

CSE 493 G1 / 599 G1
Deep Learning
Winter 2026 Exam 2

SOLUTIONS

March 5, 2026

Full Name (as on Gradescope): _____

UW Net ID: _____

| Question | Score |
|--------------------------|-------|
| True/False (30 pts) | |
| Multiple Choice (32 pts) | |
| Free Response 1 (18 pts) | |
| Free Response 2 (20 pts) | |
| Total (100 pts) | |

Welcome to the CSE 493 G1 / 599 G1 Exam!

- The exam is 80 min and is **double-sided**.
- You may use a non-graphing calculator and no other electronic devices.
- One handwritten double sided cheat sheet is allowed.

I understand and agree to uphold the University of Washington Student Conduct Code during this exam.

Mean: 76.1
Median: 77.0
Stdev: 12.6

1 True/False (30 points) - Recommended 15 Minutes

Fill in the circle next to your choice (like this: ●). No explanations are required.

Each question is worth 3 points

- 1.1 Two 3×3 convolutions with D input and D output channels have more parameters than one 5×5 convolution with D input and D output channels.

- True
 False **SOLUTION:**

False. Two 3×3 convolutions have $2 \times (3 \times 3 \times D \times D) = 18D^2$ parameters. One 5×5 convolution has $5 \times 5 \times D \times D = 25D^2$ parameters. $18D^2 < 25D^2$.

- 1.2 Two models with the same number of parameters but different architectures will require roughly the same amount of memory to perform a forward and backward pass on the same data.

- True
 False **SOLUTION:**

False. Memory during training is dominated by storing activations (intermediate values) for backpropagation, not by the number of parameters. A deeper model with fewer parameters per layer can require significantly more activation memory than a shallow, wide model with the same total parameter count.

- 1.3 Generally, larger vocabulary sizes for tokenizers will lead to shorter sequence lengths.

- True
 False **SOLUTION:**

True. With a larger vocabulary, each token can represent a longer substring of text, so fewer tokens are needed to encode the same input.

- 1.4 The number of parameters in an RNN increases linearly with the sequence length.

- True
 False **SOLUTION:**

False, the number of params remains the same

- 1.5 Increasing the number of attention heads in a multi-head attention layer increases the total number of attention maps computed, since each head produces its own separate attention map.

- True
 False **SOLUTION:**

True

1.6 Imagine you are training a language model that consists of a single stack of attention and MLP blocks, rather than an encoder-decoder. You should use unmasked (bidirectional) attention during training so the model can see the full context, and then switch to causal (masked) attention at generation time so the model only attends to previously generated tokens.

True

False **SOLUTION:**

False, you need to use the same type during training and testing

1.7 If you want a network to produce two different types of outputs simultaneously, such as a probability distribution over class labels and bounding box coordinates, you can achieve this by attaching multiple separate fully connected heads to the same shared feature representation, each with its own output dimension and loss function

True

False **SOLUTION:**

True

1.8 A key idea in contrastive self-supervised learning (e.g., SimCLR) is that two augmented views of the same image should have similar representations, while views of different images should be pushed apart.

True

False **SOLUTION:**

True

1.9 Transfer learning only works when the source task and target task have the same set of output classes, because the learned features are specific to the classes the network was originally trained on

True

False **SOLUTION:**

False, you can transfer previous layers and reinitialize the final layer

1.10 Self-supervised learning models are defined by the fact that the model's output is a similarity score between embeddings or a reconstructed input, rather than a probability distribution over discrete class labels.

True

False **SOLUTION:**

False, its defined by how the labels are obtained, not by the format of the output

2 Multiple Choice (32 points) - Recommended 25 Minutes

Fill in the circle next to the letter(s) of your choice (like this: ●). No explanations are required. Choose ALL options that apply.

Each question is worth 8 points and the answer may contain multiple correct options, or none at all. Selecting all of the correct options and none of the incorrect options will get full credit. For questions with multiple correct options, each incorrect or missing selection gets a 4-point deduction (up to 8 points).

