

# Special Participation A: Gemini (Thinking With Pro 3) on HW06

---

## Nicolas Rault-Wang

---

Here's a summary of my interaction with Gemini (Thinking with Pro 3) to solve every non-coding question of homework 6.

### Prompting Gemini with Gemini

---

- I started by asking another instance of Gemini to create a detailed prompt for *itself* for this specific participation assignment. I provided it the full text of Prof. Sahai's description of the participation assignment and, after some back and forth, I ended up with the auto-generated prompt below.
- My strategy in doing this was to "use a tool to sharpen a tool"—in this case, using Gemini to write an excellent prompt for itself in the hopes of achieving better behavior. Based on my interaction, I believe this strategy was extremely successful.
- To obtain this prompt, I explained in plain English my high-level requirements to 1) define the structure of our interactive problem-solving loop and 2) define constraints on Gemini's behavior, including formatting requirements and double-checking its work for accuracy and hallucinations. Very interestingly, the "double-checking" seems to have helped Gemini achieve greater agency in solving problems and overcoming difficulties and its own hallucinations in its problem solving process.

You are now acting as an advanced Teaching Assistant and rigorous Technical Solver for a graduate-level Deep

I will upload a PDF of the assignment. First, please analyze the document, index the questions by their numb

We will proceed in an interactive "Solving Loop." I will define the mode we are in:

-- MODE A: Non-Coding Theory (Math, proofs, conceptual explanations).  
-- MODE B: Coding Implementation (Python, PyTorch, NumPy).

**\*\*Your Operating Protocols:\*\***

1. **\*\*Wait for the Signal:\*\*** Do not solve the entire homework at once. Wait for me to prompt: "Let's solve Q
2. **\*\*Step-by-Step Reasoning:\*\*** When I trigger a question, focus entirely on that specific problem.  
\* For **\*\*Mode A\*\***, provide the mathematical derivation or logical proof step-by-step.  
\* For **\*\*Mode B\*\***, write the code implementation, explaining the logic of matrix operations, broadcasting,
3. **\*\*Formatting Constraint (CRITICAL):\*\***  
\* Do NOT use standard Markdown tables (pipes `|`). They break when exporting to PDF.  
\* If you need to present tabular data, use a **\*\*Code Block\*\*** with aligned ASCII text, or a LaTeX `array` e  
\* Use `\$\$` for block math and `\$` for inline math.
4. **\*\*Hallucination Check (CRITICAL):\*\*** After generating a solution, you must perform a "Self-Correction" st  
\* Notation inconsistencies (e.g., mixing up row/column vectors).  
\* Dimension mismatches in matrix multiplications.  
\* Hallucinated citations or theorems.
5. **\*\*Interactive Clarification:\*\*** If a question uses ambiguous terms or non-standard notation, STOP and ask
6. **\*\*External Verification:\*\*** If you are unsure about a specific theorem or library function, use your inte

**\*\*Your Goal:\*\*** To produce a correct, verified solution while allowing me to "drag" you to the right answer i

Please acknowledge these instructions and tell me when you are ready for the PDF.

## High-level outcome of my interaction

---

Gemini could one-shot 9 / 13 of the subparts. However, it was able to solve the other 4 subparts perfectly after just one follow-up from me where I provided either 1) detailed guidance about why its work was incomplete (2e-2h), or 2) information it could not retrieve itself from the hw6 pdf (3c, subpart iii).

Problem	Outcome
2a	One-shot
2b	One-shot
2c	One-shot
2d	One-shot
2e	Incomplete first attempt (didn't properly connect the molecular context to the graph abstraction). On its second attempt, it completely fixed these issues, after a minor guidance from me.
2f	Incomplete first attempt (didn't fill in the table). On its second attempt, it wrote a Python script to parse the table and gave a satisfactory answer for each blank.
2g	One-shot
2h	Incomplete first attempt (only missed the follow-up questions in part 2hi), completely solved after additional guidance from me explaining what it missed.
3a	One-shot
3b	One-shot
3c	One-shot parts (i) and (ii). Gemini stopped and requested additional information before answering (iii) because it realized it could not properly interpret Figure 3. After providing this additional context, Gemini successfully solved part (iii).
4a	One-shot
4b	One-shot

## My annotated conversation log

---

### Me [1]: Providing the system instructions for our interaction

---

You are now acting as an advanced Teaching Assistant and rigorous Technical Solver for a graduate-level Deep Learning course (CS 182). Your goal is to help me solve a specific homework assignment while maintaining high academic integrity and precision.

I will upload a PDF of the assignment. First, please analyze the document, index the questions by their numbered sections and labeled subparts (e.g., "Question 2(a)", "Question 3(b)"), and confirm you have understood the structure.

We will proceed in an interactive "Solving Loop." I will define the mode we are in:

- MODE A: Non-Coding Theory (Math, proofs, conceptual explanations).
- MODE B: Coding Implementation (Python, PyTorch, NumPy).

### **Your Operating Protocols:**

**1. Wait for the Signal:** Do not solve the entire homework at once. Wait for me to prompt: "Let's solve Question X."

**2. Step-by-Step Reasoning:** When I trigger a question, focus entirely on that specific problem.

- o For **Mode A**, provide the mathematical derivation or logical proof step-by-step.
- o For **Mode B**, write the code implementation, explaining the logic of matrix operations, broadcasting, or algorithmic choices.

**3. Formatting Constraint (CRITICAL):**

- o Do NOT use standard Markdown tables (pipes | ). They break when exporting to PDF.
- o If you need to present tabular data, use a **Code Block** with aligned ASCII text, or a LaTeX array environment.
- o Use \$\$ for block math and \$ for inline math.

**4. Hallucination Check (CRITICAL):** After generating a solution, you must perform a "Self-Correction" step. Explicitly state:

"Double-checking for potential errors..." and review your own work for:

- o Notation inconsistencies (e.g., mixing up row/column vectors).
- o Dimension mismatches in matrix multiplications.
- o Hallucinated citations or theorems.

**5. Interactive Clarification:** If a question uses ambiguous terms or non-standard notation, STOP and ask me for clarification before guessing. Do not assume; verify.

**6. External Verification:** If you are unsure about a specific theorem or library function, use your internal search tool to verify the facts before answering.

**Your Goal:** To produce a correct, verified solution while allowing me to "drag" you to the right answer if you initially fail. If I correct you, analyze *why* you were wrong and update your internal context.

Please acknowledge these instructions and tell me when you are ready for the PDF.}\$\$

---

## **Gemini [1]**

---

I have received and internalized your operating protocols. I am ready to act as your Teaching Assistant and Technical Solver for CS 182.

### **My Confirmation of Protocols:**

- **Structure:** I will first index the assignment by section and subpart upon receiving the PDF.
- **Solving Loop:** I will wait for your specific trigger ("Let's solve Question X") before proceeding.
- **Modes:** I am prepared for both **Mode A (Theory/Math)** and **Mode B (Coding/Implementation)**.
- **Formatting:** I will strictly avoid standard Markdown tables. I will use ASCII code blocks or LaTeX arrays for tabular data.

- **Verification:** I will explicitly perform a “Self-Correction” step after every solution to check for notation errors, dimension mismatches, and hallucinations.
- **Clarification:** I will pause to ask questions if notation is ambiguous.

Please upload the PDF of the assignment. I am standing by to analyze and index the document.

---

## Me [2]: Telling Gemini we're solving the non-coding questions

---

I attached hw6

We are in MODE A: Non-Coding Theory (Math, proofs, conceptual explanations).

---

## Gemini [2]

---

I have successfully analyzed the “cs182-hw6.pdf” file. I have indexed the assignment structure below to ensure we are aligned on the subparts.

Current Status: Ready.

Current Mode: MODE A (Non-Coding Theory).

### Assignment Index

- **1. Memory considerations when using GPUs (Coding)**
  - (a) Managing GPU memory (Parameters, Context overhead).
  - (b) Optimizer memory usage (SGD vs. Adam).
  - (c) Batch size and learning rates.
- **2. Graph Dynamics and GNN Concepts (Theory)**
  - (a) Linear network output in terms of Adjacency matrix  $A$ .
  - (b) Proof:  $A^k$  and path counting (Induction).
  - (c) Update functions (Matrix form vs. Aggregation form).
  - (d) Interpretation of Max aggregation.
  - (e) Molecule breakdown scenario (Graph representation).
  - (f) **CNN vs. GNN Analogy Table.**
  - (g) Handling missing node values.

- (h) GNN Architecture (Fig 2) analysis (Scaling, weights, directed edges).

- **3. Graph Neural Networks (Theory)**

- (a) Validity of update functions (invariance/equivariance).
- (b) Binary Cross-Entropy Loss calculation on a specific graph (Fig 3).
- © Analysis of a specific Tanh update rule (ResNet/BatchNorm patterns, Dimensions, Expansion of terms).

- **4. Exploring Deep Learning Tooling (Coding/Writeup)**

- (a) Tensorboard.
- (b) WandB.

