# Lecture 29: Transformers

Mark Hasegawa-Johnson

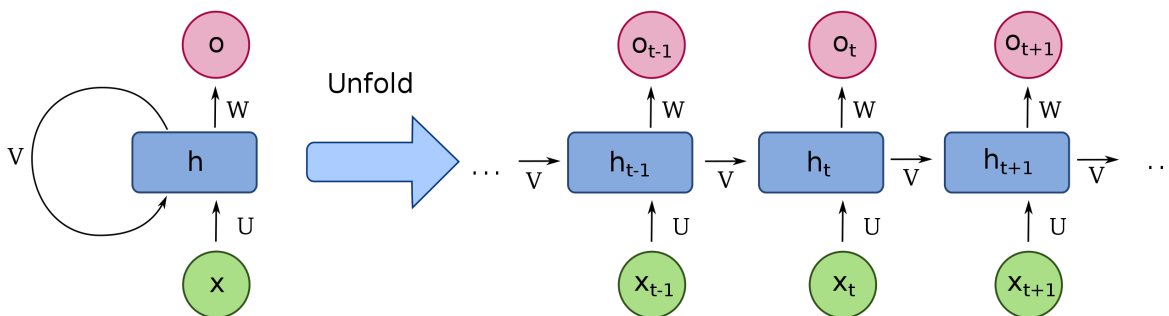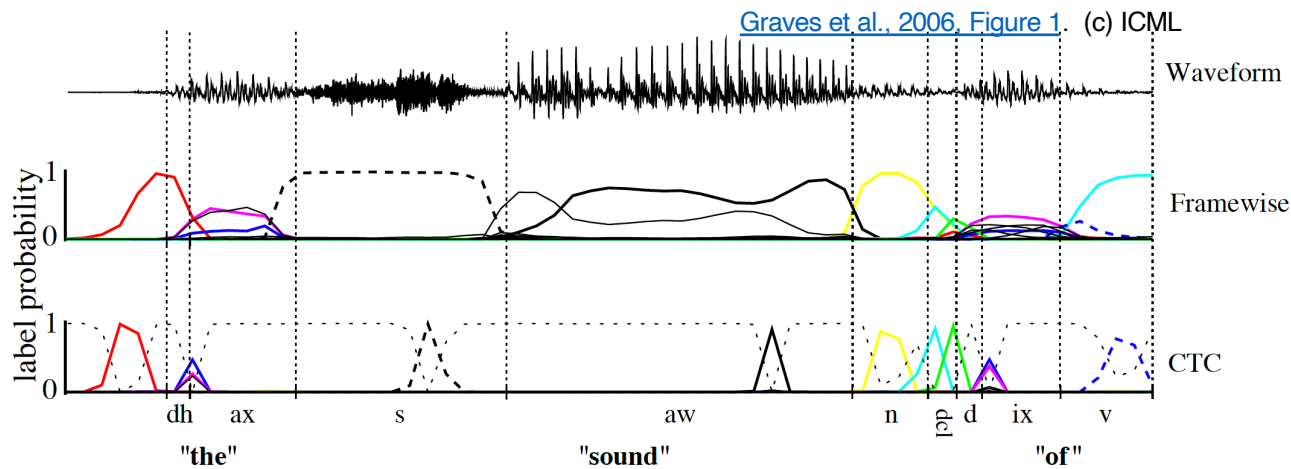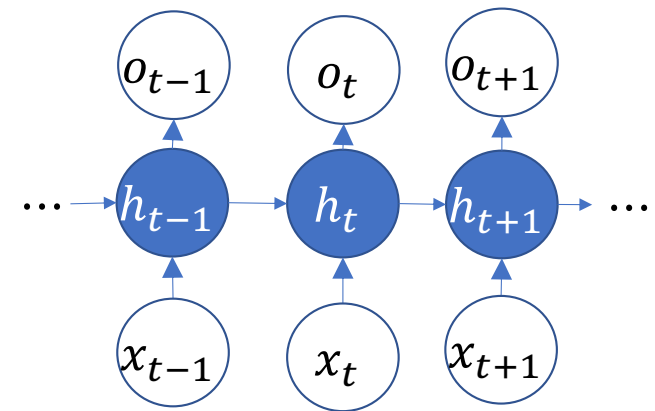4/2024

# Outline

- Recurrent neural networks
- Attention
- Self-attention, Multi-headed attention, Cross-attention, and Masked attention
- Positional embedding

