
Discussion 11

Note: Your TA will probably not cover all the problems on this worksheet. The discussion worksheets are not designed to be finished within an hour. They are deliberately made slightly longer so they can serve as resources you can use to practice, reinforce, and build upon concepts discussed in lectures, discussions, and homework.

This Week's Cool AI Demo/Video:

Anthropic credit research: <https://www.youtube.com/watch?v=Y6wiWlch5jM>

Why aren't LLMs deterministic with temperature 0?:

<https://thinkingmachines.ai/blog/defeating-nondeterminism-in-llm-inference/>

1 Positional Encoding

- (a) Why do we need positional encoding? Describe a situation where word-order information is necessary for the task performed.
- (b) What is *relative* positional encoding? How is it different from *absolute* positional encoding?
- (c) **Rotary positional embeddings (RoPE)** are a way of encoding token positions by rotating pairs of embedding dimensions by an angle that grows with the position index, so that inner products between tokens automatically depend only on their relative positions. They've become a de-facto default in many modern LLMs.

Consider a sequence of N two-dimensional vectors

$$x_n := (x_n^{(1)}, x_n^{(2)}), \quad n = 1, 2, \dots, N.$$

For a parameter $\theta \in \mathbb{R}$, the RoPE encoding at position n is

$$\text{RoPE}(x_n, n) = R(n\theta)x_n, \quad \text{where} \quad R(\alpha) := \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix}.$$

Show that the dot product between two RoPE-encoded vectors depends only on their *relative* positions. In particular, prove that for any integers m , n , and k , the following holds:

$$\text{RoPE}(x, m)^\top \text{RoPE}(y, n) = \text{RoPE}(x, m+k)^\top \text{RoPE}(y, n+k).$$

(a) Why do we need positional encoding? Describe a situation where word-order information is necessary for the task performed.

- input embeddings from input tokens don't capture position of token in sequence
- self-attention mechanism is permutation invariant
- sometimes word order changes the meaning, e.g.,
 - Only John texted me on my birthday
 - John only texted me on my birthday

(b) What is *relative* positional encoding? How is it different from *absolute* positional encoding?

Consider example (a) from previous part:

absolute → we encode "only" as first word, "John" as second word, & so on

relative → we encode "only" is one word before "John" & such

* more generalizable

In lecture, we discussed absolute encodings (like from Attention paper), where x_n is an input embedding, r_n is that token's position, & we combine them as $\tilde{x}_n = x_n + r_n$.

- (c) **Rotary positional embeddings (RoPE)** are a way of encoding token positions by rotating pairs of embedding dimensions by an angle that grows with the position index, so that inner products between tokens automatically depend only on their relative positions. They've become a de-facto default in many modern LLMs.

→ similarity score QK^T is composed of inner products

Consider a sequence of N two-dimensional vectors

$$x_n := (x_n^{(1)}, x_n^{(2)}), \quad n = 1, 2, \dots, N.$$

For a parameter $\theta \in \mathbb{R}$, the RoPE encoding at position n is rotation matrix

$$\text{RoPE}(x_n, n) = R(n\theta)x_n, \quad \text{where} \quad R(\alpha) := \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix}.$$

Show that the dot product between two RoPE-encoded vectors depends only on their *relative* positions. In particular, prove that for any integers m, n , and k , the following holds:

$$\text{RoPE}(x, m)^\top \text{RoPE}(y, n) = \text{RoPE}(x, m+k)^\top \text{RoPE}(y, n+k).$$

$$x = W_k \tilde{x} \quad y = W_k \tilde{y}$$

$$\text{RoPE}(x, m) = R(m\theta)x$$

$$\text{RoPE}(x, m+k) = R((m+k)\theta)x$$

$$\text{RoPE}(y, n) = R(n\theta)y$$

$$\text{RoPE}(y, n+k) = R((n+k)\theta)y$$

$$\text{RoPE}(x, m)^\top \text{RoPE}(y, n) = x^\top R(m\theta)^\top R(n\theta)y$$

