

MM-REACT : Prompting ChatGPT for Multimodal Reasoning and Action

Zhengyuan Yang*, Linjie Li*, Jianfeng Wang*, Kevin Lin*, Ehsan Azarnasab*, Faisal Ahmed*,
Zicheng Liu, Ce Liu, Michael Zeng, Lijuan Wang[♣]
Microsoft Azure AI

{zhengyang, lindsey.li, jianfw, keli, ehazar, fiahmed, zliu, ce.liu, nzeng, lijuanw}@microsoft.com

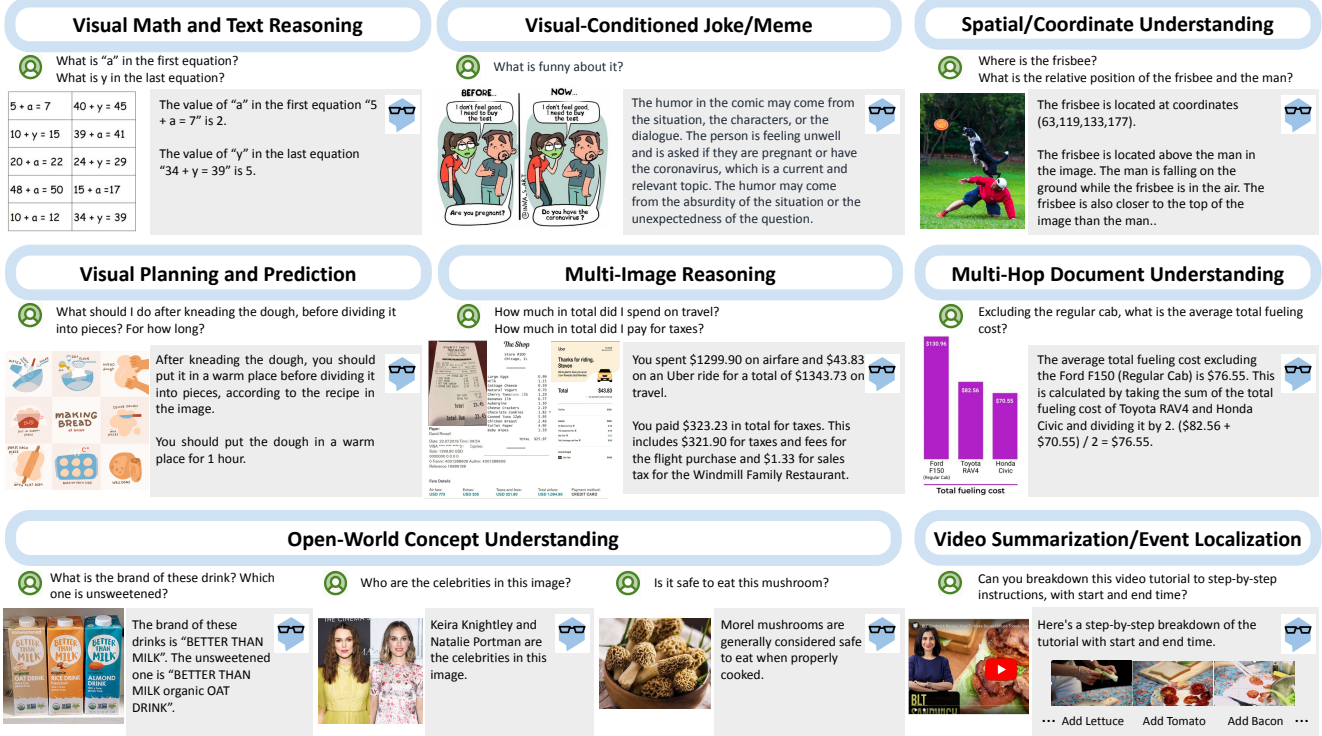


Figure 1. MM-REACT allocates specialized vision experts with ChatGPT to solve challenging visual understanding tasks through multi-modal reasoning and action. For example, the system could associate information from multiple uploaded receipts and calculate the total travel cost ("Multi-Image Reasoning"). We only highlight key information here and postpone full MM-REACT responses to Figures 4-14.

Abstract

We propose MM-REACT, a system paradigm that integrates ChatGPT with a pool of vision experts to achieve multimodal reasoning and action. In this paper, we define and explore a comprehensive list of advanced vision tasks that are intriguing to solve, but may exceed the capabilities of existing vision and vision-language models. To achieve such advanced visual intelligence, MM-REACT introduces a textual prompt design that can represent text descriptions, textualized spatial coordinates, and aligned file names for dense visual signals such as images and videos. MM-REACT's prompt design allows language models to accept,

associate, and process multimodal information, thereby facilitating the synergetic combination of ChatGPT and various vision experts. Zero-shot experiments demonstrate MM-REACT's effectiveness in addressing the specified capabilities of interests and its wide application in different scenarios that require advanced visual understanding. Furthermore, we discuss and compare MM-REACT's system paradigm with an alternative approach that extends language models for multimodal scenarios through joint fine-tuning. Code, demo, video, and visualization are available at <https://multimodal-react.github.io/>.

*Equal Contribution [♣]Project Lead

1. Introduction

Recent years have seen significant advancement for computer vision, thanks to improved network architecture [4, 9, 12, 28], large-scale model training [8, 29, 35], and other factors. However, different vision problems typically require different models, which often require manual selection and composition of individual models for each use case. For example, when determining if an image contains “people”, we may choose the image tagging model [7, 16, 19] and check if the predicted tag list contains “people”. If “people” exists, we may select the celebrity model [20] to further understand whether a celebrity appears and who he/she is.

One research direction is to combine the vision and language modules as one end-to-end model, such as Flamingo [2], PaLM-E [10], to provide a dialogue-based experience to the end user. That is, the user can use natural language to interact with the model around the image content. The vision module encodes vision signals into special text tokens or features that the language module can understand, enabling the system to utilize the language module for understanding user queries and providing responses. However, these joint finetuning approaches require a large amount of computing resources and annotated data to enable specific capabilities. In this work, we aim to combine existing individual vision models with the language model in a more flexible manner to tackle complicated visual understanding problems, *e.g.*, the ones illustrated in Figure 1.

Large language models (LLMs) [3, 6], such as ChatGPT, have shown impressive dialogue capability with text as both input and output. Recent NLP research [11, 24, 26, 34] (*e.g.*, REACT [34]) demonstrates the effectiveness of integrating external NLP tools, such as search engines and math calculators, with LLMs by proper instruction. Specifically, REACT [34] prompts an LLM to generate *reasoning* texts that break down complex problems into intermediate steps, and *action* texts that allocate NLP tools for solving these steps. One example is that the LLM can suggest a text query to a modern search engine to grab the latest internet information, and return the user with the information that is not in the pre-training corpus. Inspired by the efficacy of reasoning and acting with LLMs and NLP tools, we explore the integration of vision expert tools with LLMs.

To this end, we present MM-REACT, a system paradigm that composes numerous vision experts with ChatGPT for multimodal reasoning and action. To enable images and videos as inputs, we use their file path as the input to ChatGPT. The file path functions as a placeholder, allowing ChatGPT to treat it as a black box. Whenever a specific property such as celebrity names or box coordinates is required, ChatGPT is expected to seek help from a specific vision expert to identify the desired information. To inject the knowledge of vision experts’ usages into ChatGPT, we add instructions to ChatGPT prompts about each expert’s

capability, input argument type, and output type, along with a few in-context examples for each expert. Additionally, a special watchword is instructed such that we can use regex expression matching to invoke the expert accordingly.

We show MM-REACT’s representative visual understanding capabilities in Figure 1. For example, MM-REACT could associate information from multiple uploaded receipts and calculate the total travel cost (“Multi-Image Reasoning”), recognize and answer questions about the “morel mushrooms” (“Open-World Concept Understanding”), and condense a long video into representative thumbnails (“Video Summarization and Event Localization”). These visual intelligence features are similar to those found in recent models, such as multimodal GPT-4 [23] and PaLM-E [10], but achieved via prompting instead of additional multimodal training. The MM-REACT system may provide extra flexibility in module upgrades, and may be effective in certain visual understanding tasks by better utilizing existing specialized vision experts, such as celebrity recognition and dense captioning.

2. Related Work

LLMs Prompting Methods. Large language models (LLMs) [3, 6] demonstrate a strong chain-of-thought (CoT) capability [17, 31] that could break down complex problems into solvable intermediate steps. On the other hand, research [1, 15, 22] shows that LLMs, when equipped with a range of external NLP tools, can effectively serve as action planners to select and utilize tools for problem-solving, such as using search or mathematical tools to address knowledge or math problems.

Nevertheless, LLMs for reasoning [17, 31] and LLMs for action [1, 15, 22], when used independently, fail to solve complex tasks that require breaking down the problem via reasoning and solving sub-steps via planned actions. Recent studies [11, 24, 26, 34] have attempted to merge the action and reasoning phases to enhance LLMs’ capabilities in solving complicated tasks that require advanced planning and reasoning. One representative work, REACT [34], treats reasoning text generation as an executable action and achieves the synergetic combination of reasoning and action for NLP tasks. In this work, we explore how to extend such intriguing properties into multimodal scenarios by modeling thought and invoking vision tools as executable actions.

Vision+LLMs. Our MM-REACT is related to the previous studies that extend language models to understand visual inputs. The representative framework adds a vision module to project visual inputs into representations that the language model can understand. These representations can be either discrete text words [13, 30, 33, 36] or continuous features projected into the textual feature space [2, 10, 14, 18, 27]. Recent vision-language studies

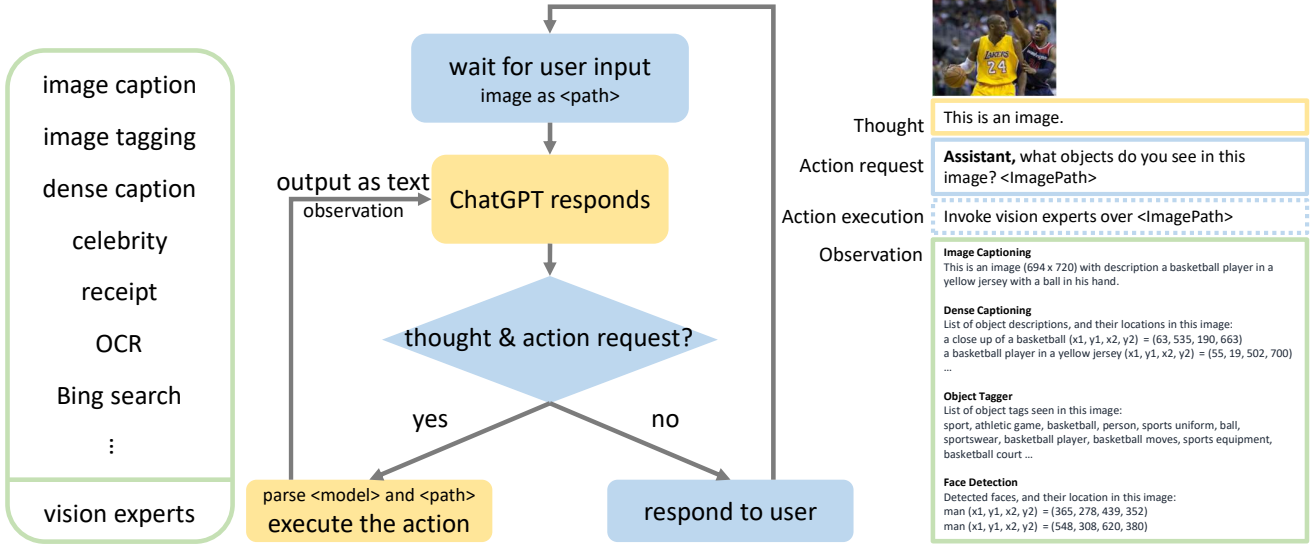


Figure 2. Flowchart of MM-REACT for enhanced visual understanding with ChatGPT. The user input can be in the form of text, images, or videos, with the latter two represented as file path strings. ChatGPT is instructed to say specific watchwords in *action request* if a vision expert is required to interpret the visual inputs. Regular expression matching is applied to parse the expert’s name and the file path, which are then used to call the vision expert (*action execution*). The expert’s output (*observation*) is serialized as text and combined with the history to further activate ChatGPT. If no extra experts are needed, MM-REACT would return the final response to the user. The right figure shows a single-round vision expert execution, which is the component that constructs the full execution flow illustrated in Figure 3.

