810 A SUPPLEMENTARY MATERIALS

A.1 ACCESS INSTRUCTIONS

Ins-MMBench and the code for the evaluation experiments will be available.

816 A.2 DATA FORMAT

814

815

817

823

824 825

826

827 828

829 830

837

The INS-MMBench dataset is formatted as a multiple-choice visual question and is stored in a .tsv file. Each entry includes the following fields: *index, question, answer, A, B, C, D, image, category, l2 category, and l3 category.* These fields represent the instance index, the question associated with the image, the correct answer, the four answer options (A, B, C, D), the base64-encoded image, the specific type of insurance task, the meta-task name, and the basic task name, respectively.

A.3 HOSTING PLAN

INS-MMBench will be hosted on Github for the foreseeable future. Should that change, we will create an issue on the GitHub repository.

B DATASHEET

B.1 MOTIVATION

For what purpose was the dataset created?

The purpose of creating INS-MMBench was to comprehensively evaluate the capabilities of LVLMs
in the insurance field. It includes 12 meta-tasks and 22 fundamental tasks. It is the first comprehensive
multimodal benchmark in the insurance field.

B.2 COMPOSITION

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?

The instance of INS-MMBench represents a multiple-choice visual question, including the instance index, the question corresponding to the image, the correct answer, four options A, B, C, and D, the base64 encoding of the image, the insurance type to which the task belongs, the meta-task name, and the basic task name.

B46 Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?

Given the expense associated with LVLMs and the need for convenient model evaluation, the image data for each of our core tasks was randomly selected from a larger dataset.

850851 What data does each instance consist of?

Each instance in the INS-MMBench dataset includes the following components: *index, question, answer, A, B, C, D, image, category, l2 category, and l3 category.* These components correspond
respectively to the instance index, the question related to the image, the correct answer, four answer
options (A, B, C, and D), the base64-encoded image, the type of insurance related to the task, the
meta-task name, and the basic task name.

Is there a label or target associated with each instance?

858 859 Yes.

860 Is any information missing from individual instances?

861 No.

Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?

N/A.	
Are t	here recommended data splits (e.g., training, development/validation, testing)?
No.	
Are t	here any errors, sources of noise, or redundancies in the dataset?
The a	uthors are not aware of any errors in the dataset.
Does tecte indiv	the dataset contain data that might be considered confidential (e.g., data that is pro- d by legal privilege or by doctor-patient confidentiality, data that includes the content of iduals' non-public communications)?
No.	
Does or mi	the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, ight otherwise cause anxiety?
No.	
Does	the dataset relate to people?
No.	
B.3	COLLECTION PROCESS
How	was the data associated with each instance acquired?
First, rando effort multi	we collected datasets from multiple public sources. Next, we performed manual filtering and om sampling of the datasets, followed by the necessary data processing. Finally, both manual and GPT-40 were utilized to construct task questions and multiple-choice options, creating a -choice visual question dataset.
Wha senso	t mechanisms or procedures were used to collect the data (e.g., hardware apparatus or or, manual human curation, software program, software API)?
We u	sed both manual human curation and the GPT-40 API.
If the prob	e dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, abilistic with specific sampling probabilities)?
We en	mployed the uniform sampling according to different label distributions.
Who how	was involved in the data collection process (e.g., students, crowdworkers, contractors) and were they compensated (e.g., how much were crowdworkers paid)?
We di	id not employ external crowdworkers or contractors to construct INS-MMBench.
Over	what timeframe was the data collected?
We co	ompleted data collection in May 2024.
Were	any ethical review processes conducted (e.g., by an institutional review board)?
No.	
B.4	PREPROCESSING/CLEANING/LABELING
Was toker of mi	any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, nization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing ssing values)?
Yes, assist	we constructed questions and options for different images using both manual and GPT-40 ance, and evenly distributed the answers.
B.5	Uses

Has the dataset been used for any tasks already?