# Recurrent neural network



Graves et al., 2006, Figure 1. (c) ICML

By fdeloche - Own work, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=60109157
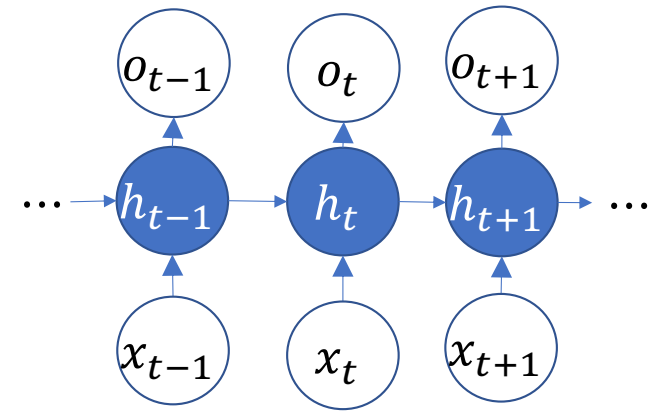
- In a recurrent neural network (RNN), the hidden node activation vector, $h_t$, depends on the value of the same vector at time $t - 1$.

- From 2014-2017, the best speech recognition and machine translation used RNNs.

- The input is $x_t$=speech or input-language text

- The output is $o_t$=text in the target language

# Example: Part of speech tagging



- $x_t$ = vector representation of the t$^{th}$ word, e.g., trained using CBOW
- $h_t$ = hidden state vector = $\tanh(\boldsymbol{U}\boldsymbol{x}_t + \boldsymbol{V}\boldsymbol{h}_{t-1})$
- $\boldsymbol{o}_t = \text{softmax}(\boldsymbol{W}\boldsymbol{h}_t) = [P(Y_t = \text{Noun}|X_1, \ldots, X_t), P(Y_t = \text{Verb}|X_1, \ldots, X_t), \ldots]$

# Training an RNN



An RNN is trained using gradient descent, just like any other neural network!

$$u_{j,i} \leftarrow u_{j,i} - \eta \frac{\partial \mathcal{L}}{\partial u_{j,i}}$$
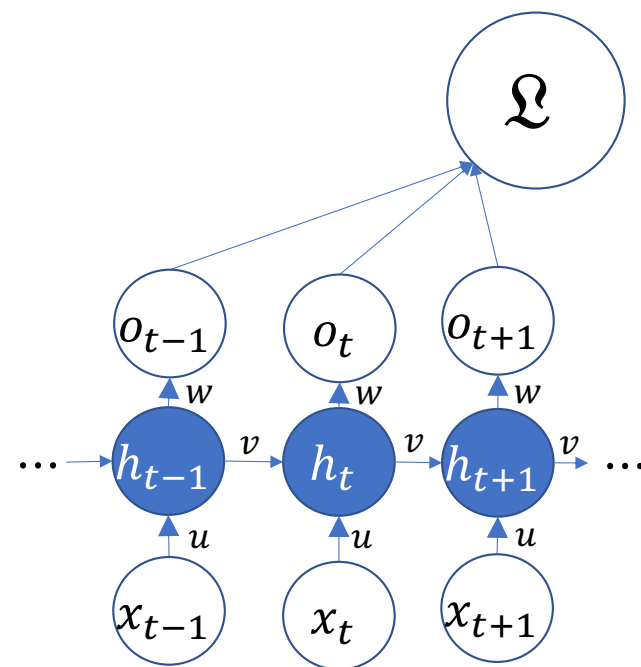
$$w_{j,k} \leftarrow w_{j,k} - \eta \frac{\partial \mathcal{L}}{\partial w_{j,k}}$$

…where $\mathcal{L}$ is the loss function, and $\eta$ is a step size.

# Training an RNN: Infinite recursion?

The big difference is that now the loss function depends on $U, V$ and $W$ in many different ways:

- The loss function depends on each of the state vectors $h_t$, which depends directly on $U$ and $V$.

- But $h_t$ also depends on $h_{t-1}$, which, in turn, depends on $U$ and $V$.

- ... and so on.

# Back-propagation through time

The solution is something called back-propagation through time:

$$\frac{d\mathfrak{L}}{dh_{i,t}} = \frac{\partial \mathfrak{L}}{\partial h_{i,t}} + \sum_{j} \frac{d\mathfrak{L}}{dh_{j,t+1}} \frac{\partial h_{j,t+1}}{\partial h_{i,t}}$$

- The first term measures losses caused directly by $h_{i,t}$, for example, if $o_{i,t}$ is wrong.

- The second term measures losses caused indirectly, for example, because $h_{i,t}$ caused $h_{j,t+1}$ to be wrong.