explore the chain-of-thought capability [17, 31] in multi-modal settings. MM-CoT [37] finetunes on the reasoning chain annotated in the ScienceQA [21] dataset to achieve the CoT capability in the science question answering task. KOSMOS-1 [14] and PaLM-E [10] demonstrate the zero-shot multimodal CoT capabilities with large-scale training.

Multimodal Reasoning and Action. A key distinction between MM-REACT and prior vision+LLM studies discussed above is that MM-REACT leverages LLMs’ high-level planning abilities to allocate various vision experts, rather than solely using LLMs for text generation conditioned on visual inputs. MM-REACT is closely related to the recent concurrent work of Visual ChatGPT [32] and ViperGPT [25]. In comparison, Visual ChatGPT [32] primarily focuses on image generation and editing, while our MM-REACT mainly focuses on visual understanding. ViperGPT [25] instructs LLMs to generate Python code for a one-round query answering. In contrast, MM-REACT is a multi-round, dialogue-based system that may integrate the strong QA model as one of its vision experts.

3. MM-REACT Prompting

The goal of MM-REACT is to compose numerous vision experts to empower ChatGPT with visual understanding. A vision expert is a computer vision model that takes an image as input and interprets the content from different perspectives.

For instance, the image captioning expert generates a natural description, the OCR expert extracts the scene text in the image, the celebrity recognition model identifies the celebrity names, and the object detection model extracts the salient object with bounding box locations. At present, one may have to manually decide which vision experts to employ for specific use cases, and manually compose them. Instead, our target is to automate this process based on the requirements presented by the user query in natural language.

ChatGPT is an artificial intelligence chatbot with text as both input and output, without visual understanding. However, ChatGPT exhibits strong instruct learning capability, which motivates us to instruct ChatGPT to properly determine which vision expert should be invoked and which image should be processed.

Figure 2 shows the flowchart of our MM-REACT system. The terms *thought* and *action request* refer to the reasoning and action-oriented texts generated by ChatGPT to break down the problem or invoke vision experts. *Observation* refers to the vision expert’s responses after the *action execution* requested in the *action request* text. Next, we detail each step in the flowchart as follows.

3.1. User Input

As ChatGPT only accepts texts as input, the first challenge is how to accommodate non-text inputs, such as multiple images and videos. Since most vision experts accept

the file path or URL, we use the path string to indicate non-text inputs. The file path itself is meaningless and is essentially a placeholder. Although no visual recognition task can be performed directly with file paths, ChatGPT may seek help from different vision experts to understand the image content from different perspectives, *e.g.*, identifying the celebrity names of the detected person. By including the provided file path in its text output, ChatGPT can indicate which image should be processed by the vision expert.

3.2. ChatGPT Response

Given the user’s input, ChatGPT is expected to provide two kinds of responses. The first is to seek help from vision experts, while the second is to respond to the user directly. A key challenge is to set up a protocol such that we know when to invoke the vision expert. Inspired by REACT [34], we instruct ChatGPT to respond with certain watchwords, such as “Assistant, what objects are there in the image? <file path>”, if a specific vision expert is required. In our implementation, we use the keyword “Assistant,” to distinguish whether a vision expert is required.

To further improve the performance, we encourage ChatGPT to show the *thought* (reasoning) process to highlight why an external tool is required. It is also shown to be beneficial in NLP studies [34] to incorporate such reasoning.

3.3. Vision Experts

Given the action request from ChatGPT, we use the regular expression matching to parse the expert name and the file path, and invoke the action (vision expert execution).

The expert’s output can be in different forms but is standardized into the text format such that ChatGPT can understand it. For certain experts, such as the captioning model or the celebrity model, it is straightforward to represent the output as text. However, the standardization is less intuitive for others. For example, the detection model outputs a list of object names with bounding box locations. In this case, we concatenate all the boxes, each of which is represented as <object name, x1, y1, x2, y2>, where (x1,y1), (x2,y2) are the coordinates of the top-left and bottom-right corners, respectively. An additional text description is added to explain the meaning of the last four numerical values. In some cases, we find ChatGPT is capable of understanding these coordinates, *e.g.*, identifying which object is on the left.

The text-formed output from vision experts can be interpreted as the *observation* resulting from ChatGPT’s action of invoking the vision expert. Combining observations with the chat history, ChatGPT can further invoke additional experts or return the final answer to the user. We provide examples of full execution flows in Section 4.2 and Figure 3.

To inject the knowledge of various vision experts’ usages, we add both instructions and in-context examples in the prefix when prompting ChatGPT. Each expert is de-

scribed with the model name, a general description of its capability, the input data format, and the output information. After describing each expert, we add a few in-context dialogue examples to enhance the performance. With the injected knowledge, ChatGPT can effectively select one or multiple vision experts to understand the images or videos from different perspectives.

3.4. Extensibility

Our scheme is motivated by REACT, which invokes different tools in the NLP field. As only the text is involved, no specific design is required to incorporate other modalities. In this work, we extend MM-REACT to the vision domain. The key is to replace the non-text modality with a path string, enabling ChatGPT to ask specific vision experts to recognize the visual content. Therefore, we could further extend MM-REACT to other modalities, such as speech and audio. Meanwhile, we can also easily incorporate more tools by formatting their outputs in a text format. While ChatGPT serves as the primary LLM in our main implementation, performance could be further enhanced through the simple upgrade to a more powerful LLM, such as GPT-4 [23] discussed in Section 4.5.

4. Experiments

4.1. Experiment Setup

We implement MM-REACT based on the LangChain codebase [5] and reference ideas from ReAct [34]. We access ChatGPT via the Azure “gpt-3.5-turbo” API that has a token length limit of 4,096, and utilize vision experts publicly available via the Azure Cognitive Services APIs¹, including the ones for image captioning, image tagging, dense captioning, optical character recognition (OCR), and specialized recognition models for celebrities, receipts, *etc.* We further expand the toolset with customized tools for spatial understanding and image editing, and tools from other modalities such as Bing search and PAL math.

4.2. MM-REACT’s Full Execution Flow

Figure 3 provides an example to illustrate MM-REACT’s full execution flow. We highlight the exact order to call different models (*i.e.*, executions) with numbered blue circles. The executions, highlighted by bold underlined text, can be either a ChatGPT call (*e.g.*, “**ChatGPT:**”) or the execution of one or multiple selected vision experts (*e.g.*, “**Image Captioning**”). We add a commentary text *action execution* in dashed boxes to help understand the vision expert execution. The *action execution* is not an actual input or output in the MM-REACT flow. ChatGPT executions can be used to generate thought (reasoning) and action texts that

¹<https://azure.microsoft.com/en-us/products/cognitive-services/vision-services>

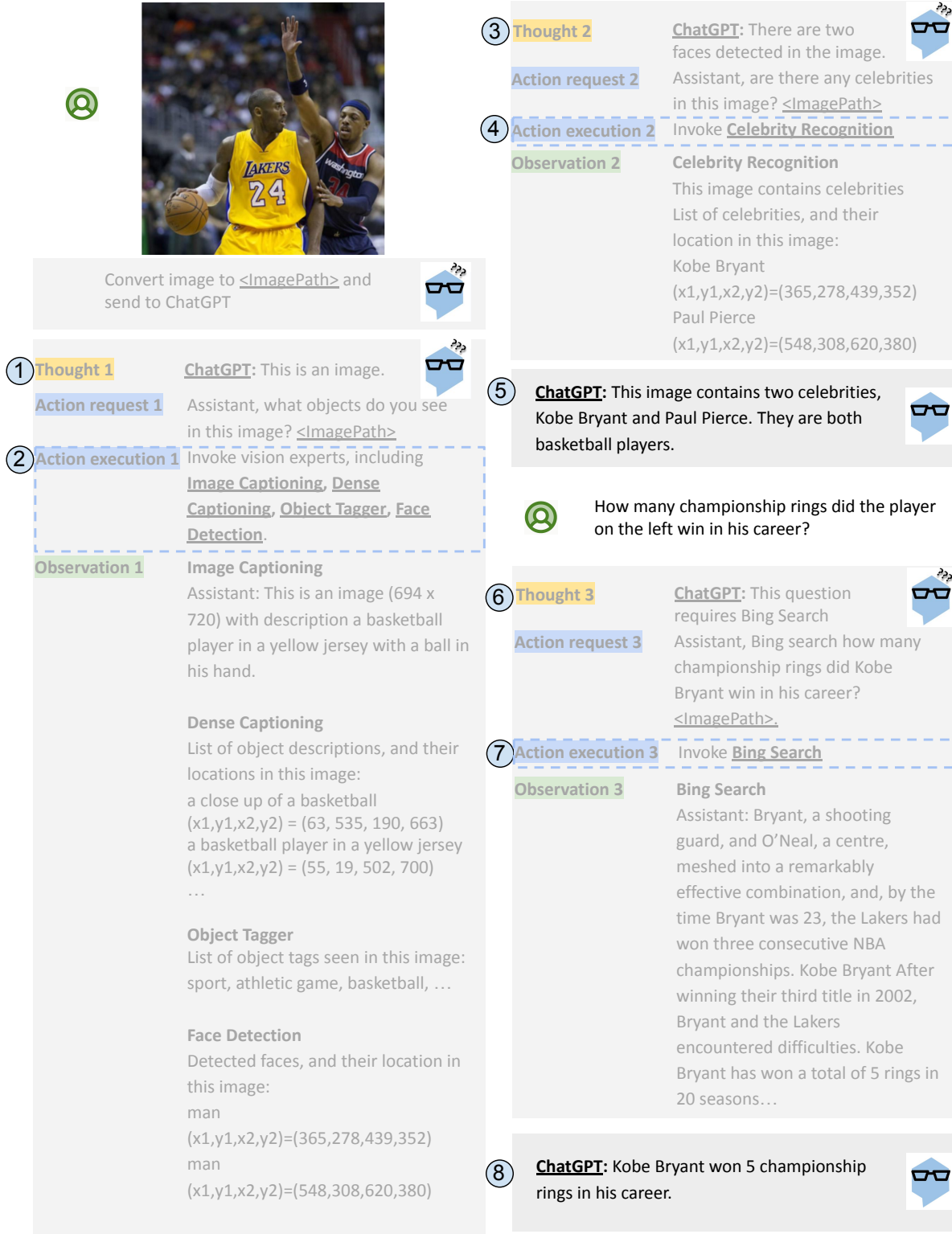




Figure 3. An example of MM-REACT’s full execution flow. The blue circles with numbered indices indicate the order in which different models are called (*i.e.*, the executions). The executions, highlighted by bold underlined text, can be either a ChatGPT call (*e.g.*, “**ChatGPT:**”) or running one or multiple selected vision experts (*e.g.*, “**Image Captioning**”). We add a commentary text *action execution* in dashed boxes to help understand the expert execution. Each ChatGPT execution takes the preceding text as input and generates the text leading up to the next execution (*e.g.*, “*This is an image. Assistant, what . . . image?* `<ImagePath>`” for Execution 1). Texts in gray represent MM-REACT’s thoughts or vision experts’ actions and outputs, which are invisible to users. This multimodal reasoning and action process occurs behind the scene to gather the necessary information for generating final responses to users, which are shown in black.