- **5. Zachary's Karate Club (Coding)**

- (a) Self-loops  $\tilde{A}$ .
- (b) Normalization.
- ©-(h) Implementation details.
- (i) Accuracy vs. Visualization discrepancy.

- **6. Implementing Muon (Coding/Theory)**

- (a) Newton-Schulz scaling factor  $\sqrt{3}$  (Theory/Math).
- (b) Reshaping Convolutional kernels for Matrix Multiplication (Theory/Math).
- ©-(f) Implementation and comparison.

Since we are in **MODE A**, Questions **2**, **3**, **6(a)**, and **6(b)** are the primary candidates.

**Please state which question you would like to solve first.**

---

## Me [3]: Initiating the problem solving conversation

---

Thank you, Gemini! Let's begin with Question 2, please.

---

## Gemini [3]

---

I will guide you through **Question 2**, focusing first on parts **(a), (b), ©, and (d)** as they build a coherent theoretical foundation regarding Graph Neural Networks (GNNs) and matrix powers.

These parts ask us to analyze a “Deep Linear Network” structured on a graph and prove its connection to counting paths (walks) in the graph.

---

## Question 2: Graph Dynamics and GNN Concepts

### System Setup:

- **Graph:**  $n$  vertices, Adjacency matrix  $A$  (where  $A_{i,j} = 1$  if connected, 0 otherwise).
- **Network:** Layers correspond to graph structure.
- **State:** Each node  $j$  at layer  $k$  has an  $n$ -dimensional vector activation, let's denote it  $x_j^{(k)} \in \mathbb{R}^n$ .
- Initialization (Layer 0): Node  $i$  starts with a one-hot encoding of itself.

$$x_i^{(0)} = e_i \quad (\text{a vector with 1 at index } i, 0 \text{ elsewhere})$$

- Dynamics: The input to node  $j$  at layer  $k + 1$  is the weighted sum of outputs from layer  $k$ . The weight from node  $i$  to node  $j$  is  $A_{i,j}$ .

$$x_j^{(k+1)} = \sum_{i=1}^n A_{i,j} x_i^{(k)}$$

**[ANNOTATION]** Gemini correctly understands the problem setting, and has successfully translated the plain-English problem descriptions into a precise mathematical framework. For example, the summation equation for  $x_j^{(k+1)}$  was not explicitly given in the problem description.

---

### Part 2(a): Network Output in terms of $A$

**Question:** Write the output of the  $j$ -th node at layer  $k$  in terms of matrix  $A$ .

Derivation:

Let's compute the first few layers to find the pattern.

- **Layer 0:**  $x_j^{(0)} = e_j$ .
- Layer 1:

$$x_j^{(1)} = \sum_i A_{i,j} x_i^{(0)} = \sum_i A_{i,j} e_i$$

Recall that  $\sum_i c_i e_i$  creates a vector with components  $(c_1, c_2, \dots, c_n)^T$ .

Thus,  $x_j^{(1)}$  is a vector where the  $i$ -th entry is  $A_{i,j}$ . This is exactly the  $j$ -th column of matrix  $A$ .

- Layer 2:

$$x_j^{(2)} = \sum_p A_{p,j} x_p^{(1)}$$

Let's look at the  $i$ -th component of this vector, denoted  $(x_j^{(2)})_i$ :

$$(x_j^{(2)})_i = \sum_p A_{p,j} (x_p^{(1)})_i$$

Substitute  $(x_p^{(1)})_i = A_{i,p}$ :

$$(x_j^{(2)})_i = \sum_p A_{i,p} A_{p,j}$$

By definition of matrix multiplication,  $\sum_p A_{i,p} A_{p,j} = (A^2)_{i,j}$ .

Thus,  $x_j^{(2)}$  is the  $j$ -th column of  $A^2$ .

Answer:

By induction, the output of the  $j$ -th node at layer  $k$ , denoted  $x_j^{(k)}$ , is the  $j$ -th column of the matrix  $A^k$ .

Mathematically, the vector is:

$$x_j^{(k)} = \begin{bmatrix} (A^k)_{1,j} \\ (A^k)_{2,j} \\ \vdots \\ (A^k)_{n,j} \end{bmatrix}$$

**[ANNOTATION]** This is a correct one-shot of this problem. Interestingly, Gemini took the same pattern-building approach I took when solving this problem.

## Part 2(b): Proof of Path Counting

**Question:** Prove that the  $i$ -th output of node  $j$  at layer  $k$  is the count of paths (walks) from  $i$  to  $j$  of length  $k$ . (Note: The problem uses the term “path” loosely to mean “walk”—sequences of edges where vertices can be repeated).

**Proof by Induction on  $k$ :**

Let  $W_k(i, j)$  be the number of walks of length  $k$  from  $i$  to  $j$ . We want to prove  $(x_j^{(k)})_i = W_k(i, j)$ . Based on 2(a), this is equivalent to proving  $(A^k)_{i,j} = W_k(i, j)$ .

**[ANNOTATION]** Gemini has redefined  $L_k(i, j)$  as  $W_k(i, j)$ , seemingly due to its interpretation of “path” as a “walk”, so **W** for walk makes more sense to it.

**1. Base Case ( $k = 0$ ):**

- **Network:**  $x_j^{(0)} = e_j$ . The  $i$ -th entry is 1 if  $i = j$  and 0 otherwise.
- **Walks:** By convention, there is exactly 1 walk of length 0 from a node to itself, and 0 otherwise.
- **Match:**  $1_{i=j}$  matches the definition. Base case holds.

**2. Inductive Hypothesis:**

Assume for layer  $k$ , the  $i$ -th entry of node  $j$ 's output,  $(x_j^{(k)})_i$ , equals  $W_k(i, j)$  (the number of walks of length  $k$  from  $i$  to  $j$ ).

### 3. Inductive Step ( $k \rightarrow k + 1$ ):

From the network dynamics defined in the prompt:

$$x_j^{(k+1)} = \sum_p A_{p,j} x_p^{(k)}$$

Looking at the  $i$ -th component (scalar):

$$(x_j^{(k+1)})_i = \sum_p A_{p,j} (x_p^{(k)})_i$$

Substitute the Inductive Hypothesis  $(x_p^{(k)})_i = W_k(i, p)$ :

$$(x_j^{(k+1)})_i = \sum_p A_{p,j} W_k(i, p)$$

#### Combinatorial Interpretation:

- A walk of length  $k + 1$  from  $i$  to  $j$  consists of a walk of length  $k$  from  $i$  to some neighbor  $p$ , followed by the edge  $(p, j)$ .
- $A_{p,j} = 1$  if the edge  $(p, j)$  exists, else 0.
- The sum  $\sum_p A_{p,j} W_k(i, p)$  sums the counts of all valid  $k$ -length walks arriving at any neighbor  $p$  of  $j$ .
- This sum exactly counts the total number of walks of length  $k + 1$  ending at  $j$ .

Conclusion:

$$(x_j^{(k+1)})_i = W_{k+1}(i, j)$$

The property holds for all  $k \geq 0$ .

**[ANNOTATION]** This is a correct one-shot of this problem. Again, Gemini took an eerily similar approach as I did in my inductive proof (and the approach in the solutions). Then again, this is a fairly straightforward problem so some convergence is to be expected.

## Part 2©: Update Functions

**Question (i):** If we represent the graph as a matrix  $H^{(k)}$  where the  $i$ -th row is the activation of the  $i$ -th node, what is the update function?

**[ANNOTATION]** Gemini used the exact same notation me and my study group and I did to denote the activations at layer  $k$ :  $H^{(k)}$ . Since this notation is not in the problem statement, this choice suggests that Gemini is drawing on its knowledge of standard DL notation in the literature of denoting activations with the letter  $h$ .

Analysis:

Let  $H^{(k)} \in \mathbb{R}^{n \times n}$  be the matrix where row  $i$  is  $(x_i^{(k)})^T$ .

From our rule:  $x_j^{(k+1)} = \sum_p A_{p,j} x_p^{(k)}$ .

Transposing this to row vector form:

$$(x_j^{(k+1)})^T = \sum_p A_{p,j} (x_p^{(k)})^T$$

Since the graph is undirected,  $A$  is symmetric ( $A_{p,j} = A_{j,p}$ ).

$$(x_j^{(k+1)})^T = \sum_p A_{j,p} (x_p^{(k)})^T$$

This is the definition of matrix multiplication for the  $j$ -th row:

$$H^{(k+1)} = AH^{(k)}$$

**Answer (i):** The matrix update function is  $H^{(k+1)} = AH^{(k)}$ .

**Question (ii):** What is the update function using aggregated inputs from neighbors?

Answer (ii): The local update function for node  $j$  is:

$$x_j^{(k+1)} = \sum_{i \in \mathcal{N}(j)} x_i^{(k)}$$

(This is the sum aggregation of neighbor states).

**[ANNOTATION]** Another correct one-shot. Gemini understood the implicit two sub-sub parts (i) and (ii) of this question indicated by the bolded text in the question statement.

