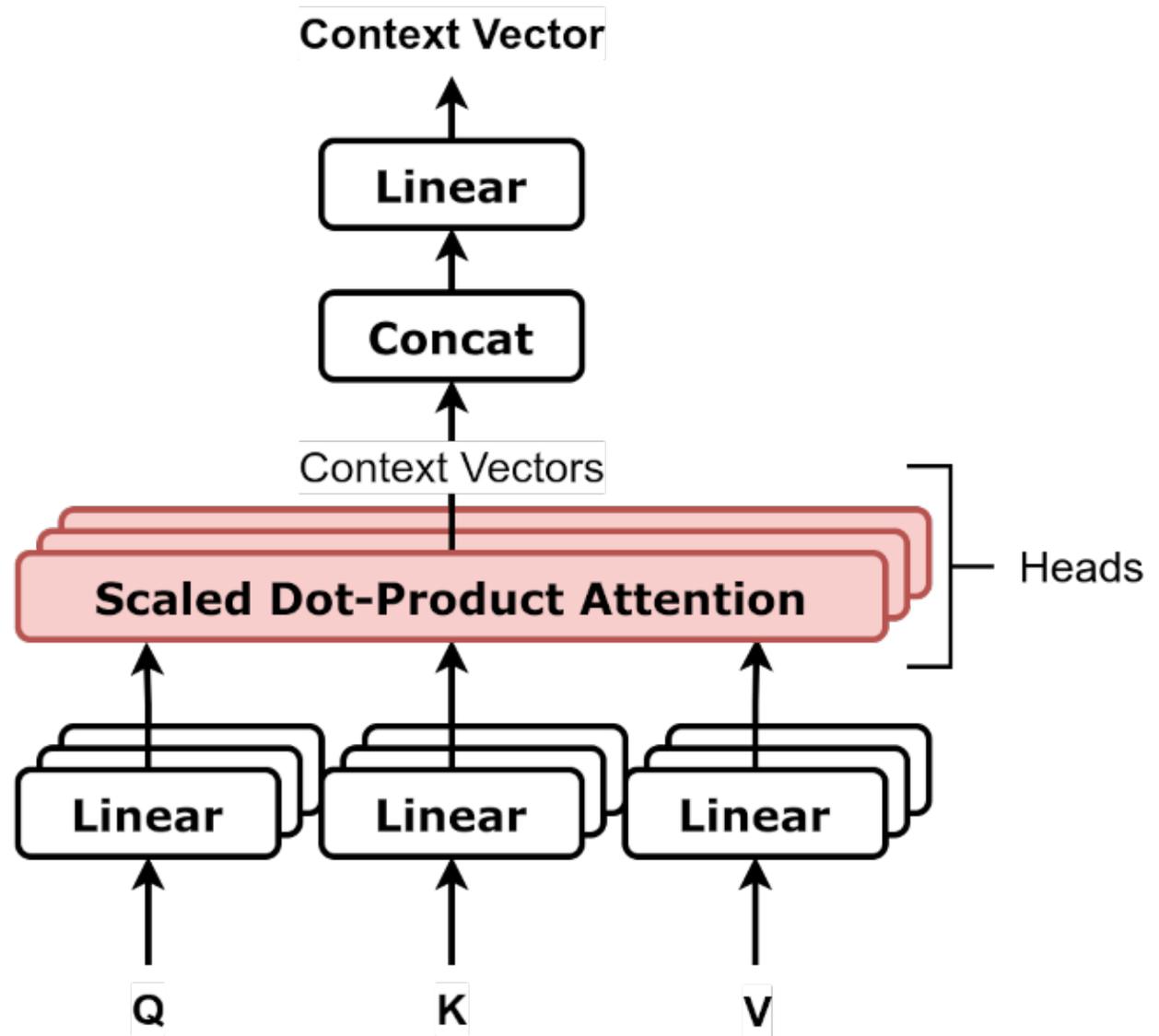
Let’s break the problem down into smaller pieces:

- “Value vector” \mathbf{v}_t = information from \mathbf{x}_t that will be used in \mathbf{c}_i
- “Key vector” \mathbf{k}_t = describes why \mathbf{v}_t might be useful
- “Query vector” \mathbf{q}_i = describes what we’re looking for

...but what if we need different types of context information?