allocate vision experts  (invisible to users), or produce the final response to users . Each ChatGPT execution takes the preceding text as input and generates the text leading up to the next execution (e.g., “This is an image. Assistant, what objects do you see in this image? <ImagePath>” for Execution 1). ChatGPT “learns” the proper text to generate based on the instructions and in-context examples in the prompt prefix, as detailed in Section 3.3. Additional examples of the reasoning and execution procedures are in Figures 18-22.

4.3. MM-REACT Capabilities and Applications

Figures 4-14 show the representative capabilities and application scenarios that MM-REACT demonstrates. Specifically, we examine MM-REACT’s capabilities in visual math and text reasoning (Figure 4), understanding visual-conditioned jokes and memes (Figure 5), spatial and coordinate understanding, visual planning and prediction (Figure 6), multi-image reasoning (Figure 7), multi-hop document understanding on bar charts (Figure 8), floorplans (Figure 9), flowcharts (Figure 10), tables (Figure 11), open-world concept understanding (Figure 12), and video analysis and summarization (Figure 13, 14). We provide an example of the unfolded steps in Figure 18.

4.4. Capability Comparison with PaLM-E

MM-REACT is a training-free scheme which composes existing vision experts with ChatGPT, while PaLM-E [10] trains a vision-language model which combines an image encoder and a text decoder with dedicated datasets. Figures 15-17 shows our MM-REACT can achieve competitive results to PaLM-E. We further illustrate the complete multimodal reasoning and action procedures in Figures 21,22.

4.5. MM-REACT Extensibility

In Figure 23 and 24, we explore the enhancement of MM-REACT’s LLM from ChatGPT (“gpt-3.5-turbo”) to GPT-4 (language-only). We access the language-only GPT-4 via the ChatGPT website and reference the multimodal GPT-4 demo provided by OpenAI for comparison. These examples demonstrate the benefit of MM-REACT’s extensibility: MM-REACT equipped with GPT-4 correctly answers the physics question (Figure 23), while the version with ChatGPT (GPT-3.5) fails. Furthermore, MM-REACT is designed with the flexibility to incorporate new tools without training. Figure 25 provides a case study of plugging an image editing tool from X-decoder [38] for multi-round, dialogue-based image editing.

4.6. Limitations

We identify the following limitations. 1). Considering the recognition capability in the wild, we find it hard to sys-

tematically evaluate the performance with concrete accuracy numbers, due to a lack of annotated benchmarks. Thus, it is worth investing efforts towards how to effectively evaluate such system’s performance. 2). The vision capability is limited by the integrated vision experts. On one hand, the integrated experts may make mistakes; on the other hand, the system may fail if the necessary experts are missing. 3). We inject the vision experts’ knowledge in the prefix, and thus the number of experts is limited by the context window (4096 tokens) of ChatGPT. 4). Visual signals are converted to text words for ChatGPT understanding, which might not be the optimal solution for certain vision tasks. 5). MM-REACT requires manual prompt engineering. We expect research work to automate this process, making the system even easier to develop.

5. Conclusion

We have presented MM-REACT, a system paradigm that synergizes multimodal reasoning and action to tackle complicated visual understanding problems. MM-REACT presents a simple and flexible way to empower LLMs with a pool of vision experts. Extensive zero-shot experiments demonstrate MM-REACT’s capabilities in solving a wide range of challenging understanding tasks, such as multi-image reasoning, multi-hop document understanding, open-world concept understanding, video summarization, and more.

Acknowledgment

We would like to express our gratitude to Jianfeng Gao for his valuable suggestions, as well as to Jianwei Yang for generously providing the image editing tool utilizing the X-Decoder framework.

References

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022. 2
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022. 2
- [3] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 2
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 2
- [5] Harrison Chase. Langchain. <https://langchain.readthedocs.io/>, 2023. 4



$5 + a = 7$	$40 + y = 45$
$10 + y = 15$	$39 + a = 41$
$20 + a = 22$	$24 + y = 29$
$48 + a = 50$	$15 + a = 17$
$10 + a = 12$	$34 + y = 39$

This image contains 10 math equations.



What is "a" in the first equation?

The value of "a" in the first equation " $5 + a = 7$ " is 2.



What is y in the last equation?

The value of "y" in the last equation " $34 + y = 39$ " is 5.



What is the answer to the third equation on the left?

The answer to the third equation on the left " $20 + a = 22$ " is 2.



What is the answer to the second last equation on the right?

The answer to the second last equation on the right " $15 + a = 17$ " is 2.



This image is likely a classroom with children writing on a chalkboard. The texts in this image are mathematical equations.



What is the answer to the first equation from the left hand side?

The answer to the first equation " $3 \times 3 =$ " is 9.



What is the answer to the first equation from right?

The answer to the first equation from right " $11 - 2 =$ " is 9.



What is the answer to the equation in the middle?

The answer to the equation in the middle " $7 \times 2 =$ " is 14.



Visual Math and Text Reasoning

Figure 4. Case studies of MM-REACT's capabilities and application scenarios: **visual math and text reasoning**.

- [6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 2
- [7] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. Nus-wide: a real-world web image database from national university of singapore. In *Proceedings of the ACM international conference on image and video retrieval*, pages 1–9, 2009. 2
- [8] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. *arXiv preprint arXiv:2302.05442*, 2023. 2
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [10] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model. In *arXiv preprint arXiv:2303.03378*, 2023. 2, 3, 6, 20, 21, 22
- [11] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. *arXiv preprint arXiv:2211.10435*, 2022. 2
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2
- [13] Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. Promptcap: Prompt-guided task-aware image captioning. *arXiv preprint arXiv:2211.09699*, 2022. 2
- [14] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*, 2023. 2, 3
- [15] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR, 2022. 2
- [16] Mark J Huiskes and Michael S Lew. The mir flickr retrieval evaluation. In *Proceedings of the 1st ACM international conference on Multimedia information retrieval*, pages 39–43, 2008. 2
- [17] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022. 2, 3
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 2
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 2
- [20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 2
- [21] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *arXiv preprint arXiv:2209.09513*, 2022. 3
- [22] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021. 2
- [23] OpenAI. Gpt-4 technical report. 2023. 2, 4
- [24] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023. 2
- [25] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *arXiv preprint arXiv:2303.08128*, 2023. 3
- [26] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*, 2022. 2
- [27] Maria Tsimpoukelli, Jacob Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *arXiv preprint arXiv:2106.13884*, 2021. 2
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2
- [29] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*, 2022. 2
- [30] Zhenhailong Wang, Manling Li, Ruochen Xu, Luwei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chengguang Zhu, Derek Hoiem, et al. Language models with image descriptors are strong few-shot video-language learners. In *Advances in Neural Information Processing Systems*. 2
- [31] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022. 2, 3

- [32] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 3
- [33] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 3081–3089, 2022. 2
- [34] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022. 2, 4
- [35] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. 2
- [36] Andy Zeng, Adrian Wong, Stefan Welker, Krzysztof Choromanski, Federico Tombari, Aavek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022. 2
- [37] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023. 3
- [38] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Jianfeng Wang, Lu Yuan, Nanyun Peng, Lijuan Wang, Yong Jae Lee, and Jianfeng Gao. Generalized decoding for pixel, image and language. 2022. 6, 30



This image is a cartoon illustration of a squirrel holding a camera



Why would a person find this image funny?

The humor in this image is subjective and can vary from person to person. Some people might find the cartoon squirrel holding a camera funny because it is a cute and quirky illustration. Others might find it funny because it is a humorous take on the idea of a squirrel taking photos. Ultimately, the humor in this image is up to the interpretation of the individual viewer. So, it depends on the person's sense of humor and personal preferences.



This image is likely a comic or cartoon illustration.



Is there text in the image? If so, can you read it and explain it to me?

The text in the image is a conversation between two people. One person says they don't feel good and need to buy a test. The other person asks if the first person is pregnant or has the coronavirus. The text is likely part of a comic or cartoon illustration. The image is tagged with "@INNA_S_ART", which may indicate the artist or creator of the image.



What is funny about it?

The humor in the comic may come from the situation, the characters, or the dialogue. The person is feeling unwell and is asked if they are pregnant or have the coronavirus, which is a current and relevant topic. The humor may come from the absurdity of the situation or the unexpectedness of the question. However, humor is subjective and what one person finds funny, another person may not. It is best to view the comic and decide for yourself if it is funny.



What is unusual about this image?



The unusual thing about this image is that a man is ironing clothes in the back of a yellow taxi. This is not a common sight and may be considered unusual or unique.



Visual-conditioned Joke/Meme

Figure 5. Case studies of MM-REACT's capabilities and application scenarios: **visual-conditioned joke/meme**. We provide an example of the unfolded multi-modal reasoning and action steps for meme understanding in Figure 18.



This image is likely a dog playing frisbee with a man.



Where is the frisbee in the image?

The frisbee is located at the top left corner of the image with coordinates (63, 119) and (133, 177).



Where is the man in the image?

The man is located at the coordinates (66, 386) to (453, 604) in the image.

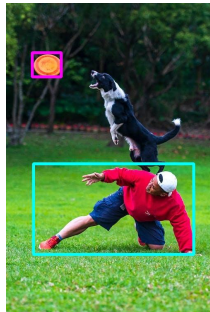


What is the spatial relation between the frisbee and the man?

The frisbee is located above the man in the image. The man is falling on the ground while the frisbee is in the air. The frisbee is also closer to the top of the image than the man..



Illustration of detected bbox from MM-React



Spatial/Coordinate Understanding



This image is a recipe for making bread at home. The ingredients include water, salt, flour, sugar, and yeast. The instructions include kneading the dough, dividing it into pieces, and baking it on each side until well done. The temperature is in Celsius.