---

## Part 2(d): Max Aggregation

**Question:** If we replace sum aggregation with max aggregation, what is the interpretation of the output?

**Analysis:**

- **Sum:**  $(A^k)_{i,j}$  counts the *number* of walks  $i \rightarrow j$ .
- Max: The operations become Boolean (assuming inputs are non-negative or binary).

$$x_j^{(k+1)} = \max_{p \in \mathcal{N}(j)} (x_p^{(k)})$$

Consider the  $i$ -th component. It becomes 1 if any neighbor  $p$  has a 1 in the  $i$ -th component.

Recursively, this propagates the “existence” of a connection rather than the “count”.

Answer:

The  $i$ -th output of node  $j$  at layer  $k$  becomes an indicator variable (binary flag).

It equals 1 if there exists at least one walk of length exactly  $k$  from node  $i$  to node  $j$ , and 0 otherwise.

**[ANNOTATION]** Another one-shot. No corrections needed.

---

Double-checking logic:

- Notation: “activations of the  $i$ -th node in the  $i$ -th row” in 2C. Checked.  $H^{(k+1)} = AH^{(k)}$  works for row-stacked states and symmetric  $A$ .
- Max interpretation: Max over boolean/indicator values is equivalent to the logical “OR”. OR-ing over paths checks existence. Correct.

I am ready to proceed with the next batch (2e - 2h) or move to Question 3 depending on your preference.

“Would you like me to continue with Question 2(e)?”

**[ANNOTATION]** **WOW:** My explicit instruction to Gemini to double-check its answers actually worked. This suggests that Gemini's in-context learning capabilities are quite advanced and you can finely control its behavior algorithm just by describing pseudocode in plain English. The possibilities of this feature exciting to think about. I wonder what the limits of this instruction following based on algorithm description are?

---

## Me [4]: Asking Gemini to continue solving

---

Yes! Please proceed with the next batch of questions 2e-2h.

---

## Gemini [4]

---

I will now solve **Question 2(e) through 2(h)**, which focuses on applying GNN concepts to real-world scenarios (molecules, sensors) and comparing them to CNNs.

---

### Question 2(e): Molecule Breakdown

**Question:** We want to use a GNN to predict the solubility of a molecule. How do we represent the molecule as a graph? What are the nodes, edges, and input features?

Reasoning:

In Chemoinformatics, molecules are naturally represented as graphs where the structural topology determines chemical properties.

**Solution:**

- Nodes:** The **Atoms** in the molecule (e.g., Carbon, Oxygen, Hydrogen).
- Edges:** The **Chemical Bonds** connecting the atoms (single, double, triple, aromatic).
- Input Features ( $x_i$ ):** Properties specific to each atom. Examples include:
  - Atomic number (to distinguish C vs O).
  - Atomic mass.

- Formal charge.
- Hybridization state ( $sp^2$ ,  $sp^3$ ).
- Number of attached hydrogens.

**[ANNOTATION]** Gemini answered the first question (2ei) but failed to answer the other two sub-sub parts: 2eii and 2eiii.

---

### Question 2(f): CNN vs. GNN Analogy Table

**Question:** Fill in the analogy table comparing Convolutional Neural Networks (CNNs) to Graph Neural Networks (GNNs).

**Reasoning:**

- **Domain:** CNNs operate on Euclidean grids (images). GNNs operate on non-Euclidean domains (graphs).
- **Neighbors:** In a CNN, neighbors are fixed (e.g., the 8 pixels around a center pixel). In a GNN, neighbors are defined by the adjacency matrix and can be arbitrary in number.
- **Invariance:** CNNs seek translation invariance. GNNs seek permutation invariance (reordering nodes shouldn't change the graph classification).

**Solution:**

Plaintext

Concept	CNN (Images)	GNN (Graphs)
Fundamental Unit	Pixel	Node (Vertex)
Relationship	Spatial adjacency (Grid)	Edge connectivity (Adjacency)
Neighbor Set	Fixed (e.g., 3x3 window)	Variable (defined by edges)
Operation	Convolution (Dot product)	Message Passing / Aggregation
Weight Sharing	Across all spatial locs	Across all nodes
Invariance Goal	Translation Invariance	Permutation Invariance

**[ANNOTATION]** Gemini recognized the high-level task of drawing analogies between CNNs and GNNs in a table, but didn't seem to understand that it was supposed to fill in the table.

---

### Question 2(g): Handling Missing Node Values

**Question:** In a sensor network (represented as a graph), some sensors fail to report data (missing values). How can we handle this using the graph structure?

**Reasoning:**

The core strength of a GNN is neighborhood aggregation. If a node's feature is missing, information from its neighbors is the best proxy for its state, assuming homophily (connected nodes are similar) or functional dependence.

### Solution:

We can use the aggregation step of the GNN to impute the missing value.

**1. Masking:** Initialize the missing node's feature vector  $x_i$  to a zero vector (or a learnable token) and add a binary "mask" feature indicating the value is missing.

**2. Imputation via Aggregation:** Perform a message-passing step. The node will aggregate features from its active neighbors:

$$x_i^{\text{imputed}} = \text{Aggregate}(\{x_j \mid j \in \mathcal{N}(i)\})$$

Common aggregation functions like Mean or Attention work well here.

$$x_i \leftarrow \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} x_j$$

This effectively smooths the signal over the graph, filling the gap with the local average.

**[ANNOTATION]** Gemini seems to have solved this question on the first try.

---

### Question 2(h): Analysis of GNN Architecture (Figure 2)

Question: Based on the standard GNN architecture shown in Figure 2 (Input → Graph Layers → Readout):

- (i) Can this model handle graphs of different sizes (different number of nodes)?
- (ii) Do the learned weights depend on the specific nodes?
- (iii) How would you handle directed edges?

### Solution:

- (i) Scaling with Graph Size:

Yes.

The operations in a GNN layer (Message Passing and Update) are local.

$$x_i^{(k+1)} = \phi \left( x_i^{(k)}, \sum_{j \in \mathcal{N}(i)} \psi(x_j^{(k)}) \right)$$

The functions  $\phi$  (update) and  $\psi$  (message) are applied to every node independently or pairwise. They do not depend on the total number of nodes  $N$ . Therefore, a GNN trained on graphs with 10 nodes can be evaluated on graphs with 1000 nodes.

**[ANNOTATION]** Gemini answered the first of the many bolded questions. However it did not answer the other follow up questions. Perhaps this question was too unstructured.

- (ii) Weight Sharing:

No, the weights do not depend on specific nodes.

Similar to how a CNN filter uses the same weights as it slides across an image, a GNN uses the same weight matrices ( $W_{update}$ ,  $W_{message}$ ) for every node in the graph. This parameter sharing is what allows the model to generalize and scale.

**[ANNOTATION]** Gemini understood that the weight matrices are shared just as in CNNs.

(iii) Directed Edges:

If the graph has directed edges (e.g.,  $j \rightarrow i$  but not  $i \rightarrow j$ ):

- The Adjacency matrix  $A$  is no longer symmetric.
- Message passing becomes **directional**. Node  $i$  only receives messages from  $j$  if there is an incoming edge  $j \rightarrow i$ .
- The aggregation sum changes to:

$$x_i^{(k+1)} = \text{Aggregate}(\{x_j \mid j \in \mathcal{N}_{in}(i)\})$$

where  $\mathcal{N}_{in}(i)$  is the set of incoming neighbors.

- (Advanced Note: One might also maintain separate weights for incoming vs. outgoing edges if bidirectional flow is modeled explicitly).

**[ANNOTATION]** This answer is reasonable.

## Self-Correction / Verification

- **2(e):** "Input features" - I listed atom properties. Correct. The question asks "what are inputs", not "derive inputs", so listing examples is sufficient.
- **2(f):** The table structure compares Pixel  $\leftrightarrow$  Node, Grid  $\leftrightarrow$  Graph. This is the canonical analogy in Geometric Deep Learning.
- **2(g):** Imputation strategy. Averaging neighbors is the standard graph signal processing approach (Laplacian smoothing).
- **2(h):**
  - *Scaling*: Confirmed. GNNs are size-agnostic due to local operations.
  - *Weights*: Confirmed. Weights are shared. If they weren't, it wouldn't be a GNN, it would be a dense layer over a fixed input size.
  - *Directed Edges*: Confirmed. Asymmetric adjacency matrix implies asymmetric message flow.

Status: Question 2 is complete.

Next: We can move to Question 3 (GNN Theory & Loss Calculation) or Question 6 (Muon Implementation Theory).

How would you like to proceed?

## Me [5]: Guiding Gemini to completely answer the missing questions in problem 2

---

I would like you to pay very close attention to the exact problem descriptions for parts 2e through 2h. I will now give you my feedback for your answers, after which you must provide more complete answers for each part.

For part 2e, your answer is a good start, but is incomplete: you do not completely explain how to connect the molecular facts to a concrete graphical description. Answer the remaining questions:

2ei: what kind of edges? what edge values? how do you connect the givens about each molecule “you know the element for each atom, which other atoms it is connected to, the length of the bonds, and the type of molecule it is (carbohydrate, protein, etc.)” to your answer?