# Back-propagation through time

Notice that this is just like training a very deep network!

- Back-propagation through time: back-propagate from time step $t + 1$ to time step $t$

- Back-propagation in a very deep network: back-propagate from layer $l + 1$ to layer $l$

Toolkits like PyTorch may use the same code in both cases.

# Outline

- Recurrent neural networks
- Attention
- Self-attention, Multi-headed attention, Cross-attention, and Masked attention
- Positional embedding

# 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.
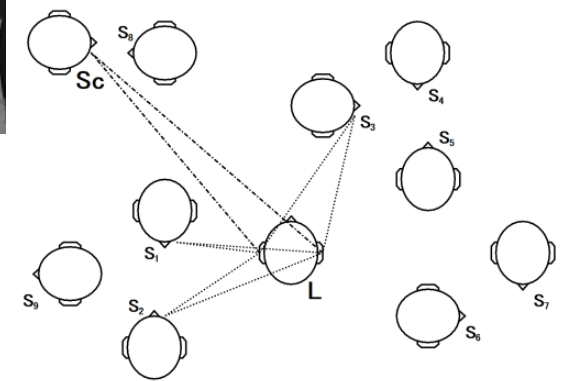
- This "cocktail-party effect" suggests a model of hearing in which all sounds are processed preconsciously. Trigger sounds in an unattended source will cause attention to re-orient to that source.

# Bottom-up attention as a strategy for machine listening

- In 2014, researchers proposed that the past 200ms of RNN state vectors should be stored in a "short-term memory buffer"
- A speech recognizer can attend to several centiseconds, all at one time, to decide what words it thinks it is hearing



FDHC0_SX209: Michael colored the bedroom wall with crayons.

Chorowski, Bahdanau, Serdyk, Cho & Bengio, Attention-Based Models for Speech Recognition, Fig. 1

# The Transformer: "Attention is all you need"

- In 2017, researchers proposed that the short-term memory buffer should contain raw signals, not processed signals.
- All processing is done using a model of bottom-up attention.

$\pi = [ \; - \quad w \quad i \quad i \quad - \quad - \quad t \quad h \quad - \; ]$

# Attention: Key concepts

- The neural net needs to make a series of decision vectors, $\boldsymbol{o}_i$
- Each decision needs to be based on some context vector, $\boldsymbol{c}_i$
- Each context vector is a weighted sum of input values, $\boldsymbol{c}_i = \sum_t \alpha_{i,t} \boldsymbol{v}_t$
- $\alpha_{i,t}$ is the amount of attention that the output decision $\boldsymbol{o}_i$ is paying to the input value $\boldsymbol{v}_t$. It is based on the similarity between a key vector, $\boldsymbol{k}_t$, that describes the type of information available in $\boldsymbol{v}_t$, and a query vector, $\boldsymbol{q}_i$, that describes the type of information necessary in order to make the output decision

# Inputs to an attention network

- Neural net inputs: a sequence of row vectors, $\boldsymbol{x}_t$
- Neural net outputs: a sequence of row vectors, $\boldsymbol{o}_i$
- Value: What type of information should $\boldsymbol{x}_t$ provide to the output? This may be just a linear transform of $\boldsymbol{x}_t$, e.g.: $\boldsymbol{v}_t = \boldsymbol{W}_V \boldsymbol{x}_t$
- Query: What type of information does $o_i$ need? This may be just a linear transform of $\boldsymbol{o}_{i-1}$, e.g.: $\boldsymbol{q}_i = \boldsymbol{W}_Q \boldsymbol{o}_{i-1}$
- Key: The dot product $\boldsymbol{q}_i^T \boldsymbol{k}_t$ should be positive if $\boldsymbol{v}_t$ is useful, and negative if $\boldsymbol{v}_t$ is useless. This may be $\boldsymbol{k}_t = \boldsymbol{W}_K \boldsymbol{x}_t$

# Attention = a probability mass over time

- Attention is like probability: You only have a fixed amount of attention, so you need to decide how to distribute it.

- $\alpha_{i,t} = P(\boldsymbol{v}_t|q_i)$ = the probability that $\boldsymbol{v}_t$ is the context that you need in order to make a decision related to the query vector $\boldsymbol{q}_i$.