What should I do after kneading the dough, before dividing it into pieces?

After kneading the dough, you should put it in a warm place before dividing it into pieces, according to the recipe in the image.



How long should I put in warm places?

You should put the dough in a warm place for 1 hour.



Visual Planning and Prediction

Figure 6. Case studies of MM-REACT's capabilities and application scenarios: **spatial/coordinate understanding** and **visual planning and prediction**.



Receipt 1:

Your receipt Invoice number: 108 24087321

Payer: David Rowell
Recipient: David Rowell

Date: 22.07.2019 Time: 09:54
VISA **** * 51 Expires: 06/
Sale: 1299.90 USD
0000000 0 0 0 0
0 Transr: 4001288606 Authnr: 4001288606
Reference 16986186

Fare Details

Air fare:	Extras:	Taxes and fees:	Total airfare:	Payment method:
USD 773	USD 205	USD 321.90	USD 1,094.90	CREDIT CARD

Receipt 3:

The Shop
Store #100
Chicago, IL

=====

Large Eggs	0.99
Milk	1.15
Cottage Cheese	0.59
Natural yogurt	0.70
Cherry Tomatoes 1lb	1.29
Bananas 1lb	0.77
Aubergine	1.50
Cheese Crackers	2.19
Chocolate Cookies	1.82
Canned Tuna 12pk	5.95
Chicken Breast	2.46
Toilet Paper	4.98
Baby wipes	1.59

TOTAL \$25.97

Receipt 4:

Windmill Family Restaurant
2630 Neth Street
Woodstock, IL 60093
(630) 266-5980

Emp1: MALISSA Chaf: 0000108
Open: May 26, 2017 01:17 PM
Station: 3 Guest: 2 Table: 0051

Qty	Description	Price	Amount
1	SOFT DRINK	2.09	2.09
1	SOFT DRINK	2.09	2.09
1	HOT TURK SAND/HP	9.19	9.19
1	L*CRAB CAKE SMD/FF	8.75	8.75

Sub Total 22.12
Sales Tax 1.33

Total 23.45

Total Due 23.45

Receipt 2:

Uber Total: **\$43.83**
Wed, Nov 06, 2019

Thanks for riding, Steven

We're glad to have you as an Uber Rewards Gold Member.

Total \$43.83
You earned 87 points on this trip

Trip Fare	\$38.81
Subtotal	\$38.81
NY Black Car Fund	\$1.03
TNC Assessment Fee	\$1.59
Wait Time	\$0.00
Tolls, Surcharges, and Fees	\$2.40
Amount Charged	
Uber Cash	\$43.83



How much did I spend on grocery?

You spent \$25.97 on groceries at The Shop Store #100 in Chicago, IL.



How much did I spend on dining out?

You spent \$23.45 on dining out at the Windmill Family Restaurant.



How much in total did I spend on travel?

You spent \$1299.90 on airfare and \$43.83 on an Uber ride for a total of \$1343.73 on travel.



How much in total did I pay for taxes?

You paid \$323.23 in total for taxes. This includes \$321.90 for taxes and fees for the flight purchase and \$1.33 for sales tax for the Windmill Family Restaurant expenses.



Multi-Image Reasoning

Figure 7. Case studies of MM-REACT's capabilities and application scenarios: **multi-image reasoning**. For the multi-image inputs shown above, we input one receipt image at a time. Once all four receipt images are provided as inputs to MM-REACT, we prompt it to answer questions that require reasoning over multiple images.



This image is a bar chart of the total fueling cost of different cars. The cars are Ford F150 (Regular Cab), Toyota RAV4, and Honda Civic. The total fueling cost of Ford F150 (Regular Cab) is \$130.96, the total fueling cost of Toyota RAV4 is \$82.56, and the total fueling cost of Honda Civic is \$70.55.



Excluding Ford F150 Lightning, what is the average total charging cost?

The average total charging cost excluding the Ford F150 Lightning (Estimate) is \$22.44. This is calculated by taking the sum of the total charging cost of Rivian R1T/S and Tesla Model 3 (Standard Range) and dividing it by 2. $(\$31.35 + \$12.54) / 2 = \$22.44$.

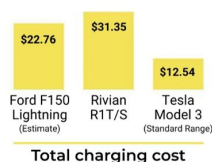


Now how much cheaper is the average charging cost over the electric cars (including Ford F150 Lightning) than the average fueling cost over the gasoline cars (including Ford F150 (Regular Cab))? Let's think step-by-step.



Excluding the regular cab, what is the average total fueling cost?

The average total fueling cost excluding the Ford F150 (Regular Cab) is \$76.55. This is calculated by taking the sum of the total fueling cost of Toyota RAV4 and Honda Civic and dividing it by 2. $(\$82.56 + \$70.55) / 2 = \$76.55$.



This image is a chart of the total charging cost of different electric vehicles. The electric vehicles are Ford F150 Lightning (Estimate), Rivian R1T/S, and Tesla Model 3 (Standard Range). The total charging cost of Ford F150 Lightning (Estimate) is \$22.76, the total charging cost of Rivian R1T/S is \$31.35, and the total charging cost of Tesla Model 3 (Standard Range) is \$12.54.



The average charging cost over the three electric cars (including Ford F150 Lightning) is \$18.66, which is calculated by taking the sum of the total charging cost of Ford F150 Lightning (Estimate), Rivian R1T/S, and Tesla Model 3 (Standard Range) and dividing it by 3. $(\$22.76 + \$31.35 + \$12.54) / 3 = \18.66 .

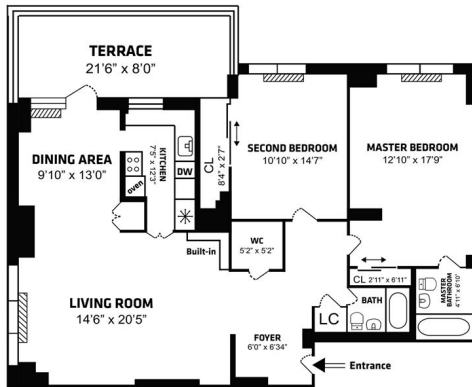


The average fueling cost over the three gasoline cars (including Ford F150 (Regular Cab)) is \$101.07, which is calculated by taking the sum of the total fueling cost of Ford F150 (Regular Cab), Toyota RAV4, and Honda Civic and dividing it by 3. $(\$130.96 + \$82.56 + \$70.55) / 3 = \101.07 .

Therefore, the average charging cost is \$82.41 cheaper than the average fueling cost. This is calculated by subtracting the average charging cost from the average fueling cost. $\$101.07 - \$18.66 = \$82.41$.

Multi-Hop Document Understanding

Figure 8. Case studies of MM-REACT's capabilities and application scenarios: **multi-hop document understanding (bar charts)**. The unfolded multi-modal reasoning and action steps of this example are detailed in Figure 19.



This image contains a floor plan of a house with dimensions of various rooms.



How many bedrooms are there in this floor plan?

There are two bedrooms in this floor plan. The first one is the second bedroom and the second one is the master bedroom.



What is the dimension of Master Bedroom?

The dimension of the Master Bedroom is 12'10" x 17'9".



Is there a bathroom in master bedroom?

Yes, there is a bathroom in the master bedroom. The dimension of the master bathroom is 4'11" x 6'10".



What is the dimension of Second Bedroom?

The dimension of the Second Bedroom is 10'10" x 14'7".



Is the terrace bigger than the kitchen?

The terrace is 21'6" x 8'0" and the kitchen is 7'5" x 12'3". Based on these dimensions, the terrace is larger than the kitchen.



What are the appliances in the kitchen in this floorplan?

The appliances in the kitchen in this floor plan are oven and DW (dishwasher). The size of the oven is not specified in the floor plan.



What will you see when you enter through the entrance? and what is the dimension of it?

When you enter through the entrance, you will see the foyer. The dimension of the foyer is 6'0" x 6'34".



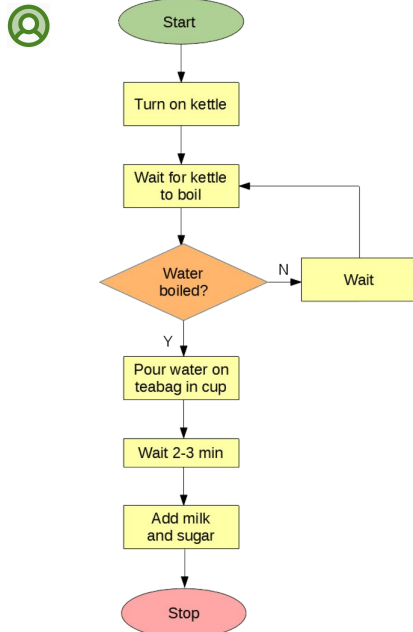
Can you summarize this floorplan based on the aspects I was interested in, in a few sentences?

Sure, this floor plan is of a house with two bedrooms, a master bedroom with a bathroom and a second bedroom. The master bedroom is 12'10" x 17'9" and the second bedroom is 10'10" x 14'7". The kitchen is 7'5" x 12'3" and the terrace is 21'6" x 8'0". The foyer is 6'0" x 6'34" and the entrance leads to it. The kitchen has an oven and a dishwasher.



Document Understanding

Figure 9. Case studies of MM-REACT's capabilities and application scenarios: **document understanding (floorplan)**.



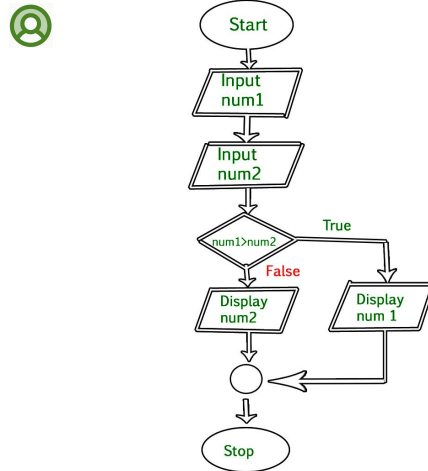
This image is a flowchart diagram with text and shapes.

Can you read the text in this flowchart and explain the process in detail?

This flowchart describes a process for making tea. The process starts with turning on the kettle and waiting for it to boil. Then, the water is poured over a teabag in a cup and the tea is allowed to steep for 2-3 minutes. Finally, milk and sugar are added to the tea. The process ends with the stop.

Which step in the flowchart is drawn with a diamond shape?

In the flowchart above, the step "Water boiled?" is drawn with a yellow diamond with black text. This step represents the decision of whether the water has boiled or not. If the water has boiled, the process continues to the next step. If not, the process waits until the water has boiled.



This image is a flowchart diagram of a computer program. The text in the image is likely the labels for the different parts of the flowchart.

Can you read the text in this flowchart and explain it to me?

The flowchart appears to be a simple program that takes two inputs, num1 and num2, and compares them. If num1 is greater than num2, it displays num1. If num2 is greater than num1, it displays num2. The program starts with the "Start" label and ends with the "Stop" label.

Can you transform this flowchart into a simple program in python?