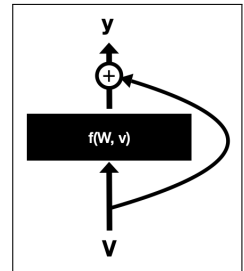
2.1

You are training a neural network. At some point in the network, there is an activation vector \mathbf{v} of shape $(1 \times d)$. This vector is fed into a linear layer f with weight matrix \mathbf{W} of shape $(d \times d)$ (i.e., $f(\mathbf{v}) = \mathbf{v}\mathbf{W}$). There is a residual connection from just before f to just after f , as shown to the right.

The output is:

$$\mathbf{y} = f(\mathbf{v}) + \mathbf{v} = \mathbf{v}\mathbf{W} + \mathbf{v}$$

Assume that $\frac{\partial \mathcal{L}}{\partial \mathbf{y}}$ is given (shape $1 \times d$). Let \mathbf{I} denote the identity matrix. What is $\frac{\partial \mathcal{L}}{\partial \mathbf{v}}$?



- A: $\frac{\partial \mathcal{L}}{\partial \mathbf{y}} \mathbf{W}^T$
- B: $\frac{\partial \mathcal{L}}{\partial \mathbf{y}} \mathbf{W}^T + \mathbf{I}$
- C: $\frac{\partial \mathcal{L}}{\partial \mathbf{y}} (\mathbf{I} + \mathbf{W}^T)$
- D: $(\mathbf{I} + \frac{\partial \mathcal{L}}{\partial \mathbf{y}}) \mathbf{W}^T$

SOLUTION:

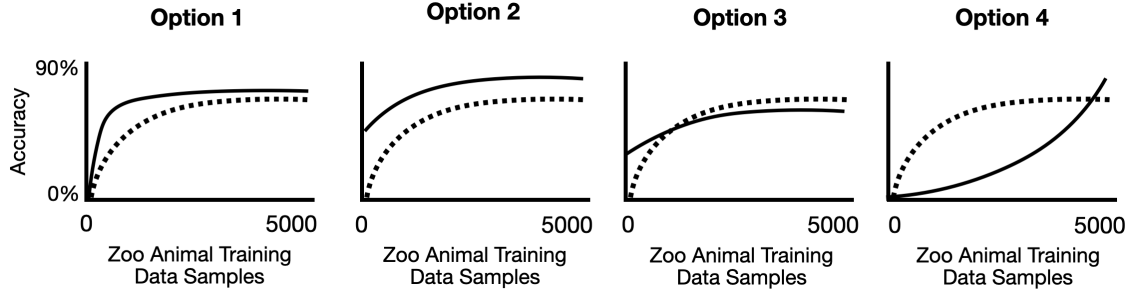
C. Since $\mathbf{y} = \mathbf{v}\mathbf{W} + \mathbf{v} = \mathbf{v}(\mathbf{W} + \mathbf{I})$, we have $\frac{\partial \mathbf{y}}{\partial \mathbf{v}} = (\mathbf{W} + \mathbf{I})^T = \mathbf{I} + \mathbf{W}^T$, so by the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} (\mathbf{I} + \mathbf{W}^T)$$

This has shape $(1 \times d)(d \times d) = (1 \times d)$ ✓ **(A)** would be correct if there were no residual connection—it is missing the gradient from the skip path. **(B)** incorrectly adds the identity matrix as a standalone term—the skip connection contributes $\frac{\partial \mathcal{L}}{\partial \mathbf{y}}$, not \mathbf{I} , to the gradient. **(D)** incorrectly adds \mathbf{I} to the upstream gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{y}}$ rather than to \mathbf{W}^T .

2.2 You pretrain a convolutional neural network on ImageNet. You then transfer it to a smaller Zoo Animal classification dataset with 500 classes by removing the final fully-connected classification layer and replacing it with a new, randomly initialized fully-connected layer for your target classes. You then fine-tune the entire network on your target data. The **dashed line** in each plot below shows the accuracy curve when training the same architecture from **random initialization** on the target dataset. The **solid line** shows the accuracy curve of the fine-tuned model.