$$\text{RoPE}(x, m+k)^\top \text{RoPE}(y, n+k) = x^\top R((m+k)\theta)^\top R((n+k)\theta)y$$

$$R(\alpha)^\top R(\beta) = \begin{bmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{bmatrix} \begin{bmatrix} \cos \beta & -\sin \beta \\ \sin \beta & \cos \beta \end{bmatrix}$$

$$= \begin{bmatrix} \cos \alpha \cos \beta + \sin \alpha \sin \beta & -\cos \alpha \sin \beta + \sin \alpha \cos \beta \\ -\sin \alpha \cos \beta + \cos \alpha \sin \beta & \sin \alpha \sin \beta + \cos \alpha \cos \beta \end{bmatrix}$$

We need the product to sum formulas from trig:

$$\sin(\alpha)\sin(\beta) = \frac{1}{2}[\cos(\alpha-\beta) - \cos(\alpha+\beta)]$$

$$\cos(\alpha)\cos(\beta) = \frac{1}{2}[\cos(\alpha-\beta) + \cos(\alpha+\beta)]$$

$$\sin(\alpha)\cos(\beta) = \frac{1}{2}[\sin(\alpha+\beta) + \sin(\alpha-\beta)]$$

$$\cos(\alpha)\sin(\beta) = \frac{1}{2}[\sin(\alpha+\beta) - \sin(\alpha-\beta)]$$

$$R(\alpha)^T R(\beta) = \begin{bmatrix} \cos(\alpha-\beta) & \sin(\alpha-\beta) \\ -\sin(\alpha-\beta) & \cos(\alpha-\beta) \end{bmatrix}$$

B/c cosine is even $\cos(-\theta) = \cos\theta$

B/c sine is odd $\sin(-\theta) = -\sin\theta$

$$R(\alpha)^T R(\beta) = \begin{bmatrix} \cos(\beta-\alpha) & -\sin(\beta-\alpha) \\ \sin(\beta-\alpha) & \cos(\beta-\alpha) \end{bmatrix} = R(\beta-\alpha)$$

$$\begin{aligned} \text{RoPE}(x, m)^T \text{RoPE}(y, n) &= x^T R(m\theta)^T R(n\theta) y \\ &= x^T R(n\theta - m\theta) y \end{aligned}$$

$$\begin{aligned} \text{RoPE}(x, m+k)^T \text{RoPE}(y, n+k) &= x^T R((m+k)\theta)^T R((n+k)\theta) y \\ &= x^T R((n+k)\theta - (m+k)\theta) y \\ &= x^T R(n\theta - m\theta) y \end{aligned}$$

∴ Shifting the position of both tokens by the same amount, k , won't change the dot product of embeddings

2 Matching Similarity Matrices to Attention Heatmaps

Recall that in self-attention mechanisms, we compute the key-query similarity matrix QK^T , scale the matrix by $\frac{1}{\sqrt{D_k}}$, and then apply the softmax function row-wise to obtain attention scores. Understanding how pre-softmax similarity scores translate to post-softmax attention weights is crucial for interpreting attention patterns. Below are four 4×4 pre-softmax similarity matrices and four corresponding post-softmax attention heatmaps. Your task is to match each pre-softmax matrix to its corresponding heatmap and explain the transformation.

Pre-softmax matrices:

QK^T

Matrix A:
$$\begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \end{bmatrix}$$