Sure, I can try to write a simple program in python based on the flowchart.

```

1 def display_greater_number(num1, num2):
2     if num1 > num2:
3         print(num1)
4     else: print(num2)
5 num1 = int(input("Input num1: "))
6 num2 = int(input("Input num2: "))
7 display_greater_number(num1, num2)
  
```

Document Understanding

Figure 10. Case studies of MM-REACT's capabilities and application scenarios: **document understanding (flowchart)**. We provide an example of the unfolded multi-modal reasoning and action steps for flowchart understanding in Figure 20.



Dataset 1

Department	Name	Salary
Administration	Emily	\$2,000
Finance	Nick	\$1,800
Marketing	John	\$2,600
Marketing	Gloria	\$1,800
Administration	Jane	\$2,500
Finance	Max	\$2,000
Administration	Tyler	\$3,000
Finance	Wilkins	\$2,500
Marketing	Miles	\$2,800
Administration	Sabine	\$1,900

This is a table of datasets with Department, Name, and Salary columns. The table has 10 rows of data. The highest salary is \$3,000 and the lowest salary is \$1,800. The most common department is Administration with 4 employees. The least common department is Marketing with 3 employees.



Sure, I can calculate the yearly salary for each employee and group them by department.



Department	Name	Monthly Salary	Yearly Salary
Administration	Emily	\$2,000	\$24,000
Administration	Jane	\$2,500	\$30,000
Administration	Tyler	\$3,000	\$36,000
Administration	Sabine	\$1,900	\$22,800
Finance	Nick	\$1,800	\$21,600
Finance	Max	\$2,000	\$24,000
Finance	Wilkins	\$2,500	\$30,000
Marketing	John	\$2,600	\$31,200
Marketing	Gloria	\$1,800	\$21,600
Marketing	Miles	\$2,800	\$33,600

The total yearly salary for Administration is \$112,800. The total yearly salary for Finance is \$75,600. The total yearly salary for Marketing is \$86,400.



How much is Wilkins earning?

Wilkins is earning \$2,500.



Given the yearly bonus is 10%, 12% and 15% for Administration, Finance, Marketing, respectively. Please provide me with the bonus for Emily, Max and Gloria.



Which department is the person who has the highest salary from ?

The person who has the highest salary is from the Administration department. They are earning \$3,000.



Sure, I can calculate the bonus for Emily, Max and Gloria.



Emily is from the Administration department and her yearly salary is \$24,000. Her bonus is \$2,400 (10% of \$24,000).

Max is from the Finance department and his yearly salary is \$24,000. His bonus is \$2,880 (12% of \$24,000).

Gloria is from the Marketing department and her yearly salary is \$21,600. Her bonus is \$3,240 (15% of \$21,600).



This table is reporting monthly salary, can you calculate the yearly salary for each employee, and summarize it into a table and group them by department?

Document Understanding

Figure 11. Case studies of MM-REACT's capabilities and application scenarios: **document understanding (table)**.



This image contains different types of milk cartons with different brands and flavors.



What is the brand of these drink? Which one is unsweetened?

The brand of these drinks is "BETTER THAN MILK". The unsweetened one is "BETTER THAN MILK organic OAT DRINK".



This image is likely a picture of a white couch with pillows.



I am trying to buy this couch. Can you find which brand it is?

The couch in the image is likely a Soderhamn Beige Viarp couch from IKEA.



Keira Knightley and Natalie Portman are the celebrities in this image.



This image contains celebrities Justin Bieber, Angelina Jolie, Kim Kardashian, Beyonce, Selena Gomez, and Taylor Swift.



This image is likely a bowl of morel mushrooms.



Is it safe to eat morel mushroom?

Morel mushrooms are generally considered safe to eat when properly cooked. However, it is important to note that some people may have an allergic reaction to morel mushrooms. It is recommended to consult a professional or do further research to determine the safety of consuming morel mushrooms.



Where and when was this image taken?

This image is taken in İzmir, Turkey during an earthquake that happened on 2023-02-16.



Open-World Concept Understanding

Figure 12. Case studies of MM-REACT's capabilities and application scenarios: **open-world concept understanding**.



Please create a summary of the video in less than 10 sentences.

The speaker is making a BLT sandwich with grilled tomatoes and fried bacon. They are using fresh bread, mayo, Worcestershire sauce, lettuce, and tomatoes. They are sprinkling salt and pepper on top. The bread can be toasted or grilled. The speaker mentions that the few ingredients used in the sandwich will give it a lot of flavor. The speaker mentions that the bread is important and can be toasted or grilled. The final product is a BLT sandwich.



Can you breakdown this tutorial of making a BLT sandwich to step-by-step instructions, with start and end time?

Here's a breakdown of the start and end times for each step in the video:

1. Heat a little olive oil in a pan and grill the tomatoes: 00:00:22 - 00:00:40
2. Cook the bacon in a frying pan: 00:00:51 - 00:01:41
3. Place the cooked bacon on a paper towel: 00:01:39 - 00:01:49
4. Spread mayo and add Worcestershire sauce: 00:01:55 - 00:02:16
5. Assemble the sandwich (lettuce, tomatoes, bacon): 00:02:22 - 00:03:11
 - a. Layer lettuce: 00:02:44 - 00:02:49
 - b. Add grilled tomatoes: 00:02:55 - 00:03:01
 - c. Place bacon on top: 00:03:15 - 00:03:27
6. Sprinkle salt and pepper: 00:03:08 - 00:03:11
7. Top the sandwich with the other slice of bread: 00:03:29 - 00:03:34



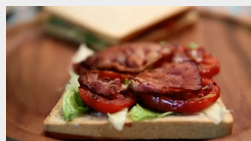
I am trying to create a short version of the video, can you suggest a shorter segment for each step?

Here's a suggested shorter segment for each step:

1. Heat a little olive oil in a pan and grill the tomatoes: 00:00:24 - 00:00:30
2. Cook the bacon in a frying pan: 00:00:53 - 00:00:59
3. Place the cooked bacon on a paper towel: 00:01:39 - 00:01:41
4. Spread mayo and add Worcestershire sauce: 00:02:00 - 00:02:08
5. Assemble the sandwich (lettuce, tomatoes, bacon): 00:02:44 - 00:03:01
 - a. Layer lettuce: 00:02:44 - 00:02:47
 - b. Add grilled tomatoes: 00:02:55 - 00:02:59
 - c. Place bacon on top: 00:03:15 - 00:03:20
6. Sprinkle salt and pepper: 00:03:08 - 00:03:11
7. Top the sandwich with the other slice of bread: 00:03:29 - 00:03:34



Which segment is best for video thumbnail to attract more audience?



Suggested Thumbnail at 00:03:29

For a video thumbnail that attracts more audience, the segment showing the assembled sandwich with all its layers would be the most appealing. This image will showcase the delicious combination of ingredients and the final result of the recipe. The suggested segment for this is:

- Assembled BLT sandwich: 00:03:29 - 00:03:34

This segment shows the completed sandwich just before adding the final slice of bread on top. You can take a screenshot of the sandwich from this segment or create a more stylized image based on it for your thumbnail.



Video Summarization/Event Localization

Figure 13. Case studies of MM-REACT's capabilities and application scenarios: **video summarization/event localization.**



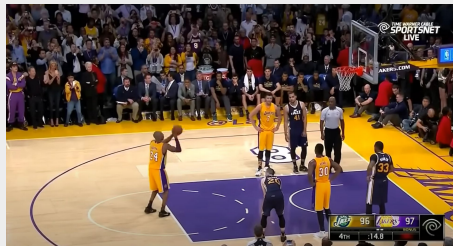
Please create a summary of the video in less than 10 sentences.

The video is about Kobe Bryant's final game in the NBA. He scored 60 points and led the Lakers to a victory over the Utah Jazz. The crowd was out of their minds and the performance was considered one of the greatest in the history of the game. Kobe was fouled and made two free throws for 60 points. The video ends with the crowd applauding and Kobe's final memory being a 60-point game and a victory over the Utah Jazz. Overall, it was a magical farewell performance for one of the greatest athletes ever.



How did Kobe win the last 2 points to get 60? Explain and locate the start and end time of relevant segment in this video.

Kobe Bryant won the last 2 points to reach 60 by making two free throws. He was fouled and went to the free-throw line. The relevant section starts at 00:04:02 with the mention of the foul, it ends at 00:04:39 after Kobe successfully converts both free throws.



First Free throw at 00:04:11 - 00:04:15



Second Free throw at 00:04:31 - 00:04:34



I want to make a short video to highlight the long deuce by Kobe if there is any in the video. Can you tell me roughly the start and end time of that.

The long deuce by Kobe Bryant is mentioned at timestamp 00:03:28. The play leading up to this shot starts around 00:03:21, and the audience's reaction to the long deuce ends at approximately 00:03:32.



Video Summarization/Event Localization

Figure 14. Case studies of MM-REACT's capabilities and application scenarios: **video summarization/event localization**.



PaLM-E

Describe the following : A dog jumping over a hurdle at a dog show.

Ours



Please describe the image

This image is likely a dog jumping over a hurdle in a dog agility competition.



Image Captioning, VQA



PaLM-E

I see (image omitted). Description: A cow sitting on a rug, Joke: Why is the cow sitting on a rug? Because it wants to be a cow-ch! .. I see . Description: A donkey is carrying a dog, cat, and rooster. Joke: What do you call a donkey with a rooster on his back. A rooster booster.

Ours



Please write a joke for this image

Why did the cat stand on the dog's back? Because the donkey was full!



Visual-Conditioned Jokes

PaLM-E

Q: Where is <img1> in this scene <img2>? A: A stove top with a pot on it.

Ours



Where is the coffee pot in the second image?

The coffee pot is located in the kitchen with a stove and cabinets. It is a silver coffee pot on a stove and its location is at (242, 230) to (295, 299).



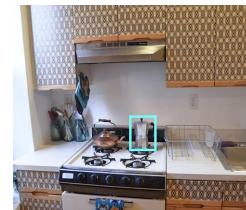
Image1



Image2



Illustration of detected bbox from MM-React



Multi-Image Relationships

Figure 15. Comparison of MM-REACT with PaLM-E [10] on illustrated capabilities. We empirically show that text prompts are as effective as expensive joint fine-tuning in solving complicated vision problems.



PaLM-E

Given **. Q: How many championship rings did the player on the left win in his career A: Let's think step by step 1. kobe bryant is on the left. 2. kobe bryant has won 5 championship rings.

Ours



How many championship rings did the player on the left win in his career?

Kobe Bryant won 5 championship rings in his career.



Multimodal Chain-of-Thought



PaLM-E

Given **. Q: I am getting just two custom pizza for me and my friend. How much should I pay in total? A: Let's think step by step. 1st step: 2 custom pizzas. 2nd step: 2 x \$9.99. 3rd step: \$19.98.

Ours



I am getting just two custom pizza for me and my friend. How much should I pay in total?

