**Paper Title:**

**Improving Zero-shot Visual Question Answering via Large Language Models with Reasoning Question Prompts**

**Paper Link:**

**https://paperswithcode.com/paper/improving-zero-shot-visual-question-answering**

## 1. Summary:

### 1.1 Motivation/purpose/aims/hypothesis:

The motivation for the paper titled "Improving Zero-shot Visual Question Answering via Large Language Models with Reasoning Question Prompts " comes from the challenges related to zero-shot Visual Question Answering (VQA) tasks using Large Language Models (LLMs). The aim of the authors is to enhance the performance of the existing methods in the context of caption generation from images. To create a bridge between the semantic gap of questions and captions by introducing reasoning question prompts in the hypothesis here. All of these will eventually improve the understanding of LLMs in zero-shot VQA scenarios.

### 1.2 Contribution:

The introduction of reasoning question prompts to facilitate zero-shot VQA tasks is the main contribution of this paper. Explicating intermediate reasoning steps in questions and effectively eliminating the semantic gap between questions and captions is the main role of the prompts. Through the experiments, it is explained that the proposed method enhances the performance of the existing zero-shot VQA methods and touches various LLM backbones.

### 1.3 Methodology:

Methodology includes creating captions from images, and subsequently training LLMs to answer questions based on these captions. The thing that really stands out here is the innovative approach of generating reasoning question prompts through an unsupervised method. The prompts are crafted carefully to meet certain criteria, like being self-contained and free of any supervision signal. With a two-step prompt, the LLMs are able to both generate and select answers. Also, enhancing their reasoning skills and ability to provide responses accurately.

**1.4 Conclusion:**

To wrap it up, this paper helps tackle the obstacles regarding zero-shot visual question answering by proposing the use of reasoning question prompts. The method establishes its prowess through experiments in enhancing current techniques and setting new standards in zero-shot VQA. The accuracy of LLMS is elevated with reasoning question prompts which is clear from the results

**2 <u>Limitations:</u>**

**2.1 First Limitation/Critique:**

The first limitation of the proposed approach is its dependence on pre-trained LLMs. The success of this method depends on the capabilities of these models which makes it susceptible to difficulties, where the pre-trained models lack necessary knowledge or demonstrate biases.

**2.2 Second Limitation/Critique:**

Secondly, the requirement for syntactic examples during the zero-shot evaluation is a restriction to be noted. The study does show improvement without explicit training data, and use of exemplars, even if artificially created, but questions the university of this approach. So dependency on these exemplars may hinder the method overall regarding stability in completely data-deficient environments.

**3 <u>Synthesis:</u>**

The paper comes up with a set of ideas that shows great potential for improving the interpretability and effectiveness of LLms. Through involving reasoning question prompts, this approach can be easily applied to other multimodal tasks. For instance, image captioning and natural language understanding. Moreover, more research can open the doors by adjusting pre-trained models in zero-shot VQA datasets with different fine-tuning techniques. Refining the reasoning question prompts approach and creating more transparent reasoning models for deeper understanding can be another potential way.