918	Yes, we have evaluated 10 LVLMs on this dataset.
919	Is there a repository that links to any or all papers or systems that use the dataset?
921	No
922	
923	What (other) tasks could the dataset be used for?
924 925	Beyond assessing the capabilities of LVLMs in the insurance sector, this dataset can also serve as fine-tuning training data to enhance LVLM performance.
926 927	Is there anything about the composition of the dataset or the way it was collected and prepro- cessed/cleaned/labeled that might impact future uses?
928 929 930	The task scope and number of instances in this dataset are limited. However, future usage and research could expand and enhance based on this foundation.
930 931	Are there tasks for which the dataset should not be used?
932	
933	No.
934 935	B.6 DISTRIBUTION
936	Will the dataset be distributed to third parties outside of the entity (e.g., company, institu-
937	tion,organization) on behalf of which the dataset was created?
938	Yes.
939	How will the detect will be distributed (e.g. terbell on website ADI CitHub)?
940 941	now will the dataset will be distributed (e.g., tarball off website, AP1, Github)?
942	Ins-MMBench will be distributed at the Github repository.
943	When will the dataset be distributed?
944	The dataset will be made available upon acceptance.
945	
946	
947	
948	
949	
951	
952	
953	
954	
955	
956	
957	
958	
960	
961	
962	
963	
964	
965	
966	
967	
960	
970	
971	

C CASES

To offer a detailed view of the task settings in INS-MMBench, we have selected sample cases for each core task and present responses from GPT-40, Gemini 1.5 Flash, and Qwen-VL-Max in this section.

978		
979		
980	License plate recognition	
981	License place recognition	
982	Prompt [.]	
983	What is the license plate number of the vehicle in	
984	the image?	
985	Options:	
986	A. NEWJOB	
987	C. N3WJ0B	
988	D. MEWJOB	
989	Correct answer [.] A	
990		
991		
992		
002		
00/	The license plate number of the vehicle in the image is:	
005	A. NEWJOB	
006		
990		
002		
000	Cemini 1 5 Flash	
1000	The ensurer is A. The license glate number is NEWIOD	
1000	The answer is A. The license plate number is NEWJOB.	
1001		
1002		
1003		
1004		
1005	Owen-VL-Max	
1007	The license plate number of the vehicle in the image is "NEWIOB"	
1007	The needse place number of the vehicle in the image is TVE wood .	
1000		
1010		
1011		
1012		
1012		
1013		
1014		
1015		
1017		
1012		
1010		
1013		
1020		
1021	Figure 5: A sample case of license plate recognition and the responses of GDT 4. Comini \mathbf{ProV}	icion
1022	and Owen-VI -Max	131011
1023		
1024		





- 1132 GeminiProVision and Qwen-VL-Max.
- 1133









Figure 10: A sample case of incar driving behavior detection and the responses of GPT-40, GeminiProVision and Qwen-VL-Max.





Figure 12: A sample case of vehicle damage type detection and the responses of GPT-40, GeminiPro-Vision and Qwen-VL-Max.

1402 1403

Vehicle damage severity detection

Prompt:

What is the damaged severity level of the vehicle in the image? Options:

A. severe

B. moderate

C. no_accident

D. minor

minor



Correct answer: D

GPT-40:

The correct answer is:

D. minor

The vehicle in the image has some visible damage, particularly around the front left fender and headlight area, but it does not appear to be severe.

Gemini 1.5 Flash:

The answer is B. moderate.

The vehicle has visible damage to the front bumper and fender, but the overall structure of the vehicle appears intact. This suggests moderate damage.

Qwen-VL-Max:

The correct answer is D. minor. Based on the image, there appears to be some damage to the front bumper and possibly the headlights of the white car. However, the overall damage seems to be relatively minor compared to more severe accidents or collisions. Therefore, option D, minor, would be the most appropriate description for the damage severity level of the vehicle in the image.

Figure 13: A sample case of vehicle damage severity detection and the responses of GPT-40,GeminiProVision and Qwen-VL-Max.





