

QWEN 2.5VL

.....

UNIVERSAL RECOGNITION



.....

What are these attractions? Please give their names in Chinese and English.

Qwen2.5-VL

These attractions are:

Top left: **The Great Pyramids of Giza**
(埃及的吉萨大金字塔)

Top right: **The Great Wall of China**
(中国的长城)

Bottom left: **The Statue of Liberty**
(美国的自由女神像)

Bottom right: **The Terracotta Army**
(中国的秦始皇兵马俑)



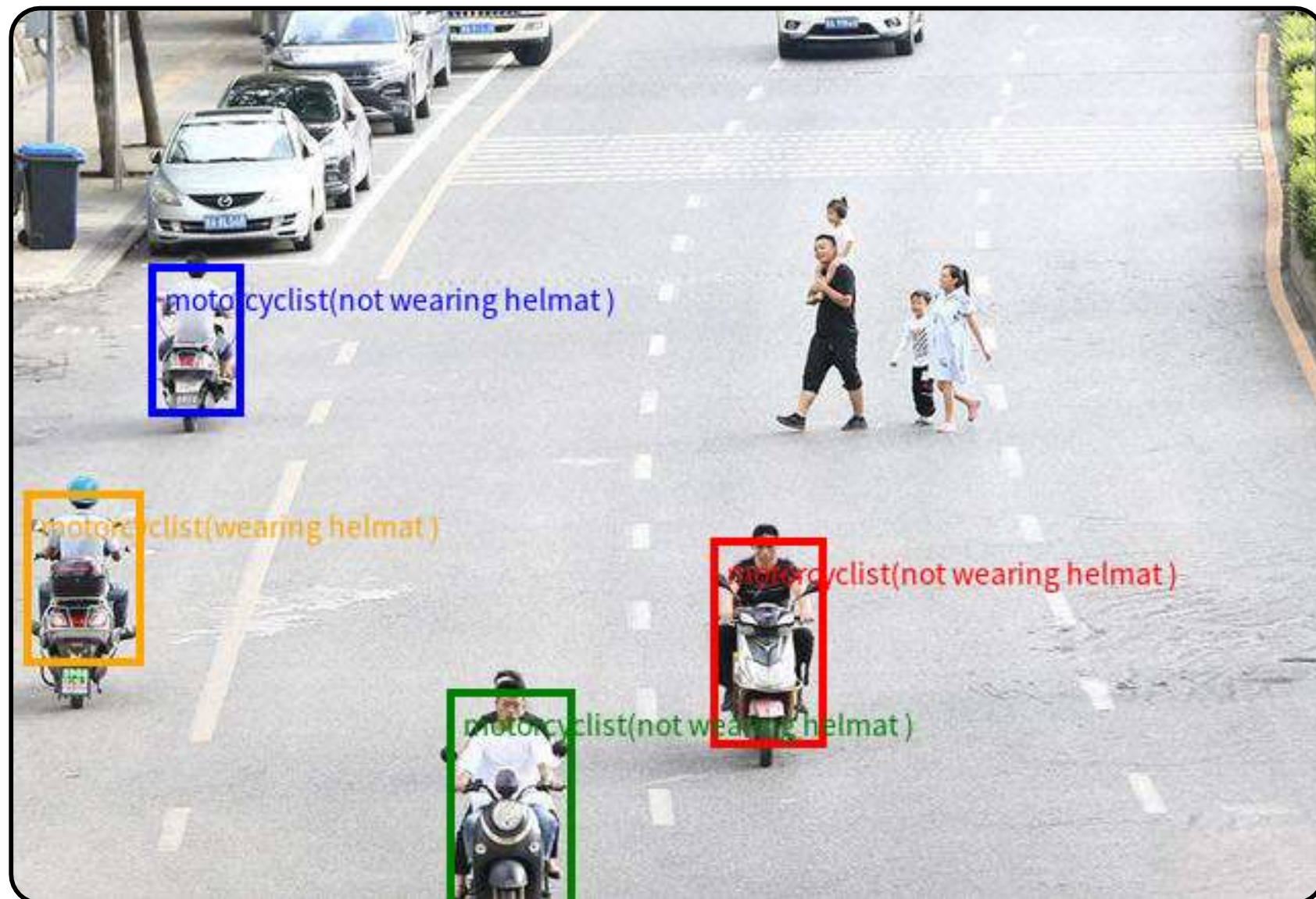
OBJECT DETECTION



.....