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

- Each output context vector ($\boldsymbol{c}_i$) is based on some input value vectors ($\boldsymbol{v}_t$). But which ones? Answer: decide which inputs to pay attention to, then pay attention.

$$\boldsymbol{c}_i = \sum_t \alpha_{i,t}\boldsymbol{v}_t$$

# 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, $\boldsymbol{k}_t$, such that $\boldsymbol{q}_i^T \boldsymbol{k}_t > 0$ if $v_t$ is relevant to $\boldsymbol{q}_i$, otherwise $\boldsymbol{q}_i^T \boldsymbol{k}_t < 0$.

2. Convert the similarity measures into a probability distribution using softmax:

$$\alpha_{i,t} = \frac{\exp\!\left(\boldsymbol{q}_i^T \boldsymbol{k}_t\right)}{\sum_\tau \exp\!\left(\boldsymbol{q}_i^T \boldsymbol{k}_\tau\right)}$$

$x$

$\theta = \arccos(x \cdot y / |x| |y|)$

$y$

# Putting it all together

- Stack up $\boldsymbol{v}_t$, $\boldsymbol{k}_t$, and $\boldsymbol{q}_i$ into matrices:

$$V = \begin{bmatrix} \boldsymbol{v}_1^T \\ \vdots \\ \boldsymbol{v}_n^T \end{bmatrix}, K = \begin{bmatrix} \boldsymbol{k}_1^T \\ \vdots \\ \boldsymbol{k}_n^T \end{bmatrix}, Q = \begin{bmatrix} \boldsymbol{q}_1^T \\ \vdots \\ \boldsymbol{q}_m^T \end{bmatrix}$$

- $\alpha_{i,t}$ is the t$^{\text{th}}$ output of a softmax whose input vector is $\boldsymbol{Kq}_i$:

$$\alpha_{i,t} = \text{softmax}_t(\boldsymbol{Kq}_i) = \frac{\exp(\boldsymbol{q}_i^T \boldsymbol{k}_t)}{\sum_\tau \exp(\boldsymbol{q}_i^T \boldsymbol{k}_\tau)}$$

- $c_i$ is the product of the vector $\text{softmax}(\boldsymbol{Kq}_i)$ times the $\boldsymbol{V}^T$ matrix:

$$\boldsymbol{c}_i = \boldsymbol{V}^T \text{softmax}(\boldsymbol{Kq}_i) = \sum_t \alpha_{i,t} \boldsymbol{v}_t$$

# Quiz!

Try the quiz!

https://us.prairielearn.com/pl/course_instance/147925/assessment/2412318
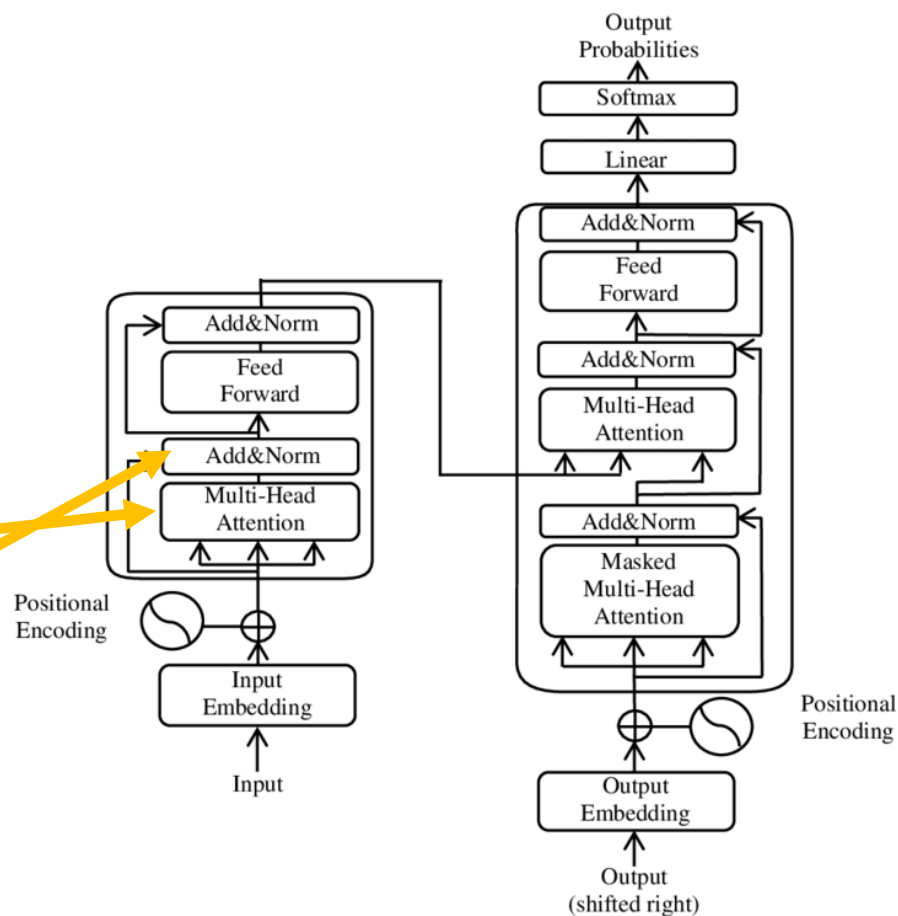
# Outline

