

# UMP-Net: Uncertainty-Aware Mixture of Prompts Network for Efficient Instruction Tuning

Anonymous authors

Paper under double-blind review

## Abstract

Instruction tuning has greatly improved how large language models (LLMs) respond to human-like instructions. However, fully fine-tuning these models is still computationally demanding, and many existing parameter-efficient methods fall short—particularly when it comes to uncertainty estimation and working effectively across different modalities. To address this, we introduce UMP-Net (Uncertainty-Aware Mixture of Prompts Network), a new approach designed to enhance the ability of LLaMA to follow instructions. UMP-Net combines a novel mixture of prompts (MoPs) technique with Latent Noise Prompting, KNN-based Heterogeneous Clustering, and Conformal Predictions to select the most reliable prompts dynamically while accounting for uncertainty. In addition, it features a CLIP-based multi-modal architecture to streamline vision-language integration. We evaluated UMP-Net on a range of benchmarks including ScienceQA, COCO Caption, and various zero-shot multi-modal tasks. The results show a strong performance: an average accuracy of 88.41% on ScienceQA and a CIDEr score of 158.3 on COCO Caption—surpassing models such as LLaVA, LLaMA-Adapter, and LLaMA-Excitor. These findings suggest that UMP-Net offers both improved multi-modal capability and computational efficiency.

## 1 Introduction

Instruction tuning has quickly become a key method for improving the capabilities of large language models (LLMs), allowing them to better interpret and follow human instructions in a wide range of tasks (Ouyang et al., 2022; Wei et al., 2022). Early successes with models like FLAN Wei et al. (2022) and InstructGPT Ouyang et al. (2022) highlighted how fine-tuning pre-trained LLMs using instruction datasets could significantly boost their zero-shot and few-shot performance. Despite these gains, most of these approaches depend on full model fine-tuning—a process that is not only resource-intensive but also becomes impractical when working with very large models such as LLaMA Touvron et al. (2023), which contain billions of parameters. Moreover, the emergence of multi-modal large language models (MMLMs) adds another layer of complexity. Combining visual and textual input often demands even more extensive pre-training or fine-tuning, intensifying the already high computational costs (Liu et al., 2023b; Li et al., 2023a).

To reduce the computational burden of full fine-tuning, researchers have developed parameter-efficient fine-tuning (PEFT) methods such as LoRA Hu et al. (2021) and prompt tuning Lester et al. (2021), which adjust only a small portion of the parameters of a model while keeping the core language model unchanged. While these approaches are notably more efficient, they often fall short when it comes to generalizing across a wide range of tasks and modalities—especially in multi-modal contexts, where aligning visual and textual data plays a crucial role. Another key limitation is the lack of tools to quantify and manage uncertainty in model predictions. This becomes particularly important when dealing with ambiguous or noisy input, where clear guidance is essential. These challenges are even more apparent in multi-modal applications, where models such as Flamingo Alayrac et al. (2022) and LLaVA Liu et al. (2023b) depend heavily on large-scale datasets for vision-language alignment—making them less scalable and harder to adopt in resource-constrained environments.

In this work, we present UMP-Net (Uncertainty-Aware Mixture of Prompts Network), a new framework aimed at addressing the limitations of existing instruction-tuned and multi-modal systems. UMP-Net combines uncertainty-aware prompt tuning with an efficient strategy for multi-modal adaptation. At its core is a mixture of prompts (MoPs) mechanism, which brings together Latent Noise Prompting, KNN-based Heterogeneous Clustering (HeteroGraphPrompt), and Cluster-Wise Uncertainty Estimation (CUE) to dynamically tailor prompts for the LLaMA model. To improve reliability, the system incorporates Conformal Predictions, allowing it to measure uncertainty at both the prompt and cluster levels and guide the selection process accordingly. On the multi-modal front, UMP-Net uses CLIP-based embeddings [Radford et al. \(2021\)](#) to integrate visual data, enabling effective cross-modal reasoning without the need for costly pre-training. This approach not only strengthens LLaMA’s ability to follow instructions in both language-only and multi-modal settings, but also keeps computational demands low—making it well-suited for use in environments with limited resources. To evaluate the efficacy of our proposed UMP-Net, we conducted a comparative analysis with existing models, LLaMA-Adapter [Zhang et al. \(2024\)](#) and LLaMA-Excitor [Zou et al. \(2024\)](#), across diverse tasks blending visual and textual inputs. Figure 1 illustrates this comparison, highlighting UMP-Net’s superior performance in tasks such as identifying solution concentrations, recognizing botanical features, listing medical specialties, and generating functional code.

This paper makes three key contributions. (1) We introduce UMP-Net, a parameter-efficient framework that combines uncertainty-aware prompt tuning with multi-modal adaptation. It achieves state-of-the-art results on both text-only and vision-language benchmarks. (2) We propose a new clustering and uncertainty estimation pipeline that leverages KNN-based prompt categorization and Conformal Predictions. It helps to enhance prompt reliability while minimizing redundancy. (3) We validate UMP-Net’s performance through comprehensive experiments on benchmarks such as ScienceQA [Lu et al. \(2022a\)](#), COCO Caption [Chen et al. \(2015\)](#), and a range of zero-shot multi-modal tasks. Across these settings, UMP-Net consistently outperforms leading models including LLaVA [Liu et al. \(2023b\)](#), LLaMA-Adapter [Zhang et al. \(2024\)](#), and LLaMA-Excitor [Zou et al. \(2024\)](#).

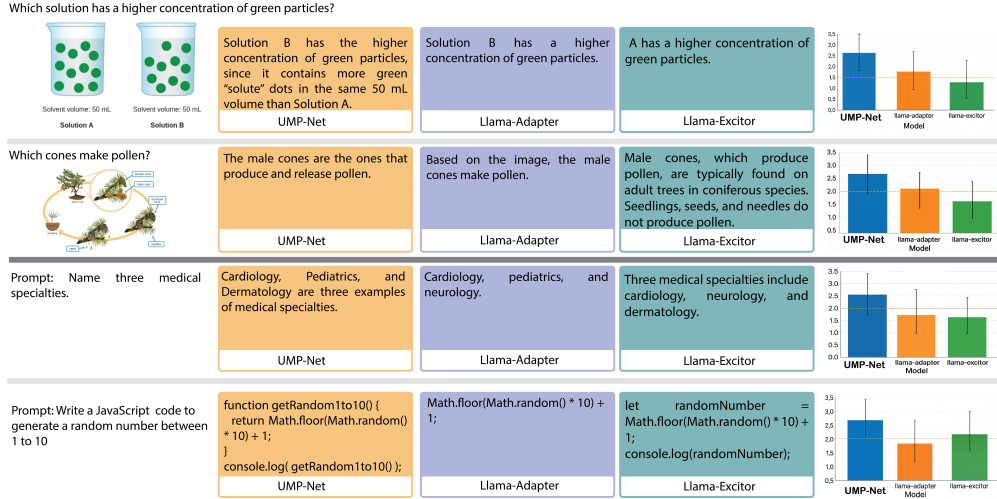


Figure 1: Comparison of UMP-Net, LLaMA-Adapter [Zhang et al. \(2024\)](#), and LLaMA-Excitor [Zou et al. \(2024\)](#) across four mixed visual-text tasks: identifying solution concentration, recognizing pollen-producing cones, listing medical specialties, and generating JavaScript code. Bar charts to the right show UMP-Net’s superior mean human evaluation scores with standard deviation error bars.

## 2 Proposed Method

In this section, we present UMP-Net, a novel framework developed to adapt the LLaMA model by enhancing its performance through a learnable adaptation prompt based on a mixture of prompts strategy. Our method brings together Bayesian reasoning, Conformal Prediction, and KNN-based heterogeneous clustering

to build a more robust and uncertainty-aware prompting mechanism. This combination allows UMP-Net to dynamically tailor prompts while maintaining a high degree of reliability in its predictions.

## 2.1 Overview of UMP-Net

UMP-Net uses a modular architecture designed to dynamically generate and weight prompts based on their associated uncertainty scores. As illustrated in Figure 2, the model is composed of three core components. (1) Latent Noise Prompting combined with MoPs strategy, (2) KNN-based Heterogeneous Clustering for selecting and aggregating relevant prompts, and (3) Conformal Predictions for estimating uncertainty across prompt candidates. These modules are interconnected through Attention Gates and Softmax layers, enabling the system to compute a single, reliable prompt that adapts effectively to the LLaMA model’s needs.

## 2.2 Latent Noise Prompting with MoPs

The Latent Noise Prompting module is a core component of UMP-Net, designed to introduce controlled variability into the prompt generation process for the adaptation of LLaMA. This module begins by sampling latent noise  $Z$  from a standard Gaussian distribution  $N(0, I)$ , where  $I$  denotes the identity matrix, ensuring isotropic noise with zero mean and unit variance. The dimensionality of  $Z$  is denoted by  $d_z$ , representing the dimension of the latent space, typically aligned with the input embedding size of the LLaMA model.

The sampled latent noise  $Z \in \mathbb{R}^{d_z}$  is then processed through a Multi-Layer Perceptron (MLP) to generate an MoP, denoted as  $P_{1:n}$ , where  $n$  represents the number of prompts in the mixture. Each prompt  $P_i \in \mathbb{R}^{d_p}$  (where  $d_p$  is the embedding dimension of the prompt) is a vector representation that captures the various semantic and syntactic characteristics of potential inputs. The MLP, parameterized by weights  $W^{(l)}$  and biases  $b^{(l)}$  across  $L$  layers, transforms the latent noise as follows:

$$H^{(l)} = \sigma(W^{(l)}H^{(l-1)} + b^{(l)}), \quad l = 1, 2, \dots, L, \quad (1)$$

where  $H^{(0)} = Z$ ,  $H^{(L)} = P_{1:n}$ , and  $\sigma$  is a non-linear activation function (ReLU). The resulting  $P_{1:n}$  serves as the input for subsequent modules, such as Conformal Predictions and Heterogeneous Clustering, to further refine prompt selection and weighting.

## 2.3 Heterogeneous Clustering by KNN

To refine prompt selection and enhance the adaptability of the UMP-Net for LLaMA, we introduce a KNN based heterogeneous clustering approach. This module organizes the MoPs  $P_{1:n}$  into distinct clusters—textual, visual, and cross-modal—based on their feature representations, enabling effective handling of diverse input modalities and ensuring robustness across tasks.