Detect all motorcyclists in the image and return their locations in the form of coordinates.

The format of output should be like {"bbox_2d": [x1, y1, x2, y2], "label": "motorcyclist", "sub_label": "wearing helmet" # or "not wearing helmet"}.



```
[  
    {"bbox_2d": [341. 258. 397. 360], "label": "motorcyclist",  
     "sub_label": "not wearing helmet"},  
    {"bbox_2d": [212. 332. 274. 448], "label": "motorcyclist",  
     "sub_label": "not wearing helmet"},  
    {"bbox_2d": [66. 124. 112. 198], "label": "motorcyclist",  
     "sub_label": "not wearing helmet"},  
    {"bbox_2d": [5. 235. 63. 320], "label": "motorcyclist",  
     "sub_label": "wearing helmet"}  
]
```

OBJECT DETECTION



.....



Locate the person who act bravely, report the bbox coordinates in JSON format.

```
[  
    {"bbox_2d": [74, 58, 526, 619], "label": "person who act  
bravely"}  
]
```


● ● ● ● ● ● ● ● ● ●

arXiv:2412.15115v2 [cs.CL] 3 Jan 2025

```
<address class="author"><p>Qwen Team</p></address>
```

<p data-bbox="885 491 1460 541"><https://huggingface.co/Qwen></p>

<p data-bbox="885 541 1708 589"><https://modelscope.cn/organization/qwen></p>

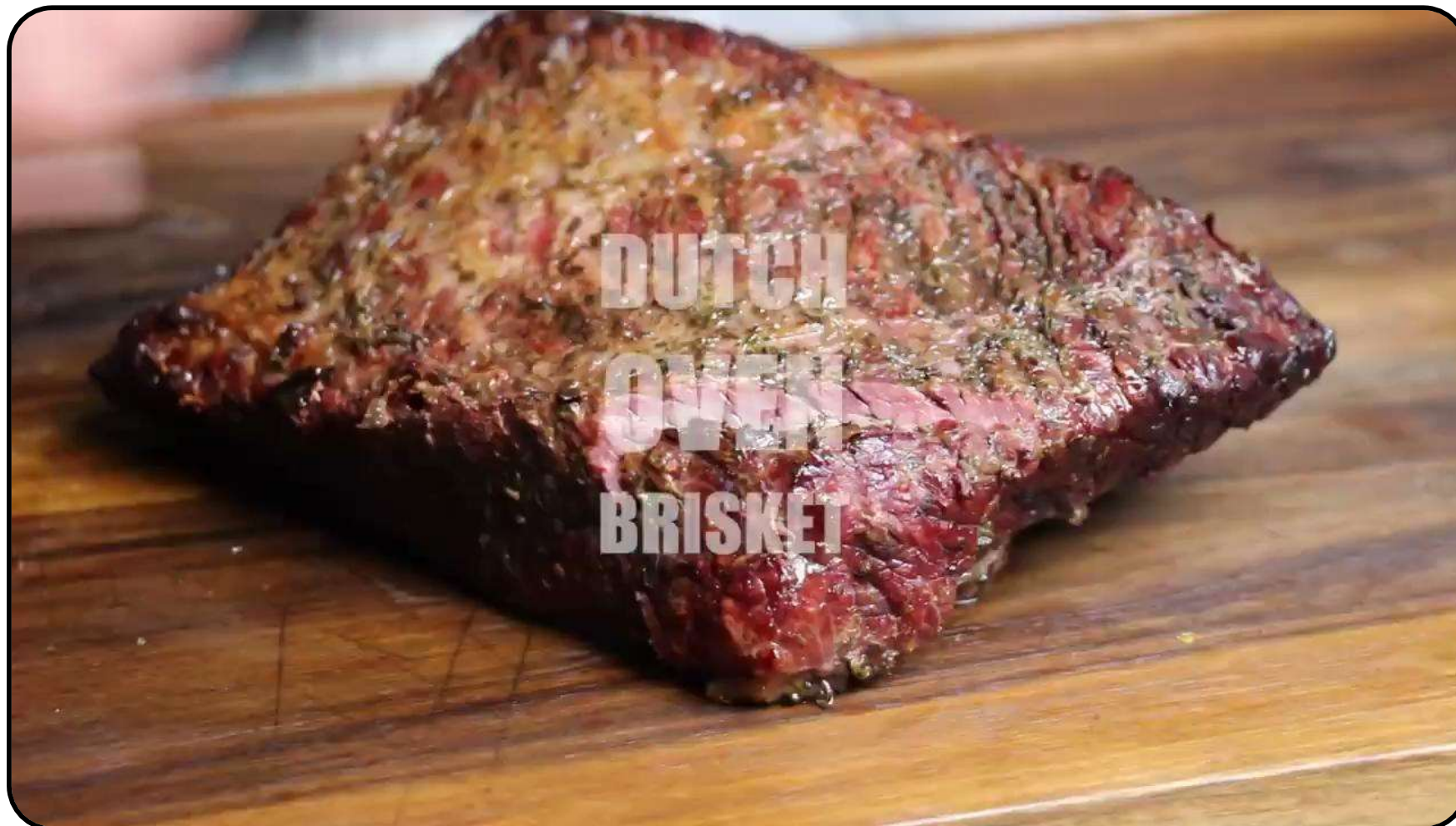
<p data-bbox="885 589 1584 639"><https://github.com/QwenLM/Qwen2.5></p>