- Recurrent neural networks
- Attention
- Self-attention, Multi-headed attention, Cross-attention, and Masked attention
- Positional embedding

# Self-attention

Self-attention (literally!) adds context to each input vector:

$$q_i = W_Q x_i$$
$$k_t = W_K x_t$$
$$v_t = W_V x_t$$
$$c_i = V^T \text{softmax}(K q_i)$$
$$y_i = \frac{x_i + c_i - E[x_i + c_i]}{\sqrt{Var(x_i + c_i)}}$$



Output Probabilities

Softmax

Linear

Add&Norm

Feed Forward

Add&Norm

Feed Forward

Add&Norm

Multi-Head Attention

Add&Norm

Multi-Head Attention

Add&Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Input

Output (shifted right)

# Multi-headed-attention

Multi-headed-attention uses 8 different $w_Q$, $w_K$, and $w_V$ matrices, in order to get 8 different views of the input data:

$$q_{j,i} = W_{j,Q} x_{j,i}, \qquad 1 \leq j \leq 8$$

$$k_{j,t} = W_{j,K} x_{j,t}, \qquad 1 \leq j \leq 8$$

$$v_{j,t} = W_{j,V} x_{j,t}, \qquad 1 \leq j \leq 8$$

$$h_{j,i} = V_j^T \operatorname{softmax}(K_j q_{j,i})$$

$$c_i = W_{j,O} \begin{bmatrix} h_{1,i} \\ \vdots \\ h_{8,i} \end{bmatrix}$$



Multi-Head Attention

# Cross-attention

Cross-attention: query depends on preceding output, key and value depend on input:

$$q_{j,i} = W_{j,Q} o_{j,i-1}$$
$$k_{j,t} = W_{j,K} x_{j,t}$$
$$v_{j,t} = W_{j,V} x_{j,t}$$
$$h_{j,i} = V_j^T \text{softmax}(K_j q_{j,i})$$

# Masked attention

Masked attention forces $c_i$ to pay attention to value vectors $v_t$ only if $t < i$:

$$s(\boldsymbol{q}_i, \boldsymbol{k}_t) = \begin{cases} \boldsymbol{q}_i^T \boldsymbol{k}_t & t < i \\ -\infty & t \geq i \end{cases}$$

$$\alpha_{i,t} = \frac{\exp\big(s(\boldsymbol{q}_i, \boldsymbol{k}_t)\big)}{\sum_\tau \exp\big(s(\boldsymbol{q}_i, \boldsymbol{k}_\tau)\big)}$$

$$= \begin{cases} \text{softmax}(\boldsymbol{q}_i^T \boldsymbol{k}_t) & t < i \\ 0 & t \geq i \end{cases}$$

# 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

# Word Error Rates using Transformers

By 9/2020, transformers had error rates of:

- 2%: English, quiet recording conditions
- 4%: Chinese or Japanese, quiet recording conditions
- 5-7%: if the reference transcript has errors
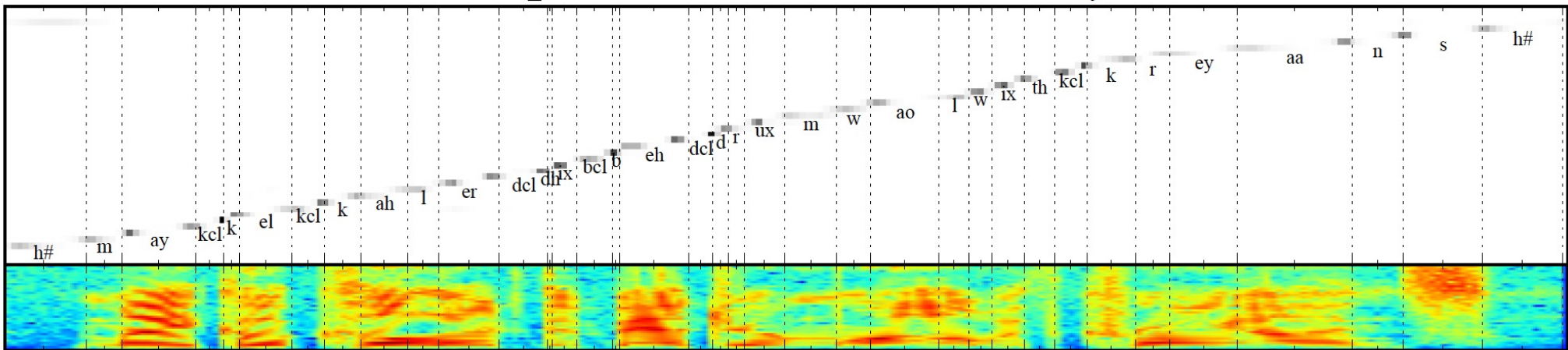- 14%: 2-talker mixtures, synthetic reverberation
- 38%: actual in-home recordings in noisy households