- The general topic of the discussion
- The most recent noun that the human mentioned
- Is next word part of subject, direct object, indirect object, or other?

Multi-headed attention



By dvgodoy -
<https://github.com/dvgodoy/dl-visuals/?tab=readme-ov-file>,
CC BY 4.0,
<https://commons.wikimedia.org/w/index.php?curid=1512160>
29

Multi-Headed Attention

- In the j^{th} head, multiply by $\mathbf{W}_{j,Q}$, $\mathbf{W}_{j,K}$, and $\mathbf{W}_{j,V}$:

$$\mathbf{Q}_j = \mathbf{XW}_{j,Q}, \quad \mathbf{K}_j = \mathbf{XW}_{j,K}, \quad \mathbf{V}_j = \mathbf{XW}_{j,V}$$

- Compute context vectors using dot-product attention:

$$\mathbf{C}_j = \text{softmax}\left(\frac{\mathbf{Q}_j\mathbf{K}_j^T}{\sqrt{d}}\right)\mathbf{V}_j$$

- Concatenate all heads, and linearly transform one more time:

$$\mathbf{O} = [\mathbf{C}_1, \dots, \mathbf{C}_h]\mathbf{W}_O$$

Multi-headed attention gives multiple contexts

- The general topic of the discussion
- The most recent noun that the human mentioned
- Is next word part of subject, direct object, indirect object, or other?
- What is the head word of the current phrase?
- What is the word that started this sentence?
- What is the most recent word?

.... wait those last two need information about time. Where do we get that information?

Outline

- Text generation
- Attention
- Multi-headed attention
- Positional embedding
- Layer norm
- Residual connections

What we have lost...

- With a transformer, each state vector pays attention to the input that is most similar, regardless of what time it happened:

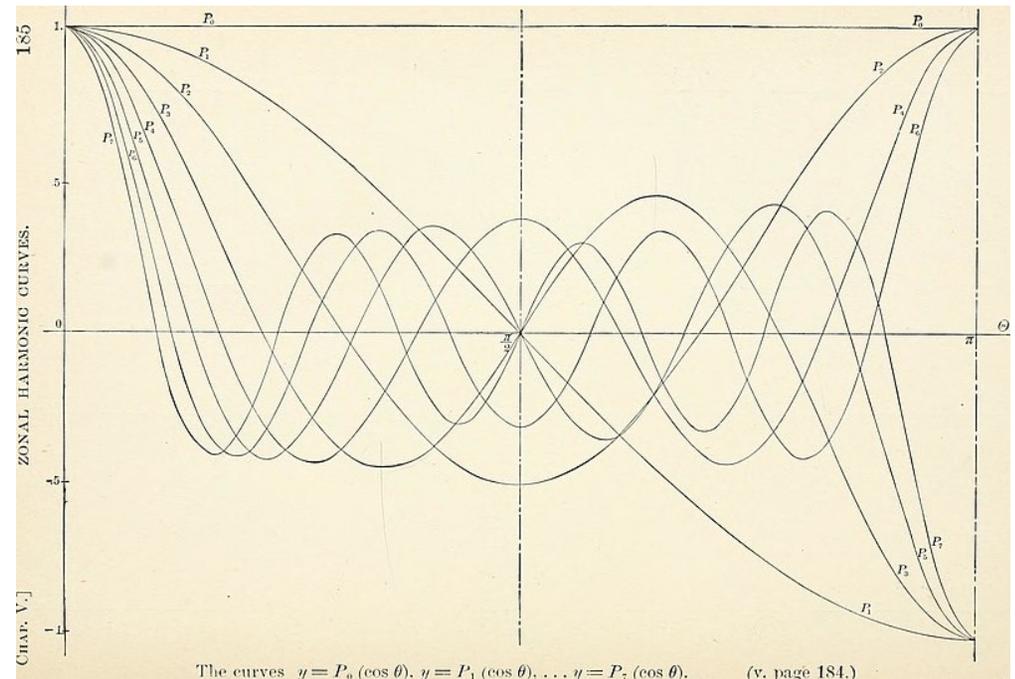
$$\mathbf{c}_i = \sum_t \alpha_{i,t} \mathbf{v}_t, \quad \alpha_{i,t} = \frac{\exp(\mathbf{q}_i^T \mathbf{k}_t)}{\sum_\tau \exp(\mathbf{q}_i^T \mathbf{k}_\tau)}$$

- What if we always want \mathbf{h}_t to pay special attention to \mathbf{v}_{t-1} ? Is that possible?

