# Left Panel

**Input files**



AgentX_181.jpg

**Query**

Which model performs the best in Visual Knowledge Acquisition, and what is the total number of distinct models shown in the radar chart?

## Reasoning steps: 5

**Tools used**

**STEP 1**
- TASK: Describe the image to understand the content and context.
- THOUGHT: Understanding the image structure helps guide the following steps to identify which model performs best and count the number of distinct models.

SCENE DESCRIBER

**STEP 2**
- TASK: Identify the model performing best in Visual Knowledge Acquisition.
- THOUGHT: To find the best-performing model, I need to locate the highest value in the 'Visual Knowledge Acquisition' axis among all models.

LOCATE OBJECTBYTEXT

**STEP 3**
- TASK: Overlay output to identify the best performer clearly.
- THOUGHT: Using coordinates, I can locate the peak performance on the specific axis and identify the top model.

OVERLAY TEXT

**STEP 4**
- TASK: Count the total number of distinct models represented in the chart.
- THOUGHT: Recognizing and listing out each model name ensures accurate counting.

OCR

**STEP 5**
- TASK: Compute the total number of distinct models.
- THOUGHT: Calculating the length of the list from OCR results gives the number of distinct models.

CALCULATOR

**GT Answer**

{'best_model': 'Bard', 'total_models': 12}

**Justification**

Bard was identified as the model performing best in Visual Knowledge Acquisition based on its position in the radar chart. The OCR tool was used to enumerate and count all distinct models, resulting in a total of 12.

# Right Panel

**Input files**



AgentX_629.mp4

**Query**

What place is shown in the video, and how many people are performing duties?

## Reasoning steps: 3

**Tools used**

**STEP 1**
- TASK: Identify the scene using the video.
- THOUGHT: SceneDescriber provides a general overview of the setting to identify that it's a courtroom.

SCENE DESCRIBER

**STEP 2**
- TASK: Determine the location of the video scene from text in the image.
- THOUGHT: The OCR tool scanned the image for text and identified the location as Lorain Municipal Court.

OCR

**STEP 3**
- TASK: Count the number of people performing duties in the video.
- THOUGHT: The ObjectCounter tool detects how many people are present who appear to be performing duties in the courtroom scene.

OBJECT COUNTER

**GT Answer**

Lorain Municipal Court, 3

**Justification**

Using the SceneDescriber and OCR tools, I identified the scene as Lorain Municipal Court. The ObjectCounter tool confirmed three people performing duties: a judge, security officer, and lawyer, while the rest appear to be seated as the audience.