**Table 1**. CER/WER results on various open source ASR corpora. Both Transformer and Conformer models are implemented based on ESPnet toolkit. ∗ marks ESPnet2 results. † and ‡ indicate only w/ speed or only w/ SpecAugment, respectively. § denotes w/o any data augmentation.

| Dataset | Vocab | Metric | Evaluation Sets | Transformer | Conformer |
|---|---|---|---|---|---|
| AIDATATANG | Char | CER | dev / test | (†) 5.9 / 6.7 | **4.3 / 5.0** |
| AISHELL-1 | Char | CER | dev / test | (†) 6.0 / 6.7 | (∗) **4.4 / 4.7** |
| AISHELL-2 | Char | CER | android / ios / mic | (†) 8.9 / 7.5 / 8.6 | **7.6 / 6.8 / 7.4** |
| AURORA4 | Char | WER | dev_0330 (A / B / C / D) | **3.3 / 6.0 / 4.5** / 10.6 | 4.3 / 6.0 / 5.4 / **9.3** |
| CSJ | Char | CER | eval{1, 2, 3} | (∗) 4.7 / 3.7 / 3.9 | (∗) **4.5 / 3.3 / 3.6** |
| CHiME4 | Char | WER | {dt05, et05}_{simu, real} | (†) 9.6 / 8.2 / 15.7 / 14.5 | **9.1 / 7.9 / 14.2 / 13.4** |
| Fisher-CallHome | BPE | WER | dev / dev2 / test / devtest / evltest | 22.1 / 21.5 / 19.9 / 38.1 / 38.2 | **21.5 / 21.1 / 19.4 / 37.4 / 37.5** |
| HKUST | Char | CER | dev | (†) 23.5 | (†) **22.2** |
| JSUT | Char | CER | our split | (†) 18.7 | **14.5** |
| LibriSpeech | BPE | WER | {dev, test}_{clean, other} | 2.1 / 5.3 / 2.5 / 5.5 | **1.9 / 4.9 / 2.1 / 4.9** |
| REVERB | Char | WER | et_{near, far} | (†) 13.1 / 15.4 | (†) **10.5 / 13.9** |
| Switchboard | BPE | WER | eval2000 (callhm / swbd) | 17.2 / 8.2 | **14.0 / 6.8** |
| TEDLIUM2 | BPE | WER | dev / test | 9.3 / 8.1 | **8.6 / 7.2** |
| TEDLIUM3 | BPE | WER | dev / test | 10.8 / 8.4 | **9.6 / 7.6** |
| VoxForge | Char | CER | our split | (§) 9.4 / 9.1 | (§) **8.7 / 8.2** |
| WSJ | BPE | WER | dev93/ eval92 | (‡) **7.4 / 4.9** | (‡) 7.7 / 5.3 |
| WSJ-2mix | Char | WER | tt | (§) 12.6 | (§) **11.7** |

Guo, Boyer, Chang, Hayashi, Higuchi et al., ICASSP 2021, © IEEE

# Outline

- Recurrent neural networks
- Attention
- Self-attention, Multi-headed attention, Cross-attention, and Masked attention
- **Positional Encoding**

# What we have lost…

- With the recurrent neural net, each state vector paid attention to the one that preceded it:

$$\boldsymbol{h}_t = \tanh(\boldsymbol{U}\boldsymbol{v}_t + \boldsymbol{V}\boldsymbol{h}_{t-1})$$