Abstract

LONG VIDEO ANALYSIS



.....



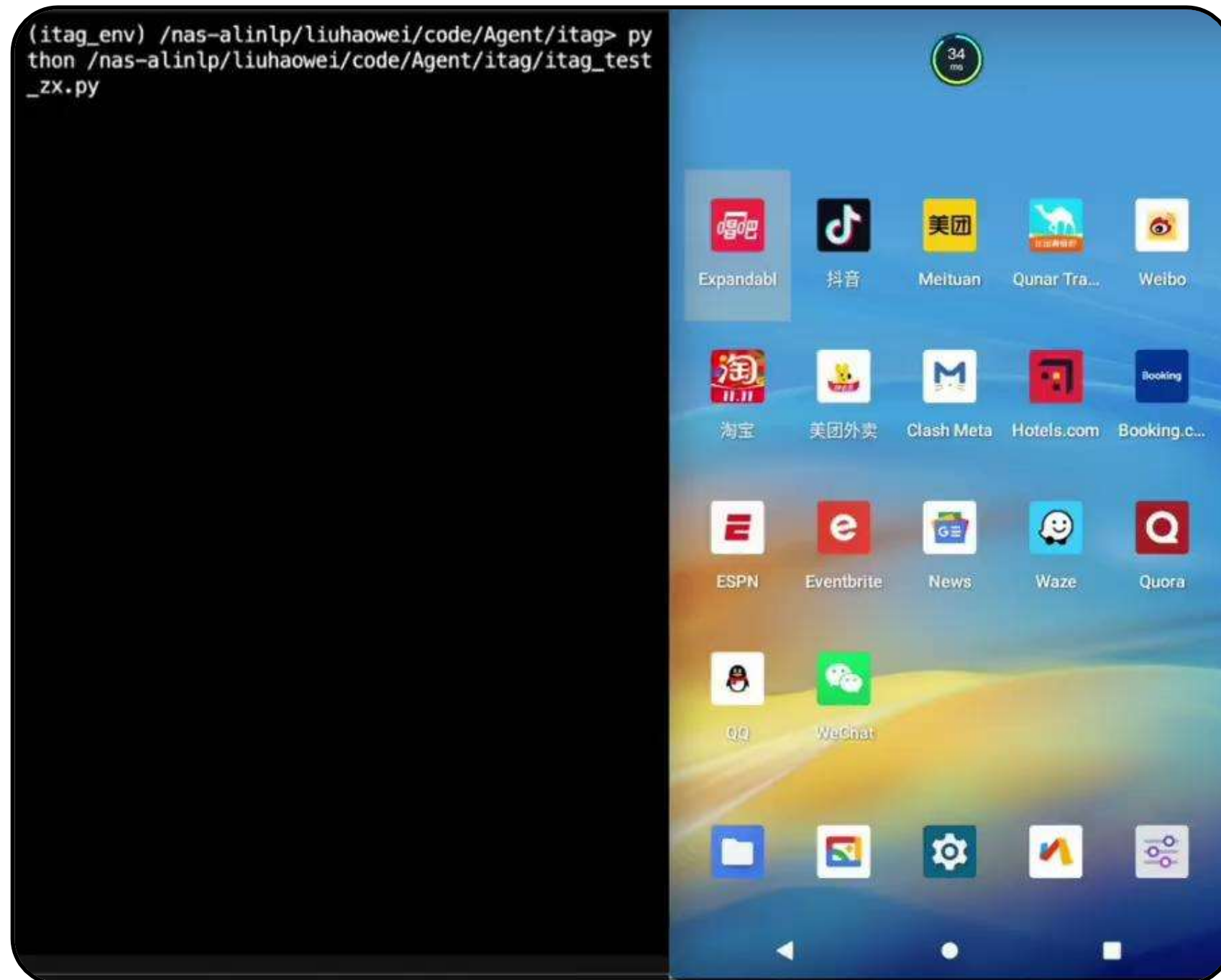
Localize a series of activity events in the video, output the start and end timestamp for each event, and describe each event with sentences. Provide the result in JSON.

```
[
  {
    "start_time": "00:21.00",
    "end_time": "00:49.00",
    "description": "A person removes a piece of meat from its packaging and cuts off the fat."
  },
  {
    "start_time": "00:50.00",
    "end_time": "01:27.00",
    "description": "The person seasons the meat with salt and pepper on both sides."
  },
  {
    "start_time": "01:28.00",
    "end_time": "01:36.00",
    "description": "The seasoned meat is placed on a grill."
  }
]
```