We leverage the KNN algorithm to identify structurally similar prompts, grouping them into  $K$  clusters  $C_{1:K}$ , where each cluster  $C_k$  contains  $m_k + 1$  prompts. Here,  $m_k$  represents the number of prompts in cluster  $k$ , and the "+1" accounts for a representative or centroid prompt. The KNN clustering is performed as follows:

$$\text{Distance}(P_i, P_j) = \|P_i - P_j\|_2, \quad (2)$$

where  $\|\cdot\|_2$  denotes the Euclidean distance. For each prompt  $P_i$ , we identify its  $k$  nearest neighbors based on  $\text{Distance}(P_i, P_j)$ , and group prompts into clusters  $C_k$ . This approach offers several benefits, particularly when managing a large number of prompts to avoid redundancy, improve performance, and enhance uncertainty quantification in downstream modules like Conformal Predictions. The KNN-based clustering organizes prompts into coherent, modality-specific clusters, each tailored to distinct functional roles: Visual Prompts (V-Prompts Cluster,  $C_k^V$ ): Includes prompts specialized in processing visual inputs, such as object recognition, spatial reasoning, or image understanding. These prompts are clustered based on visual feature similarities, e.g., embeddings from a vision transformer. Textual Prompts (T-Prompts Cluster,  $C_k^T$ ): Groups prompts focused on textual reasoning, such as sentence embeddings, text completion, or semantic parsing. Clustering is based on linguistic feature similarities, derived from a language model encoder. Unified Cross-Modal Prompts (VL-Prompts Cluster,  $C_k^{VL}$ ): Combines prompts that handle tasks involving both visual and

linguistic modalities, such as visual question answering or image captioning. These prompts are clustered based on joint embeddings that integrate both visual and textual features.

Moreover without clustering, a large collection of prompts could lead to redundancy or conflicting predictions, making it difficult to calibrate and fuse them effectively within UMP-Net. The KNN-based clustering approach addresses this issue by grouping structurally similar prompts into functional units, which reduces redundancy and simplifies prompt management. For example, with  $K = 3$  (a typical choice for modality-specific clustering), prompts are divided into distinct clusters— $C_k^V$ ,  $C_k^T$ , and  $C_k^{VL}$ —minimizing overlap and ensuring that each cluster serves a unique, well-defined purpose.

## 2.4 Conformal Predictions for Uncertainty Quantification

To quantify uncertainty in the UMP-Net, we employ Conformal Predictions, a distribution-free statistical framework that provides reliable uncertainty estimates for model outputs. This module assesses the reliability of each prompt in the mixture  $P_{1:n}$  by computing nonconformity scores, which measure how well a given prompt aligns with the expected output for a specific input. These scores are used to derive confidence levels, enabling the selection of the most reliable prompts for LLaMA adaptation. For each prompt  $P_i$ , we compute a nonconformity score  $S(P_i, x, y)$  based on the input  $x$  and corresponding label or output  $y$ , where  $x$  represents the task input (e.g., text, image, or multimodal data), and  $y$  is the predicted or target output. The nonconformity scores are computed differently for each prompt type (visual, textual, or cross-modal) within their respective clusters, as detailed below.

**Visual Prompts (V-Prompts).** For visual prompt  $P_i$ , the nonconformity score reflects how well the input image aligns with the visual features expected by prompt  $P_i$  in the visual cluster. We define the V-Prompt nonconformity score as:

$$S(P_i, x, y) = \|f(x) - g_i(y)\|_2^2, \quad (3)$$

where  $f(x) \in \mathbb{R}^{d_v}$  represents the visual feature embedding of the input image  $x$ , extracted using a pre-trained vision model (e.g., a convolutional neural network or vision transformer).  $g_i(y) \in \mathbb{R}^{d_v}$  is the label embedding generated by the  $i$ -th visual prompt  $P_i$  in the cluster, mapping the output  $y$  (e.g., a predicted class or description) into the visual feature space.

The intuition behind this formulation is that aggregating the Euclidean distances across this prompt ensures the nonconformity measure reflects the collective alignment of the input image with the prompt’s expected visual features. A lower  $S(P_i, x, y)$  indicates a higher conformity (i.e., the input aligns well with the prompt), corresponding to a lower uncertainty.

**Textual Prompts (T-Prompts).** For the cluster of textual prompts, where predictions are based on linguistic input, the nonconformity score reflects the negative log-likelihood of the label  $y$  given the input  $x$ . We define the prompt-level nonconformity score as:

$$S(P_i, x, y) = -\log P_i(y|x), \quad (4)$$

where  $P_i(y|x)$  is the probability of the label  $y$  given the input  $x$ , as predicted by the  $i$ -th textual prompt  $P_i$  in the cluster.

This aggregated score captures the collective confidence of the textual prompts in the cluster. A higher  $S(P_i, x, y)$  indicates lower conformity (i.e., greater uncertainty), as it reflects a lower likelihood of the predicted label  $y$  aligning with the input  $x$ . This formulation leverages the probabilistic nature of language models, ensuring that uncertainty is quantified in terms of predictive confidence.

**Unified Cross-Modal Prompts (VL-Prompts Cluster).** For cross-modal prompts, which rely on both visual and textual modalities, we define a weighted hybrid nonconformity score that balances contributions from both domains. The cluster-level nonconformity score is given by:

$$S(P_i, x, y) = [\lambda \|f(x) - g_i(y)\|_2^2 - (1 - \lambda) \log P_i(y|x)], \quad (5)$$

where  $\lambda \in [0, 1]$  is a hyperparameter that balances the contributions of the visual ( $\|f(x) - g_i(y)\|_2^2$ ) and textual ( $-\log P_i(y|x)$ ) components. Moreover,  $f(x)$ ,  $g_i(y)$ , and  $P_i(y|x)$  are defined as in the visual and textual cases, respectively.



This formulation ensures that the nonconformity score reflects the collective judgment of the cluster across modalities. The parameter  $\lambda$  is tuned based on the task requirements, allowing flexibility to emphasize either visual or textual information. A lower  $S(P_i, x, y)$  indicates higher conformity and lower uncertainty, enabling UMP-Net to adapt LLaMA effectively to multimodal inputs. The nonconformity scores of all prompts guides the selection of the best confident prompt from each cluster,  $P_{\text{best}}$ , for LLaMA adaptation, as described in the Attention Gate and Weighted Prompt Creation module.

## 2.5 Attention Gate and Weighted Prompt Creation

The best confident prompts from each cluster are passed through an Attention Gate, which dynamically weights each prompt based on its relevance. The Attention Gate employs a Softmax layer to normalize attention scores, producing a weighted prompt  $P_{\text{weighted}}$ . The weighting process is guided by:

$$P_{\text{weighted}} = \sum_{k=1}^K \alpha_k P_{\text{best},k}, \quad (6)$$

where  $\alpha_k$  represents the attention weight for prompt  $P_{\text{best},k}$ , learn through the training time. The output weighted prompt is then selected for LLaMA adaptation.

The final weighted prompt is integrated into the LLaMA model, enabling it to adapt to inputs effectively. By incorporating uncertainty-aware prompts, UMP-Net improves LLaMA’s ability to generate coherent and contextually appropriate responses, particularly in scenarios with limited or noisy data (see Figure 2). The UMP-Net pipeline can be summarized as Algorithm 1.

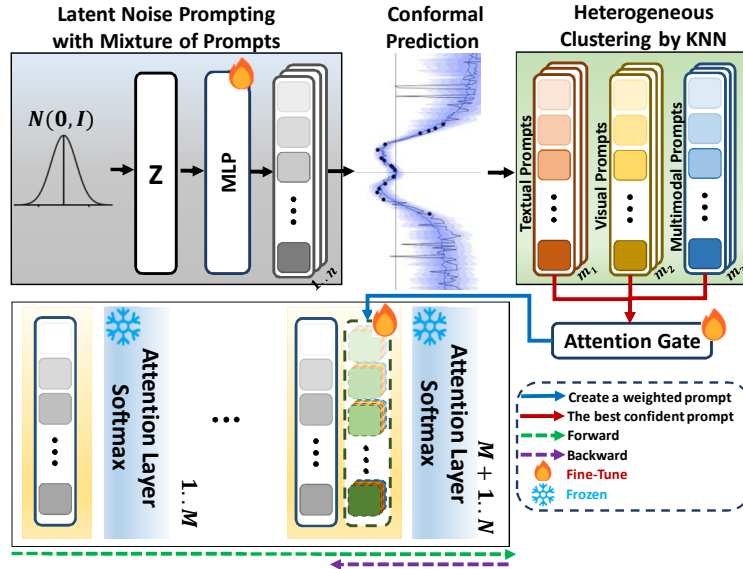


Figure 2: Diagram illustrating the architecture of UMP-Net for LLaMA adaptation. The process begins with Latent Noise Prompting, where noise from a normal distribution  $N(0, I)$  is processed through an MLP to generate initial prompts. These prompts undergo Conformal Prediction to assess uncertainty, followed by Heterogeneous Clustering using KNN to categorize them into textual, visual, and multi-modal clusters. An Attention Gate then selects the best confident prompt from each cluster via a Softmax Layer, creating a weighted prompt. The final weighted prompt is integrated into LLaMA for enhanced instruction-following, with frozen layers ensuring efficiency.

## 2.6 Multi-modal Architecture

The multi-modal architecture of UMP-Net enhances its ability to process various input modalities by integrating image embeddings in multiple stages of the pipeline, as illustrated in Figure 3. This architecture

**Algorithm 1** UMP-Net Algorithm for LLaMA Adaptation

---

**Require:** Input  $x$  (task input),  $d_p$  (latent dimension),  $n$  (number of prompts),  $K$  (number of clusters),  $k$  (KNN neighbors),  $L$  (MLP layers),  $\lambda$  (cross-modal weight), pre-trained LLaMA model

**Ensure:** Weighted prompt  $P_{\text{weighted}}$ , predicted output  $y_{\text{pred}}$

- 1: **1. Latent Noise Prompting:** Sample  $Z \sim N(0, I)$  with dimension  $d_z$
- 2: Process  $Z$  through MLP with  $L$  layers ( $H^{(l)} = \sigma(W^{(l)}H^{(l-1)} + b^{(l)})$ ) to generate  $P_{1:n} \in \mathbb{R}^{d_p}$
- 3: **2. Heterogeneous Clustering by KNN:**
- 4: Partition into  $K$  clusters  $C_{1:K}$  ( $C_k^V, C_k^T, C_k^{VL}$ )
- 5: **3. Conformal Predictions:**
- 6: **for** each  $P_i \in P_{1:n}$  **do**
- 7:   **if**  $P_i \in C_k^V$  **then**
- 8:      $S(P_i, x, y) = \|f(x) - g_i(y)\|_2^2$  ▷ Visual
- 9:   **else if**  $P_i \in C_k^T$  **then**
- 10:      $S(P_i, x, y) = -\log P_i(y|x)$  ▷ Textual
- 11:   **else if**  $P_i \in C_k^{VL}$  **then**
- 12:      $S(P_i, x, y) = \lambda\|f(x) - g_i(y)\|_2^2 - (1 - \lambda)\log P_i(y|x)$  ▷ Cross-modal
- 13:   **end if**
- 14: **end for**
- 15: Select best prompt  $P_{\text{best},k}$  per cluster with lowest  $S(P_i, x, y)$
- 16: **4. Attention Gate:**
- 17: Compute  $P_{\text{weighted}} = \sum_{k=1}^K \alpha_k P_{\text{best},k}$  through learning
- 18: **5. LLaMA Integration:**
- 19: Feed  $P_{\text{weighted}}$  into LLaMA to get  $y_{\text{pred}} = \text{LLaMA}(x, P_{\text{weighted}})$
- 20: **return**  $P_{\text{weighted}}, y_{\text{pred}}$

---

leverages the CLIP model [Radford et al. \(2021\)](#) to project multimodal features and incorporates image embeddings into the prompt embeddings after clustering and across all attention layers, ensuring robust multimodal integration. The architecture includes the following key components and processes:

**CLIP-based Image Embedding.** The input image  $x_{\text{img}}$  (e.g., the cat image in Figure 3) is processed through CLIP to extract a visual embedding:

$$e_{\text{img}} = \text{CLIP}_{\text{visual}}(x_{\text{img}}) \in \mathbb{R}^{d_c}, \quad (7)$$

where  $d_c$  is the CLIP embedding dimension, aligning with the prompt embedding dimension  $d_p$ .