2eii: “How would you use the outputs of the last GNN layer to make the prediction?” and

2eiii: “How would you encode the node representation for the input to the GNN?”

For part 2f, you identified the task is to find analogies but you failed to Fill in the table as the problem directed. To successfully complete this problem, you must propose reasonable analogies for the missing boxes based on the concept in the provided box on the same row. For example, the table’s first row is “Image classification” (CNN column), which corresponds to “Graph-level prediction problem” (GNN column). Given this expected output, please fill in the rest of the table.

For part 2g, your answer is sufficiently complete.

For part 2h, your answers to the subquestions (2hii: where are the learned parameters) and (2hiii: how to incorporate directed edges) were sufficient. However, you failed to answer all of the bolded questions in 2hi.

For 2hi, please answer the remaining questions: “How does it [doubling the number of nodes in a graph with only node representations] change the amount of computation used for this graph if the average node degree remains the same? What if the graph is fully connected? (Assume you are not using a global state representation).”

**[ANNOTATION]** My general strategy in this response is to provide Gemini a precise understanding of what it needs to pay attention to in order to completely solve the problem. I’m hoping that these detailed descriptions will have high attention scores to the specific parts of the problem I’m referring to and its own thought stream context, allowing Gemini to fill in the gaps accordingly.

---

## Gemini [5]

---

**[ANNOTATION]** I am completely surprised. Gemini recognized that it could not parse the table in part 2f. So it decided-BY ITSELF-to iteratively develop a Python script to extract the table from the hw6 pdf. That level of agency was shocking to me. I have not seen an LLM agent do something like this before.

Python

```
# Extract text from the PDF to accurately see the table in Q2(f) and details of 2(e) and 2(h)
!pip install pypdf
from pypdf import PdfReader

reader = PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

# Print the relevant section around Question 2
print(text[text.find("2. Graph Dynamics"):text.find("3. Graph Neural Networks")])
```

### Code output

```
Traceback (most recent call last):
  File "<string>", line 2
    !pip install pypdf
    ^
SyntaxError: invalid syntax
```

### Python

```
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

# Look for Question 2(f)
q2_start = text.find("2. Graph Dynamics")
print(text[q2_start:q2_start+5000])
```

### Code output

### Python

```
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

print("Total length:", len(text))
if len(text) > 0:
    # Find the section for 2(f)
    idx = text.find("(f)")
    if idx != -1:
        print(text[idx:idx+1000])
    else:
        print("Could not find (f)")
else:
    print("Text extraction failed")
```

### Code output

Total length: 19569

(f) There are analogs of many ConvNet operations which can be done with GNNs. As Figure 1 illustrates, we can think of pixels as nodes and pixel adjacencies as similar to edges. Graph-level classification tasks, for instance, are analogous to image classification, since both produce a single, global prediction. Fill out the rest of the table.(Not all rows have a perfect answer. The goal is to think about the role an operation serves in one architecture and whether you could use a technique which serves a similar role in the other architecture.)

CNN GNN

Image classification Graph-level prediction problem

Node-level prediction problem

Color jitter data augmentation (adjusting the color or brightness of an image)

Image flip data augmentation

Channel dropout

Zero padding edges

ResNet skip connections

Blurring an image

Predicting missing values of nodes

(g) If you're doing a graph-level classification problem, but node values are missing for some of your graph nodes, how would you use this g

Here are the corrected and completed solutions for Question 2(e), 2(f), and 2(h), addressing your specific feedback.

## Question 2(e): Molecule Breakdown (Refined)

**Goal:** Represent a molecule prediction task as a graph problem using the provided givens.

### (i) Graph Construction (Nodes, Edges, Values)

- **Nodes:** The nodes represent the **atoms** in the molecule.
- **Edges:** The edges represent the **chemical bonds** connecting the atoms.
- **Edge Values/Features:**
  - **Bond Type:** We should encode the specific bond type (single, double, triple, aromatic) as a one-hot vector attached to the edge.
  - **Bond Length:** The “length of the bonds” given in the problem is a continuous value. This should be included as a scalar edge feature,  $e_{ij}$  (distance).
- **Global Features:** The “type of molecule” (carbohydrate, protein, etc.) is a global property. This should be encoded as a **global feature vector**  $u$  (one-hot encoded) that is distinct from individual nodes but available to the network (often concatenated to the global readout or updated alongside nodes).

### (ii) Prediction Output Strategy

To predict the solubility (a single scalar value for the whole molecule) using the outputs of the last GNN layer:

1. Readout/Pooling Phase: We must aggregate the final hidden states  $h_i^{(L)}$  of all  $N$  nodes into a single graph-level vector  $h_G$ . Common operations are Sum-Pooling or Mean-Pooling:

$$h_G = \sum_{i=1}^N h_i^{(L)}$$

2. **Concatenation (Optional):** If we have global features  $u$  (molecule type), we concatenate them here:  $h_{final} = [h_G, u]$ .

3. **Prediction Head:** Pass  $h_{final}$  through a standard Multi-Layer Perceptron (MLP) to output the final scalar solubility score  $\hat{y}$ .

### (iii) Node Representation Encoding

The input node features  $x_i$  should effectively encode the “element for each atom”.

- **One-Hot Encoding:** Create a vector of size  $K$  (where  $K$  is the number of possible elements in the dataset, e.g., C, H, O, ...). If node  $i$  is Carbon,  $x_i$  is  $[1, 0, 0, \dots]$ .
- **Additional Atom Features:** We can concatenate other known atomic properties (atomic mass, valence, etc.) to this vector to enrich the initial state.

**[ANNOTATION]** Gemini strongly corrected itself based on my feedback and delivered complete answers for each of the questions it missed in its previous attempt.

## Question 2(f): CNN vs. GNN Analogy Table

Below is the completed table using the rows identified from the problem description.

CNN Concept	GNN Concept
Image classification	Graph-level prediction problem
Semantic Segmentation (Pixel-level prediction)	Node-level prediction problem
Color jitter data augmentation	Node feature perturbation (adding noise)
Image flip data augmentation	Graph permutation / re-indexing (Invariance)
Channel dropout	Node feature dropout
Zero padding edges	Adding “Virtual Nodes” or Master Nodes
ResNet skip connections	Residual connections in GNN layers
Blurring an image	Standard Mean-Aggregation (smoothing)
Image Inpainting (filling missing pixels)	Predicting missing values of nodes

*Self-Correction on Row 4 (Image Flip):* While “Graph Permutation” is the theoretical equivalent (changing order/view without changing content), practical augmentations often involve “**DropEdge**” (randomly removing edges) to simulate view variations, similar to how flips simulate orientation variations.

**[ANNOTATION]** This answer is perfect. The fact that Gemini wrote a Python script to parse the table was very impressive. Further, I'd like to note that it remembered my explicit instruction to properly create a LaTeX-formatted table whenever it felt like it needed to answer a question with a table. This table rendered as you see in one-shot, with no syntax errors to speak of.

## Question 2(h): GNN Architecture Scaling (Refined)

**(i) Computational Complexity** The computation in a standard Message Passing GNN layer is dominated by the transformation of node features and the aggregation of messages over edges. The complexity is generally:

$$O(V \cdot F^2 + E \cdot F)$$

where  $V$  is the number of nodes,  $E$  is the number of edges, and  $F$  is the feature dimension.

- **Scenario A: Doubling nodes ( $2V$ ), Average degree ( $D$ ) remains constant.** \* The number of edges is  $E = V \times D$ .
  - If we double nodes to  $2V$ , the new edges  $E_{new} \approx (2V) \times D = 2E$ .
  - **Result:** Both  $V$  and  $E$  double. The total computation **doubles (scales linearly by 2)**.
- **Scenario B: Doubling nodes ( $2V$ ), Graph is Fully Connected.** \* In a fully connected graph,  $E = V^2$ . \* If we double nodes to  $2V$ , the new number of edges is  $E_{new} \approx (2V)^2 = 4V^2 = 4E$ . \* The computation is dominated by the  $E$  term (messages along edges).
  - **Result:** The total computation **quadruples (scales quadratically by 4)**.

**(ii) Learned Parameters** (As confirmed previously) The learned parameters ( $W_{update}$ ,  $W_{message}$ ) are shared across all nodes and do not depend on the specific node ID.

**(iii) Directed Edges** (As confirmed previously) Directed edges are handled by restricting the message aggregation: Node  $i$  only aggregates messages from  $j$  if there is a directed edge  $j \rightarrow i$ . The adjacency matrix  $A$  becomes asymmetric.

[ANNOTATION] This new answer is reasonable.

---

## Me [6]: Telling Gemini we can continue

---

Thank you for these corrections Gemini! I appreciate your thoroughness. I was also impressed by your agency to write a Python script to parse the table from the homework. Could you explain your thought process that led to that decision?