Rotary Positional Embedding (RoPE) = Fourier series

The solution is to encode the relative position of each input, \mathbf{v}_t , using a Fourier basis \mathbf{e}_t :

$$\mathbf{e}_t = \begin{bmatrix} \cos\left(\frac{\pi t}{T}\right) \\ \sin\left(\frac{\pi t}{T}\right) \\ \vdots \\ \cos\left(\frac{\pi D t}{2T}\right) \\ \sin\left(\frac{\pi D t}{2T}\right) \end{bmatrix}$$



Public domain image,

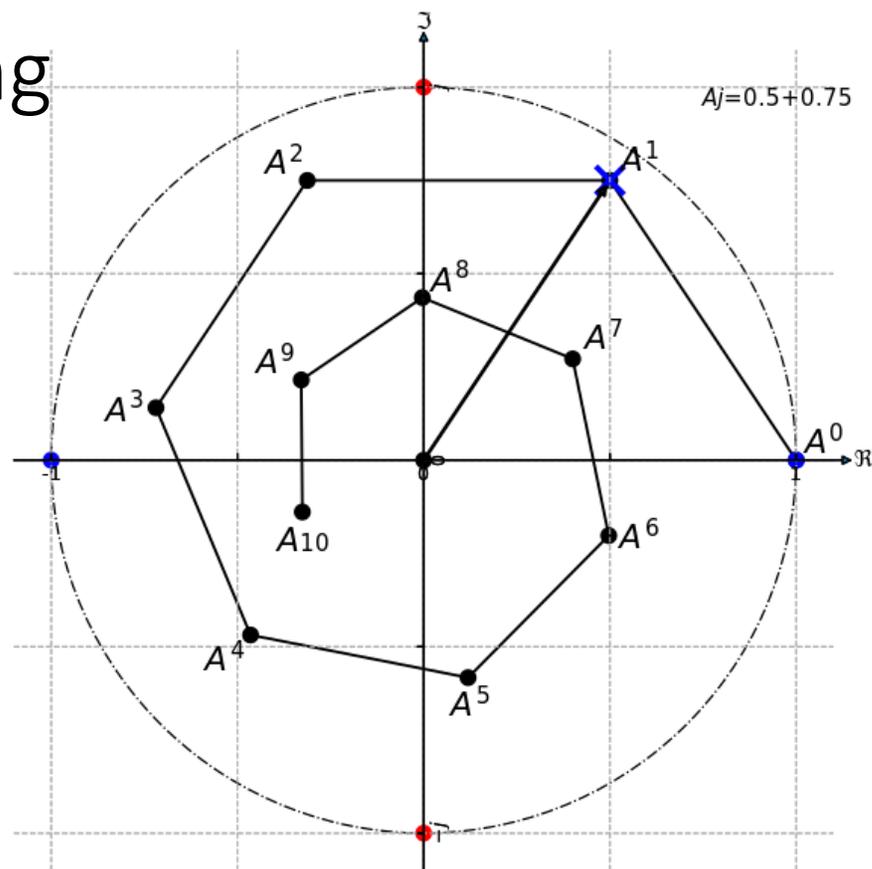
[https://commons.wikimedia.org/wiki/File:An elementary treatise on Fourier's series and spherical, cylindrical, and ellipsoidal harmonics, with applications to problems in mathematical physics \(1893\) \(14780364665\).jpg](https://commons.wikimedia.org/wiki/File:An_elementary_treatise_on_Fourier%27s_series_and_spherical,_cylindrical,_and_ellipsoidal_harmonics,_with_applications_to_problems_in_mathematical_physics_(1893)_%2814780364665%29.jpg)