**Embedding Addition to Each Prompt.** After Heterogeneous Clustering by KNN, the image embedding  $e_{\text{img}}$  is added to each prompt  $P_i \in C_k$  within clusters  $C_k^V, C_k^T$ , and  $C_k^{VL}$ . For each prompt  $P_i \in \mathbb{R}^{d_p}$ , the augmented embedding is computed as:

$$P_i^{\text{aug}} = P_i + W_{\text{proj}} e_{\text{img}}, \quad (8)$$

where  $W_{\text{proj}} \in \mathbb{R}^{d_p \times d_c}$  is a learnable projection matrix ensuring dimensional compatibility ( $d_p = d_c$  after projection).

**Confidence Score Computation for Each Prompt.** Using the augmented prompts  $P_i^{\text{aug}}$ , we recompute the nonconformity scores as described in Section 2.4. For each prompt  $P_i$ , the nonconformity score  $S(P_i^{\text{aug}}, x, y)$  is calculated based on its cluster type. The confidence score  $\text{conf}(P_i^{\text{aug}})$  is then derived as the inverse of the nonconformity score:

$$\text{conf}(P_i^{\text{aug}}) = \frac{1}{1 + S(P_i^{\text{aug}}, x, y)}, \quad (9)$$

ensuring that lower nonconformity (higher conformity) corresponds to higher confidence.

**Selection of Best Prompts from Each Cluster.** For each cluster  $C_k$ , we select the prompt with the highest confidence score as the best confident prompt:

$$P_{\text{best},k} = \arg \max_{P_i^{\text{aug}} \in C_k} \text{conf}(P_i^{\text{aug}}). \quad (10)$$

This results in  $K$  best prompts  $P_{\text{best},1:K}$ , one from each cluster.

The selected best prompts  $P_{\text{best},k}$  are passed to the Attention Gate, which computes attention weights  $\alpha_k$  using a softmax layer.

**Multimodal Integration:** The integration of  $e_{\text{img}}$  into each prompt enhances UMP-Net’s ability to handle tasks such as visual question answering (e.g., processing the cat image in Figure 3). The confidence-based selection and attention mechanism ensure that the most reliable prompts are prioritized, improving the quality of the final weighted prompt for the LLaMA adaptation.

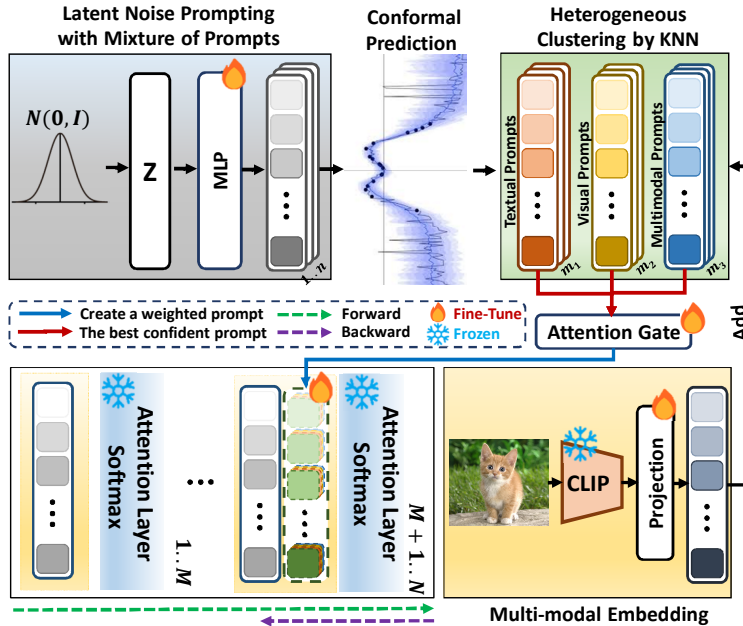


Figure 3: Illustration of the multi-modal architecture of UMP-Net, emphasizing the integration of visual and textual embeddings for enhanced LLaMA adaptation. A key multi-modal step involves adding CLIP-based image embeddings to each prompt to create augmented multi-modal embeddings.

This multi-modal architecture strengthens UMP-Net’s capability to process diverse data types, leveraging CLIP’s pre-trained visual representations and the systematic integration of image embeddings to optimize performance for LLaMA adaptation. This proposed method significantly improves LLaMA’s robustness and adaptability, as demonstrated in subsequent experimental sections.

### 3 Experiments

#### 3.1 Language Only Performance Assessment

**Experimental Setup.** Following the Stanford Alpaca [Taori et al. \(2023a\)](#), we employ a data set of 52K instruction-following examples for training purposes. The UMP-Net model is fine-tuned using 2 RTX 4090 GPUs over 4 epochs. We configure the training with two warmup epochs, a batch size of 8, a learning rate of 0.009, and a weight decay of 0.02. By default, we utilize the LLaMA-Adapter [Zhang et al. \(2024\)](#) pre-trained for version LLaMA2 7B and the foundation pre-trained LLaMA model with 8B version LLaMA3 parameters and  $N = 32$  transformer layers. The prompt length is set to  $d_p = 40$ , and the adaptation prompts are integrated into the final  $M = 30$  layers of the model. For quantitative evaluation, we compare

our approach against methods trained on the same 52K instruction dataset, specifically Alpaca [Taori et al. \(2023a\)](#), Alpaca-LoRA [Tloen \(2023\)](#), LLaMA-Adapter [Zhang et al. \(2024\)](#) and LLaMA-Excitor [Zou et al. \(2024\)](#), using the widely recognized GPT-4 evaluation benchmark ([Chiang et al., 2023b](#)).

**Instruction-Following Performance Assessment.** Table 1 provides a comprehensive evaluation of the ability to follow instruction in the proposed UMP-Net model and competing approaches.

This question assesses UMP-Net’s capability to interpret abstract philosophical concepts, integrate historical perspectives, and provide a reasoned preference, evaluating its critical thinking and language comprehension skills. Table 1 compares responses from various models to the prompt on defining intelligence across philosophical perspectives. UMP-Net excels by delivering a comprehensive definition of intelligence that synthesizes cognitive, practical, moral, and emotional dimensions, drawing on the philosophies of Plato, Aristotle, and Confucius, and offering a balanced preference. In contrast, LLaMA7B provides a narrow definition focused on problem-solving, lacking philosophical depth. Alpaca-LoRA [Tloen \(2023\)](#) and LLaMA-Adapter [Zhang et al. \(2024\)](#) discuss Plato and Aristotle but fail to incorporate modern or emotional aspects, limiting their responses’ breadth. LLaMA-Excitor [Zou et al. \(2024\)](#) covers multiple theorists but presents a fragmented perspective without a cohesive synthesis. UMP-Net’s response is distinguished by its holistic integration and contextual relevance, making it the most robust solution for this philosophical inquiry.

Table 1: Instruction-following Comparison for "What’s the intelligence? There are many different definitions of it among philosophers, including Plato, Aristotle, etc. Which do you prefer? Explain in two sentences only".

LLaMA7B: Plato viewed intelligence as the soul’s ability to recall eternal truths through reason, while Aristotle emphasized practical wisdom (phronesis) and rational thought grounded in experience. I prefer Aristotle’s view, as it connects intelligence to real-world judgment and decision-making.
Alpaca-LoRA <a href="#">Tloen (2023)</a> : Intelligence is the capacity to learn, understand, and apply knowledge, as Aristotle emphasized, which I prefer for its practical focus. Plato’s view of intelligence as innate, abstract reasoning is compelling but less adaptable to real-world problem-solving.
LLaMA-Adapter <a href="#">Zhang et al. (2024)</a> : As an AI, I don’t possess personal preferences like humans do regarding philosophical definitions. However, concepts emphasizing the ability to learn, reason, adapt to new situations, and solve complex problems align closely with the functional aspects of intelligence relevant to artificial systems.
LLaMA-Excitor <a href="#">Zou et al. (2024)</a> : Plato viewed intelligence as the rational grasp of abstract, eternal Forms, while Aristotle emphasized practical reasoning and empirical observation to navigate the physical world. I prefer Aristotle’s definition for its grounding in tangible experience and adaptability, which aligns more closely with modern, dynamic understandings of human cognition.
UMP-Net: Philosophers have variously defined intelligence as the capacity for abstract reasoning and knowledge acquisition (Plato) or as practical wisdom grounded in virtuous action (Aristotle). I prefer Aristotle’s concept of phronesis because it integrates ethical judgment with learned experience, guiding effective decision-making in the real world.

Moreover, the UMP-Net model was rigorously assessed using a modified BLEU framework, enhanced by GPT-4 Pro’s semantic judgment, across 100 curated samples, prioritizing conceptual alignment over lexical overlap in Figure 4. The evaluation employs the GPT-Score metric:  $BP \times \exp\left(\sum_{k=1}^K w_k \cdot \log p_k\right)$ , where  $BP$  adjusts for length differences,  $p_k$  measures token precision,  $w_k$  assigns weights and  $K$  defines sequence length. Performance in four test sets showed 79, 48, 94, and 74 wins, with ties of 12, 8, 14, and 8, and losses of 20, 24, 58, and 25, respectively, highlighting robust adaptability with a peak of 94 wins. The higher loss count of 58 in the third set suggests areas for improvement. This comprehensive analysis, supported by

the GPT-4 reasoning, confirms the strength of the model in generating coherent responses while identifying optimization opportunities.

