Figure 1: Illustration of a 5-iteration GenCeption procedure run on artwork images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.

Original/Seed Input Image $\mathbf{X}^{(0)}$

**Visual-Intensive** Group

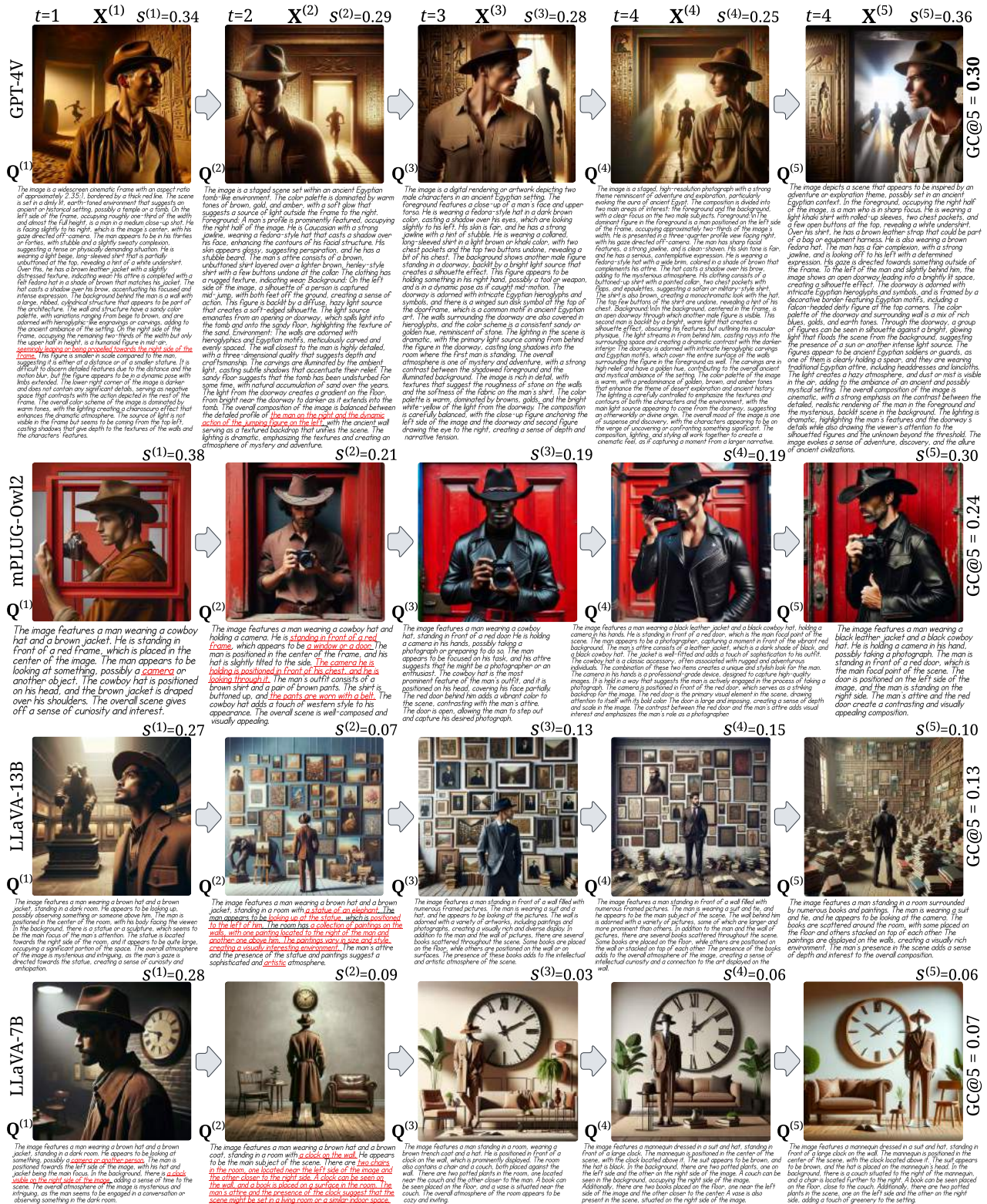**Celebrity** category: tt0082971_shot_0831_img_0.jpg

Figure 2: Illustration of a 5-iteration GenCeption procedure run on celebrity images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.