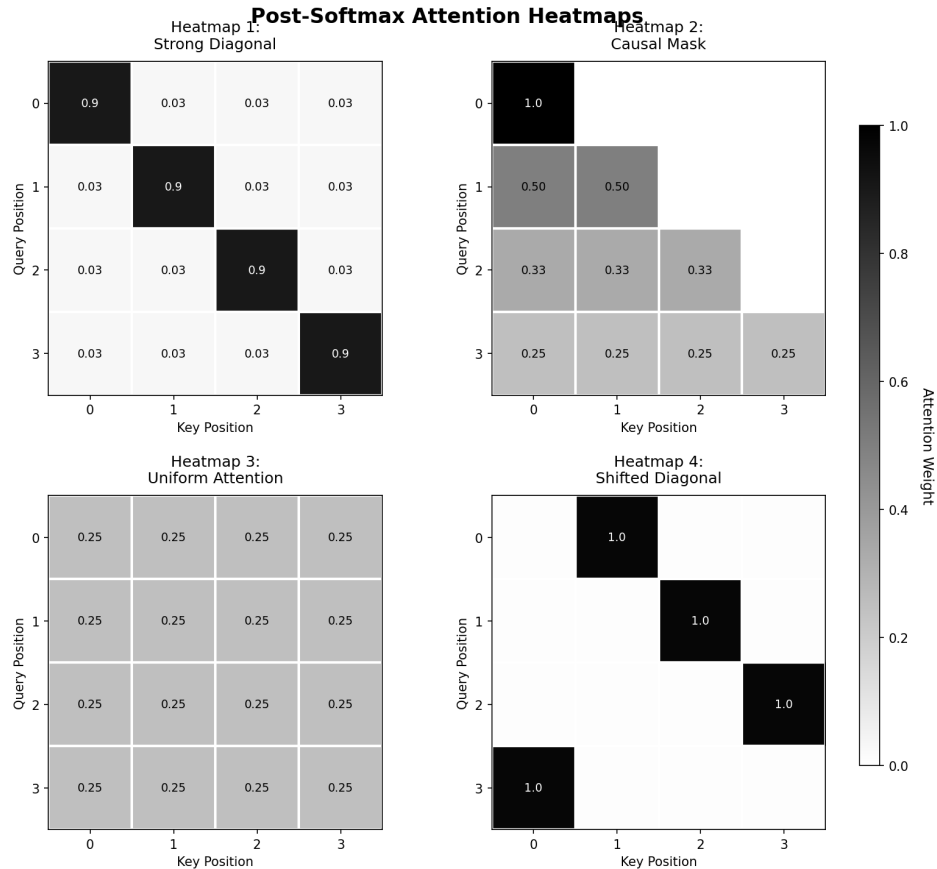
Matrix B:
$$\begin{bmatrix} 10 & 1 & 1 & 1 \\ 1 & 10 & 1 & 1 \\ 1 & 1 & 10 & 1 \\ 1 & 1 & 1 & 10 \end{bmatrix}$$

Matrix C:
$$\begin{bmatrix} 5 & -\infty & -\infty & -\infty \\ 5 & 5 & -\infty & -\infty \\ 5 & 5 & 5 & -\infty \\ 5 & 5 & 5 & 5 \end{bmatrix}$$

Matrix D:
$$\begin{bmatrix} -\infty & 5 & -\infty & -\infty \\ -\infty & -\infty & 5 & -\infty \\ -\infty & -\infty & -\infty & 5 \\ 5 & -\infty & -\infty & -\infty \end{bmatrix}$$

Post-softmax heat maps (darker = higher attention weight):

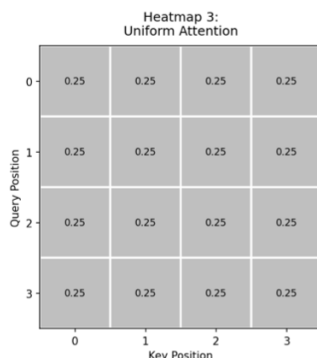
$\text{SoftMax}_{\text{row}}\left(\frac{QK^T}{\sqrt{D_k}}\right)$
 $D_k = \text{key / query dimension}$



- (a) Match each pre-softmax matrix (A-D) to its corresponding heatmap (1-4). Provide brief reasoning for each match.
- (b) Explain why Matrix C represents a causal attention mask. In what type of model would you typically see this pattern?
- (c) Consider a temperature-scaled softmax: $\text{softmax}(x/T)$ where T is temperature. How would changing T affect the attention distribution for Matrix B? What happens as $T \rightarrow 0$ and $T \rightarrow \infty$?

