

The background features several large, stylized geometric shapes in the corners. These shapes are composed of nested triangles and polygons in shades of dark gray and light gray, creating a modern, abstract design.

# **VISUAL SPATIAL REASONING OF LARGE LANGUAGE MODELS**

Presentation: Jinlong Liu

# LIST OF CONTENT

**01.**

## **PROJECT INTRODUCTION**

A brief description of this project

**02.**

## **AIMS AND OBJECTIVES**

Research aims and objectives

**03.**

## **FEASIBILITY ESTIMATE**

Talk about feasibility according to the martial arts

**04.**

## **EXPERIMENT DESIGN**

Introduction the key aspects of experiment design

**05.**

## **PROJECT PLAN**

Show detailed project plan

**06.**

## **REFERENCES**

List of References



# **01.**

# **PROJECT INTRODUCTION**

In recent years, Large Language Models (LLMs) have transformed the field of Natural Language Processing (NLP) with their impressive achievements. Such as ChatGPT, Bert and Bard. As they updated in a raptly speed, there are a lot of test work also starting.




Source of all images:

<https://chat.openai.com>

<https://bard.google.com/>

<https://ivet360.com/google-bert-for-vets/>

- 
- The testing works are focused on the professional domain, like the ability of logic, answer test questions in various expert domain.
  - This project is based on the work of testing ability of spatial reasoning, which is a subject of testing ability of commonsense reasoning.
  - But it is different with existed work with adding visual input into prompt not just text input. Which used to be works of multi-model problem.




# **02.**

## **AIMS AND OBJECTIVES**



# AIMS:

1. Investigate performance of Visual Question Answering (VQA) tasks.
  2. Make comparison between different versions of LLMs.
  3. Find challenges of LLM's on Visual Spatial Reasoning.
- 



# MAIN OBJECTIVES

1. Design a diverse dataset with corresponding questions.
2. Conduct systematic experiments by using prepared dataset.
3. Analyse the experimental results to identify the challenges of LLMs
4. Provide the insight of strong points and weakness of current LLMs.





# **03.**

# **FEASIBILITY ESTIMATE**

# Former works

Name	Year	Source	Feature
VSR	2022	[1]	66 distinct types of spatial relations
GQA	2019	[2]	22M diverse reasoning questions
AGQA	2021	[3]	Add videos paired with question
CVR	2022	[4]	Measures of sample efficiency(Train dataset)
VALUE	2021	[5]	Test the linguistic phenomena in visual modality
Spatial commonsense benchmark	2022	[6]	Focus on positional relationship between people and objects
SpartQA	2021	[7]	Focus on generate spatial description by limited grammar rules

## **SAME POINTS:**

- The new version of GPT models now can consider as multi-models, and it can support get in visual content and text prompt at same time. It can work with this kind of tasks now.

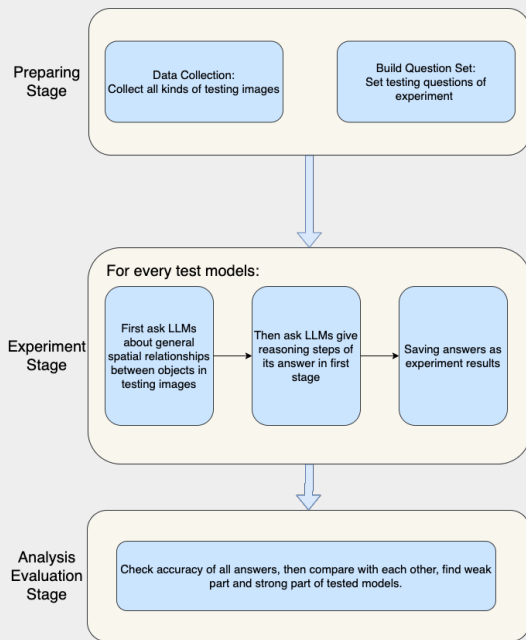
## **DIFFERENT POINTS:**

- This VAQ is mostly like testing dataset for VLMs, LLM's are usually tested by prompt.
- All the benchmarks give specific words to describe spatial reasoning or give choice or label to select. This project is to let LLMs generate by its self.