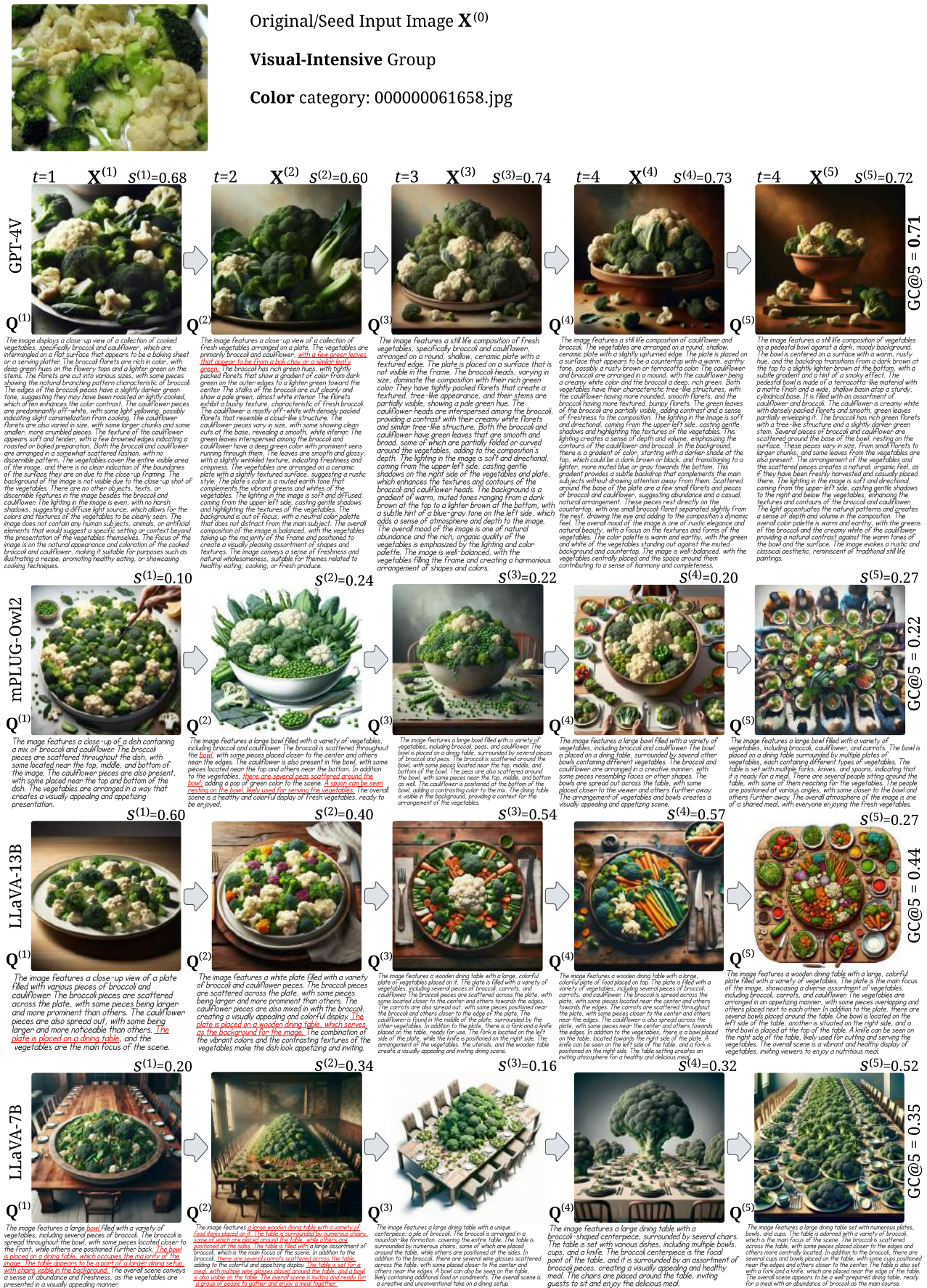Original/Seed Input Image $\mathbf{X}^{(0)}$

**Visual-Intensive** Group

**Color** category: 000000061658.jpg

$t=1$   $\mathbf{X}^{(1)}$   $s^{(1)}=0.68$    $t=2$   $\mathbf{X}^{(2)}$   $s^{(2)}=0.60$    $t=3$   $\mathbf{X}^{(3)}$   $s^{(3)}=0.74$    $t=4$   $\mathbf{X}^{(4)}$   $s^{(4)}=0.73$    $t=4$   $\mathbf{X}^{(5)}$   $s^{(5)}=0.72$

**GPT-4V** — GC@5 = 0.71

$\mathbf{Q}^{(1)}$ $\mathbf{Q}^{(2)}$ $\mathbf{Q}^{(3)}$ $\mathbf{Q}^{(4)}$ $\mathbf{Q}^{(5)}$

**mPLUG-Owl2** — $s^{(1)}=0.10$, $s^{(2)}=0.24$, $s^{(3)}=0.22$, $s^{(4)}=0.20$, $s^{(5)}=0.27$, GC@5 = 0.22

$\mathbf{Q}^{(1)}$ $\mathbf{Q}^{(2)}$ $\mathbf{Q}^{(3)}$ $\mathbf{Q}^{(4)}$ $\mathbf{Q}^{(5)}$