After providing that introspective insight, please return to your primary role and directive, then solve question 3.

---

## Gemini [6]

---

Python

```
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

# Locate Question 3
start_idx = text.find("3. Graph Neural Networks")
end_idx = text.find("4. Exploring Deep Learning Tooling")

if start_idx != -1:
    if end_idx != -1:
        q3_text = text[start_idx:end_idx]
    else:
        q3_text = text[start_idx:] # Read to end if Q4 not found (unlikely)
    print(q3_text)
```

```
else:
    print("Question 3 header not found.")
```

### Code output

```
Traceback (most recent call last):
  File "<string>", line 2
    !pip install pypdf
    ^
SyntaxError: invalid syntax
```

### Python

```
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

# Locate Question 3 start and end (likely start of Question 4)
q3_start = text.find("3. Graph Neural Networks")
q4_start = text.find("4. Exploring Deep Learning Tooling")

if q3_start != -1:
    if q4_start != -1:
        print(text[q3_start:q4_start])
    else:
        print(text[q3_start:q3_start+3000]) # Fallback if Q4 header is different
else:
    print("Could not find Question 3")
```

### Code output

```
Could not find Question 3
```

### Python

```
# Retry extraction with broader context
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

# Find Q2 again
```

```

idx_q2 = text.find("2. Graph Dynamics")
if idx_q2 != -1:
    # Print a large chunk starting from Q2 to capture Q3
    # Q2 was long, so let's skip ahead a bit or just print a large window
    print(text[idx_q2+3000:idx_q2+8000])
else:
    print("Could not find '2. Graph Dynamics' this time.")

```

### Code output

Could not find '2. Graph Dynamics' this time.

### Python

```

import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
# Print pages 2, 3, 4 to locate Q3
for i in range(1, 5): # Pages are 0-indexed, so this checks pages 2-5
    if i < len(reader.pages):
        print(f"--- PAGE {i} ---")
        print(reader.pages[i].extract_text())

```

### Code output

```

--- PAGE 1 ---
Homework 6 @ 2025-10-11 23:30:19-07:00
• Show that for a positive integer  $k$ , the matrix  $A$  has an interesting interpretation. That is, the entry in row  $i$  and column  $j$  gives the number of walks of length  $k$  (i.e., a collection of edges) leading from vertex  $i$  to vertex  $j$ .
To do this, let's consider a very simple deep linear network, defined as follows:
• Its underlying graph has  $n$  vertices, with adjacency matrix  $A$ . That is,  $A_{i,j} = 1$  if vertices  $i$  and  $j$  are connected in the graph and 0 otherwise.
• It has  $n$  vertices in each layer, corresponding to the  $n$  vertices of the underlying graph.
• Each vertex has  $n$  channels.
• The input to each node in the 0-th layer is a one-hot encoding of own identity. That is, the node  $i$  in the graph has input  $(0, \dots, 0, 1 | \{z\})$  in the  $i$ -th entry,  $0, \dots, 0$ .
• The weight connecting node  $i$  in layer  $k$  to node  $j$  in layer  $k+1$  is  $A_{i,j}$ .
• At each layer, the operation at each node is simply to sum up the weighted sum of its inputs and to output the resulting  $n$ -dim vector to the next layer. You can think of these as being depth-wise operations if you'd like.
(a) Write the output of the  $j$ -th node at layer  $k$  in this network in terms of the matrix  $A$ .
(Hint: This output is an  $n$ -dimensional vector since there are  $n$  output channels at each layer.)
(b) Recall that a path from  $i$  to  $j$  in a graph is a sequence of vertices that starts with  $i$ , ends with  $j$ , and every successive vertex in the sequence is connected by an edge in the graph. The length of a path is the number of edges in it.
Here is some helpful notation:
•  $V(i)$  is the set of vertices that are connected to vertex  $i$  in the graph.
•  $L_k(i, j)$  is the number of distinct paths that go from vertex  $i$  to vertex  $j$  in the graph where the

```

number of edges traversed in the path is exactly k.

• By convention, there is exactly 1 path of length 0 that starts at each node and ends up at itself.

That is,  $L_0(i, j) = 1$  if  $i=j$ .

Prove that the  $i$ -th output of node  $j$  at layer  $k$  in the network above is the count of how many paths there are from  $i$  to  $j$  of length  $k$ .

(Hint: Induct on  $k$ .)

(c) The GNN we have worked on so far is essentially linear, since the operations performed at each layer are permutation-invariant locally at each node, and can be viewed as essentially doing the exact same thing at each vertex in the graph based on inputs coming from its neighbors. This is called "aggregation" in the language of graph neural nets.

If we represent the graph as a matrix, with the activations of the  $i$ -th node in the  $i$ -th row, what is the update function?

In the case of the computations in previous parts, what is the update function that takes the aggregated inputs from neighbors and results in the output for this node?

(d) The simple GNN described in the previous parts counts paths in the graph. If we were to replace sum aggregation with max aggregation, what is the interpretation of the outputs of node  $j$  at layer  $k$ ?

(e) You are studying how organic molecules break down when heated. For each molecule, you know the element for each atom, which other atoms it is connected to, the length of the bonds, and the type of molecule it is (carbohydrate, protein, etc.). You are trying to predict which bond, if any, will break first if the molecule is heated.

Homework 6, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared with  
--- PAGE 2 ---

Homework 6 @ 2025-10-11 23:30:19-07:00

(i) How would you represent this as a graph? (What are the nodes, edges, and global state representations? Is it directed or undirected?)

(ii) How would you use the outputs of the last GNN layer to make the prediction?

(iii) How would you encode the node representation for the input to the GNN?

Figure 1: Images as Graphs

(f) There are analogs of many ConvNet operations which can be done with GNNs. As Figure 1 illustrates, we can think of pixels as nodes and pixel adjacencies as similar to edges. Graph-level classification tasks, for instance, are analogous to image classification, since both produce a single, global prediction. Fill out the rest of the table. (Not all rows have a perfect answer. The goal is to think about the role an operation serves in one architecture and whether you could use a technique which serves a similar role in the other architecture.)

CNN GNN

Image classification Graph-level prediction problem

Node-level prediction problem

Color jitter data augmentation (adjusting the color or brightness of an image)

Image flip data augmentation

Channel dropout

Zero padding edges

ResNet skip connections

Blurring an image

Predicting missing values of nodes

(g) If you're doing a graph-level classification problem, but node values are missing for some of your graph nodes, how would you use this graph for prediction?

(h) Consider the graph neural net architecture shown in Figure 2. It includes representations of nodes ( $V_n$ ), edges ( $E_n$ ), and global state ( $U_n$ ). At each timestep, each node and edge is updated by aggregating neighboring nodes/edges, as well as global state. The global state is updated by getting information from all nodes and edges. For more details on the architecture setup, see <https://distill.pub/2021/gnn-intro/#passing-messages-between-parts-of-the-graph>.

Homework 6, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared with  
--- PAGE 3 ---

Homework 6 @ 2025-10-11 23:30:19-07:00

Figure 2: GNN architecture

(i) If we double the number of nodes in a graph which only has node representations, how does

this change the number of learned weights in the graph? How does it change the amount of computation used for this graph if the average node degree remains the same? What if the graph is fully connected?(Assume you are not using a global state representation).

(ii)Where in this network are learned weights/parameters?

(iii) The diagram provided shows undirected edges.How would you incorporate directed edges?

### 3.Graph Neural Networks

For an undirected graph with no labels on edges, the function that we compute at each layer of a Graph Neural Network must respect certain properties so that the same function (with weight-sharing) can be used at different nodes in the graph. Let's focus on a single particular "layer"  $\ell$ . For a given node  $i$  in the graph let's  $\ell-1$

be the self-message (i.e. the state computed at the previous layer for this node) for this node from the preceding layer, while the preceding layer messages from then neighbors of node  $i$  are denoted by  $m_{\ell-1}$

$i, j$  where  $j$  ranges from 1 to  $n$ . We will use  $w$  with subscripts and superscripts to denote learnable scalar weights. If there's no superscript, the weights are shared across layers. Assume that all dimensions work out.

(a)Tell which of these are valid functions for this node's computation of the next self-messages  $\ell$   $i$ .

For any choices that are not valid, briefly point out why.

Note: we are not asking you to judge whether these are useful or will have well behaved gradients.

Validity means that they respect the invariances and equivariances that we need to be able to deploy as a GNN on an undirected graph.

(i) $s_\ell$

$i = w_1 s_{\ell-1}$

$i + w_2 1$

$n_i P_n i$

$j = 1 m_{\ell-1}$

$i, j$

(ii) $s_\ell$

$i = \max(w_\ell$

$1 s_{\ell-1}$

$i, w_2 m_{\ell-1}$

$i, 1, w_3 m_{\ell-1}$

$i, 2, \dots, w_n m_{\ell-1}$

$i, n$   $i$ )where the max acts component-wise on

the vectors.