Rotary positional embedding

The Fourier basis is useful because shifting by a fixed time offset, to $t - \tau$, can be accomplished by a matrix multiplication:

$$\mathbf{e}_{t-\tau} = \begin{bmatrix} \cos\left(\frac{\pi\tau}{T}\right) & \sin\left(\frac{\pi\tau}{T}\right) & \dots \\ -\sin\left(\frac{\pi\tau}{T}\right) & \cos\left(\frac{\pi\tau}{T}\right) & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \mathbf{e}_t$$

...so if we want a particular query to pay attention to vectors with a time delay of τ , we just set $\mathbf{W}_{j,Q}$ to the matrix shown above.



Public domain image,
https://commons.wikimedia.org/wiki/File:Exponentials_of_complex_number_within_unit_circle-2.svg

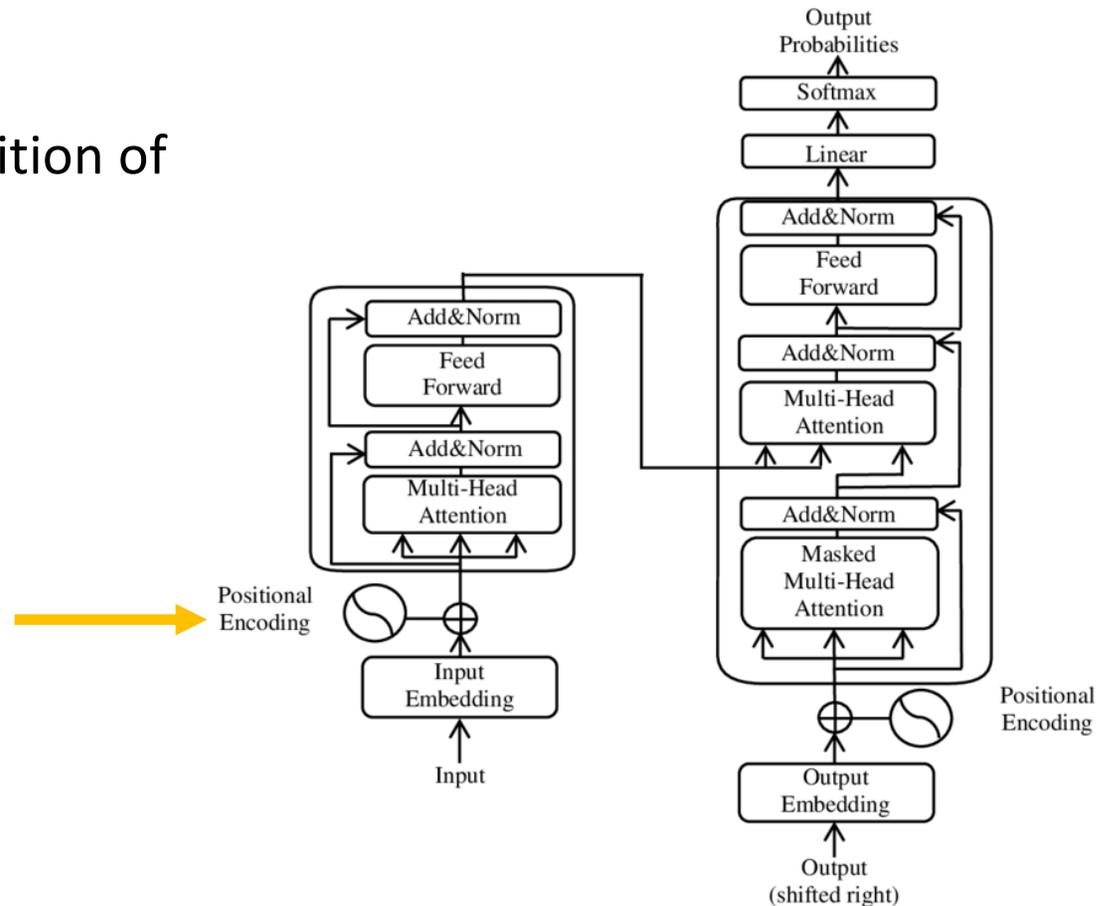
Where do we put the positional encoding?