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

$$\boldsymbol{h}_i = \sum_t \alpha_{i,t}\, \boldsymbol{v}_t, \qquad \alpha_{i,t} = \frac{\exp\left(\boldsymbol{q}_i^T \boldsymbol{k}_t\right)}{\sum_\tau \exp\left(\boldsymbol{q}_i^T \boldsymbol{k}_\tau\right)}$$
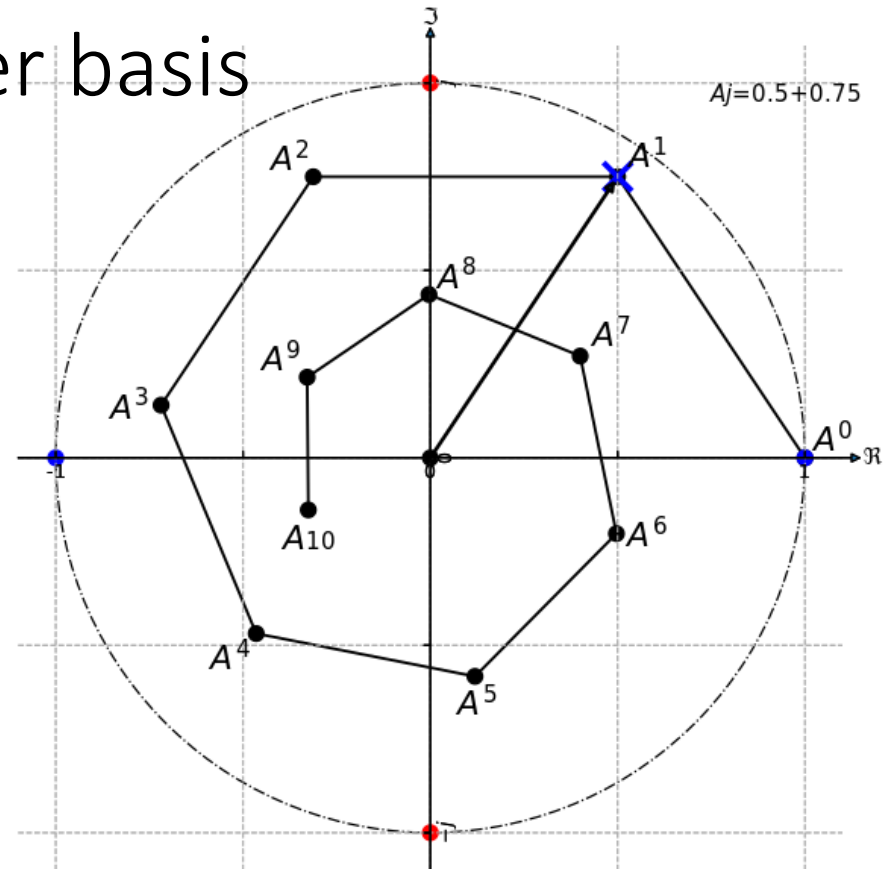
- What if we always want $\boldsymbol{h}_t$ to pay special attention to $\boldsymbol{v}_{t-1}$? Is that possible?

# Position encoding by Fourier basis

The solution is to encode the relative position of each input, $v_t$, using a Fourier basis $e_t$:

$$e_t = \begin{bmatrix} \cos\left(\dfrac{\pi t}{T}\right) \\ \sin\left(\dfrac{\pi t}{T}\right) \\ \vdots \\ \cos\left(\dfrac{\pi D t}{2T}\right) \\ \sin\left(\dfrac{\pi D t}{2T}\right) \end{bmatrix}$$



The curves $y = P_0(\cos\theta), y = P_1(\cos\theta), \ldots y = P_7(\cos\theta)$.  (v. page 184.)

# Position encoding by Fourier basis

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

$$\boldsymbol{e}_{t-d} = \begin{bmatrix} \cos\left(\dfrac{\pi d}{T}\right) & \sin\left(\dfrac{\pi d}{T}\right) & \cdots \\ -\sin\left(\dfrac{\pi d}{T}\right) & \cos\left(\dfrac{\pi d}{T}\right) & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \boldsymbol{e}_t$$