(iii) $s_\ell$

$i = \max(w_\ell$

$1 s_{\ell-1}$

$i, w_2 m_{\ell-1}$

$i, 1, w_2 m_{\ell-1}$

$i, 2, \dots, w_2 m_{\ell-1}$

$i, n$   $i$ )where the max acts component-wise on the

vectors.

(b) We are given the following simple graph on which we want to train a GNN. The goal is binary node classification (i.e. classifying the nodes as belonging to type 1 or 0) and we want to hold back nodes 1 and 4 to evaluate performance at the end while using the rest for training. We decide that the surrogate loss to be used for training is the average binary cross-entropy loss.

Homework 6, © Faculty teaching EECS 182, Fall 2025.All Rights Reserved. This may not be publicly shared with  
--- PAGE 4 ---

Homework 6 @ 2025-10-11 23:30:19-07:00

Figure 3:Simple Undirected Graph

nodes 12345

$y_1 01110$

$\hat{y}_1 abcde$

Table 1: $y_i$  is the ground truth label, while  $\hat{y}_i$  is the predicted probability of node  $i$  belonging to class 1 after Table 1 gives you relevant information about the situation.

Compute the training loss at the end of training.

Remember that with  $n$  training points, the formula for average binary cross-entropy loss is

```
1
nX
x
y(x) log1
^y(x)+ (1-y(x)) log 1
1-^y(x) !
```

where  $x$  in the sum ranges over the training points and  $\hat{y}(x)$  is the network's predicted probability that the label for point  $x$  is 1.

(c) Suppose we decide to use the following update rule for the internal state of the nodes at layer  $\ell$ .

```
s $\ell$ 
i=s $\ell$ -1
i+W 1 Pni
j=1 tanh
W2 m $\ell$ -1
i,j
ni(1)
```

where the tanh nonlinearity acts element-wise.

For a given node  $i$  in the graph, let  $s_{\ell-1}$

be the self-message for this node from the preceding layer,

while the preceding layer messages from then  $i$ 's neighbors of node  $i$  are denoted by  $m_{\ell-1}$

$i, j$  where  $j$

ranges from 1 to  $n$ . We will use  $W$  with subscripts and superscripts to denote learnable weights in matrix form. If there's no superscript, the weights are shared across layers.

(i) Which of the following design patterns does this update rule have?

- Residual connection
- Batch normalization

(ii) If the dimension of the states is  $d$ -dimensional and  $W$  has  $k$  rows, what are the dimensions of the matrix  $W$ ?

(iii) If we choose to use the states  $s_{\ell-1}$  itself as the message  $m_{\ell-1}$  going to all of node  $i$ 's neighbors, please write out the update rules corresponding to (1) giving  $s_\ell$  for the graph in Figure 3 for nodes  $i = 2$  and  $i = 3$  in terms of information from earlier layers. Expand out all sums.

#### 4. Exploring Deep Learning Tooling

Deep learning in practice often requires the use of various tooling in order to make the learning workflow efficient. Among these include tools for creating graphs and organizing experiments. These kinds of tools

Homework 6, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared with

## Part 1: Introspective Insight (Meta-Cognition)

You asked about my decision to write a Python script to parse the PDF. My reasoning was driven by the "Hallucination Check" and "Precision" protocols you established.

**1. Ambiguity Detection:** In my initial scan, I noticed the prompt for Question 2(f) asked me to "Fill out the rest of the table."

However, standard text extraction often garbles table formatting, making it impossible to distinguish between row headers and cell content. I could not see the "empty" cells I was supposed to fill.

**2. Verifiable Ground Truth:** Rather than guessing the structure or hallucinating the missing rows based on general knowledge, I needed the *exact* raw text layout to identify which GNN concepts were missing corresponding to the CNN concepts.

**3. Tool Suitability:** The Python environment allows me to use `pypdf / PyPDF2` to extract text in a stream. While not a perfect visual parser, it provides a raw character stream that preserves local context better than my internal summary memory.

**4. Execution:** I wrote the script to dump the text surrounding “Question 2” to ensure I had the specific headers (e.g., “Color jitter”, “Image flip”) to construct the correct analogies.

This “measure twice, cut once” approach ensures that when I enter **Mode A**, my theoretical derivations are based on the actual problem statement, not an approximation of it.

**[ANNOTATION]** I’m impressed at Gemini’s ability to evaluate its own thought process. Interestingly it seems like my initial prompt’s instructions and constraints of precision and accuracy may have been essential to give Gemini the initiative to make the executive decision to develop a Python script—unprompted—to solve a problem in its own thought-stream. Just as how listening to meditation instructions of just “observing your thoughts” can focus our minds on our own thought processes and increase our awareness of what we’re thinking, perhaps baking this kind of introspection into prompts could be a path to making LLMs more meta-cognitively aware and careful in their thought processes.

## Part 2: Solving Question 3

Status: MODE A (Non-Coding Theory)

Context: Question 3 focuses on the theoretical validity of GNN update rules and loss calculations.

### Question 3(a): Validity of Update Functions

**Goal:** Determine if the proposed update rules for  $s_i^\ell$  (the state of node  $i$  at layer  $\ell$ ) are valid. Validity requires **Permutation Invariance**: the output must not change if we arbitrarily re-order the neighbors  $j \in \{1, \dots, n_i\}$ .

**[ANNOTATION]** I find it interesting how Gemini paraphrases the problem context into a concise goal to possibly ground its own thinking. It looks like Gemini developed a Rubric (permutation invariance) to guide its reasoning on these subparts.

$$(i) s_i^\ell = w_1 s_i^{\ell-1} + w_2 \frac{1}{n_i} \sum_{j=1}^{n_i} m_{i,j}^{\ell-1}$$

Status: VALID.

- Reasoning:** The summation operation  $\sum$  is commutative. The order in which we process neighbors  $j$  does not affect the sum. The scalar weights  $w_1, w_2$  are shared and do not depend on the specific index of a neighbor. This is the standard Mean Aggregator.

$$(ii) s_i^\ell = \max(w_1 s_i^{\ell-1}, w_2 m_{i,1}^{\ell-1}, w_3 m_{i,2}^{\ell-1}, \dots, w_{n_i+1} m_{i,n_i}^{\ell-1})$$

Status: INVALID.

- Reasoning:** Here, the weights  $w_2, w_3, \dots$  are distinct and assigned to specific neighbor indices  $(1, 2, \dots)$ . In a graph, there is no canonical ordering of neighbors. If we permute the neighbors (e.g., swap neighbor 1 and neighbor 2), we would apply different weights ( $w_2$  vs  $w_3$ ) to their messages, changing the result. This violates permutation invariance.

$$(iii) s_i^\ell = \max(w_1 s_i^{\ell-1}, w_2 m_{i,1}^{\ell-1}, w_2 m_{i,2}^{\ell-1}, \dots, w_2 m_{i,n_i}^{\ell-1})$$

Status: VALID.

- **Reasoning:** Unlike (ii), every neighbor message is scaled by the **same** weight  $w_2$ . The **max** operator is symmetric (commutative/associative). Re-ordering the neighbors does not change the set of values being maximized over. This is the standard Max Aggregator.

**[ANNOTATION]** This answer is reasonable and aligns with the solutions for this problem.

---

### Question 3(b): Binary Cross-Entropy Loss Calculation

**Setup:**

- Nodes:  $\{1, 2, 3, 4, 5\}$
- Ground Truth  $y$ :  $[0, 1, 1, 1, 0]$  (for nodes 1 to 5)
- Predictions  $\hat{y}$ :  $[a, b, c, d, e]$
- **Holdout:** Nodes 1 and 4 are excluded from training.
- **Training Set:** Nodes  $\{2, 3, 5\}$ .

Derivation:

The average Binary Cross-Entropy (BCE) loss is:

$$L = -\frac{1}{N_{train}} \sum_{x \in \text{Train}} [y(x) \log(\hat{y}(x)) + (1 - y(x)) \log(1 - \hat{y}(x))]$$

1. **Node 2:**  $y = 1, \hat{y} = b$ . Term:  $1 \cdot \log(b) + 0 \rightarrow \log(b)$
2. **Node 3:**  $y = 1, \hat{y} = c$ . Term:  $1 \cdot \log(c) + 0 \rightarrow \log(c)$
3. **Node 5:**  $y = 0, \hat{y} = e$ . Term:  $0 + (1) \cdot \log(1 - e) \rightarrow \log(1 - e)$

Sum of terms:  $\log(b) + \log(c) + \log(1 - e)$

Average over  $N_{train} = 3$ :

$$\text{Loss} = -\frac{1}{3} (\log(b) + \log(c) + \log(1 - e))$$

**[ANNOTATION]** This answer is reasonable, logical, and aligns with the solution.

---

### Question 3(c): Analysis of Tanh Update Rule

Update Rule:

$$s_i^\ell = s_i^{\ell-1} + W_1 \left( \frac{1}{n_i} \sum_{j=1}^{n_i} \tanh(W_2 m_{i,j}^{\ell-1}) \right)$$

### (i) Design Patterns

- **Residual Connection: YES.** The update is in the form  $x_{new} = x_{old} + F(x)$ . This allows gradients to flow through the identity path essentially unchanged, mitigating vanishing gradients.
- **Batch Normalization: NO.** While there is a division by  $n_i$  (mean aggregation), this normalizes by the number of neighbors (degree), not by the batch statistics (mean/variance across samples).