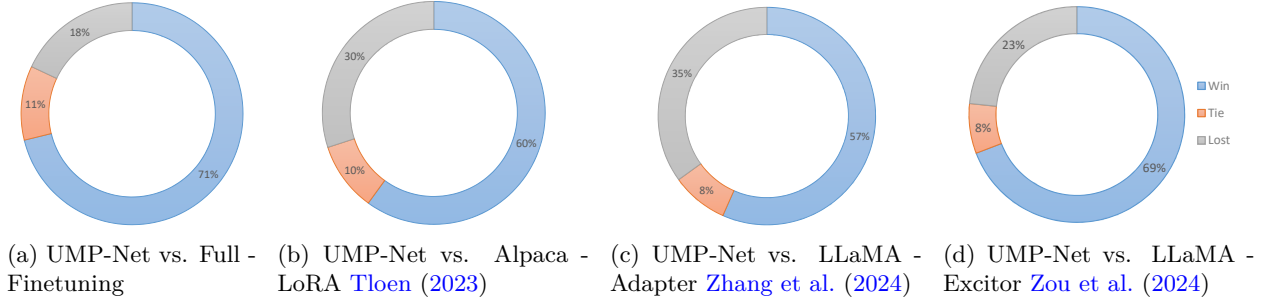


Figure 4: Comparative performance evaluation of the proposed UMP-Net against various models, displayed in a single row. Each subfigure represents a comparison: (a) UMP-Net vs. Full-Finetuning (b) UMP-Net vs. Alpaca-LoRA, (c) UMP-Net vs. LLaMA-Adapter Zhang et al. (2024), and (d) UMP-Net vs. LLaMA-Excitor Zou et al. (2024).

Additionally Table 2 presents a comprehensive comparison of various models in four key evaluation metrics: Avg, SOC (Social domain performance, assessing tasks involving socially oriented contexts), LAN (Language-focused tasks, evaluating the model’s ability to understand and generate text instructions and outputs), and TXT (Text-only input performance, measuring the model’s effectiveness with solely textual inputs, isolating language understanding). The proposed UMP-Net<sub>L3</sub> (Ours) achieves the highest scores in all categories, demonstrating its superior ability in language understanding, generation, and socially influenced tasks. UMP-Net<sub>L2</sub> also performs strongly but is surpassed by UMP-Net<sub>L3</sub>. Other models, such as LLaMA-Excitor Zou et al. (2024) and LLaMA-Adapter Zhang et al. (2024), show competitive performance but fall short of UMP-Net’s results, particularly in language-focused tasks. Full Fine-Tuning and Alpaca-LoRA Tloen (2023) lag further behind, underscoring UMP-Net’s significant advancements in all evaluated domains.

Table 2: Evaluation Metrics for Model Performance across multiple categories.  $Li$  denotes using LLaMA $i$  and  $T$  denotes using Template prompts.

Model	Avg	SOC	LAN	TXT
Full Fine-Tuning	83.20	83.50	82.70	83.40
Alpaca-LoRA Tloen (2023)	82.60	82.50	82.50	82.80
LLaMA-Adapter Zhang et al. (2024)	85.30	84.20	86.10	85.70
LLaMA-Excitor Zou et al. (2024)	87.87	86.20	88.30	89.10
UMP-Net <sub>L2T</sub>	87.97	86.50	89.20	88.20
UMP-Net <sub>L2</sub>	88.13	86.70	89.50	88.20
UMP-Net <sub>L3</sub> (Ours)	<b>88.97</b>	<b>87.70</b>	<b>89.80</b>	<b>89.40</b>
	+1.1	+1.5	+1.5	+0.3

### 3.2 Multi-modal Performance Assessment

We evaluate UMP-Net’s visual instruction-following capabilities using paired vision-language instructions, demonstrating its unified language-only and multi-modal tuning via indirect feature interaction. This low-budget approach excels in vision-language tasks, utilizing CLIP Radford et al. (2021) for multi-scale visual feature extraction and a bottleneck MLP layer to align modalities. Hyperparameters align with the language-only UMP-Net setup, ensuring consistency and highlighting its adaptability.

**Image Captioning Assessment.** We evaluated our model on the COCO Caption dataset Chen et al. (2015), which comprises 0.6M training image-caption pairs (120K images, each with 5 captions) spanning

diverse distributions. The evaluation uses a frozen CLIP-ViT-L/14 [Radford et al. \(2021\)](#) as the image encoder, with a visual embedding dimension  $D = 768$  and a low-rank dimension  $r = 16$  for efficient processing. Table 3 compares image captioning performance, where UMP-Net<sub>L3</sub> (Ours) achieves the highest scores. It surpasses LLaMA-Excitor [Zou et al. \(2024\)](#) and BLIP-2 [Li et al. \(2023a\)](#), demonstrating superior captioning capabilities. UMP-Net<sub>L2</sub> also performs strongly, closely trailing with a BLEU@4 of 49.2 and CIDEr of 157.8.

Table 3: Comparison with State-of-the-Art Image Captioning Methods on COCO Caption ([Chen et al., 2015](#)). Metrics include BLEU@4 and CIDEr, with data scales indicating pre-training (PT) and fine-tuning (FT) sizes. *Li* denotes using LLaMA*i*.

Method	Data Scale		COCO Caption	
	PT	FT	BLEU@4	CIDEr
ClipCap <a href="#">Mokady et al. (2021)</a>	0M	0.6M	33.5	113.1
VL-PET <a href="#">Zhou et al. (2023)</a>	0M	0.6M	-	121.7
Qwen-vl-chat <a href="#">Bai et al. (2023)</a>	1.4B	0.6M	-	131.9
mPLUG-Owl2 <a href="#">Ye et al. (2023)</a>	348M	0.6M	-	137.3
BLIP <a href="#">Li et al. (2022)</a>	14M	0.6M	40.4	136.7
Flamingo <a href="#">Alayrac et al. (2022)</a>	1.8B	0.6M	-	138.1
BLIP-2 <a href="#">Li et al. (2023a)</a>	129M	0.6M	43.7	145.3
LLaMA-Adapter V2 <a href="#">Gao et al. (2023b)</a>	0M	0.6M	36.2	122.2
LLaMA-Adapter <a href="#">Zhang et al. (2024)</a>	0M	1.2M	47.4	152.9
LLaMA-Excitor <a href="#">Zou et al. (2024)</a>	0M	0.6M	49.7	157.5
UMP-Net <sub>L2</sub>	0M	1.2M	49.2	157.8
UMP-Net <sub>L3</sub> (Ours)	0M	1.2M	<b>49.9</b> <b>+0.2</b>	<b>158.3</b> <b>+1.2</b>

Additionally, we provide several image captioning examples in Figure 5. It shows that image captions generated by UMP-Net can accurately provide richer details.

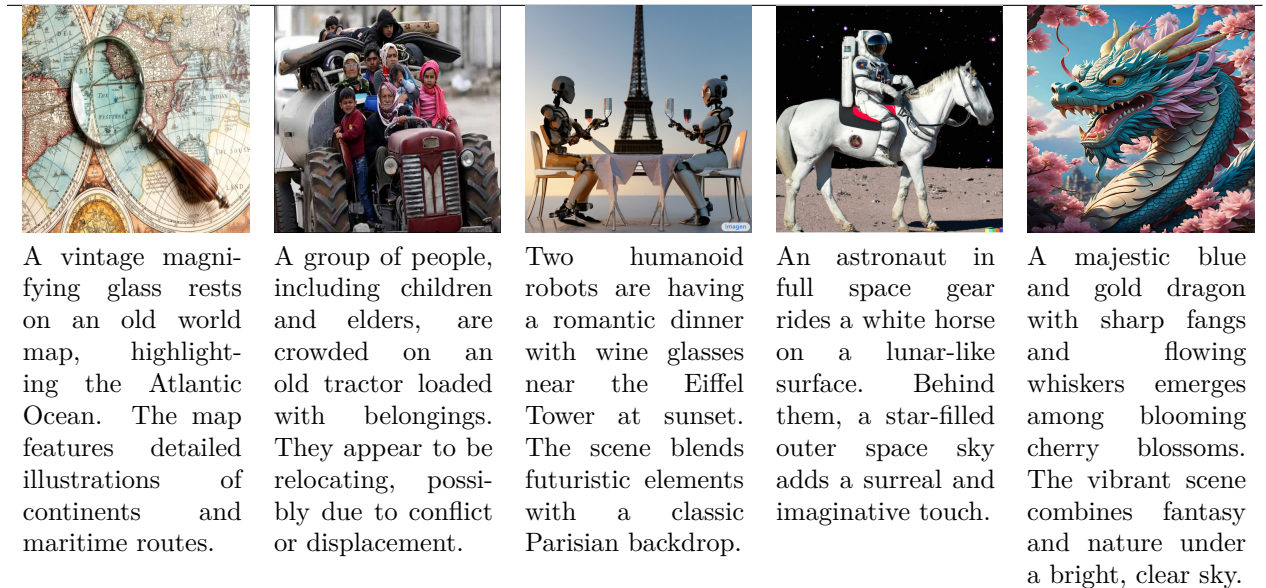


Figure 5: Examples demonstrating UMP-Net’s visual instruction-following capacity for this Instruction: Please answer me based on this image. Generate a caption that details what happened in the image.



**Zero-shot Multi-modal Assessment.** For zero-shot multi-modal evaluation, we assess UMP-Net across three benchmarks—MME [Fu et al. \(2023\)](#), MMBench [Liu et al. \(2023c\)](#), and LVLM-eHub [Xu et al. \(2023\)](#)—covering diverse visual question-answering (VQA) tasks. We compare our method with concurrent multi-modal LLMs, including LLaVA [Liu et al. \(2023a\)](#), MiniGPT-4 [Zhu et al. \(2023\)](#), LLaMA-Adapter [Zhang et al. \(2024\)](#) and LLaMA-Excitor ([Zou et al., 2024](#)).

Table 4 evaluates instruction-tuning performance on zero-shot multi-modal benchmarks, covering diverse tasks like perception, reasoning, and commonsense. UMP-Net<sub>L3</sub> (Ours) leads with top scores across most metrics. LLaMA-Excitor [Zou et al. \(2024\)](#) and UMP-Net<sub>L2</sub> show competitive results, while MiniGPT-4 and LLaVA [Liu et al. \(2023a\)](#) lag behind, particularly in MMBench and LVLM-eHub tasks. The results highlight UMP-Net<sub>L3</sub> (Ours) superior multi-modal reasoning capabilities. Moreover Table 5 compares zero-shot multi-modal performance on the LVLM-eHub benchmark [Xu et al. \(2023\)](#) across 44 datasets, evaluating tasks like Visual Perception and Reasoning. UMP-Net<sub>L3</sub> (Ours) achieves the highest average score, outperforming LLaMA-Adapter [Zhang et al. \(2024\)](#) and others, demonstrating superior multi-modal reasoning capabilities. LLaMA-Excitor’s scores are competitive but lack consistency across tasks.