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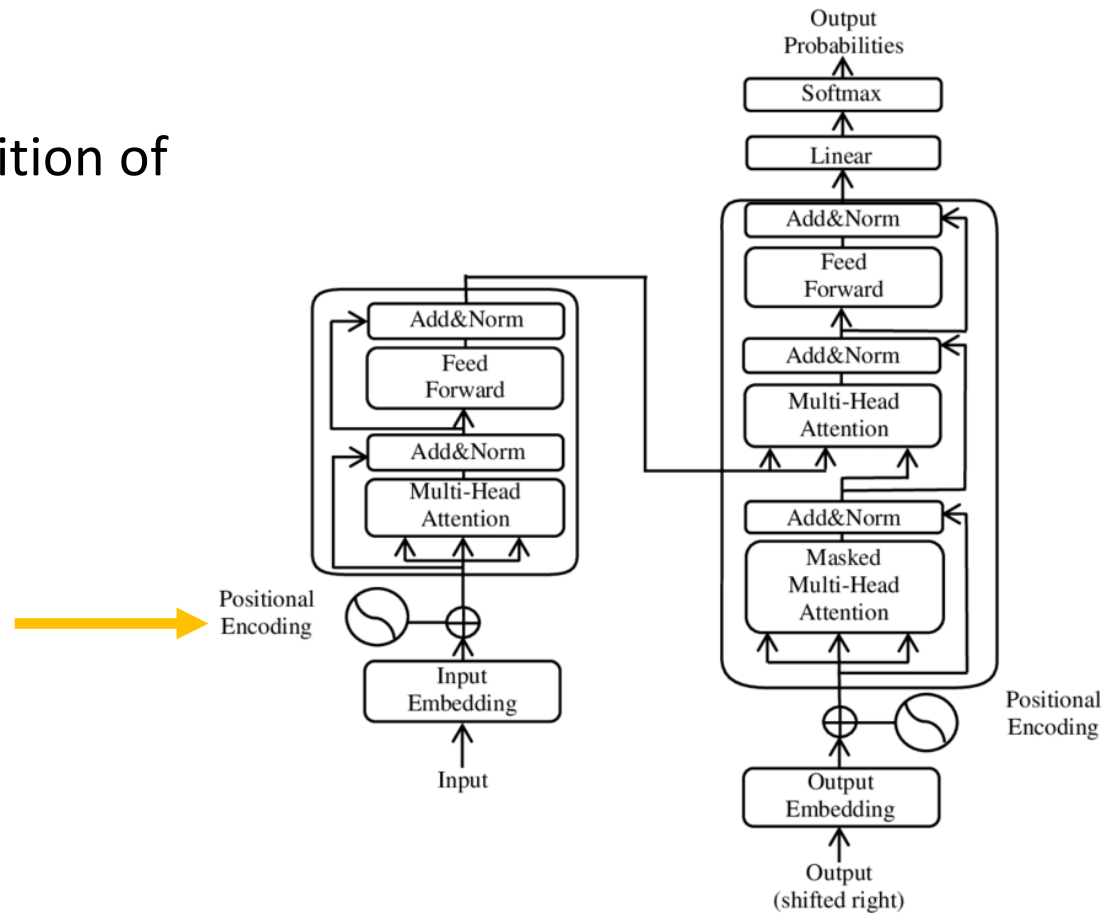


$Aj = 0.5 + 0.75$

# Where do we put the positional encoding?

- Possibility #1: Concatenate it, i.e., $x_t^T = [x_t^T, e_t^T]$
  - Advantage: $W_{j,Q}$ can learn to operate separately on the content $x_t^T$ and the positional encoding $e_t^T$
  - Disadvantage: every vector is twice as large, and every matrix is four times as large

- Possibility #2: Add it, i.e., $x_t = x_t + e_t$
  - Advantage: fewer parameters to learn
  - Disadvantage: $W_{j,Q}$ can only operate directly on $e_t$ if $x_t$ is mostly zero
  - Surprise: this works well in practice. Apparently, the positional encoding can learn to ignore local fluctuations in $x_t$, and pretend that it's mostly 0 on average

# Positional encoding

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

$$\boldsymbol{x}_t = \boldsymbol{x}_t + \begin{bmatrix} \cos\left(\dfrac{\pi t}{T}\right) \\ \sin\left(\dfrac{\pi t}{T}\right) \\ \vdots \\ \cos\left(\dfrac{\pi D t}{2T}\right) \\ \sin\left(\dfrac{\pi D t}{2T}\right) \end{bmatrix}$$

# Summary

- Recurrent neural networks

$$h_t = \tanh(Uv_t + Vh_{t-1})$$

- Attention

$$c_i = V^T \text{softmax}(Kq_i) = \sum_t \frac{\exp(q_i^T k_t)}{\sum_\tau \exp(q_i^T k_\tau)} v_t$$

- Self-attention, Multi-headed attention, Cross-attention, and Masked attention
- Positional encoding

$$x_t \mathrel{+}= \begin{bmatrix} \cos\left(\dfrac{\pi t}{T}\right) \\ \sin\left(\dfrac{\pi t}{T}\right) \\ \vdots \end{bmatrix}$$