Based on the information you provided, you are ordering two custom pizzas. The price of each custom pizza is \$9.99. So, the total cost of two custom pizzas would be $2 \times \$9.99 = \19.98



Scene Text Reasoning

Photo1 at 10:30am

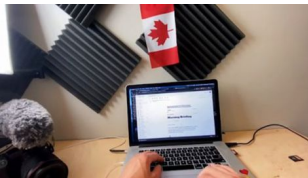


Photo2 at 12:45pm



Photo3 at 4:45pm



PaLM-E

Photo 1, at 10:30 am: *<img1>*. Photo 2, at 12:45 pm: *<img2>*. Photo 3, at 3:45 pm: *<img3>*. Q: I forget, what did I have for lunch, and what time was it? A: Let's think step by step. 1. you had a sandwich for lunch. 2. it was 12:45 pm

Ours



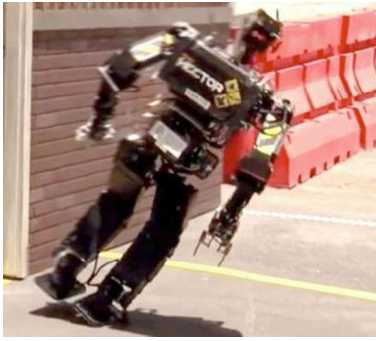
I forget, what did I have for lunch, and what time was it?

You had a sandwich for lunch at 12:45pm.



End-to-End Egocentric Q&A, Dialogue

Figure 16. Comparison of MM-REACT with PaLM-E [10] on illustrated capabilities. We empirically show that text prompts are as effective as expensive joint fine-tuning in solving complicated vision problems.



PaLM-E

Given . Q: What will the robot do next?
A: fall.

Ours



What will the robot do next?

I'm sorry, I cannot predict the future actions of the robot. The image only shows its current state.



Physical Prediction

PaLM-E

When a human asks me to do a task, I will respond with the sequence of actions I would do to accomplish the task with only the items I see.

Human: Use all of the ingredients you see to make a cake batter. 1. crack egg. 2. put egg in bowl. 3. put flour in bowl. 4. put sugar in bowl. 5. mix. 6. put in pan. 7. bake. 8. eat. 9. clean up.

Ours



When a human asks me to do a task, I will respond with the sequence of actions I would do to accomplish the task with only the items I see. Human: Use all of the ingredients you see to make a cake batter.

To make a cake batter using the ingredients I see in the image, I would do the following:

1. Gather all of the bags of flour and empty them into a large mixing bowl.
2. Add any other dry ingredients, such as sugar or baking powder, to the bowl.
3. Mix the dry ingredients together until they are well combined.
4. Add any liquids, such as milk or eggs, to the bowl and mix until the batter is smooth.
5. Pour the batter into a greased cake pan and bake according to the recipe instructions.



Robot Perception, Dialogue, Planning

Figure 17. Comparison of MM-REACT with PaLM-E [10] on illustrated capabilities. We empirically show that text prompts are as effective as expensive joint fine-tuning in solving complicated vision problems.

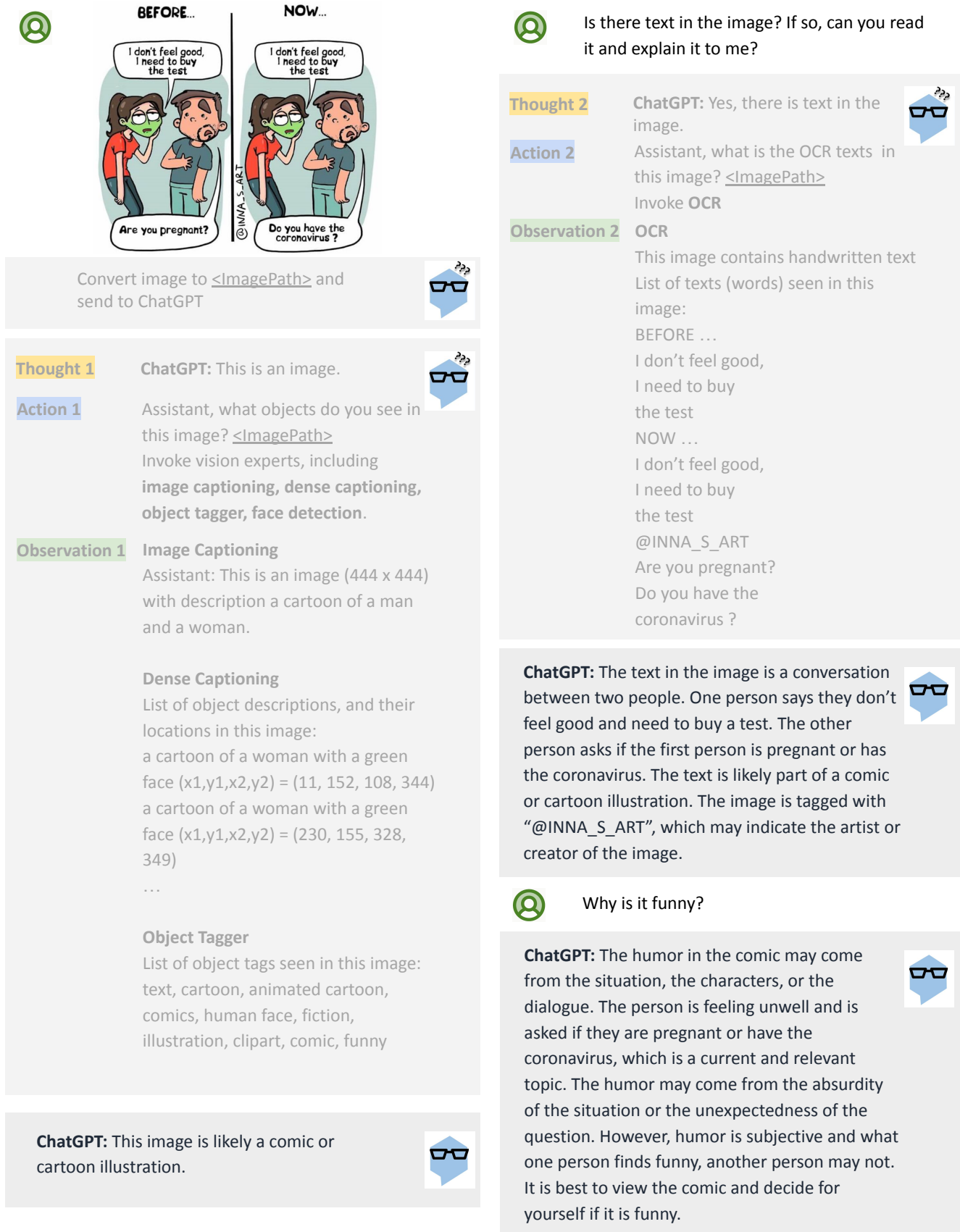





Figure 18. Unfolded multimodal reasoning and action steps for example of visual-conditioned joke/meme in Figure 5.

Total fueling cost

Vehicle	Total fueling cost
Ford F150 (Regular Cab)	\$130.96
Toyota RAV4	\$82.56
Honda Civic	\$70.55

Convert image to `<ImagePath>` and send to ChatGPT



Thought 1 ChatGPT: This is an image.

Action 1 Assistant, what objects do you see in this image? `<ImagePath>`
Invoke vision experts, including image captioning, dense captioning, object tagger.

Observation 1 **Image Captioning**
This is an image (410 x 593) with description a graph of a number of cars.

Dense Captioning
List of object descriptions, and their locations in this image:
a purple rectangular object with white border (x1,y1,x2,y2) = (152, 169, 245, 446)
a purple rectangular bar chart (x1,y1,x2,y2) = (25, 135, 365, 456)
...


Object Tagger
List of object tags seen in this image: text, screenshot, purple, violet, graphic design, ...


Thought 2 ChatGPT: This image is likely a bar chart or graph.

Action 2 Assistant, what is the OCR texts in this image? `<ImagePath>`
Invoke OCR

Observation 2 **OCR**
This image contains text
List of texts (words) seen in this image:
\$130.96
Ford
F150
(Regular Cab)
\$82.56
Toyota
RAV4
\$70.55
Honda
Civic
Total fueling cost


ChatGPT: This image is a bar chart of the total fueling cost of different cars. The cars are Ford F150 (Regular Cab), Toyota RAV4, and Honda Civic. The total fueling cost of Ford F150 (Regular Cab) is \$130.96, the total fueling cost of Toyota RAV4 is \$82.56, and the total fueling cost of Honda Civic is \$70.55.


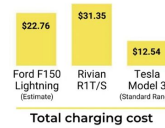




Excluding the regular cab, what is the average total fueling cost?

ChatGPT: The average total fueling cost excluding the Ford F150 (Regular Cab) is \$76.55. This is calculated by taking the sum of the total fueling cost of Toyota RAV4 and Honda Civic and dividing it by 2. $(\$82.56 + \$70.55) / 2 = \$76.55$.




Total charging cost

Vehicle	Total charging cost
Ford F150 Lightning (Estimate)	\$22.76
Rivian R1T/S	\$31.35
Tesla Model 3 (Standard Range)	\$12.54

Convert image to `<ImagePath>` and send to ChatGPT



Thought 3 ChatGPT: This is an image.

Action 3 Assistant, what objects do you see in this image? `<ImagePath2>`
Invoke vision experts, including image captioning, dense captioning, object tagger.

Observation 3 **Image Captioning**
This is an image (400 x 297) with description a chart of cost comparison.

Dense Captioning
List of object descriptions, and their locations in this image:
a yellow rectangle with black text (x1,y1,x2,y2) = (154, 12, 245, 150)
a chart of cost per charge (x1,y1,x2,y2) = (1, 1, 391, 288)
...


Object Tagger
List of object tags seen in this image: text, screenshot, line, yellow, number, ...


Thought 4 ChatGPT: This image is likely a chart of cost comparison.

Action 4 Assistant, what is the OCR texts in this image? `<ImagePath>`
Invoke OCR

Observation 4 **OCR**
This image contains text
List of texts (words) seen in this image:
\$22.76
Ford F150
Lightning
(Estimate)
\$31.35
Rivian
R1T/S
\$12.54
Tesla
Model 3
(Standard Range)
Total charging cost


ChatGPT: This image is a chart of the total charging cost of different electric vehicles. The electric vehicles are Ford F150 Lightning (Estimate), Rivian R1T/S, and Tesla Model 3 (Standard Range). The total charging cost of Ford F150 Lightning (Estimate) is \$22.76, the total charging cost of Rivian R1T/S is \$31.35, and the total charging cost of Tesla Model 3 (Standard Range) is \$12.54.






Excluding Ford F150 Lightning, what is the average total charging cost?

ChatGPT: The average total charging cost excluding the Ford F150 Lightning (Estimate) is \$22.44. This is calculated by taking the sum of the total charging cost of Rivian R1T/S and Tesla Model 3 (Standard Range) and dividing it by 2. $(\$31.35 + \$12.54) / 2 = \$22.44$.





Now how much cheaper is the average charging cost over the electric cars (including Ford F150 Lightning) than the average fueling cost over the gasoline cars (including Ford F150 (Regular Cab))? Let's think step-by-step.