Which plot best represents the **expected training curve** of the fine-tuned model compared to training from scratch? **Select one.**

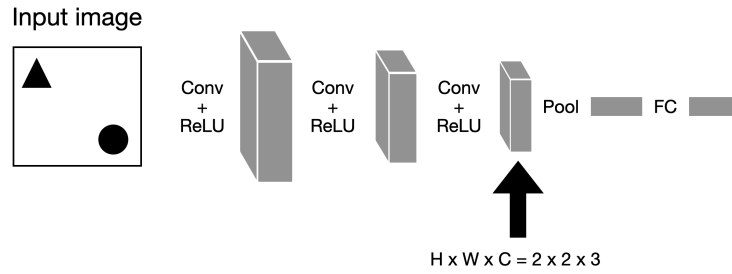


- A: Option 1
- B: Option 2
- C: Option 3
- D: Option 4

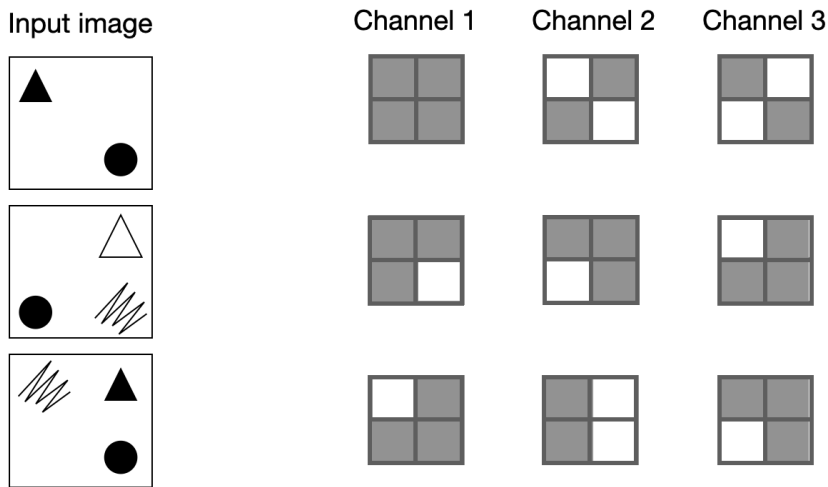
SOLUTION:

A. Since the final fully-connected layer is replaced with a new, randomly initialized one, both models start at a similarly low accuracy. However, the fine-tuned model benefits from pretrained features in the earlier layers, leading to two key differences: **Faster learning:** The pretrained convolutional features are already useful for visual recognition—especially since ImageNet contains many animal classes. The new classification head only needs to learn a mapping from these good features to the target classes, so accuracy rises much more quickly. **Higher final accuracy:** With a smaller target dataset, training from scratch is more prone to underfitting or learning weaker features. The pretrained initialization converges to a better solution. Option 1 correctly shows both effects: same starting point, faster rise, and higher final accuracy. **(B)** Incorrectly shows the fine-tuned model starting at high accuracy. Since the final FC layer is randomly initialized, accuracy at step 0 should be near chance regardless of pretraining. **(C)** Shows the from-scratch model learning faster initially and outperforming the fine-tuned model early on. This is unlikely when transferring from a closely related domain—pretrained features should help immediately, not hinder early training. **(D)** Shows the fine-tuned model learning slower than from-scratch early on, only overtaking it late. Pretrained features should accelerate learning from the start, not slow it down.

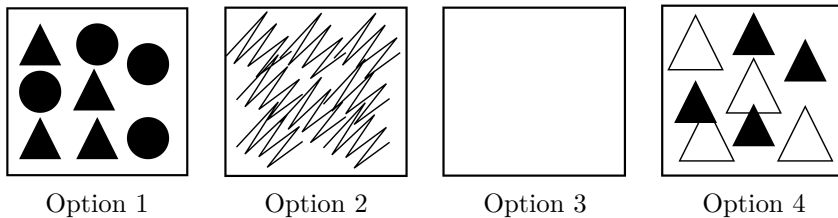
2.3 You have trained the CNN shown below and want to understand what features it has learned at a specific layer. You examine the activations at the layer indicated by the arrow, which has shape $H \times W \times C = 2 \times 2 \times 3$.