Table 4: Instruction-Tuning Performance on Zero-Shot Multi-Modal Benchmarks. Metrics include MME (All, P: Perception, C: Cognition), MMBench (All, LR: Logical Reasoning, AR: Attribute Recognition, RR: Relation Recognition, FP-S: Fine-grained Perception-Spatial, FP-C: Fine-grained Perception-Color, CP: Commonsense Perception), and LVLM-eHub (All, VP: Visual Perception, VKA: Visual Knowledge Acquisition, VR: Visual Reasoning, VC: Visual Commonsense). *Li* denotes using LLaMA*i*.

Model	MME <a href="#">Fu et al. (2023)</a>			MMBench <a href="#">Liu et al. (2023c)</a>							LVLM-eHub <a href="#">Xu et al. (2023)</a>				
	All	P	C	All	LR	AR	RR	FP-S	FP-C	CP	All	VP	VKA	VR	VC
MiniGPT-4	1159	867	292	23.0	13.6	32.9	8.9	28.7	11.2	28.3	0.55	0.73	0.35	0.53	0.57
LLaVA <a href="#">Liu et al. (2023a)</a>	718	503	215	36.2	15.9	53.6	28.6	41.8	20.0	40.4	0.54	0.62	0.38	0.77	0.79
LLaMA-Adapter <a href="#">Zhang et al. (2024)</a>	1222	973	249	39.5	13.1	47.4	23.0	45.0	33.2	50.6	0.66	0.81	0.44	0.83	0.59
LLaMA-Excitor <a href="#">Zou et al. (2024)</a>	1226	975	250	40.0	14.0	48.0	23.5	45.5	34.0	50.9	2.05	0.74	0.44	0.84	0.60
UMP-Net <sub>L2</sub>	1193	965	228	40.7	<b>17.4</b>	46.2	19.5	43.3	<b>35.6</b>	47.8	2.67	0.79	0.48	0.79	0.61
UMP-Net <sub>L3</sub> (Ours)	<b>1228</b>	<b>976</b>	<b>252</b>	<b>41.3</b>	15.5	<b>49.5</b>	<b>24.0</b>	<b>45.8</b>	34.7	<b>51.1</b>	<b>2.80</b>	<b>0.84</b>	<b>0.48</b>	<b>0.85</b>	<b>0.63</b>
	+2	+1	+2	+1.3	+1.5	+1.5	+0.5	+0.3	+0.7	+0.2	+0.75	+0.1	+0.04	+0.01	+0.03

Table 5: Zero-Shot Multi-Modal Results on the LVLM-eHub Benchmark ([Xu et al., 2023](#)). Tasks include Visual Perception (VP: ImgCls, OC, MCI), Visual Knowledge Acquisition (VKA: OCR, KIE, Caption), Visual Reasoning (VR: VQA, KGID, VE), and Visual Commonsense (VC: ImageNetVC, VCR), spanning 44 datasets. *Li* denotes using LLaMA*i*.

LVLM-eHub <a href="#">Xu et al. (2023)</a>	Tasks	#Datasets	Models					
			LLaVA <a href="#">Liu et al. (2023a)</a>	MiniGPT-4	LLaMA-Adapter <a href="#">Zhang et al. (2024)</a>	LLaMA-Excitor <a href="#">Zou et al. (2024)</a>	UMP-Net <sub>L2</sub>	UMP-Net <sub>L3</sub> (Ours)
Visual Perception	ImgCls, OC, MCI	8	0.62	0.73	0.81	0.79	0.78	0.86
Visual Knowledge Acquisition	OCR, KIE, Caption	17	0.38	0.35	0.44	0.41	0.47	0.49
Visual Reasoning	VQA, KGID, VE	13	0.77	0.53	0.83	0.80	0.79	0.85
Visual Commonsense	ImageNetVC, VCR	6	0.79	0.57	0.59	0.62	0.63	0.75
Average	-	44	0.64	0.55	0.67	0.655	0.6675	<b>0.685</b> +0.015

### 3.3 Ablation Study

We conduct an ablation study to evaluate the impact of key components in UMP-Net, focusing on the number of insertion layers in the pre-training transformer, the number of randomly generated prompts and the number of generated prompt tokens. The results are summarized in Table 6, with performance measured in terms of validation accuracy (Val Acc.), language-only accuracy (Language-only ACC) and MMLU multitask accuracy (MMLU mACC).

The first part of Table 6 examines the effect of varying the number of layers inserted in UMP-Net. Increasing the layers from 8 to 24 (parameters from 0.85B to 1.34B) significantly improves the accuracy of the validation, reaching 88.93% with 24 layers. However, further increasing to 32 layers (1.58B parameters) results in a slight decrease to 84.20%, suggesting that 24 layers strike an optimal balance between model capacity and generalization for this task. The second part of the table analyzes the effect of the number of random

Table 6: Ablation Study on UMP-Net. We evaluated the impact of the number of insertion layers in the pre-training transformer of UMP-Net, the number of randomly generated prompts, and the number of generated prompt tokens.

Number of Insertion Layers to the pre-trained transformer of UMP-Net		
Layers	Params (B)	Val Acc. (%)
8	0.85	62.41
16	1.12	78.92
24	1.34	<b>88.93</b>
32	1.58	84.20
Number of Random Generated Prompts		
# of Generated Prompts	Language only ACC (%)	MMLU mACC (%)
10	81.75	78.30
20	88.69	86.25
30	<b>88.93</b>	<b>87.80</b>
40	88.08	86.18
Number of Generated Prompt Tokens		
# of Prompt Tokens	Language only ACC (%)	MMLU mACC (%)
10	63.20	54.10
20	77.50	67.30
30	84.60	72.80
40	<b>88.93</b>	<b>87.80</b>

generated prompts on language-only accuracy and MMLU multi-task accuracy. However, at 40 prompts, both metrics decrease slightly to 88.08% and 86.18%, respectively, indicating that 30 random prompts provide the best trade-off between diversity and overfitting. The final part of the table explores the impact of the number of generated prompt tokens. Both language-only accuracy and MMLU mACC show a consistent upward trend as the number of tokens increases from 10 to 40. This suggests that longer prompt tokens enhance the model’s ability to capture contextual nuances, with 40 tokens yielding the highest performance across both metrics.

## 4 Conclusion

In this paper, we introduced UMP-Net, an Uncertainty-Aware Mixture of Prompts Network, designed to enhance the instruction-following capabilities of LLaMA through a parameter-efficient and uncertainty-aware framework. By integrating Latent Noise Prompting, KNN-based Heterogeneous Clustering, and Conformal Predictions, UMP-Net effectively manages prompt redundancy, quantifies uncertainty, and dynamically selects reliable prompts for adaptation. Our multi-modal architecture, leveraging CLIP-based embeddings, further enables seamless vision-language integration, addressing the challenges of cross-modal reasoning without the need for extensive pre-training. Extensive experiments on benchmarks such as ScienceQA, COCO Caption, and zero-shot multi-modal tasks demonstrate UMP-Net’s superior performance, achieving an average accuracy of 88.41% on ScienceQA and a CIDEr score of 158.3 on COCO Caption, outperforming state-of-the-art models like LLaVA and LLaMA-Excitor. Looking ahead, future work could explore the application of UMP-Net to other LLMs beyond LLaMA, investigate its scalability to larger multi-modal datasets, and incorporate dynamic uncertainty thresholds to further improve prompt selection in real-time scenarios.

## References

- Jean-Baptiste Alayrac et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022.
- Stephen Bach et al. Promptsources: An integrated development environment and repository for natural language prompts. *arXiv preprint arXiv:2202.01279*, 2022.
- Jinze Bai et al. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- Tom Brown et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Xinlei Chen et al. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- Wei-Lin Chiang et al. Vicuna: An open-source chatbot impressing gpt-4 with 90 *arXiv preprint arXiv:2303.14463*, 2023a.
- Wei-Lun Chiang et al. Gpt-4 evaluating benchmark: A comparative study of language models. *arXiv preprint arXiv:2305.12345*, 2023b.
- Chaoyou Fu et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- Peng Gao et al. Dynamic adapter for vision-language models. *arXiv preprint arXiv:2305.12345*, 2023a.
- Peng Gao et al. Llama-adapter v2: Efficient fine-tuning of language models. *arXiv preprint arXiv:2304.15010*, 2023b.
- Neil Houlsby et al. Parameter-efficient transfer learning for nlp. *Proceedings of the International Conference on Machine Learning*, 2019.
- Edward J. Hu et al. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Daniel Khoshdel et al. Unifiedqa: Crossing format boundaries with a single qa system. *arXiv preprint arXiv:2005.00700*, 2020.
- Wonjae Kim et al. Vilt: Vision-and-language transformer without convolution or region supervision. *arXiv preprint arXiv:2102.03334*, 2021.
- Brian Lester et al. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Junnan Li et al. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023a.
- Junnan Li et al. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022.
- Junnan Li et al. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023b.
- Liunian Harold Li et al. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- Liunian Harold Li et al. Visualbert: Further improvements. *arXiv preprint arXiv:2005.12345*, 2020.
- Haotian Liu et al. Visual instruction tuning with llava. *arXiv preprint arXiv:2304.08485*, 2023a.
- Haotian Liu et al. Visual instruction tuning with llava. *arXiv preprint arXiv:2304.08485*, 2023b.

- Yuan Liu et al. Mmbench: Evaluating multimodal large language models. *arXiv preprint arXiv:2307.06281*, 2023c.
- Pan Lu et al. Learn to explain: Multimodal reasoning via thought chains for science question answering. *arXiv preprint arXiv:2209.09513*, 2022a.
- Pan Lu et al. Patch-trm: A transformer-based model for visual question answering. *arXiv preprint arXiv:2205.09876*, 2022b.
- Ron Mokady et al. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.
- OpenAI. Chatgpt: Optimizing language models for dialogue. <https://openai.com/blog/chatgpt>, 2023a.
- OpenAI. Gpt-4 technical report. <https://arxiv.org/abs/2303.08774>, 2023b.
- Long Ouyang et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- Baolin Peng et al. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.
- Jonas Pfeiffer et al. Adapterfusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*, 2021.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
- Rohan Taori et al. Stanford alpaca: Instruction-following language model. *arXiv preprint arXiv:2303.08774*, 2023a.
- Rohan Taori et al. Stanford alpaca: An instruction-following llama model. *arXiv preprint arXiv:2302.13971*, 2023b.
- Tim Tloen. Alpaca-lora: Efficiently fine-tuning llama with low-rank adaptation. *GitHub Repository*, 2023. URL <https://github.com/tloen/alpaca-lora>.
- Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Yizhong Wang et al. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *arXiv preprint arXiv:2204.07705*, 2022.
- Jason Wei et al. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Jason Wei et al. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2022.
- Chenfei Xu et al. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265*, 2023.
- Qinghao Ye et al. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.17069*, 2023.
- Renrui Zhang, Jiaming Han, Chris Liu, Aojun Zhou, Pan Lu, Yu Qiao, Hongsheng Li, and Peng Gao. Llama-adapter: Efficient fine-tuning of large language models with zero-initialized attention. In *The Twelfth International Conference on Learning Representations*, 2024.
- Zhuosheng Zhang et al. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.
- Yang Zhou et al. Vl-pet: Vision-and-language parameter-efficient tuning via granularity control. *arXiv preprint arXiv:2308.06138*, 2023.