- Possibility #1: Concatenate it, i.e., $\mathbf{x}_t^T = [\mathbf{x}_t^T, \mathbf{e}_t^T]$
 - Advantage: $\mathbf{W}_{j,Q}$ can learn to operate separately on the content \mathbf{x}_t^T and the positional encoding \mathbf{e}_t^T
 - Disadvantage: every vector is twice as large, and every matrix is four times as large
- Possibility #2: Add it, i.e., $\mathbf{x}_t = \mathbf{x}_t + \mathbf{e}_t$
 - Advantage: fewer parameters to learn
 - Disadvantage: $\mathbf{W}_{j,Q}$ can only operate directly on \mathbf{e}_t if \mathbf{x}_t is mostly zero
 - Surprise: this works well in practice. Apparently, the positional encoding can learn to ignore local fluctuations in \mathbf{x}_t , and pretend that it's mostly 0 on average

Positional encoding

In the standard transformer, position of the input is encoded using

$$\mathbf{x}_t = \mathbf{x}_t + \begin{bmatrix} \cos\left(\frac{\pi t}{T}\right) \\ \sin\left(\frac{\pi t}{T}\right) \\ \vdots \\ \cos\left(\frac{\pi Dt}{2T}\right) \\ \sin\left(\frac{\pi Dt}{2T}\right) \end{bmatrix}$$



Outline

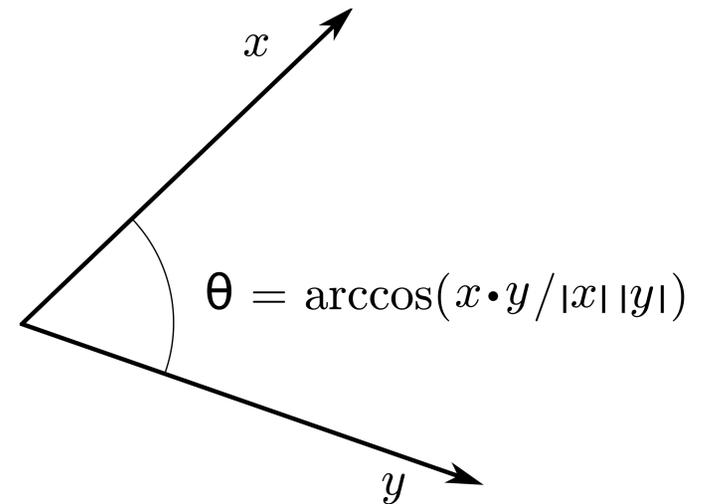
- Text generation
- Attention
- Multi-headed attention
- Positional embedding
- Layer norm
- Residual connections

Dot product = cosine if vectors are unit-norm

$$\mathbf{q}_i^T \mathbf{k}_t = \cos(\mathbf{q}_i, \mathbf{k}_t) \text{ iff } |\mathbf{q}_i| = 1 \text{ and } |\mathbf{k}_t| = 1.$$

Incorrect normalization is sometimes useful, e.g.: to include \mathbf{v}_t regardless of the query, you can make $|\mathbf{k}_t|$ large.

Often, though, allowing very large or very small vectors has bad effects.



By BenFrantzDale at the English Wikipedia, CC BY-SA 3.0,
<https://commons.wikimedia.org/w/index.php?curid=49972362>

How can we normalize the vectors?

- Normalization: subtract the mean, divide by standard deviation
- Batch normalization: Each neuron is normalized over the batch
- Instance normalization: Each instance is normalized over neurons
- Layer normalization: Normalize over the whole input

How can we normalize the vectors?

$$\mathbf{O} = \begin{bmatrix} O_{1,1} & \cdots & O_{1,d} \\ \vdots & \ddots & \vdots \\ O_{n,1} & \cdots & O_{n,d} \end{bmatrix}$$

- Normalization: subtract the mean, divide by standard deviation

$$O_{i,k} \leftarrow \frac{O_{i,k} - \mu}{\sigma}$$

- Batch normalization: Each neuron is normalized over the batch
- Instance normalization: Each instance is normalized over neurons
- Layer normalization: Normalize over the whole input

How can we normalize the vectors?

