**Text Guidance:** Are there **six** people?

Reference [Ans: Yes]     Generated [Ans: Yes]

(a)  Numerical Attribute

Is there a **red** hat?

Reference [Ans: No]     Generated [Ans: No]

(b)  Physical Attribute (Color)

Is there a **bottle**?

Reference [Ans: Yes]     Generated [Ans: Yes]

(c)  Object Existence

**Text Guidance:** Is the dog <u>under</u> the pool?

Reference [Ans: No]     Generated [Ans: No]

(d)  Spatial Relationship

Is this a **general store outdoor**?

Reference [Ans: Yes]     Generated [Ans: Yes]

(e)  Visual Attribute (Scene)

Is person **swinging** or **holding** a tennis racket?

Reference [Ans: swing]     Generated [Ans: swing]

(f)  Human-Object Interaction Reasoning

(g)  Diversity/Fidelity Comparison

Reference Image     I2T2I SDXL     Image Variation     Image Translation     Textual Inversion     **Hummingbird**

✓ Diversity ✗ Fidelity     ✓ Diversity ✗ Fidelity     ✗ Diversity ✓ Fidelity     ✓ Diversity ✗ Fidelity     ✓ Diversity ✓ Fidelity

<u>Text Guidance</u>: Is the clock above people?

<u>Context Description from MLLM</u>: The image shows a city street scene with a clock mounted on a building. The clock is positioned above the street level, where pedestrians are walking. The street is lined with shops and buildings, and there are people walking on the sidewalk. The clock is clearly visible and serves as a prominent feature in the image.