JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT

Make-up Examination-Nov-2025

COURSE CODE (CREDITS):25B1WCI511(2)　　　　　MAX. MARKS: 25

COURSE NAME: PROMPT ENGINEERING

COURSE INSTRUCTORS:VANI SHARMA　　　　　MAX. TIME: 1 Hour 30 Minutes

Note:Note:(a)All questions are compulsory.

(b) The candidate is allowed to make Suitable numeric assumptions wherever required

for solving problems

| Q.No | Question | CO | Marks |
|---|---|---|---|
| Q1 | A language model predicts the following next-token probabilities: <br><br> | Token | Probability | <br>|---|---|<br>| A | 0.28 |<br>| B | 0.22 |<br>| C | 0.18 |<br>| D | 0.15 |<br>| E | 0.10 |<br>| F | 0.07 |<br><br>(a) If greedy decoding is used, which token will be selected? <br>(b) If top-k sampling with $k = 3$ is used, identify the candidate set and compute the renormalized probabilities. <br>(c) Explain how temperature scaling could affect the diversity of the generated output. | [2] | [7] |
| Q2 | (a) Explain in detail how positional encoding enables Transformers to understand and process sequential data.Discuss why it is necessary despite the model's parallel processing nature, and describe the mathematical or conceptual mechanism used to represent positional information. <br><br>(b) Describe thoroughly how cross-attention facilitates encoder–decoder interaction in models such as T5 or BERT2BERT. Explain the role of query, key, and value vectors in this process and how this interaction allows the decoder to generate context-aware outputs. | [2] | [3+3] |
| Q3 | (a) What are system prompts, and how do they influence model behaviour compared to user prompts? <br><br>(b)An LLM is tasked with solving a math word problem: <br>"If a shop sells 8 pencils for ₹40, how much will 15 pencils cost?" <br>You notice that the model sometimes gives direct, wrong answers. <br>To improve accuracy, you modify the prompt: | [3] | [2+3] |

| | | | |
|---|---|---|---|
| | "Think step by step before answering."<br><br>Explain how *chain-of-thought prompting* changes the reasoning behavior of the model. What does this approach reveal about the model's internal reasoning process? | | |
| Q4 | (a) Compare the architecture and functionality of LSTM and GRU networks in detail. Explain how their gating mechanisms control information flow and address the vanishing gradient problem.<br><br>(b) Write the mathematical equations of all gates used in LSTM (input, forget, output) and GRU (update, reset). Highlight the differences in how each model updates its hidden state.<br><br>(c) Which of the two — LSTM or GRU — is generally considered more efficient, and why? | [2] | [3+3+1] |