- Batch normalization: Each neuron is normalized over the batch

$$\mu = \frac{1}{n} \sum_{i=1}^n o_{i,k}, \quad \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (o_{i,k} - \mu)^2}$$

- Instance normalization: Each instance is normalized over neurons

$$\mu = \frac{1}{d} \sum_{k=1}^d o_{i,k}, \quad \sigma = \sqrt{\frac{1}{d} \sum_{k=1}^d (o_{i,k} - \mu)^2}$$

- Layer normalization: Normalize over the whole input

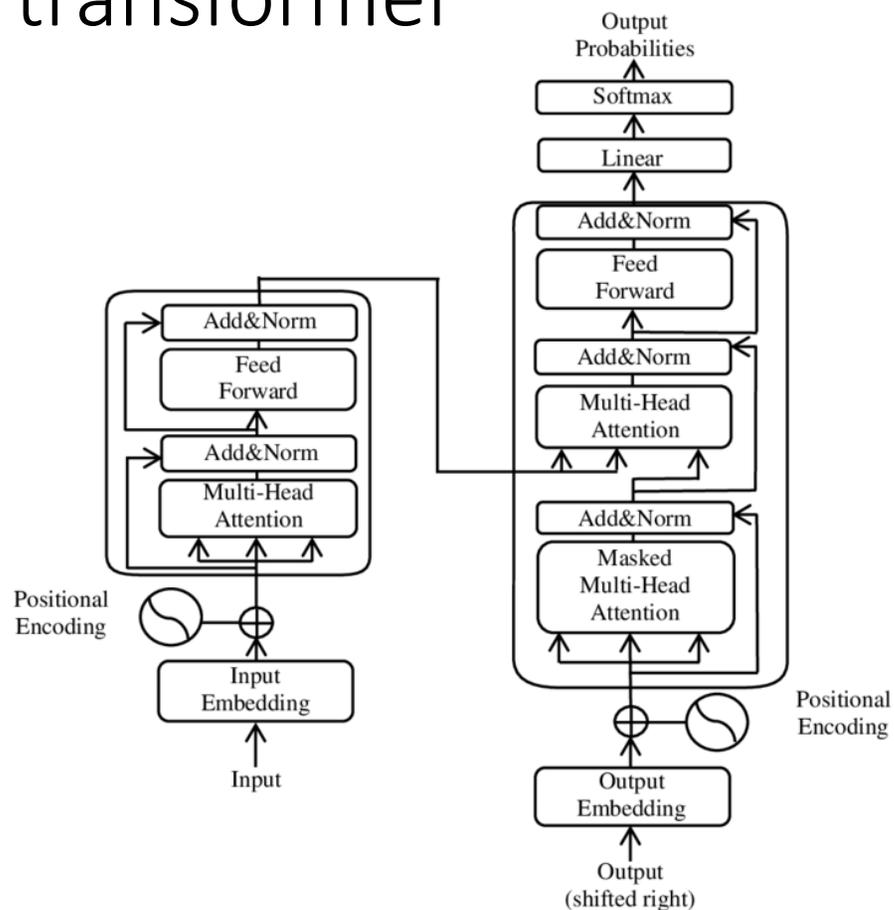
$$\mu = \frac{1}{nd} \sum_{i=1}^n \sum_{k=1}^d o_{i,k}, \quad \sigma = \sqrt{\frac{1}{nd} \sum_{i=1}^n \sum_{k=1}^d (o_{i,k} - \mu)^2}$$

How can we normalize the vectors?