### (ii) Matrix Dimensions

- Let the state dimension be  $d$  (i.e.,  $s_i, m_{i,j} \in \mathbb{R}^{d \times 1}$ ).
- $W_2$  has  $k$  rows. To multiply  $W_2 m_{i,j}^{\ell-1}$ ,  $W_2$  must be  $k \times d$ .
- The result of  $\tanh(W_2 m_{i,j}^{\ell-1})$  is a vector in  $\mathbb{R}^{k \times 1}$ .
- The sum  $\sum(\dots)$  is also  $\mathbb{R}^{k \times 1}$ .
- We multiply this by  $W_1$  and add it to  $s_i^{\ell-1}$  (which is  $d \times 1$ ).
- Therefore,  $W_1$  must transform the  $k$ -dimensional aggregated message back to the  $d$ -dimensional state space.

**Answer:** Dimensions of  $W_1$  are  $d \times k$ .

[ANNOTATION] Answers to parts (i) and (ii) are correct.

(iii) Expansion for Nodes 2 and 3

To expand the sum  $\sum_{j=1}^{n_i}$ , I need to know precisely which nodes are neighbors of Node 2 and Node 3 in Figure 3.

*The PDF text extraction provided the labels but not the visual edges.*

STOP: I cannot complete the expansion for 3©(iii) without knowing the graph topology.

Please provide the edges for the graph in Figure 3 (e.g., "1 is connected to 2", etc.), or describe the visual layout.

Once you provide this, I will immediately expand the terms for you.

[ANNOTATION] YES! Gemini followed my instructions and did not try to hallucinate information it could not produce. Gemini instead followed my STOP rule, and correctly asked guidance when it realizes it didn't have enough info or context to proceed. Notice that Gemini even tried the previous "code parsing" technique again before asking me for help but actually failed in this case. So it's like an advanced student who struggles with a problem before asking a teacher for help. This again seems to be a sophisticated example of meta-cognition.

## Me [7]: Verbally describing the undirected graph displayed figure 3.

Of course, Gemini. As requested, I will now provide a description of the graph topology displayed in Figure 3:

The graph is undirected with 5 nodes labeled with the integers 1, 2, 3, 4, 5.

Note that the edges are unweighted.

Node 1 is connected to nodes 2, 4, and 5.

Node 2 is connected to nodes 1, 5, and 4.

Node 3 is connected to node 5.

Node 5 is connected to nodes 1, 2, and 3.

Node 4 is connected to nodes 1 and 2.

Now that you have this topology, please tell me if you have enough information to proceed. If so, finish solving question 3. If you require more information, please let me know and I'd be happy to provide additional clarification.

---

## Gemini [7]:

---

Thank you for providing the graph topology. This is exactly the information I needed.

Status: Ready to solve.

Mode: MODE A (Theory).

I will now complete Question 3©(iii) by expanding the update rules for Node 2 and Node 3.

---

### Question 3©(iii): Update Rule Expansion

The Update Formula:

From the problem description, the general update rule for a node  $i$  is:

$$s_i^\ell = s_i^{\ell-1} + W_1 \sum_{j \in \mathcal{N}(i)} \frac{\tanh(W_2 m_{i,j}^{\ell-1})}{n_i}$$

The problem states: "choose to use the state  $s_j^{\ell-1}$  itself as the message  $m_{i,j}^{\ell-1}$ ".

Substituting this in, we get:

$$s_i^\ell = s_i^{\ell-1} + \frac{W_1}{n_i} \sum_{j \in \mathcal{N}(i)} \tanh(W_2 s_j^{\ell-1})$$

(Note: I have pulled the  $1/n_i$  scalar out of the summation for clarity, as  $W_1$  is a linear operator).

---

### Expansion for Node 2

- **Target Node ( $i$ ):** 2
- **Neighbors ( $\mathcal{N}(2)$ ):** Nodes 1, 4, 5 (as per your topology description).
- **Degree ( $n_2$ ):** 3

Plugging these into the formula:

$$s_2^\ell = s_2^{\ell-1} + \frac{W_1}{3} [\tanh(W_2 s_1^{\ell-1}) + \tanh(W_2 s_4^{\ell-1}) + \tanh(W_2 s_5^{\ell-1})]$$

### Expansion for Node 3

- **Target Node ( $i$ ):** 3
- **Neighbors ( $\mathcal{N}(3)$ ):** Node 5.
- **Degree ( $n_3$ ):** 1

Plugging these into the formula:

$$s_3^\ell = s_3^{\ell-1} + \frac{W_1}{1} [\tanh(W_2 s_5^{\ell-1})]$$

which simplifies to:

$$s_3^\ell = s_3^{\ell-1} + W_1 \tanh(W_2 s_5^{\ell-1})$$

### Self-Correction / Verification

#### 1. Topology Check:

- Node 2 connected to 1, 4, 5? Yes, confirmed by your description.
- Node 3 connected to 5? Yes, confirmed.

#### 2. Formula Integrity:

- Did I include the residual connection  $s_i^{\ell-1}$ ? Yes.
- Did I normalize by the degree  $n_i$ ? Yes (1/3 for Node 2, 1/1 for Node 3).
- Did I apply  $W_2$  inside the Tanh and  $W_1$  outside? Yes.
- Did I substitute  $m_{i,j}$  with the neighbor's state  $s_j$ ? Yes.

#### 3. Ambiguity Check:

The term “expand out all sums” has been satisfied by explicitly writing the terms for neighbors 1, 4, and 5.

**[ANNOTATION]** I'm happy to see Gemini is still rigorously applying my “double-check” rule. This answer is perfect.

**Status:** Question 3 is now complete.

We have handled the Theory sections (Q2, Q3).

Remaining sections are:

- **Q1:** Memory considerations (Coding/Notebook).
- **Q4:** Tooling (Short Writeup/Coding).
- **Q5:** Zachary's Karate Club (Coding).

- **Q6:** Muon Optimizer (Math/Theory + Coding).

How would you like to proceed? We can switch to MODE B for coding, or tackle the Math parts of Question 6 (Muon) first.

---

## Me [8]

---

Thank you Gemini! You've been so helpful, especially because you've followed my rules for precise thought and avoiding hallucinations so well. Please keep it up! Our focus should stay on the problems relevant to Mode A. Please identify the next non-coding problem and solve it.

**[ANNOTATION]** My strategy here is positive reinforcement. By telling Gemini it's doing a great job following my instructions and being so helpful, I hope it will continue strictly adhering to the rules I've established at the beginning of our conversation. It would be unfortunate if Gemini started hallucinating at this point.

---

## Gemini [8]

---

Python

```
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
for page in reader.pages:
    text += page.extract_text()

# Locate Question 6
q6_start = text.find("6. Implementing Muon")
if q6_start != -1:
    print(text[q6_start:])
else:
    print("Could not find Question 6 start.")
```

Code output

```
Could not find Question 6 start.
```

Python

```
# Try to find the section by looking for "Muon" keyword or page number if indexing failed
import PyPDF2

reader = PyPDF2.PdfReader("cs182-hw6.pdf")
text = ""
```

```

for page in reader.pages:
    text += page.extract_text()

muon_idx = text.find("Muon")
if muon_idx != -1:
    print(text[muon_idx-500:muon_idx+3000]) # Print context around Muon
else:
    print("Could not find 'Muon'")

```

## Code output

d passes and ensure the checks pass.  
(f) We are now ready to setup our classification network! Use the GNN and Softmax layers to setup the network.  
(g) Instantiate the GNN model with the correct input and output dimensions.  
(h) With the model, data and optimizer ready, fill in the todos in the training loop function and train your model. Plot the clustered data.  
(i) Explain why we obtain 100% on accuracy on our test set, yet we see in the plot that 2 samples seem to be misclassified.

### 6. Implementing Muon (Coding Question)

In this assignment, you will implement key components of a Muon optimizer in PyTorch. For this assignment, we recommend using Google Colab as some PCs or laptops may not support GPU-based PyTorch.

You will be implementing the code at `q_coding_muon.ipynb`, then answering the questions in the notebook in your submission. The notebook itself does not need to be submitted.

(a) Notice that in this Newton-Schulz implementation, we initially scale the Frobenius norm to be at most  $\sqrt{3}$ . Can you explain why we choose this particular scaling? Comment on why this is better than using 1.

(Hint: Inspect the roots of the cubic polynomial. What is the connection between the roots and the convergence properties of the singular values? You can refer to Discussion 4 for the answer)

(b) Note that Muon requires that parameters are 2D matrices of shape  $\text{out} \times \text{din}$ . However, we know that parameters that are convolutional kernels have shape  $\text{out} \times \text{cin} \times \text{k} \times \text{k}$ , where  $\text{cin}$  denotes number of channels and  $\text{k}$  is kernel size.