ChatGPT: The average charging cost over the three electric cars (including Ford F150 Lightning) is \$18.66, which is calculated by taking the sum of the total charging cost of Ford F150 Lightning (Estimate), Rivian R1T/S, and Tesla Model 3 (Standard Range) and dividing it by 3. $(\$22.76 + \$31.35 + \$12.54) / 3 = \18.66 .

The average fueling cost over the three gasoline cars (including Ford F150 (Regular Cab)) is \$101.07, which is calculated by taking the sum of the total fueling cost of Ford F150 (Regular Cab), Toyota RAV4, and Honda Civic and dividing it by 3. $(\$130.96 + \$82.56 + \$70.55) / 3 = \101.07 .

Therefore, the average charging cost is \$82.41 cheaper than the average fueling cost. This is calculated by subtracting the average charging cost from the average fueling cost. $\$101.07 - \$18.66 = \$82.41$.


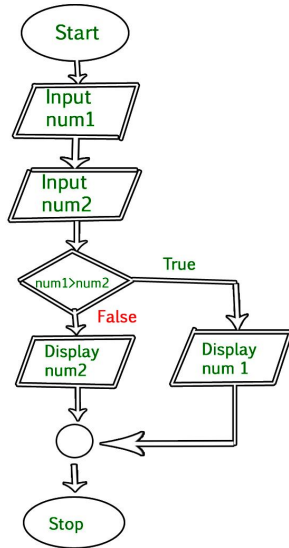


Figure 19. Unfolded multimodal reasoning and action steps for multi-hop document understanding (bar charts) in Figure 8.



Convert image to [ImagePath](#) and send to ChatGPT



Thought 1

ChatGPT: This is an image.



Action 1

Assistant, what objects do you see in this image? [ImagePath](#)
Invoke vision experts, including **image captioning, dense captioning, object tagger**.

Observation 1

Image Captioning

This is an image (682 x 1080) with description a diagram of a computer program.

Dense Captioning

List of object descriptions, and their locations in this image:

a white sign with green text
(x1,y1,x2,y2) = (104, 641, 329, 748)
a green text on a white background
(x1,y1,x2,y2) = (115, 948, 314, 1052)
...

Object Tagger

List of object tags seen in this image:
text, diagram, sketch, font, design

ChatGPT: This image is a flowchart diagram of a computer program. The text in the image is likely the labels for the different parts of the flowchart.



Can you read the text in this flowchart and explain it to me?

Thought 2

ChatGPT: Let me ask for more information about the text.



Action 2

Assistant, what is the OCR texts in this image? [ImagePath](#)
Invoke **OCR**

Observation 2

OCR

This image contains text
List of texts (words) seen in this image:
Start
Input
num1
Input
num2
num1 > num2
False
Display
num2
Stop
True
Display
num 1

ChatGPT: The flowchart appears to be a simple program that takes two inputs, num1 and num2, and compares them. If num1 is greater than num2, it displays num1. If num2 is greater than num1, it displays num2. The program starts with the "Start" label and ends with the "Stop" label.



Can you transform this flowchart into a simple program in python?

ChatGPT: Sure, I can try to write a simple program in python based on the flowchart.



```

1 def display_greater_number(num1, num2):
2     if num1 > num2:
3         print(num1)
4     else: print(num2)
5 num1 = int(input("Input num1: "))
6 num2 = int(input("Input num2: "))
7 display_greater_number(num1, num2)
  
```

Figure 20. Unfolded multimodal reasoning and action steps for document understanding (flowchart) in Figure 10.



Convert image to <ImagePath> and send to ChatGPT



Thought 1

ChatGPT: This is an image.



Action 1

Assistant, what objects do you see in this image? <ImagePath>
Invoke vision experts, including **image captioning, dense captioning, object tagger.**

Observation 1

Image Captioning
Assistant: This is an image (604 x 720) with description a silver and black coffee pot.

Dense Captioning

List of object descriptions, and their locations in this image:
a silver and black coffee pot
(x1,y1,x2,y2) = (36, 30, 573, 692)
...

Object Tagger

List of object tags seen in this image:
coffee maker, pot, kitchenware, small appliance, coffee percolator, home appliance, kitchen appliance ...

ChatGPT: This image is likely a coffee pot



Where is the coffee pot in the second image?

ChatGPT: The coffee pot is located in the kitchen with a stove and cabinets. It is a silver coffee pot on a stove and its location is at (242, 230) to (295, 299).



Convert image to <ImagePath2> and send to ChatGPT



Thought 2

ChatGPT: This is an image.



Action 2

Assistant, what objects do you see in this image? <ImagePath2>
Invoke vision experts, including **image captioning, dense captioning, object tagger.**

Observation 2

Image Captioning
Assistant: This is an image (500 x 430) with description a kitchen with a kettle on top of a stove.

Dense Captioning

List of object descriptions, and their locations in this image:
a tea kettle on a stove (x1,y1,x2,y2) = (88, 237, 317, 426)
a close up of a dish rack (x1,y1,x2,y2) = (311, 242, 432, 324)
a silver coffee pot on a stove
(x1,y1,x2,y2) = (242, 230, 295, 299)
...

Object Tagger

List of object tags seen in this image:
indoor, appliance, kitchen appliance, home appliance, countertop, wall, cabinetry, kitchen, major appliance, gas stove, cabinet, kitchen stove ...

ChatGPT: This image is likely a kitchen with a stove and cabinets.



Figure 21. Unfolded multimodal reasoning and action steps for example of multi-image relationships in Figure 15.

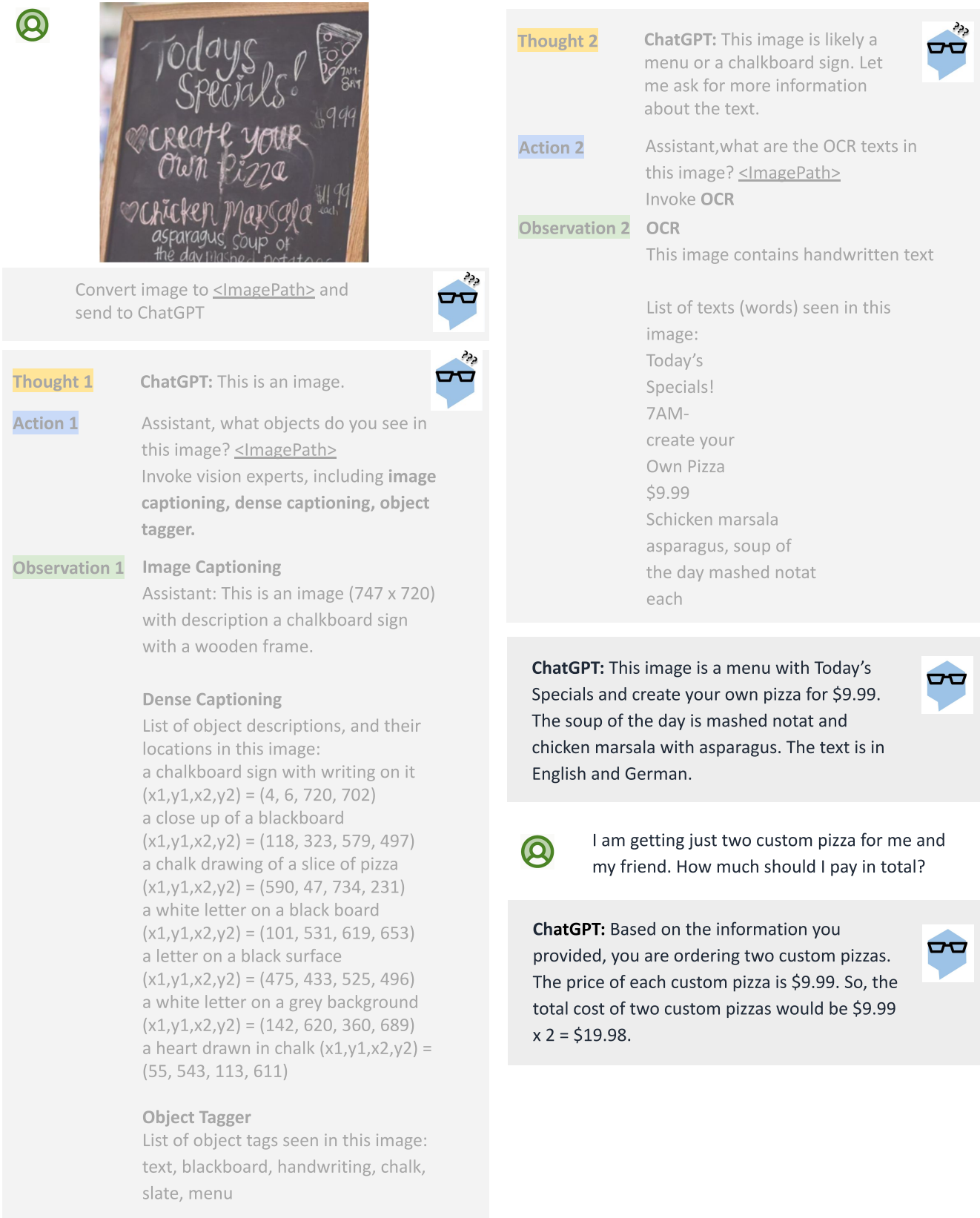


Figure 22. Unfolded multimodal reasoning and action steps for example of scene text reasoning in Figure 16.

Comme illustré sur la figure 1 un bolomètre est constitué d'un absorbeur qui reçoit le rayonnement que l'on désire détecter. Sa température T , supposée uniforme, est mesurée à l'aide d'un thermomètre incorporé, constitué d'un matériau conducteur dont la résistance $R(T)$ varie avec la température T ; cette variation est caractérisée par le coefficient $\alpha = \frac{1}{R} \frac{dR}{dT}$. L'ensemble possède la capacité thermique C_{th} .

Un barreau, conducteur thermique, homogène, de longueur L , de section S et de conductivité thermique λ et sans échanges thermiques latéraux, relie le bolomètre à un thermostat de température T_b fixe.

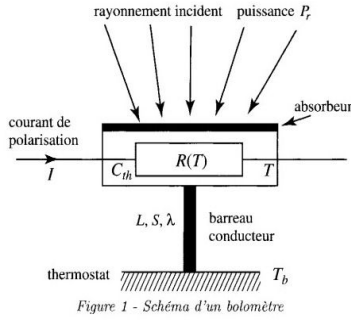


Figure 1 - Schéma d'un bolomètre

I.1. Dans cette première partie, on étudie le comportement du bolomètre en régime stationnaire, à la température de fonctionnement $T = T_0$.

I.1.a) Exprimer la température $T(x)$ en tout point du barreau conducteur repéré par son abscisse x ; on prendra $T(0) = T_0$ et $T(L) = T_b$.

GPT-4 (Multimodal)

User Answer question I.1.a. Think step-by-step.

GPT-4 To answer question I.1.a, we need to express the temperature $T(x)$ in every point of the conductive bar, which is characterized by the abscissa x . We are given that $T(0) = T_0$ and $T(L) = T_b$.

