**Step 1: Image-Instruction Pair Collection**

Instruction

Analyze: What is unusual about this picture?

Images

Random Sample

What is the Area of CHD?

Describe this image.

Someone said this man is an angel, why?

WORKING FROM HOME IS WORKING

77% of employees surveyed are currently working from home.

63% now realize their job is doable outside the office.
60% say their work-life balance has improved without a commute.
43% are more comfortable using technology.
20% report stronger relationships with colleagues.

What percentage of workers are not working from home?

Image-Instruction Pairs

**Step 2: MLLM Response Collection**

MLLM response

The image depicts a rainy night in a bustling city, with people …

To determine the number of people who lived in ….

…

The largest bar in the figure has a value of 90. …

LLaVA    CogVLM

KEG

Gemini    GPT-4V

Response MLLM

**Step 3: MLLM Judge v.s. Human Annotation**

Scoring evaluation

**Assistant A:**
The answer is….

**Judgement: 4**

Pair comparison

**Assistant A:**
The number is….

**Assistant B:**
As for the number….

**Judgement: B**

Batch ranking

**Assistant A:**
The year is….

**Assistant B:**
As for the year….

**Assistant C:**
The answer is...

**Assistant D:**
The year is ….

**Judgement: CBAD**

Judge MLLM

Human Annotation