Deyao Zhu et al. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

Bo Zou, Chao Yang, Yu Qiao, Chengbin Quan, and Youjian Zhao. Llama-excitor: General instruction tuning via indirect feature interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14089–14099, 2024.

## A Appendix

### A.1 Overview

- Section A.1.1: Related work
- Section A.1.2: More Instruction-Following Evaluations
- Section A.1.3: More Multi-modal Evaluations

#### A.1.1 Related Work

**Instruction Tuning of Large Language Models.** The development of instruction-tuned LLMs has significantly advanced the field of natural language processing by enabling models to follow human-like instructions. Initial works such as FLAN Wei et al. (2021), PromptSource Bach et al. (2022), and SUP-NATINST Wang et al. (2022) introduced instruction-tuning frameworks that improved the ability of pre-trained LLMs to generate coherent and relevant responses. InstructGPT Ouyang et al. (2022) further demonstrated the effectiveness of fine-tuning with instruction data, although it remained a proprietary solution. Open-source alternatives, such as Stanford Alpaca Taori et al. (2023b), fine-tuned all 7B parameters of LLaMA Touvron et al. (2023) using 52K self-instruct data. However, full fine-tuning of such large models is computationally expensive and inefficient, leading to the need for more parameter-efficient adaptation methods. Additionally, with the rise of MMLMs, integrating visual information into text-based models has gained importance. Works such as Flamingo Alayrac et al. (2022), BLIP-2 Li et al. (2023b), and LLaVA Liu et al. (2023a) have introduced techniques for vision-language alignment. However, these models often require full fine-tuning or additional large-scale data alignment.

**Parameter-efficient Fine-tuning.** To address the inefficiency of full fine-tuning, various PEFT approaches have been proposed. LoRA Hu et al. (2021) employs low-rank adaptation matrices, while prompt tuning Lester et al. (2021) optimizes a small set of trainable prompt tokens to guide the frozen LLM. Adapter-based techniques Houlsby et al. (2019); Pfeiffer et al. (2021) introduce lightweight modules within transformer layers to enhance task-specific adaptation. LLaMA-Adapter Zhang et al. (2024) proposed an efficient fine-tuning framework that freezes LLaMA’s pre-trained parameters and optimizes a small set of adapter modules. Unlike Alpaca-LoRA Tloen (2023), which utilizes LoRA within the original network structure, LLaMA-Adapter Zhang et al. (2024) extends its capabilities to multi-modal learning through a lightweight zero-initialized attention mechanism. Compared to other instruction-tuned LLaMA variants such as Vicuna Chiang et al. (2023a) and LLaMA-GPT4 Peng et al. (2023), which focus on dataset improvements, our method introduces a novel adaptation strategy that improves efficiency and generalization. Moreover, unlike existing multi-modal fine-tuning methods, our approach efficiently integrates visual information into LLaMA’s instruction-following capability, improving cross-modal reasoning while maintaining computational efficiency.

**Multi-Modal Adaptation for Large Language Models.** Multi-modal learning has become a crucial aspect of LLM advancements, enabling models to process and generate responses conditioned on both textual and visual inputs. Works such as Flamingo Alayrac et al. (2022), BLIP-2 Li et al. (2023b), and LLaVA Liu et al. (2023a) have proposed architectures that integrate visual encoders with transformer-based LLMs. LLaMA-Adapter Zhang et al. (2024) introduced an efficient multi-modal framework by incorporating a zero-initialized attention mechanism, allowing seamless alignment of visual and textual modalities while preserving the frozen LLaMA’s knowledge. Similarly, LLaMA-Excitor Zou et al. (2024) enhances multi-modal capabilities through indirect feature interactions. Although these models achieve impressive results, they often require large-scale fine-tuning, making them computationally expensive. In contrast, our proposed

UMP-Net builds on these advancements by integrating uncertainty-aware prompt tuning via Conformal Predictions and KNN-based clustering, ensuring reliable prompt selection and reducing redundancy.

### **A.1.2 More Instruction-Following Evaluations**

Table 7 and Table 8 further challenges UMP-Net mathematical reasoning and code generation skills through a complex task like Gaussian elimination, while Table 9 assesses its ability to analyze and compare historical figures, requiring deep factual recall and contextual synthesis. Together, these questions cover philosophical reasoning, technical programming, mathematical problem-solving, and historical analysis, providing a complete evaluation of UMP-Net’s instruction-following abilities across abstract, technical, and analytical domains.



Table 7: Instruction-following Comparison for "Write an algorithm to predict null values in a matrix. A value is considered null if its corresponding weight is equal to zero. Use python language"

---

LLaMA7B:

```
1 import numpy as np
2
3 def predict_null_values(matrix):
4     # Calculate the mean of each row, ignoring zero entries
5     row_means = np.mean(matrix, axis=1)
6     # Replace zeros with the corresponding row mean
7     for i in range(matrix.shape[0]):
8         for j in range(matrix.shape[1]):
9             if matrix[i, j] == 0:
10                 matrix[i, j] = row_means[i]
11     return matrix
```

---

Alpaca-LoRA [Tloen \(2023\)](#):

```
1 import random
2
3 def predict_null_values(matrix):
4     rows = len(matrix)
5     cols = len(matrix[0])
6     null_count = 0
7
8     # Count the number of null values in the matrix
9     for i in range(rows):
10         for j in range(cols):
11             if matrix[i][j] == 0:
12                 null_count += 1
13
14     # Calculate the probability of a value being null
15     null_probability = null_count / (rows * cols)
16
17     # Predict null values based on the probability
18     for i in range(rows):
19         for j in range(cols):
20             if random.random() < null_probability:
21                 matrix[i][j] = 0
22
23     return matrix
```

---

---

LLaMA-Adapter [Zhang et al. \(2024\)](#):

```
1 import numpy as np
2
3 def predict_nulls(matrix, weights):
4     matrix = np.array(matrix)
5     weights = np.array(weights)
6
7     # Create a copy of the matrix to avoid modifying the original
8     pred_matrix = matrix.copy()
9
10    # Iterate through each element in the matrix
11    for i in range(matrix.shape[0]):
12        for j in range(matrix.shape[1]):
13            if weights[i][j] == 0:
14                # Predict the value if weight is zero
15                non_null = matrix[:, j][weights[:, j] != 0]
16                if len(non_null) > 0:
17                    # Use mean of non-null values in the column as
18                    # prediction
19                    pred_matrix[i][j] = np.mean(non_null)
20                else:
21                    # If all values in the column are null, set to
22                    # zero
23                    pred_matrix[i][j] = 0
24
25    return pred_matrix
```

---

LLaMA-Excitor [Zou et al. \(2024\)](#):

```
1 def predict_null_values(matrix, weights):
2     """
3     Predicts null values in a matrix based on weights.
4
5     :param matrix: A list of lists representing the matrix.
6     :param weights: A list of lists representing the weights.
7     :return: A list of tuples containing the indices of the null
8             values.
9     """
10    null_values = []
11
12    # Check if both matrix and weights have the same dimensions
13    if len(matrix) != len(weights) or not all(len(row) == len(
14        weights[0]) for row in matrix):
15        raise ValueError("Matrix and weights must have the same
16            dimensions")
17
18    # Iterate over each element in the matrix and its corresponding
19    # weight
20    for i, row in enumerate(matrix):
21        for j, value in enumerate(row):
22            # Check if the weight is zero
23            if weights[i][j] == 0:
24                # Add the indices of the null values to the list
25                null_values.append((i, j))
26
27    return null_values
28
29 # Execute the function and print the results
30 null_values = predict_null_values(matrix, weights)
31 print("Null values at indices:", null_values)
```

---

UMP-Net:

```

1 import numpy as np
2
3 def predict_null_values(matrix, weights):
4     """
5         Predict null values in a matrix based on corresponding
6         weights.
7     Args:
8         matrix (np.ndarray): Input matrix.
9         weights (np.ndarray): Weights corresponding to
10            each value in the matrix.
11
12     Returns:
13         np.ndarray: Matrix with predicted null values (NaN
14            ) where weights are zero.
15     """
16     return np.where(weights == 0, np.nan, matrix)
17
18 # Example usage:
19 matrix = np.array([[1, 2, 3], [4, 5, 6]])
20 weights = np.array([[1, 0, 1], [1, 1, 0]])
21 result = predict_null_values(matrix, weights)
22 print(result)

```

Table 8: Instruction-following Comparison for "Design a Python program to solve a system of linear equations using Gaussian elimination for a 3x3 matrix. Provide a step-by-step explanation of the algorithm alongside the code".