You pass three different input images through the network and visualize the resulting activations across all 3 channels. In the activation grids below, **white squares indicate high activation values** and **gray squares indicate low activation values**.



You then perform **gradient ascent** on a randomly initialized input image to maximize the sum of all 4 spatial values in **Channel 3** of this activation. Which of the following images would you most expect gradient ascent to produce? **Select one.**



- A: Option 1
- B: Option 2
- C: Option 3
- D: Option 4

SOLUTION:

C. To determine what Channel 3 detects, we examine where it has high activation (white) across the three input images.

In every case, Channel 3 activates exactly where there is *no object*—it has learned to detect empty space. To maximize all 4 spatial values of Channel 3 via gradient ascent, we want an input with nothing in any quadrant. Option 3 (a blank image) achieves this.

2.4 For this question, use the following notation and assumptions:

- Let $T(f(x))$ denote the time it takes to compute $f(x)$
- Let $\text{Encode}_m(n)$ denote the process of encoding a sequence of length n into a vector representation using model m (e.g., converting a user's input question into a hidden representation for later processing)
- Let $\text{Decode}_m(n)$ denote the process of autoregressively generating n new tokens using model m (e.g., generating a response to a user's question, or captioning an unseen image)
- Let tf denote a Transformer and let rnn denote an RNN

Assume we have **infinite compute and memory**, so any operations that can be parallelized are fully parallelized. Under this assumption:

- Independent operations run in parallel: $T(f(x) + g(x)) = \max(T(f(x)), T(g(x)))$
- Sequential dependencies still require waiting: $T(g(f(x))) \approx T(f(x)) + T(g(y))$ where $y = f(x)$

Select all statements that are TRUE:

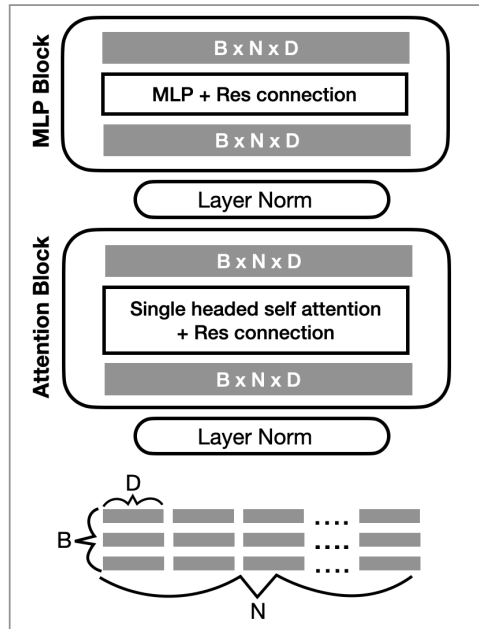
- A: $T(\text{Encode}_{tf}(n)) \approx T(\text{Encode}_{tf}(2n))$
- B: $T(\text{Encode}_{rnn}(n)) \approx T(\text{Encode}_{rnn}(2n))$
- C: $T(\text{Decode}_{tf}(n)) \approx T(\text{Decode}_{tf}(2n))$
- D: $T(\text{Decode}_{rnn}(n)) \approx T(\text{Decode}_{rnn}(2n))$

SOLUTION:

A. **(A) TRUE:** Transformer encoding processes all tokens in parallel via self-attention. With infinite compute, doubling the sequence length does not significantly increase wall-clock time. **(B) FALSE:** RNN encoding is inherently sequential—each hidden state depends on the previous one. Doubling the sequence length approximately doubles the encoding time: $T(\text{Encode}_r(2n)) \approx 2 \cdot T(\text{Encode}_r(n))$. **(C) FALSE:** Autoregressive decoding in a Transformer generates tokens one at a time, where each new token depends on all previously generated tokens. Doubling the number of generated tokens approximately doubles the decoding time. **(D) FALSE:** RNN decoding is also autoregressive—each generated token depends on the previous hidden state. Doubling the output length approximately doubles the decoding time.