# **04.**

# **EXPERIMENT DESIGN**

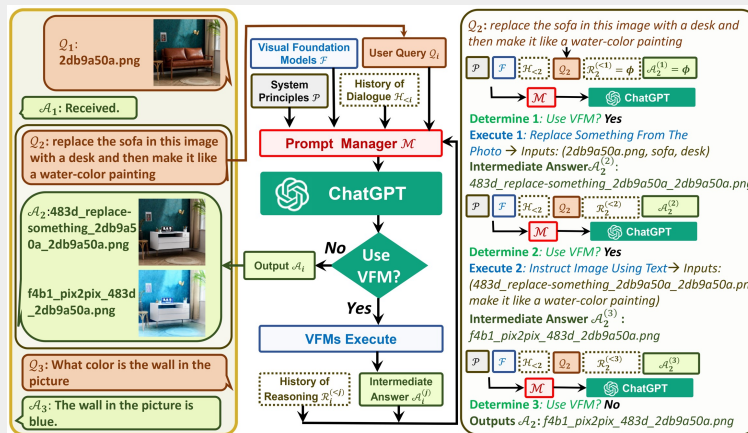


The detailed experiment shows as figure 1, it has three stage:


1. **Preparing Stage(Design Dataset):** Combine existed dataset and design questions with it.
2. **Experiment Stage:** test every selected model with image paired question one by one.
3. **Analysis Evaluation Stage:** Analysis generated answer with accuracy of answers and make judgement about reasoning steps

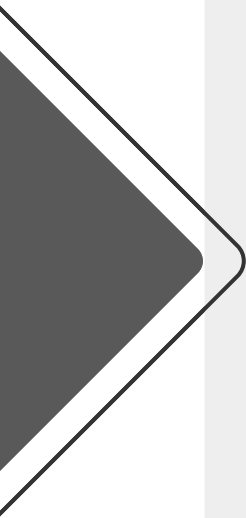
Selected Model: TaskMatrix(AKA Visual ChatGPT)  
connects ChatGPT and a series of Visual Foundation Models to  
enable sending and receiving images during chatting.[8]

Here is its architecture



The question prompt shows as below, it may contain 3-4 questions include objects relevant positions and explain reasoning steps.

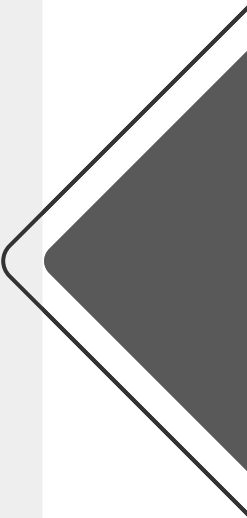
Input Image	Question
	<ol style="list-style-type: none"><li>1) Can you describe the spatial relationship between the tea table and the sofa?</li><li>2) Can you describe the relevant spatial relationships between the objects on the tea table?</li><li>3) Show the Reasoning Steps.</li></ol>



### Advantages:

1. Questions can be designed more specific follow former works.
2. The test dataset can be much more comprehensive. Which can contain much more situations.
3. Can find detailed differences between different models.

### Disadvantages:

1. The experiment may cost much more time.
  2. Analysis work is low efficient in comparing with former work.
  3. Cost more time in responding.
  4. Can't change to specific version of GPT-models.
- 