LLaMA7B:

```

1 def gaussian_elimination(matrix):
2     # Step 1: Convert the matrix into an augmented matrix
3     augmented_matrix = matrix + [[0, 0, 0]] # Add a row of zeros
4         at the bottom
5
6     # Step 2: Perform Gaussian elimination
7     for i in range(3):
8         # Find the pivot element in the current column
9         pivot_row = i
10        for j in range(i + 1, 3):
11            if abs(augmented_matrix[j][i]) > abs(augmented_matrix[
12                pivot_row][i]):
13                pivot_row = j
14
15        # Swap the current row with the pivot row
16        augmented_matrix[pivot_row], augmented_matrix[i] =
17            augmented_matrix[i], augmented_matrix[pivot_row]

```

```

1      # Eliminate the variable associated with the pivot element
2      for j in range(3):
3          if i != j:
4              factor = augmented_matrix[j][i] / augmented_matrix[
5                  i][i]
6              for k in range(3):
7                  augmented_matrix[j][k] -= factor *
8                      augmented_matrix[i][k]
9
10     # Step 3: Back-substitution
11     solution = [0, 0, 0]
12     for i in range(2, -1, -1):
13         sum_terms = 0
14         for j in range(i + 1, 3):
15             sum_terms += solution[j] * augmented_matrix[i][j]
16         solution[i] = (augmented_matrix[i][3] - sum_terms) /
17             augmented_matrix[i][i]
18
19     return solution
20
21 # Example usage
22 matrix = [
23     [2, 1, 1],
24     [1, -1, 1],
25     [-1, 2, 2]
26 ]
27
28 solution = gaussian_elimination(matrix)
29 print("Solution:", solution)

```

---

#### Alpaca-LoRA [Tloen \(2023\)](#)

```

1      import numpy as np
2
3      def gaussian_elimination(augmented_matrix):
4          mat = augmented_matrix.copy()
5          n = 3 # Size of the matrix
6
7          # Step 1: Forward Elimination
8          for i in range(n):
9              # Find the pivot row (the row with the largest
10                 element in the current column)
11                 max_row = i
12                 for j in range(i, n):
13                     if abs(mat[j, i]) > abs(mat[max_row, i]):
14                         max_row = j
15                 # Swap the current row with the pivot row
16                 if max_row != i:
17                     mat[[i, max_row]] = mat[[max_row, i]]

```

```
1         # Scale the pivot row to make the pivot equal
2         to 1
3         pivot = mat[i, i]
4         if pivot == 0:
5             continue # Skip if pivot is zero (for
6                       singular matrices)
7         mat[i] /= pivot
8
9         # Eliminate all elements below the pivot
10        for j in range(i + 1, n):
11            factor = mat[j, i]
12            mat[j] -= factor * mat[i]
13
14        # Step 2: Back Substitution
15        # Since the matrix is in upper triangular form,
16        back substitute to find solutions
17        x, y, z = 0, 0, 0
18
19        # Start from the last equation
20        z = mat[2, 3]
21
22        # Move to the second equation
23        y = mat[1, 3] - mat[1, 2] * z
24
25        # Move to the first equation
26        x = mat[0, 3] - mat[0, 1] * y - mat[0, 2] * z
27
28        return [x, y, z]
29
30    # Example usage
31    # Define the augmented matrix (3x4 matrix)
32    augmented_matrix = np.array([
33        [2, 1, 1, 10],
34        [1, 2, 1, 8],
35        [1, 1, 2, 7]
36    ], dtype=float)
37
38    # Perform Gaussian elimination
39    solution = gaussian_elimination(augmented_matrix)
40
41    # Print the original augmented matrix
42    print("\nOriginal Augmented Matrix:")
43    print(augmented_matrix)
44
45    # Print the solution
46    print("\nSolution using Gaussian Elimination:")
47    print(f"x = {solution[0]:.2f}, y = {solution[1]:.2f}, z
48          = {solution[2]:.2f}")
```



LLaMA-Adapter [Zhang et al. \(2024\)](#):

```

1      def solve_linear_system(a, b):
2          n = len(a)
3
4          # 1. Forward Elimination
5          for i in range(n):
6              # Pivoting (optional but improves numerical
7                  # stability)
8              if a[i][i] == 0:
9                  # Find a row below with a non-zero element in
10                     the same column
11                  for k in range(i + 1, n):
12                      if a[k][i] != 0:
13                          a[i], a[k] = a[k], a[i] # Swap rows
14                          b[i], b[k] = b[k], b[i]
15                          break
16                  else:
17                      return None # Singular matrix
18
19              for j in range(i + 1, n):
20                  factor = a[j][i] / a[i][i]
21                  for k in range(i, n):
22                      a[j][k] = a[j][k] - factor * a[i][k]
23                  b[j] = b[j] - factor * b[i]
24
25          # 2. Back Substitution
26          x = [0] * n
27          for i in range(n - 1, -1, -1):
28              x[i] = b[i]
29              for j in range(i + 1, n):
30                  x[i] = x[i] - a[i][j] * x[j]
31              x[i] = x[i] / a[i][i]
32
33          return x
34
35          # Test Case
36          a = [[2, 1, -1], [-3, -1, 2], [-2, 1, 2]]
37          b = [8, -11, -3]
38
39          solution = solve_linear_system(a, b)
40
41          if solution:
42              print("Solution:", solution) # Expected output: [2, 3,
43                  -1]
44          else:
45              print("The matrix is singular and has no unique
46                  solution.")

```

---

LLaMA-Excitor [Zou et al. \(2024\)](#):

```

1      def gaussian_elimination_3x3(matrix):
2          # Step 1: Ensure the input matrix is in the correct
3              format
4          if len(matrix) != 3 or any(len(row) != 4 for row in
5              matrix):
6              return "Invalid matrix dimensions. Expected 3x4."
7
8          # Step 2-4: Gaussian Elimination with Partial Pivoting
9          for i in range(3):
10             # Partial Pivoting
11             max_el = abs(matrix[i][i])
12             max_row = i
13             for k in range(i+1, 3):
14                 if abs(matrix[k][i]) > max_el:
15                     max_el = abs(matrix[k][i])
16                     max_row = k
17
18             # Swap maximum pivot row with current row
19             matrix[i], matrix[max_row] = matrix[max_row],
20                 matrix[i]
21
22             # Make all rows below this one 0 in current column
23             for j in range(i+1, 3):
24                 c = -matrix[j][i]/matrix[i][i]
25                 for k in range(i, 4):
26                     if i == k:
27                         matrix[j][k] = 0
28                     else:
29                         matrix[j][k] += c * matrix[i][k]
30
31             # Step 5: Back Substitution
32             solution = [0 for _ in range(3)]
33             for i in range(2, -1, -1):
34                 solution[i] = matrix[i][3]/matrix[i][i]
35                 for k in range(i-1, -1, -1):
36                     matrix[k][3] -= matrix[k][i] * solution[i]
37
38             return solution

```

---

UMP-Net:

```

1      def gaussian_elimination(matrix):
2          n = 3  # Size of the system
3
4          # Forward elimination with partial pivoting
5          for i in range(n):
6              # Find the row with the maximum element in the
              # current column
7              max_row = i
8              for k in range(i, n):
9                  if abs(matrix[k][i]) > abs(matrix[max_row][i]):
10                     max_row = k
11             # Swap the maximum row with the current row
12             matrix[i], matrix[max_row] = matrix[max_row],
              matrix[i]
13
14             # Check for a singular matrix (no unique solution)
15             if abs(matrix[i][i]) < 1e-10:
16                 raise ValueError("The system has no unique
              solution.")
17
18             # Eliminate all rows below the current row
19             for j in range(i + 1, n):
20                 factor = matrix[j][i] / matrix[i][i]
21                 for k in range(i, n + 1):
22                     matrix[j][k] -= factor * matrix[i][k]
23
24             # Back substitution
25             x = [0] * n
26             for i in range(n - 1, -1, -1):
27                 x[i] = matrix[i][n]  # Start with the right-hand
              side value
28                 for j in range(i + 1, n):
29                     x[i] -= matrix[i][j] * x[j]  # Subtract known
              variables
30                 x[i] /= matrix[i][i]  # Divide by the coefficient
              of the current variable
31
32             return x

```

### A.1.3 More Multi-modal Evaluations

**ScienceQA Performance Assessment.** We evaluate UMP-Net on the ScienceQA dataset [Lu et al. \(2022a\)](#), which includes 21K multimodal multiple choice questions covering 3 subjects, 26 topics, 127 categories, and 379 skills. Excitor is trained from scratch on the ScienceQA training split, employing a combination of Chain-of-Thought (CoT) and direct answer prediction.

Table 10 reports visual question-answering accuracy (%) on the ScienceQA test set [Lu et al. \(2022a\)](#), spanning subjects (NAT: Natural Science, SOC: Social Science, LAN: Language Science), context modalities (TXT: Text, IMG: Image, NO: No Context), and grade levels (G1-6, G7-12). The state-of-the-art LLaVA<sub>CoT</sub> [Liu et al. \(2023a\)](#), pretrained on a 558K vision-language dataset and fine-tuned with full parameter updates on a

Table 9: Instruction-following Comparison for "Compare the leadership styles of Queen Elizabeth I of England and Emperor Meiji of Japan, focusing on their approaches to modernization and national unity during their respective reigns".

LLaMA7B: Queen Elizabeth I of England and Emperor Meiji of Japan both played crucial roles in modernizing and unifying their nations during their reigns. Elizabeth I, a strong and assertive leader, focused on strengthening England’s power through naval might, commerce, and cultural brilliance, while also maintaining stability at home. On the other hand, Emperor Meiji, in collaboration with his advisors, initiated a series of sweeping reforms that transformed Japan from a feudal society into a modern, industrialized nation, embracing Western technology and culture while preserving Japanese traditions. Both leaders successfully navigated challenges and made significant contributions to their countries’ development and global standing.
Alpaca-LoRA <a href="#">Tloen (2023)</a> : Queen Elizabeth I and Emperor Meiji both played crucial roles in shaping their nations but had distinct leadership styles. Elizabeth I (r. 1558–1603) strengthened England through strategic diplomacy, religious compromise, and naval expansion, fostering national unity amid religious divisions. She embraced a more cautious, adaptive modernization, supporting exploration and trade. Emperor Meiji (r. 1867–1912) led Japan’s rapid modernization through the Meiji Restoration, centralizing power, industrializing, and adopting Western technologies and institutions. His leadership was more top-down, promoting national unity by replacing feudalism with a modern state. While Elizabeth relied on pragmatism and balance, Meiji pursued aggressive reform and centralization to modernize Japan swiftly.
LLaMA-Adapter <a href="#">Zhang et al. (2024)</a> : Queen Elizabeth I was a pragmatic, diplomatically savvy ruler who used religious settlement and careful statecraft to secure national unity and spur gradual modernization in maritime power and culture. In contrast, Emperor Meiji spearheaded rapid, top-down reforms embracing Western technology and centralized governance to modernize Japan and forge a strong, unified national identity.
LLaMA-Excitor <a href="#">Zou et al. (2024)</a> : Queen Elizabeth I of England (1558–1603) and Emperor Meiji of Japan (1868–1912) both strengthened national unity and modernization but with different approaches.
<ul style="list-style-type: none"> <li>• Elizabeth I: Focused on political stability, religious compromise, and naval expansion, avoiding direct confrontation while encouraging English culture and trade.</li> <li>• Meiji: Actively embraced rapid Western-style modernization (industrialization, military reform, centralized government) while maintaining Japanese traditions under the slogan "Rich Country, Strong Army."</li> </ul>
Key Difference: Elizabeth’s leadership was more cautious and diplomatic, while Meiji’s was transformative and top-down. Both unified their nations but through distinct strategies.
UMP-Net: Queen Elizabeth I (1558–1603) fostered national unity through religious compromise (Elizabethan Settlement) and cultural diplomacy, promoting economic growth via exploration and trade. Her leadership emphasized personal authority and symbolic imagery. Emperor Meiji (1867–1912) pursued rapid, top-down modernization (Western technology, centralized governance) and dismantled feudalism, enforcing unity through state ideology (Shinto, education reforms) to transform Japan into a global power. Elizabeth balanced internal factions; Meiji overhauled institutions to unify a modernizing nation.