AGENT ACTIONS

.....



Please help me book a one-way ticket in the bookings app. The starting point is Chongqing Jiangbei Airport and the ending point is Beijing Capital Airport on January 28th



Usage

.....

Qwen 2.5VL Chat Interface



.....

```
chat_response = client.chat.completions.create(  
    model="Qwen/Qwen2.5-VL-7B-Instruct",  
    messages=[  
        {"role": "system",  
         "content": "You are a helpful assistant."},  
        {"role": "user",  
         "content": [  
             {  
                 "type": "image_url",  
                 "image_url": {"url": "https://example.com/image.png"}  
             },  
             {  
                 "type": "text",  
                 "text": "What is the text in the image?"  
             }  
         ]  
        },  
    ],  
)
```

QWEN 2.5VL: ROLE

.....



SYSTEM: Sets global behaviour

```
{  
  "role": "system",  
  "content": "You are an assistant that answers briefly, returns currency values in USD, "  
    + "and uses JSON when a user asks for structured data."  
}
```


QWEN 2.5VL: ROLE

.....



USER: Sends an image and text

```
{
  "role": "user",
  "content": [
    { "type": "image",
      "image": "https://example.com/receipt-coffee.jpg" },
    { "type": "text",
      "text": "How much did the latte cost? Please give the answer as JSON." }
  ]
}
```



MODEL CARD

.....

Name	Qwen/Qwen2.5-VL-3B-Instruct
Summary	Multi-modal LLM capable of image, video and tool-augmented chat.
Parameters	3 billion
Release Date	April 2025
Developer	Alibaba Cloud (Qwen team)
License	Qwen license (research/non-commercial); 3 B & 72 B sizes are not Apache 2.0.

WILL IT FIT?

.....



Model size	Precision	VRAM needed*	Consumer cards that clear it
Qwen 2.5-VL-3B	8-bit / BF16	6 – 8 GB	RTX 3060 12 GB, RTX 4060 8–16 GB, laptop 4070
Qwen 2.5-VL-7B	8-bit	12 – 14 GB	RTX 4070 12 GB (tight), RTX 4070 Ti 16 GB
	BF16	18 – 20 GB	RTX 3090 / 4090 24 GB
Qwen 2.5-VL-32B / 72B	Any	≥ 32 GB (32B) / ≥ 64 GB (72B) or multi-GPU	Prosumer/enterprise (RTX 6000 Ada 48 GB, dual 4090s, etc.)



IMAGE CAPTIONING

.....