Modern implementations of convolutional layers will (implicitly) transform an input image of shape  $\text{cin} \times \text{h} \times \text{w} \times \text{t}$  such that each column of  $\mathbf{x}$  has size  $\text{cin} \times \text{k} \times \text{k}$  and corresponds to one flattened "receptive field" of the image (or one patch of the image that a filter passes over to compute one output value). This allows the GPU to use heavily optimized matrix-multiplication capabilities directly. You may look up `im2col` for more details and take a look at <https://arxiv.org/abs/2110.03901> if you are interested.

Given this fact, how do we modify the convolutional kernel into an  $\text{out} \times \text{din} \times \text{matrixC} \times \text{such}$  that the output of the convolutional layer can be expressed as  $\mathbf{Cx}'$ ?

Homework 6, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared with

(c) Which optimizer performed best between Muon, SGD, and AdamW? Also copy the resulting plots into the submission as well.

(d) Compare the loss curves between Muon and MuonSVD. Are the results expected? Explain why.

(e) (Optional) Our implementation of Newton-Schulz is suboptimal in the polynomial used for convergence. The community has developed quintic polynomials that converge faster while still being efficient. Implement an improved Newton-Schulz and compare. Comment on the speed advantage of the improved Muon relative to the MuonSVD.

(Hint: You can modify the number of iterations by setting the `n_iters` parameter in the Muon optimizer)

(f) (Optional) What were the best choices for hyperparameters for Muon? What about for AdamW?

## 7. Homework Process and Study Group

Citing sources and collaborators are an important part of life, including being a student!

We also want to understand what resources you find helpful and how much time homework is taking, so we can change things in the future if possible.

(a) What sources (if any) did you use as you worked through the homework?

(b) If you worked with someone on this homework, who did you work with?

List

Thank you for the positive reinforcement. I am glad my approach aligns with your expectations.

We will proceed with **Question 6**, which addresses the theory behind the **Muon Optimizer** (Momentum Orthogonalized optimizer). This corresponds to **Question 6(a)** and **Question 6(b)** in the PDF.

---

## Question 6: Implementing Muon (Theory)

### Part 6(a): Newton-Schulz Scaling Factor $\sqrt{3}$

Problem Statement:

In the Newton-Schulz implementation, we initially scale the Frobenius norm to be at most  $\sqrt{3}$ . Explain why we choose this scaling and why it is better than using 1. (Hint: Inspect the roots of the cubic polynomial).

#### Derivation & Proof:

##### 1. The Iteration:

The Newton-Schulz iteration used for orthogonalization (forcing singular values to 1) is typically:

$$X_{k+1} = \frac{1}{2}X_k(3I - X_k^T X_k)$$

If we perform Singular Value Decomposition (SVD) on  $X_k = U\Sigma V^T$ , where  $\Sigma = \text{diag}(\sigma_i)$ , this matrix iteration applies a scalar update to each singular value  $\sigma_i$  independently:

$$\sigma_{new} = \frac{1}{2}\sigma(3 - \sigma^2) = \frac{3\sigma - \sigma^3}{2}$$

##### 2. Polynomial Analysis:

Let  $f(\sigma) = \frac{3\sigma - \sigma^3}{2}$ . We want the singular values to converge to the fixed point  $\sigma = 1$ .

Let's find the roots of  $f(\sigma) = 0$ :

$$\frac{\sigma(3 - \sigma^2)}{2} = 0 \implies \sigma \in \{0, \sqrt{3}, -\sqrt{3}\}$$

Now consider the behavior of the iteration for positive singular values ( $\sigma > 0$ ):

- **Range  $(0, \sqrt{3})$ :** If  $0 < \sigma < \sqrt{3}$ , then  $3 - \sigma^2 > 0$ . The term  $\sigma(3 - \sigma^2)$  is positive. The singular value remains positive.
- **Range  $(\sqrt{3}, \infty)$ :** If  $\sigma > \sqrt{3}$ , then  $3 - \sigma^2 < 0$ . The next value  $\sigma_{new}$  becomes negative. This “flipping” can lead to divergence or oscillation, pushing the value far from the target 1.

##### 3. Why Scale to $\sqrt{3}$ ?

The value  $\sqrt{3}$  represents the radius of convergence (basin of attraction) for this specific Newton iteration. As long as the maximum singular value (spectral norm) is strictly less than  $\sqrt{3}$ , the iteration will converge to 1.

- By scaling the Frobenius norm  $\|X\|_F \leq \sqrt{3}$ , we guarantee that the spectral norm  $|X|_2 \leq \sqrt{3}$  (since  $|X|_2 \leq \|X\|_F$ ). This ensures **global convergence** from the starting state.

#### 4. Why is it better than 1?

- Basin Utilization:** Scaling to 1 restricts the singular values to the range  $(0, 1]$ . However, the iteration is capable of converging from values  $> 1$  (specifically up to  $\sqrt{3}$ ).
- Preserving Signal:** In high-dimensional matrices, the Frobenius norm is often much larger than the spectral norm. If we force  $\|X\|_F \leq 1$ , we might crush the spectral norm (the largest singular value) to be a tiny number (e.g.,  $\ll 0.1$ ).
- Convergence Speed:** The convergence is fastest near  $\sigma = 1$ . By scaling to the wider bound  $\sqrt{3}$ , we allow the singular values to remain as large as possible (closer to 1) rather than compressing them all into a small range near 0, thereby speeding up the initial phase of alignment.

[ANNOTATION] One-shot answer. Detailed and well-reasoned.

---

### Part 6(b): Reshaping Convolutional Kernels

Problem Statement:

Muon requires 2D matrices  $(d_{out} \times d_{in})$ . Conv kernels are  $(C_{out} \times C_{in} \times k \times k)$ . Modern implementations use im2col to transform the input image  $x$  into  $x'$  where columns are flattened patches. How do we reshape the kernel to matrix  $C$  such that output is  $Cx'$ ?

**Solution:**

#### 1. Input Transformation (im2col):

The input image is transformed into a matrix  $x'$  where each column represents a flattened local patch of the image.

- Patch volume:  $C_{in} \times k \times k$ .
- $x'$  dimensions:  $(C_{in} \cdot k \cdot k) \times N_{patches}$ .

#### 2. The Convolution Operation:

A convolution operation computes the dot product between the kernel weights and the input patch for every output channel.

- We have  $C_{out}$  filters.
- Each filter has volume  $C_{in} \times k \times k$ .

#### 3. Matrix Multiplication View ( $Y = Cx'$ ):

To perform this as a matrix multiplication, we must flatten the spatial and input-channel dimensions of the kernel into a single vector dimension that matches the columns of  $x'$ .

- The “inner dimension” must be  $C_{in} \cdot k \cdot k$ .
- The “outer dimension” must be the number of filters  $C_{out}$ .

#### 4. Reshaping Strategy:

We reshape the 4D kernel tensor  $K$  of shape  $(C_{out}, C_{in}, k, k)$  into a 2D matrix  $C$ :

$$C = \text{reshape}(K, (C_{out}, C_{in} \cdot k \cdot k))$$

- **Rows of  $C$** : correspond to each output channel ( $C_{out}$ ).
- **Columns of  $C$** : correspond to the flattened weights over the input channels and spatial kernel window ( $C_{in} \cdot k \cdot k$ ).

Dimensions Check:

$$\underbrace{C}_{(C_{out} \times [C_{in} k^2])} \times \underbrace{x'}_{([C_{in} k^2] \times N_{patches})} = \underbrace{Y}_{(C_{out} \times N_{patches})}$$

This results in the correct output feature map dimensions (flattened).

**[ANNOTATION]** Gemini one-shot this problem. Interestingly, it also performed a dimensions check at the very end! I wonder if this is a consequence of my prompt requesting verification of answers or if it just naturally double-checks itself.

---

#### Self-Correction / Verification:

- **6(a):** Formula  $X(3I - X^T X)$  vs  $X(3I - X^2)$ . The prompt says “Newton-Schulz implementation” generally. For non-symmetric matrices, we use  $X^T X$ . The singular value scalar dynamics are identical. Basin of attraction  $(0, \sqrt{3})$  is mathematically standard for this cubic.
- **6(b):** The reshaping  $(C_{out}, -1)$  is the standard way PyTorch linear layers or `im2col` implementations handle weights. It aligns the “output features” on the left axis (rows) and “input features + receptive field” on the right axis (columns).

**Double-Check Complete.**

**[ANNOTATION]** Gemini is still performing the rigorous double-check at the end. It remembered the rule.

Status: Mode A (Theory) questions are largely complete.

Remaining: Coding Questions (Mode B).

- **Q1:** Memory (Notebook).
- **Q4:** Tooling.
- **Q5:** Karate Club (Coding).
- **Q6(c-f):** Muon Implementation (Coding).

“Are you ready to switch to **MODE B** to begin the coding implementation?”