- Batch normalization: $E[o_{i,k}] = 0$ and $E[o_{i,k}^2] = 1$ for each k , separately, over the whole batch
- Instance normalization: Each vector $\mathbf{o}_i = [o_{i,1}, \dots, o_{i,d}]$ is guaranteed to have a constant norm, $|\mathbf{o}_i|^2 = d$
- Layer normalization: Each vector has a constant norm on average, i.e., $E[|\mathbf{o}_i|^2] = d$
- If \mathbf{W}_K and \mathbf{W}_V are also pretty well normalized, then $E[|\mathbf{k}_t|^2] \approx E[|\mathbf{v}_t|^2] \approx d$
- The query vector often comes from a different source (the “query”), so it’s reasonable to assume that it’s normalized differently, $E[|\mathbf{q}_i|^2] \approx 1$
- Hence $\cos(\mathbf{k}_t, \mathbf{q}_i) \approx \mathbf{q}_i^T \mathbf{k}_t / \sqrt{d}$

Layer norm in the standard transformer

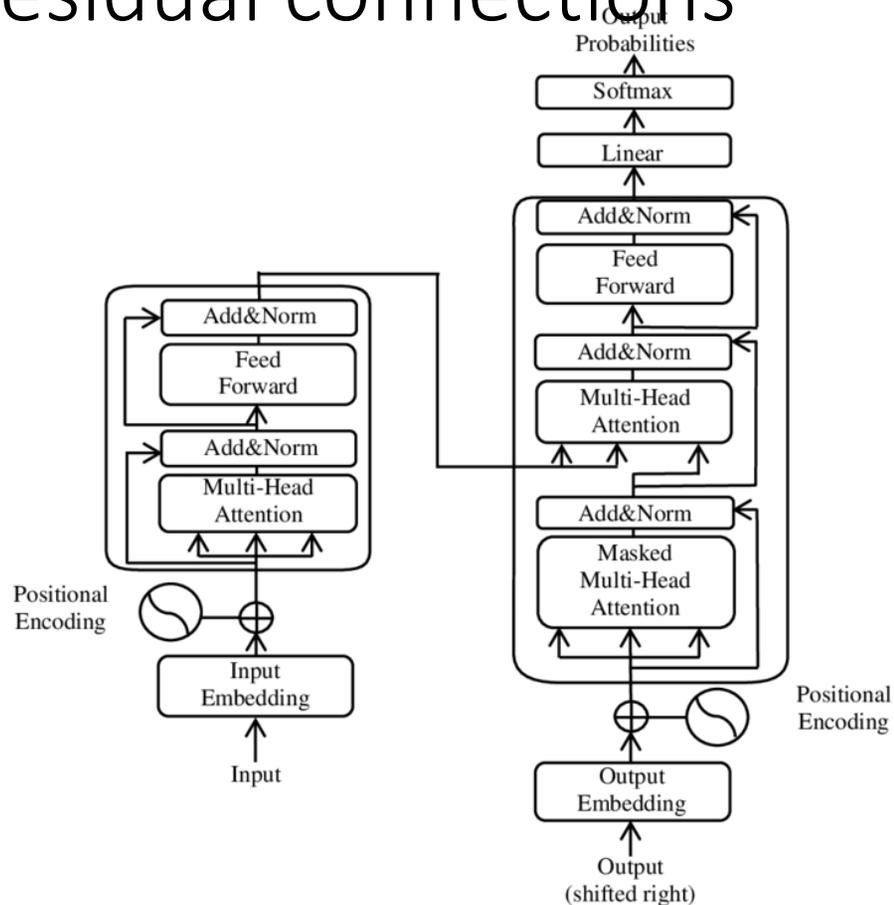
- In the “add & norm” layers, the “norm” is usually a layer norm
- What is the “add”?



Outline

- Text generation
- Attention
- Multi-headed attention
- Positional embedding
- Layer norm
- **Residual connections**