3 Transformer Memory (18 points) - Recommended 15 Minutes

You are training a Transformer that consists of n layers, where each layer contains a single-head self-attention block followed by an MLP block, as shown below. The MLP block consists of two linear layers with a nonlinearity in between. The input and output dimensions of both linear layers are equal and of size D . You perform a forward pass on a batch of sequences. A single element of a sequence within a batch is a vector of dimension D . So the input/output of each attention and MLP block has shape $B \times N \times D$, where B is the batch size, N is the sequence length, and D is the embedding dimension.



Assume that N is approximately $100\times$ larger than both B and D (i.e., $N \gg B$ and $N \gg D$). Ignore the memory required to store model weights—consider only the memory for intermediate values stored during the forward pass (e.g., activations, attention maps). For each of the following scenarios, write the **approximate multiplicative factor** by which the total memory stored during the forward pass changes. Your answer should reflect the **dominant** memory term—you may ignore terms that are negligibly small relative to the largest term given our assumptions. Write **1x** if the memory does not meaningfully change. You are not using any efficient attention implementations, such as Flash Attention.

3.1 Memory in the Attention Block (3 pts each)

a) The batch size B doubles. Attention block memory increases by approximately x

b) The sequence length N doubles. Attention block memory increases by approximately x

c) The embedding dimension D doubles. Attention block memory increases by approximately x

3.2 Memory in the MLP Block (3 pts each)

a) The batch size B doubles. MLP block memory increases by approximately x

b) The sequence length N doubles. MLP block memory increases by approximately x

c) The embedding dimension D doubles. MLP block memory increases by approximately x

SOLUTION:

The key insight is identifying the dominant memory term in each block given that $N \gg B$ and $N \gg D$.
Attention block: The attention map has shape $B \times N \times N$. The input/output activations have shape $B \times N \times D$. Since $N \gg D$, the $B \times N \times N$ attention map dominates.

a) **2x.** The attention map $B \times N \times N$ is linear in B .

b) **4x.** The attention map $B \times N \times N$ is quadratic in N , so doubling N quadruples memory.

c) **1x.** The dominant term $B \times N \times N$ does not depend on D . While the input/output activations ($B \times N \times D$) do scale with D , they are negligible compared to the attention map since $N \gg D$.

MLP block: The MLP consists of two linear layers with a nonlinearity in between. The stored activations include the input ($B \times N \times D$), the intermediate representation after the first linear layer ($B \times N \times kD$ for some constant k), and the output ($B \times N \times D$). All terms are linear in B , N , and D .

a) **2x.** Linear in B .

b) **2x.** Linear in N .

c) **2x.** Linear in D .

4 Transformer Arithmetic (20 points) - Recommended 20 Minutes

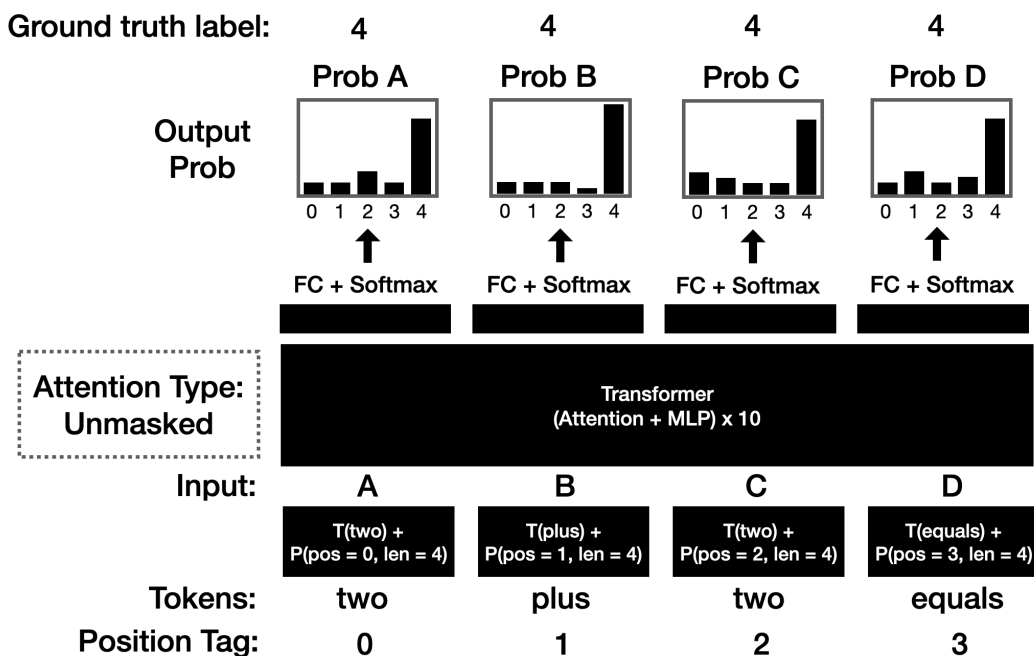
You have trained a Transformer to perform simple arithmetic in a supervised fashion. The training data consists of text sequences paired with numerical ground truth labels. For example:

| Input sequence | Ground truth label |
|------------------------------|--------------------|
| One Plus Two Equals | 3 |
| Two Minus One Equals | 1 |
| One Plus One Plus One Equals | 3 |

The input vocabulary is {One, Two, Plus, Minus, Equals} and the output classes are {0, 1, 2, 3, 4}. The final token in every sequence is always “Equals”.

The model consists of 10 stacked blocks of (Attention + MLP). Each input token is represented as the sum of a token embedding and a positional embedding. The positional encoding encodes both the position of each token and the total length of the sequence.

At each input position, the model produces a probability distribution over the output classes using an FC layer + Softmax. The same ground truth label is used as the target at **every** output position, and the model is trained with cross-entropy loss. The setup is illustrated below for the input “Two Plus Two Equals” (ground truth label: 4):



For each of the following scenarios, you are given an input sequence and an attention type (either unmasked or masked causal). Assume the model was trained and is evaluated using the attention type shown. Your task is to sketch the approximate probability distribution you would expect at the specified output position. If the model can confidently determine the answer, draw a single tall bar at that value. If the model is uncertain, draw bars reflecting the plausible outputs. **Read each table carefully, pay close attention to the attention type, the input tokens, their positions, and which output position you are asked about.**

Here is an example of a table specifying a certain scenario and what we want from you.

Read this table carefully

| Attention | Unmasked | | | |
|--------------|----------|------|-----|--------|
| Input | A | B | C | D |
| Token | Two | Plus | Two | Equals |
| Position | 0 | 1 | 2 | 3 |
| Seq Len | 4 | 4 | 4 | 4 |
| Which output | ★ | | | |

Example of what we want from you

Pay attention to which output we want!

Prob A

The graph shows a horizontal axis with tick marks at 0, 1, 2, 3, and 4. A vertical bar is drawn at position 4, filled with a diagonal hatching pattern. The label 'Prob A' is centered above the graph.

4.1 Prob F given the below input and settings (4 pts)

| Attention | Unmasked | | | | | |
|--------------|----------|------|-----|------|-----|--------|
| Input | A | B | C | D | E | F |
| Token | One | Plus | One | Plus | One | Equals |
| Position | 0 | 1 | 2 | 3 | 4 | 5 |
| Seq Len | 6 | 6 | 6 | 6 | 6 | 6 |
| Which output | ★ | | | | | |

Prob F

The graph shows a horizontal axis with tick marks at 0, 1, 2, 3, and 4. No bars are present. The label 'Prob F' is centered above the graph.

4.2 Prob A given the below input and settings (4 pts)

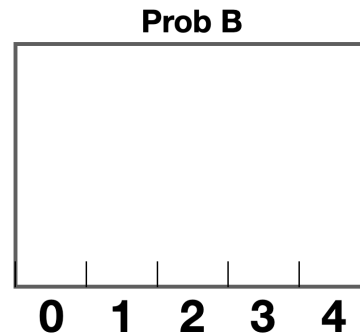
| Attention | Unmasked | | | | | |
|--------------|----------|------|-----|------|-----|--------|
| Input | A | B | C | D | E | F |
| Token | One | Plus | One | Plus | One | Equals |
| Position | 0 | 1 | 2 | 3 | 4 | 5 |
| Seq Len | 6 | 6 | 6 | 6 | 6 | 6 |
| Which output | ★ | | | | | |

Prob A

The graph shows a horizontal axis with tick marks at 0, 1, 2, 3, and 4. No bars are present. The label 'Prob A' is centered above the graph.