Prompt:

What's the risk condition of the workplace in the

image? Options:

A. There is risk of not wearing helmet in the image

B. There is risk of not wearing helmet and vest in the imageC. There is risk of not wearing vest in the image

D. There is no risk in the image



Correct answer: C

GPT-4o:

A. There is risk of not wearing helmet in the image

Gemini 1.5 Flash:

There is risk of not wearing vest in the image

Qwen-VL-Max:

C. There is risk of not wearing vest in the image. In the image, one worker is seen wearing a hard hat and helping another worker who appears to be injured. However, neither of them is wearing a safety vest, which poses a risk on the construction site. Therefore, option C is the correct answer.

Figure 15: A sample case of workplace risk assessment and the responses of GPT-40, GeminiProVision and Qwen-VL-Max.



House damage type detection

Prompt:

 What type of damage is done to the house in the image? Options: A. Combined Damage B. Flexural C. Shear



Correct answer: A

GPT-4o:

The damage shown in the image appears to be characterized by diagonal cracks and a pattern that suggests a sliding or shearing motion. This type of damage is typically associated with shear forces. The correct answer is:

C. Shear

Gemini 1.5 Flash:

The answer is A. Combined Damage.

The image shows cracks in the wall, which are a sign of both flexural and shear damage. Flexural damage occurs when a material is bent or deformed, while shear damage occurs when a material is subjected to a force that is parallel to its surface. In this case, the cracks are likely caused by a combination of these two types of damage.

Qwen-VL-Max:

The correct answer is A. Combined Damage.\n\nBased on the image provided, it appears that there is combined damage to the house's exterior wall. The wall shows signs of cracking and peeling, which could be due to a combination of factors such as age, weathering, or structural issues. The presence of a metal mesh suggests that repairs have been attempted, but the damage still remains visible. Therefore, option A, Combined Damage, best describes the type of damage seen in the image.

Figure 17: A sample case of hosue damage type detection and the responses of GPT-40, GeminiPro Vision and Qwen-VL-Max.

House damage level detection

Prompt:

What is the level of damage to the house in the image? Options: A. Undamage

- B. Minor Damage
- C. Moderate Damage
- D. Heavy Damage
- J. meary Dumage



GPT-40:

The correct answer is:

C. Moderate Damage

The image shows significant cracks and peeling on the exterior wall, indicating structural issues that are more than minor but not extensive enough to be classified as heavy damage.

Gemini 1.5 Flash:

The answer is C. Moderate Damage.

The image shows a large crack running across the wall of the house, indicating a significant structural issue. However, the house is still standing and the crack doesn't appear to have caused a complete collapse. Therefore, it's not heavy damage, but it's more than just minor damage.

Qwen-VL-Max:

The correct answer based on the image provided is:\n\nC. Moderate Damage\n\nThe image shows a pink building with visible cracks and peeling paint, indicating some degree of structural damage or wear over time. The extent of the damage appears to be moderate rather than minor or heavy.

Figure 18: A sample case of hosue damage level detection and the responses of GPT-40, GeminiPro Vision and Qwen-VL-Max.





Medical image abnormality recognition

Prompt:

what is the primary abnormality in this image? Options:

- A. Osteomyelitis
- B. Osteogenesis Imperfecta
- C. Osteoarthritis
- D. Osteosarcoma



Correct answer: B

GPT-40:

The primary abnormality in the image is consistent with Osteogenesis Imperfecta (Option B). This condition is characterized by bones that break easily, often with little or no apparent cause. The X-ray shows multiple fractures and bone deformities, which are typical features of Osteogenesis Imperfecta.

Gemini 1.5 Flash:

The answer is D. Osteosarcoma.

Osteosarcoma is a type of bone cancer that typically affects children and young adults. It is characterized by the formation of malignant tumors in the bone. The image shows a large, destructive tumor in the humerus, which is a classic finding of osteosarcoma.