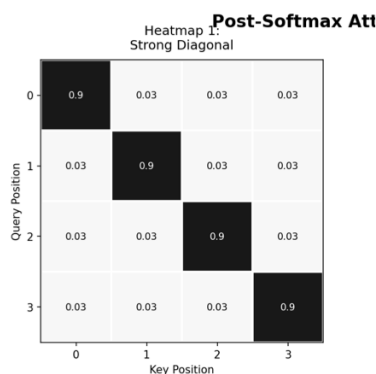
- (a) Match each pre-softmax matrix (A-D) to its corresponding heatmap (1-4). Provide brief reasoning for each match.

Matrix A:
$$\begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \end{bmatrix}$$



Why are all the values 0.25?
 we apply softmax row-wise, so each row sums to 1.

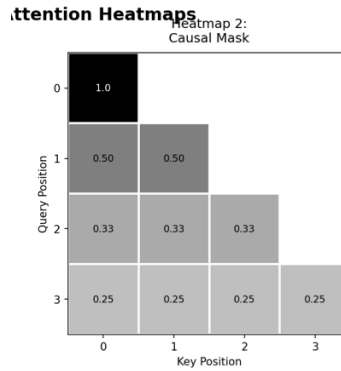
Matrix B:
$$\begin{bmatrix} 10 & 1 & 1 & 1 \\ 1 & 10 & 1 & 1 \\ 1 & 1 & 10 & 1 \\ 1 & 1 & 1 & 10 \end{bmatrix}$$



How did we get 0.9 along diagonal?
 we scale QK^T (matrix B) by $\frac{1}{\sqrt{D_k}}$ (which we don't provide) + then apply softmax row-wise.

Matrix C:

$$\begin{bmatrix} 5 & -\infty & -\infty & -\infty \\ 5 & 5 & -\infty & -\infty \\ 5 & 5 & 5 & -\infty \\ 5 & 5 & 5 & 5 \end{bmatrix}$$



How did we get zeros in upper triangle?

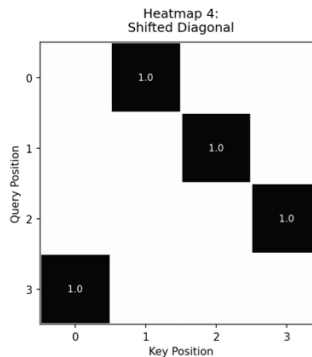
$$\frac{e^{-\infty}}{ne^{-\infty} + (4-n)e^5} = 0 \text{ for any finite integer } n$$

What about the lower triangle?

After setting upper triangle to 0, rows must sum to 1.

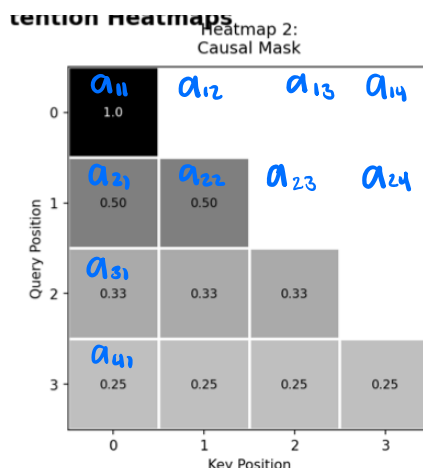
Matrix D:

$$\begin{bmatrix} -\infty & 5 & -\infty & -\infty \\ -\infty & -\infty & 5 & -\infty \\ -\infty & -\infty & -\infty & 5 \\ 5 & -\infty & -\infty & -\infty \end{bmatrix}$$



(b) Explain why Matrix C represents a causal attention mask. In what type of model would you typically see this pattern?

$$\text{Matrix C: } \begin{bmatrix} 5 & -\infty & -\infty & -\infty \\ 5 & 5 & -\infty & -\infty \\ 5 & 5 & 5 & -\infty \\ 5 & 5 & 5 & 5 \end{bmatrix}$$



By setting the upper triangle to $-\infty$, we force the attention values in those positions to be 0.