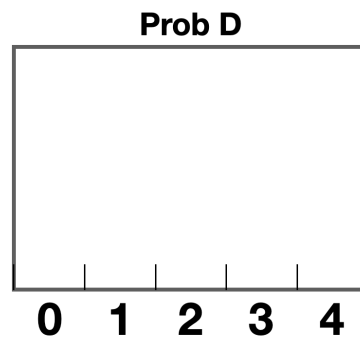
4.3 Prob B given the below input and settings (4 pts)

| Attention | Masked Causal | | | |
|--------------|---------------|------|-----|--------|
| Input | A | B | C | D |
| Token | Two | Plus | One | Equals |
| Position | 0 | 1 | 2 | 3 |
| Seq Len | 4 | 4 | 4 | 4 |
| Which output | ★ | | | |



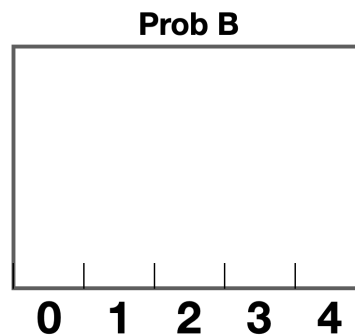
4.4 Prob D given the below input and settings (4 pts)

| Attention | Masked Causal | | | |
|--------------|---------------|------|-----|--------|
| Input | A | B | C | D |
| Token | Two | Plus | One | Equals |
| Position | 0 | 1 | 2 | 3 |
| Seq Len | 4 | 4 | 4 | 4 |
| Which output | | | | ★ |



4.5 Prob B given the below input and settings (4 pts)

| Attention | Unmasked | | | | | |
|--------------|----------|------|-----|-------|-----|--------|
| Input | A | B | C | D | E | F |
| Token | One | Plus | One | Minus | Two | Equals |
| Position | 4 | 1 | 2 | 3 | 0 | 5 |
| Seq Len | 6 | 6 | 6 | 6 | 6 | 6 |
| Which output | ★ | | | | | |



SOLUTION:

The key concepts tested are: (1) unmasked attention allows every position to see all tokens, (2) masked causal attention restricts each position to only see itself and earlier positions, and (3) the model interprets token order based on **positional encodings**, not the physical order in the sequence.

Q1: “One Plus One Plus One Equals” — unmasked, Prob A. With unmasked attention, position A sees all tokens. The expression is $1 + 1 + 1 = 3$. **Peak at 3.**

Q2: “One Plus One Plus One Equals” — unmasked, Prob F. Same input and unmasked attention. Position F also sees all tokens. Since every position has access to the same information with unmasked attention, Prob F is the same as Prob A. **Peak at 3.**

Q3: “Two Plus One Equals” — masked causal, Prob B. With masked causal attention, position B can only see positions A and B (tokens “Two” and “Plus”). However, from the positional encoding, the model knows the total sequence length is 4, meaning the expression has the form “X op Y Equals”. The model knows $X = \text{Two}$ and $\text{op} = \text{Plus}$, but cannot see Y. The possible completions are “Two Plus One = 3” and “Two Plus Two = 4”. **Spread across 3 and 4** (roughly equal).

Q4: “Two Plus One Equals” — masked causal, Prob D. Position D sees all positions A through D (“Two”, “Plus”, “One”, “Equals”). The full expression $2 + 1 = 3$ is visible. **Peak at 3.**

Q5: “One Plus One Minus Two Equals” — unmasked, Prob B. With unmasked attention, position B sees all tokens. However, the **positional encodings** are scrambled (positions are 4, 1, 2, 3, 0, 5). The model interprets token order based on positional encodings, not physical order. Re-ordering by position:

| | | | | | | |
|----------|-----|------|-----|-------|-----|--------|
| Position | 0 | 1 | 2 | 3 | 4 | 5 |
| Token | Two | Plus | One | Minus | One | Equals |

The model interprets the expression as $2 + 1 - 1 = 2$. **Peak at 2.**