For a stationary regime, the heat equation for the one-dimensional heat conduction in the bar is:

$$d^2T/dx^2 = 0$$

Integrating once, we get:

$$dT/dx = C_1$$

Integrating a second time:

$$T(x) = C_1x + C_2$$

To determine the constants C_1 and C_2 , we use the boundary conditions:

$$T(0) = T_0, \text{ which gives } C_2 = T_0$$

$$T(L) = T_b, \text{ which gives } C_1 = (T_b - T_0) / L$$

Finally, the expression for $T(x)$ is:

$$T(x) = (T_b - T_0) * (x / L) + T_0$$

This represents a linear temperature profile along the conductive bar, with a slope of $(T_b - T_0) / L$.



Answer question I.1.a. Think step-by-step.

Ours w/ ChatGPT

To answer question I.1.a, one would need to use the principles of thermodynamics and heat transfer to determine the temperature at any point in the conductor bar. This would likely involve solving a differential equation to find the temperature distribution along the bar. The solution would depend on the specific properties of the bar, such as its length, cross-sectional area, and thermal conductivity, as well as the temperature of the bolometer and the temperature of the thermostat. The solution would also depend on the power of the incident radiation and the current of polarization. Once the temperature distribution is found, the temperature at any point in the bar can be determined by evaluating the solution at that point.



Ours w/ GPT-4 (Language-only)

To answer question I.1.a, we'll consider the heat conduction equation along the thermal conductor bar. We're working under the assumptions of stationary mode and one-dimensional heat conduction with no lateral heat exchanges.

Start with Fourier's Law of heat conduction:

$$q = -k * (dT/dx)$$

where q is the heat flow (W/m^2), k is the thermal conductivity ($W/m \cdot K$), and dT/dx is the temperature gradient along the x -axis.

In a stationary regime, the heat flow q is constant throughout the bar. Thus, we have:

$$-k * (dT/dx) = \text{constant}$$

Integrate the equation with respect to x :

$$\int (dT/dx) dx = \int (\text{constant}/k) dx$$

After integration, we get:

$$T(x) = (\text{constant}/k) * x + C$$

where C is the integration constant.

Apply the boundary conditions to find the constants:

$$T(0) = T_0 \Rightarrow C = T_0$$

$$T(L) = T_b \Rightarrow (\text{constant}/k) * L + T_0 = T_b$$

Solve for the constant:

$$\text{constant} = k * (T_b - T_0) / L$$

Substitute the constant back into the $T(x)$ equation:

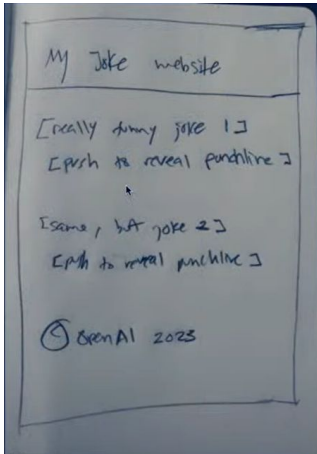
$$T(x) = (k * (T_b - T_0) / L) * x / k + T_0$$

$$T(x) = ((T_b - T_0) / L) * x + T_0$$

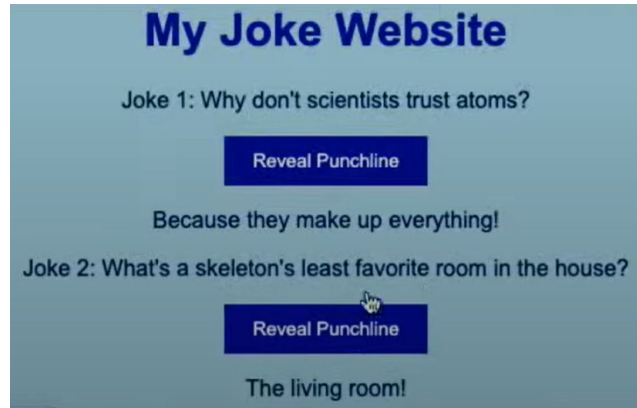
So, the temperature $T(x)$ at any point x along the thermal conductor bar is given by:

$$T(x) = ((T_b - T_0) / L) * x + T_0$$

Figure 23. Case studies of MM-REACT's **extensibility**. With the new release of GPT-4 (Language-only), we can upgrade ChatGPT in MM-REACT to GPT-4 (Language-only) and further improve the results of MM-REACT. In Figure 26, we provide the unfolded multimodal reasoning and action steps based on Ours w/ ChatGPT.

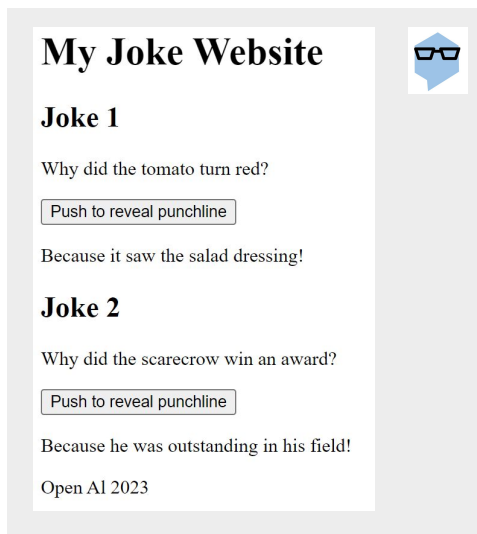


GPT-4 (Multimodal)



Can you write a template for this website in HTML? Replace the joke and punchline with a real one.

Ours w/ ChatGPT



Ours w/ GPT-4 (Language-only)

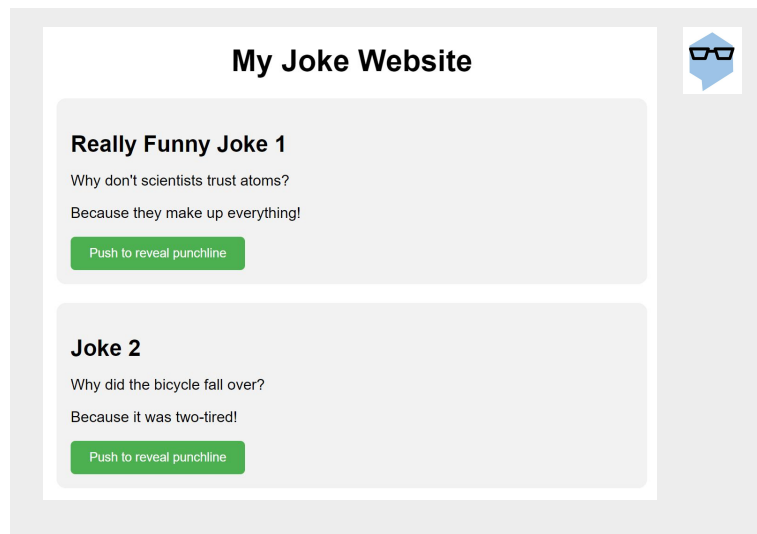


Figure 24. Case studies of MM-REACT’s **extensibility**. With the new release of GPT-4 (Language-only), we can upgrade ChatGPT in MM-REACT to GPT-4 (Language-only) and further improve the results of MM-REACT.



This image is a cartoon illustration of a squirrel holding a camera



Can you change the squirrel in the image to a mouse?



Can you change the camera the mouse is holding in the image to a phone?



There is a pine cone on the ground. Can you change it to a cupcake instead?



Image Editing

Figure 25. Case studies of MM-REACT’s **extensibility**. MM-REACT is designed with the flexibility to plug in new experts. In this example, MM-REACT is enhanced with the image editing model from X-decoder [38] and can perform dialogue-based image editing.



I. Principe de la détection de rayonnement avec un bolomètre

Comme illustré sur la figure 1 un bolomètre est constitué d'un absorbeur qui reçoit le rayonnement que l'on désire détecter. Sa température T , supposée uniforme, est mesurée à l'aide d'un thermomètre incorporé, constitué d'un matériau conducteur dont la résistance $R(T)$ varie avec la température T ; cette variation est caractérisée par le coefficient $\alpha = \frac{1}{R} \frac{dR}{dT}$. L'ensemble possède la capacité thermique C_B .

Un barreau, conducteur thermique, homogène, de longueur L , de section S et de conductivité thermique λ et sans échanges thermiques latéraux, relie le bolomètre à un thermostat de température T_b fixe.

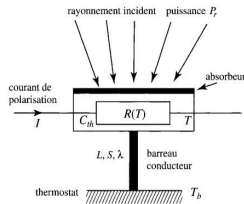


Figure 1 - Schéma d'un bolomètre

I.1. Dans cette première partie, on étudie le comportement du bolomètre en régime stationnaire, à la température de fonctionnement $T = T_b$.

I.1.a) Exprimer la température $T(x)$ en tout point du barreau conducteur repéré par son abscisse x ; on prendra $T(0) = T_b$ et $T(L) = T_b$.

Convert image to [ImagePath](#) and send to ChatGPT



Thought 1

ChatGPT: This is an image.



Action 1

Assistant, what objects do you see in this image? [ImagePath](#)
Invoke vision experts, including image captioning, dense captioning, object tagger.

Observation 1

Image Captioning

This is an image (1716 x 1449) with description a black and white text with text and a diagram.

Dense Captioning

List of object descriptions, and their locations in this image:

a diagram of different directions
(x1,y1,x2,y2) = (656, 553, 1107, 770)
diagram of a diagram of a heat exchanger (x1,y1,x2,y2) = (457, 539, 1241, 1107)
...

Object Tagger

List of object tags seen in this image:
text, diagram, screenshot, font, line, parallel, black and white, design

Thought 2

ChatGPT: This image is likely a diagram or a screenshot with text and a diagram. Let me ask for more information about the text.



Action 2

Assistant, what is the OCR texts in this image? [ImagePath](#)
Invoke OCR

Observation 2

OCR

This image contains text
List of texts (words) seen in this image:

I. Principe de la détection de rayonnement avec un

.

.

.

Figure 1 - Schéma d'un bolomètre

I.1. Dans cette première partie, on étudie le comportement du bolomètre en régime stationnaire, à la température de fonctionnement $T = T_b$.

I.1.a) Exprimer la température $T(x)$ en tout point du barreau conducteur repéré par son abscisse x ; on prendra $T(0) = T_b$ et $T(L) = T_b$.



Answer question I.1.a. Think step-by-step.

ChatGPT: To answer question I.1.a, one would need to use the principles of thermodynamics and heat transfer to determine the temperature at any point in the conductor bar. This would likely involve solving a differential equation to find the temperature distribution along the bar. The solution would depend on the specific properties of the bar, such as its length, cross-sectional area, and thermal conductivity, as well as the temperature of the bolometer and the temperature of the thermostat. The solution would also depend on the power of the incident radiation and the current of polarization. Once the temperature distribution is found, the temperature at any point in the bar can be determined by evaluating the solution at that point.



ChatGPT: This image contains text in French and English. The text appears to be a description of a bolometer and its components. The text is written in a technical style and includes terms such as "absorber," "thermometer," "resistance," "thermal capacity," "thermal conductor," and "thermostat." The text also mentions a diagram (Figure 1) and refers to a temperature profile along a conductor. The text appears to be part of a technical document or scientific paper.



Figure 26. Unfolded multimodal reasoning and action steps with ChatGPT to tackle physics problem in Figure 23.