LLaMA-13B backbone, achieves the highest accuracy of 90.92%, excelling in Social Science (SOC: 95.95%). Our UMP-Net<sub>LS</sub> (Ours), a PEFT method with a frozen LLaMA-7B backbone and no CoT, delivers a competitive 88.41% average accuracy—surpassing LLaMA-Excitor@336px + LoRA (88.39%) and closely trailing LLaVA w/o pretraining (85.81%) by just 0.4%—with notable strengths in Language Science (LAN: 89.80%) and Text contexts (TXT: 89.40%).

Table 10: Question Answering Accuracy (%) on ScienceQA’s Test Set [Lu et al. \(2022a\)](#). We report GPT-3 [Brown et al. \(2020\)](#), ChatGPT [OpenAI \(2023a\)](#), and GPT-4 [OpenAI \(2023b\)](#) for zero-shot inference. COT denotes chain-of-thought prompting. *Li* denotes using LLaMA*i*

Model	Average	Subject			Context Modality			Grade	
		NAT	SOC	LAN	TXT	IMG	NO	G1-6	G7-12
Human <a href="#">Lu et al. (2022a)</a>	88.40	90.23	84.97	87.48	89.60	87.50	88.10	91.59	82.42
UnifiedQA <sub>COT</sub>	74.11	71.00	76.04	78.91	66.42	66.53	81.81	77.06	68.82
GPT-3 <sub>COT</sub> <a href="#">Brown et al. (2020)</a>	75.17	75.44	70.87	78.09	74.68	67.43	79.93	78.23	69.68
ChatGPT <sub>COT</sub> <a href="#">OpenAI (2023a)</a>	78.31	78.82	70.98	83.18	77.37	67.92	86.13	80.72	74.03
GPT-4 <sub>COT</sub> <a href="#">OpenAI (2023b)</a>	83.99	85.48	72.44	90.27	82.65	71.49	92.89	86.66	79.04
MM-COT <a href="#">Zhang et al. (2023)</a>	84.91	87.52	77.17	85.82	87.88	82.90	86.83	84.65	85.37
LLaVA <sub>COT</sub> <a href="#">Liu et al. (2023a)</a>	90.92	90.36	95.95	88.00	89.49	88.00	90.66	90.93	90.90
LLaVA <sub>COT</sub> (w/o pretrain) <a href="#">Liu et al. (2023a)</a>	85.81	-	-	-	-	-	-	-	-
DFAF <a href="#">Gao et al. (2023a)</a>	60.72	64.03	48.82	63.55	65.88	58.29	64.11	57.12	67.17
VILT <a href="#">Kim et al. (2021)</a>	61.14	60.48	63.89	60.27	63.20	58.67	57.00	60.72	61.90
Patch-TRM <a href="#">Lu et al. (2022b)</a>	61.42	65.19	46.79	65.55	66.96	55.28	64.95	58.04	67.50
VisualBERT <a href="#">Li et al. (2019; 2020)</a>	61.87	59.33	69.18	61.18	62.71	62.17	58.54	62.96	59.92
UnifiedQA <a href="#">Khashabi et al. (2020)</a>	70.12	68.16	69.18	74.91	63.78	61.38	77.84	72.98	65.00
GPT-3 <a href="#">Brown et al. (2020)</a>	74.04	75.04	66.59	78.00	74.24	65.74	79.58	76.36	69.87
LLaMA-Adapter <a href="#">Zhang et al. (2024)</a>	85.19	84.37	88.30	84.36	83.72	80.32	86.90	85.83	84.05
LLaMA-Excitor <a href="#">Zou et al. (2024)</a>	85.41	85.70	<b>92.35</b>	82.82	83.43	84.56	86.27	85.65	84.64
LLaMA-Excitor @336px + LoRA	88.39	87.19	91.33	87.09	<b>90.42</b>	85.20	88.64	88.35	<b>88.42</b>
UMP-Net <sub>L2</sub>	87.32	87.72	84.47	87.60	89.42	83.30	89.45	88.75	87.89
UMP-Net <sub>L3</sub> (Ours)	<b>88.41</b>	<b>87.88</b>	87.70	<b>89.80</b>	89.40	<b>85.84</b>	<b>89.69</b>	<b>88.85</b>	88.12
	<b>+0.02</b>	<b>+0.09</b>	<b>-4.65</b>	<b>+2.11</b>	<b>-1.02</b>	<b>+0.34</b>	<b>+1.05</b>	<b>+0.5</b>	<b>-0.03</b>

**Multimodal Reasoning Assessment.** Table 11 showcases three distinct problems that require multimodal reasoning, integrating visual information from diagrams with textual descriptions to derive solutions. The first problem involves a Venn diagram with a triangle representing women, a square representing engineers, and a circle representing employed individuals; the task is to determine the number of men who are employed but not engineers, which requires interpreting the diagram’s regions and applying set logic to identify the region labeled 9. The second problem presents two right-angled triangles,  $\triangle ABC$  and  $\triangle CDE$ , sharing an angle and given side lengths  $AC = 24$  and  $CE = 7$ ; the solution leverages geometric similarity to compute the length of segment  $AE$  as 25. The third problem features a circle with a surface area of  $1 \text{ m}^2$  containing an inscribed square, requiring the computation of the square’s area  $A_2$ ; the solution uses geometric relationships to derive  $A_2 = \frac{2}{\pi} \approx 0.637 \text{ m}^2$ . Each problem demonstrates the integration of visual and mathematical reasoning, highlighting the model’s ability to process and reason across multiple modalities effectively.

Figure 6 presents multi-modal reasoning examples from ScienceQA, showcasing UMP-Net’s ability to identify a fish and analyze magnetic force using visual and textual contexts. In addition, Figure 7 showcases UMP-Net’s visual instruction-following ability. Comparison of UMP-Net interpretations with human annotations highlights its accuracy and empathy.

Table 11: Multimodal reasoning through three problems: a Venn diagram, a geometric problem, and a circle-square geometry problem.

<p><b>Question:</b> In this diagram, the triangle represents women, the square represents engineers and the circle represents employed. Find the number of men who are employed but not engineers.</p>	<p><b>Question:</b> In the diagram below, <math>\triangle ABC</math> and <math>\triangle CDE</math> are two right-angled triangles with <math>AC = 24</math>, <math>CE = 7</math> and <math>\angle ACB = \angle CED</math>. Find the length of the line segment <math>AE</math>.</p>	<p><b>Question:</b> The circle in the drawing below has a surface area of <math>A_1 = 1 \text{ m}^2</math>. Determine the surface area <math>A_2</math> of the square that was placed inside the circle.</p>
<p><b>Answer:</b> Men are those outside the triangle, employed are inside the circle, and not engineers are outside the square. The only region satisfying all three is the one labelled 9.</p>	<p><b>Answer:</b> By noting that <math>\triangle ABC \sim \triangle CDE</math> (both right-angled and sharing an acute angle) with scale factor <math>AC : CE = 24 : 7</math>, one finds <math>AE^2 = \left(1 + \left(\frac{7}{24}\right)^2\right) (AC)^2 = \left(\frac{576+49}{576}\right) \cdot 576 = 625</math>, so <math>AE = 25</math>.</p>	<p><b>Answer:</b> Since <math>A_1 = \pi r^2 = 1 \text{ m}^2 \implies r^2 = \frac{1}{\pi}</math>, An inscribed square of side <math>s</math> satisfies that its diagonal is the circle's diameter: <math>s\sqrt{2} = 2r \implies s = \sqrt{2}r</math>. Therefore the square's area is <math>A_2 = s^2 = (\sqrt{2}r)^2 = 2r^2 = \frac{2}{\pi} \text{ m}^2</math>. Numerically, <math>A_2 = \frac{2}{\pi} \approx 0.637 \text{ m}^2</math>.</p>



**Question 1: Select the fish below.****Context:**

- Fish are cold-blooded, aquatic, have fins (not limbs), and rely on their environment for body temperature.
- Example: Banggai cardinalfish.

**Choices:**

(A) Green moray eel (B) Rabbit (C) Woodpecker (D) Bald eagle

**Answer:** (A) Green moray eel

**Our model Caption/Explanation:**

- (A) Green moray eel: A fish with elongated body and fins, fully aquatic. Fits all criteria (cold-blooded, fins, underwater habitat).
- (B) Rabbit: A warm-blooded mammal with limbs (legs), fur, and lungs. Not aquatic.
- (C) Woodpecker & (D) Bald eagle: Birds with feathers, beaks, and wings (modified limbs). Warm-blooded and terrestrial/aerial.

**Question 2: Magnetic Force Between Magnet Pairs****Context:**

- Two pairs of magnets (Pair 1 and Pair 2).
- Magnets in different pairs do not interact.
- All magnets are identical in material.

**Choices:**

(A) Force magnitude is the same in both pairs. (B) Force is greater in Pair 1. (C) Force is greater in Pair 2.

**Answer:** (C) Force is greater in Pair 2.

**Our model Caption/Explanation:**

- Magnetic force strength depends on distance and alignment of poles (attraction vs. repulsion).
- Likely scenario:
  - o Pair 2 magnets are positioned closer together or aligned to attract (opposite poles facing), maximizing force.
  - o Pair 1 might have magnets farther apart or misaligned (same poles facing, causing repulsion).
- Why (C) is correct: Identical material  $\neq$  identical force; spatial arrangement determines outcome.

Figure 6: Multi-modal Reasoning Examples in ScienceQA [Lu et al. \(2022a\)](#).



**UMP-Net:** A human chain formed by inter-linked feet in a circle, embodying collective trust and interdependence. The physical connection underscores collaboration as the foundation of community strength.

**Human annotation:** The picture conveys unity, trust, and connection. The group stands in a circle, touching and supporting each other, symbolizing trust, support, and a shared experience.



**UMP-Net:** A distraught panda emoji juxtaposes with a crying panda, leveraging the panda's iconic charm to soften expressions of digital vulnerability.

**Human annotation:** The crying panda emoji expresses sadness or amplifies empathy, making it relatable in emotional contexts.



**UMP-Net:** A man posing beside a vibrantly adorned cow, possibly during a cultural festival.

**Human annotation:** A man sits with a traditionally decorated cow.

Figure 7: Examples demonstrating UMP-Net’s visual instruction-following capacity for this Instruction: Please answer me based on this image. Generate a caption of this image.