Qwen-VL-Max: B. Osteogenesis Imperfecta

Figure 21: A sample case of medical image abnormality detection and the responses of GPT-40, GeminiProVision and Qwen-VL-Max.

Medical image part recognition

Prompt:

what is the organ system in this image? Options:

A. Kidney, Ureter, Bladder

B. Liver, Gallbladder, Pancreas

- C. Stomach, Duodenum, Jejunum
- D. Lung, Mediastinum, Pleura



Correct answer: D

GPT-40:

The organ system in this image is the brain, which is part of the central nervous system. Therefore, none of the provided options are correct.

Gemini 1.5 Flash:

The answer is D. Lung, Mediastinum, Pleura. This is a brain scan, which is part of the nervous system.

Qwen-VL-Max:

The correct answer is B. Liver, Gallbladder, Pancreas. In the provided image, we can see an MRI scan of the liver, gallbladder, and pancreas. The image shows these organs in detail, allowing for evaluation of their structure and potential abnormalities.

Figure 22: A sample case of medical image part recognition and the responses of GPT-40, GeminiProVision and Qwen-VL-Max.



Satellite image crop type identification

Prompt:

What type of crops are in the image? Options: A. Barley B. Sorghum

- C. Wheat
- D. Maize



Correct answer: D

GPT-40:

The image appears to show a close-up of a crop, but the details are not clear enough to definitively identify the type of crop. However, based on the general appearance, it resembles the structure of wheat. So, the correct answer is:

C. Wheat

Gemini 1.5 Flash:

The answer is C. Wheat.

Wheat crops are characterized by their long, slender leaves and tall stalks. The image shows these features, making it likely that the crops are wheat.

Qwen-VL-Max:

Based on the image provided, it is difficult to determine the exact type of crop as the image is blurry and lacks detail. However, based on the given options, the most likely answer would be C. Wheat, as it is a common cereal grain that is often grown in large fields and has a similar appearance to the作物 in the image. However, without more information or a clearer image, it is impossible to确定ly identify the specific type of crop.

Figure 24: A sample case of satellite image crop type identification and the responses of GPT-40, GeminiProVision and Qwen-VL-Max.

Crop growth stage recognition

Prompt:

What is the growth stage of the crops in the image? Options:

A. Booting

- B. Ripening
- C. Tillering
- D. Flowering

Correct answer: A

GPT-40:

The growth stage of the crops in the image appears to be "C. Tillering." This stage is characterized by the development of multiple shoots from the base of the plant, which is evident in the dense and lush green appearance of the crop.

Gemini 1.5 Flash:

The answer is C. Tillering.

The crops in the image are in the tillering stage. This is the stage where the plant is producing new shoots from the base of the stem. The shoots are called tillers. The tillering stage is important for the plant because it allows it to produce more leaves and stems, which will eventually produce more grain.

Qwen-VL-Max:

Based on the image, it appears that the crops are at the tillering stage, which is when the plants are developing side shoots and branching out to form a dense canopy. Therefore, the correct answer is C. Tillering.

Figure 25: A sample case of crop growth stage recognition and the responses of GPT-40, GeminiPro-Vision and Qwen-VL-Max.





Prompt:

What type of damage has occurred on the
farmland?
Options:
A. Double plant
B. Lodging
C. Erosion
D. Pest infestation

Correct answer: A

GPT-40:

B. Lodging

Gemini 1.5 Flash:

The answer is C. Erosion.

Erosion is the process of soil and rock being moved from one place to another by wind, water, or ice. The image shows a pattern of vertical lines, which are likely caused by erosion.

Qwen-VL-Max:

Based on the image provided, it appears that there is lodging damage on the farmland. The plants appear to be leaning or fallen over, which is characteristic of lodging damage caused by strong winds or other factors. Therefore, the correct answer is B. Lodging.

Figure 26: A sample case of farmland damage detection and the responses of GPT-40, GeminiProVision and Qwen-VL-Max.