INPUT



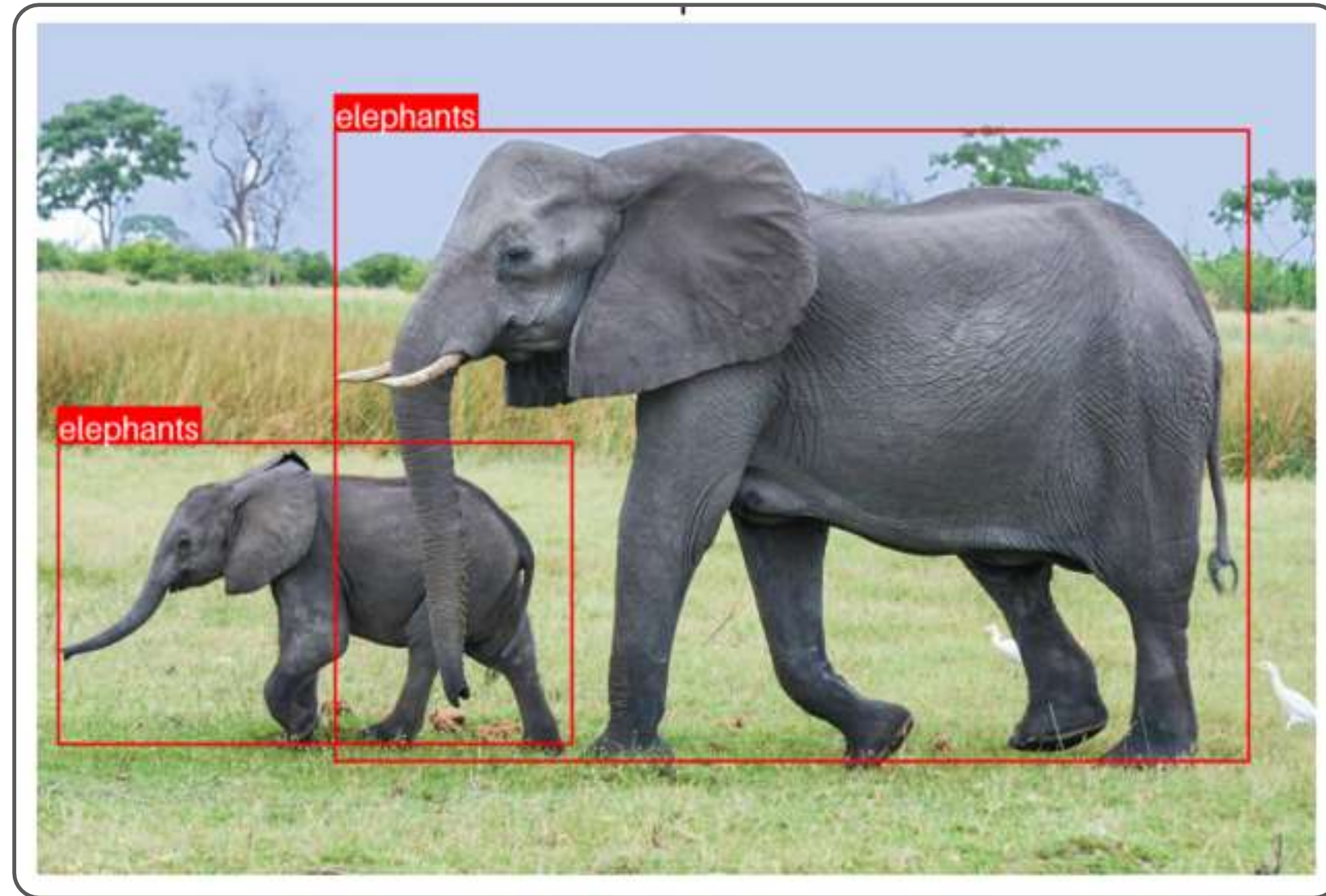
OUTPUT

The image depicts a serene and picturesque scene of a white dog sitting on a stone pathway near a stunning lake. The lake has crystal-clear turquoise water that reflects the surrounding landscape, including the mountains in the background. The mountains are covered with lush greenery, and the sky above is partly cloudy, adding to the tranquil atmosphere.

OBJECT DETECTION & SPATIAL REASONING APP



.....





Thank You

.....