Recall a_{ni} is attention that token n places on token i . Why would we want $a_{ni} = 0$ for $n < i$?

Consider the example from discussion mini lecture.

During training, we have the following:

Inputs: where is the library <eos>

Outputs: <sos> Dónde está la biblioteca

Targets: Dónde está la biblioteca <eos>

During inference, we start w/ the following:

Inputs: where is the library <eos>

Outputs: <sos>

And we build the outputs one token at a time from our trained transformer's predictions.

To train a good (i.e., generalizable) model that can perform well during inference time, we don't want it to "cheat" during training. Token 1 should not be able to attend to token 2 b/c it won't be able to do this during inference. But token 4 can look at tokens 1-3 b/c it will have this info available during inference.

row-wise

(c) Consider a temperature-scaled softmax: $\text{softmax}(x/T)$ where T is temperature. How would changing T affect the attention distribution for Matrix B? What happens as $T \rightarrow 0$ and $T \rightarrow \infty$?

$$\text{Matrix B: } \begin{matrix} & \begin{matrix} 10 & 1 & 1 & 1 \end{matrix} \\ \begin{matrix} 1 \\ 1 \\ 1 \\ 1 \end{matrix} & \begin{matrix} 10 & 1 & 1 & 1 \\ 1 & 10 & 1 & 1 \\ 1 & 1 & 10 & 1 \\ 1 & 1 & 1 & 10 \end{matrix} \end{matrix} * \frac{1}{T}, \text{ where } T = \sqrt{D_k}$$

$$\text{Along the diagonal: } a_{\text{diag}} = \frac{e^{10/T}}{e^{10/T} + 3e^{1/T}}$$

$$\text{Off of diagonal: } a_{\text{off}} = \frac{e^{1/T}}{e^{10/T} + 3e^{1/T}}$$

Low temperature ($T \rightarrow 0^+$)

i.e., low key/query dimension, D_k :

$$\lim_{T \rightarrow 0^+} \frac{e^{10/T}}{e^{10/T} + 3e^{1/T}} \cdot \frac{e^{-10/T}}{e^{-10/T}} = \lim_{T \rightarrow 0^+} \frac{1}{1 + 3e^{-9/T}} = \frac{1}{1 + 3(0)} = 1$$

$$a_{\text{diag}} + 3a_{\text{off}} = 1$$

$$1 + 3a_{\text{off}} = 1 \rightarrow a_{\text{off}} = 0$$

$$\lim_{T \rightarrow 0^+} \frac{e^{1/T}}{e^{10/T} + 3e^{1/T}} \cdot \frac{e^{-10/T}}{e^{-10/T}} = \lim_{T \rightarrow 0^+} \frac{e^{-9/T}}{1 + 3e^{-9/T}} = \frac{0}{1 + 3(0)} = 0$$

\therefore The diagonal terms dominate off-diagonal ones.

High temperature ($T \rightarrow +\infty$)

i.e., high key/query dimension, D_k :

$$\lim_{T \rightarrow \infty} \frac{e^{10/T}}{e^{10/T} + 3e^{1/T}} = \frac{1}{1+3(1)} = \frac{1}{4}$$

$$a_{diag} + 3a_{off} = 1$$

$$\frac{1}{4} + 3a_{off} = 1 \rightarrow a_{off} = \frac{1}{4}$$

$$\lim_{T \rightarrow \infty} \frac{e^{1/T}}{e^{10/T} + 3e^{1/T}} = \frac{1}{1+3(1)} = \frac{1}{4}$$

\therefore All of the terms get equal weight (i.e., attention becomes uniform)

3 KV Caching in Autoregressive Generation

In autoregressive language models like GPT, we generate tokens one at a time, conditioning on all previously generated tokens. This problem explores how KV caching dramatically improves the runtime of this process.