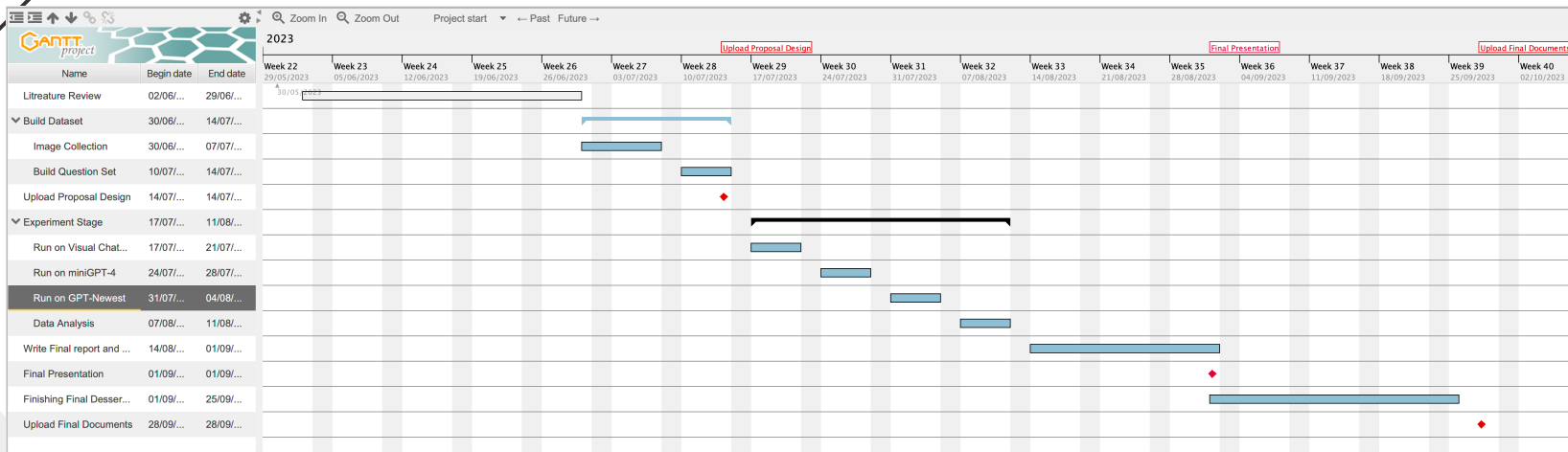


The background features four large, stylized geometric shapes in the corners, resembling arrowheads or chevrons. The top-left and bottom-right shapes are light gray, while the top-right and bottom-left shapes are dark gray. The central text is positioned in the white space between these shapes.

# **05.**

# **PROJECT PLAN**

# GANTT CHART





# **06.**

# **REFERENCES**

# PAPERS MENTIONED IN THIS SLIDES

- [1]. F. Liu, G. Emerson, and N. Collier, “Visual spatial reasoning,” arXiv preprint arXiv:2205.00363, 2022.
- [2]. D. A. Hudson and C. D. Manning, “Gqa: A new dataset for real-world visual reasoning and compositional question answering,” in Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 6700–6709.
- [3]. M. Grunde-McLaughlin, R. Krishna, and M. Agrawala, “Agqa: A benchmark for compositional spatio-temporal reasoning,” in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 11 287–11 297.
- [4]. X. Liu, D. Yin, Y. Feng, and D. Zhao, “Things not written in text: Exploring spatial commonsense from visual signals,” arXiv preprint arXiv:2203.08075, 2022.
- [5]. A. Zerroug, M. Vaishnav, J. Colin, S. Musslick, and T. Serre, “A benchmark for compositional visual reasoning,” arXiv preprint arXiv:2206.05379, 2022.
- [6]. L. Parcalabescu, M. Cafagna, L. Muradjan, A. Frank, I. Calixto, and A. Gatt, “Valse: A task-independent benchmark for vision and language models centered on linguistic phenomena,” arXiv preprint arXiv:2112.07566, 2021.
- [7]. R. Mirzaee, H. R. Faghihi, Q. Ning, and P. Kordjmeshidi, “Spartqa: A textual question answering benchmark for spatial reasoning,” arXiv preprint arXiv:2104.05832, 2021.
- [8]. C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, “Visual chatgpt: Talking, drawing and editing with visual foundation models,” 2023.



**THANK  
YOU**