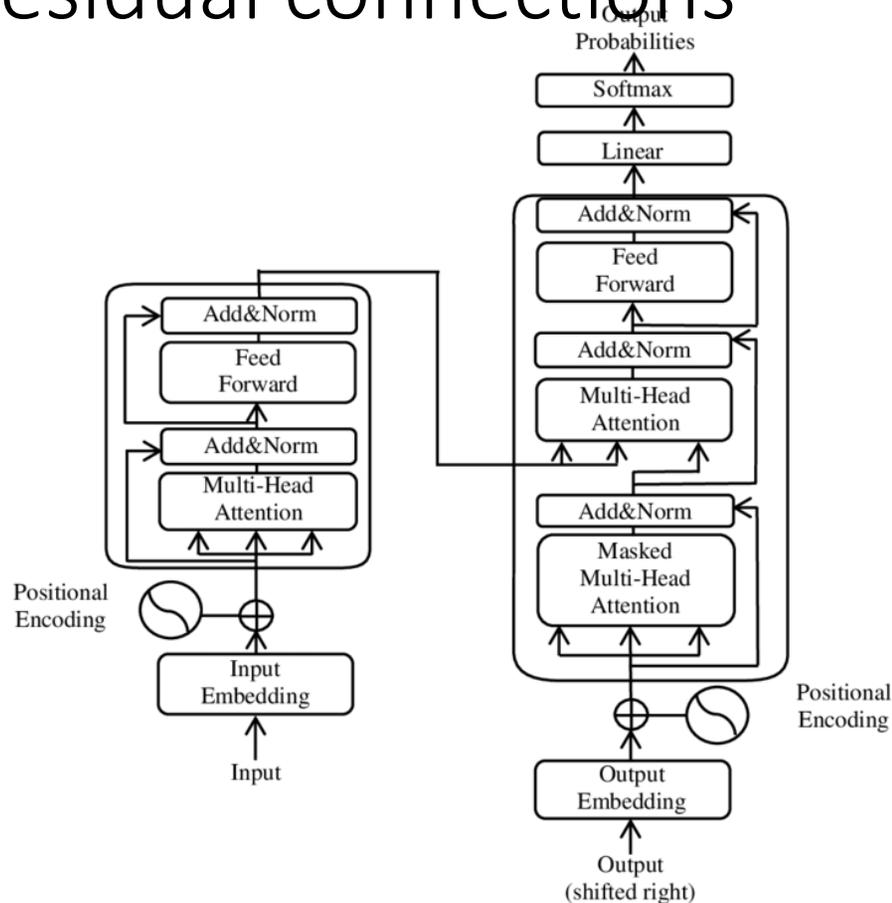
Residual connections



By Yuening Jia - DOI:10.1088/1742-6596/1314/1/012186, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=121340680>

- Most transformers have multiple layers
- For example, this diagram shows:
 - 2 layers of self-attention in the encoder (self-attention: Q, K, and V all come from the same input)
 - 1 self-attention layer + 1 cross-attention layer in the decoder (cross-attention: Q is from the previous output word, K and V are from the input words)
- More typical would be 12 layers in the encoder and 24 layers in the decoder

Residual connections



- In a residual connection, the input to the l^{th} layer, $\mathbf{X}^{(l)}$, is

$$\mathbf{X}^{(l)} = \mathbf{X}^{(l-1)} + \mathbf{O}^{(l-1)}$$

- Thus the output of each layer is only a small modification of its input (a “residual”)
- The output of the whole 24-layer transformer is just the input words, plus residuals computed by all layers
 - Residuals can be small, if the input words are all that you need
 - Residuals can be large if the input words are not enough

The standard transformer neural net

- Multi-headed dot-product attention: $\mathbf{C}_j = \text{softmax}\left(\frac{\mathbf{Q}_j \mathbf{K}_j^T}{\sqrt{d}}\right) \mathbf{V}_j$
- Concatenate and transform: $\mathbf{O} = [\mathbf{C}_1, \dots, \mathbf{C}_h] \mathbf{W}_O$
- Rotary Positional Encoding (RoPE):
 $\mathbf{x}_t += [\cos(\pi t/T), \sin(\pi t/T), \cos(2\pi t/T), \dots]^T$
- LayerNorm:

$$\mu = \frac{1}{nd} \sum_{i=1}^n \sum_{k=1}^d o_{i,k}, \quad \sigma = \sqrt{\frac{1}{nd} \sum_{i=1}^n \sum_{k=1}^d (o_{i,k} - \mu)^2}$$

- Residual Connections: $\mathbf{X}^{(l)} = \mathbf{X}^{(l-1)} + \mathbf{O}^{(l-1)}$