The Naive Approach:

Consider generating a sequence of length N with a transformer that has:

- Model (hidden) dimension, D_{model}
- Number of attention heads, H
- Number of attention layers, L

Without caching, at iteration n when generating token y_n , we compute attention as:

$$\text{Attention}(Q_n, K_{1:n}, V_{1:n}) = \text{SoftMax}_{\text{row}} \left(\frac{Q_n K_{1:n}^T}{\sqrt{D_{model}/H}} \right) V_{1:n}$$

where $K_{1:n}$ and $V_{1:n}$ are the keys and values for all tokens from position 1 to n . Notice that we recompute K_i and V_i for all $i < n$ at every step, even though these values don't change.

The KV Cache Solution: We store the computed K and V matrices from previous iterations and reuse them, only computing the new K_n and V_n for the current token.

- Without KV caching**, how many times do we compute the key and value projections for the **first** token y_1 throughout the entire generation of a sequence of length N ?
- For a single layer, what is the total number of matrix multiplications needed to compute **all** key projections $K_{1:N}$ throughout the entire generation process **without KV caching**? (Assume we compute each K_i individually at each iteration, i.e. $K_i = y_i W_k$, instead of stacking tokens i up to n into an $n \times D_{model}$ matrix to compute $K = X W_k$ at each iteration.) Express your answer in terms of N .
- With KV caching**, how many matrix multiplications are needed to compute all key projections throughout the **entire** generation?
- Briefly explain why KV caching is particularly important for applications like chatbots or code assistants that involve multiple rounds of interaction.

Without caching, at iteration n when generating token y_n , we compute attention as:

$$\text{Attention}(Q_n, K_{1:n}, V_{1:n}) = \text{SoftMax}_{\text{row}} \left(\frac{Q_n K_{1:n}^T}{\sqrt{D_{\text{model}}/H}} \right) V_{1:n}$$

where $K_{1:n}$ and $V_{1:n}$ are the keys and values for all tokens from position 1 to n . Notice that we recompute K_i and V_i for all $i < n$ at every step, even though these values don't change.

- (a) **Without KV caching**, how many times do we compute the key and value projections for the **first** token y_1 throughout the entire generation of a sequence of length N ?

Every token in our sequence of length N will use the key/value of y_1 in the attention mechanism, so $k_1 + v_1$ would need to be computed N times if we don't cache them.

- (b) For a single layer, what is the total number of matrix multiplications needed to compute **all** key projections $K_{1:N}$ throughout the entire generation process **without KV caching**? (Assume we compute each K_i individually at each iteration, i.e. $K_i = y_i W_k$, instead of stacking tokens i up to n into an $n \times D_{\text{model}}$ matrix to compute $K = XW_k$ at each iteration.) Express your answer in terms of N .

At iteration 1, we need k_1 , so we have 1 mat. mul.

At iteration 2, we need $k_1 + k_2$, so we have 2 m.m.

⋮

At iteration N , we need k_1, k_2, \dots, k_N , so we have N

In total, we have $\sum_{n=1}^N n = \frac{1}{2} N(N+1) = O(N^2)$

matrix multiplications

- (c) **With KV caching**, how many matrix multiplications are needed to compute all key projections throughout the **entire** generation?

At iteration 1, we need k_1 , so we have 1 mat. mul.

At iteration 2, we need $k_1 + k_2$, but we cached

k_1 , so we only have 1 mat. mul. for k_2

⋮

At iteration N , we need k_1, k_2, \dots, k_N , but we

cached k_1, k_2, \dots, k_{N-1} , so we only have 1
mat. mul. for k_N

In total, we have exactly N matrix multiplications.

- (d) Briefly explain why KV caching is particularly important for applications like chatbots or code assistants that involve multiple rounds of interaction.