**LLaVA-13B** — $s^{(1)}=0.60$, $s^{(2)}=0.40$, $s^{(3)}=0.54$, $s^{(4)}=0.57$, $s^{(5)}=0.27$, GC@5 = 0.44

$\mathbf{Q}^{(1)}$ $\mathbf{Q}^{(2)}$ $\mathbf{Q}^{(3)}$ $\mathbf{Q}^{(4)}$ $\mathbf{Q}^{(5)}$

**LLaVA-7B** — $s^{(1)}=0.20$, $s^{(2)}=0.34$, $s^{(3)}=0.16$, $s^{(4)}=0.32$, $s^{(5)}=0.52$, GC@5 = 0.35

$\mathbf{Q}^{(1)}$ $\mathbf{Q}^{(2)}$ $\mathbf{Q}^{(3)}$ $\mathbf{Q}^{(4)}$ $\mathbf{Q}^{(5)}$

Figure 3: Illustration of a 5-iteration GenCeption procedure run on color images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.
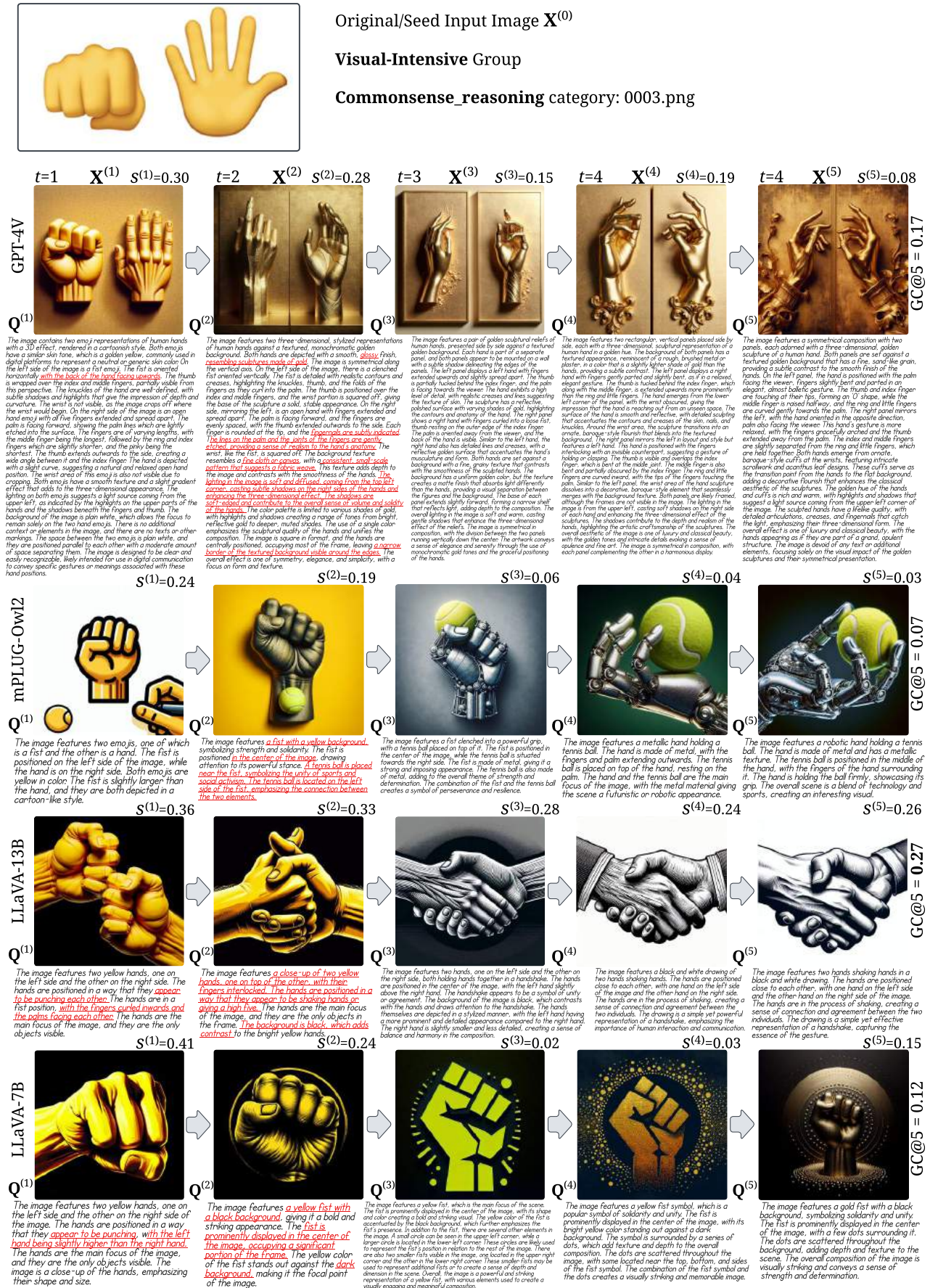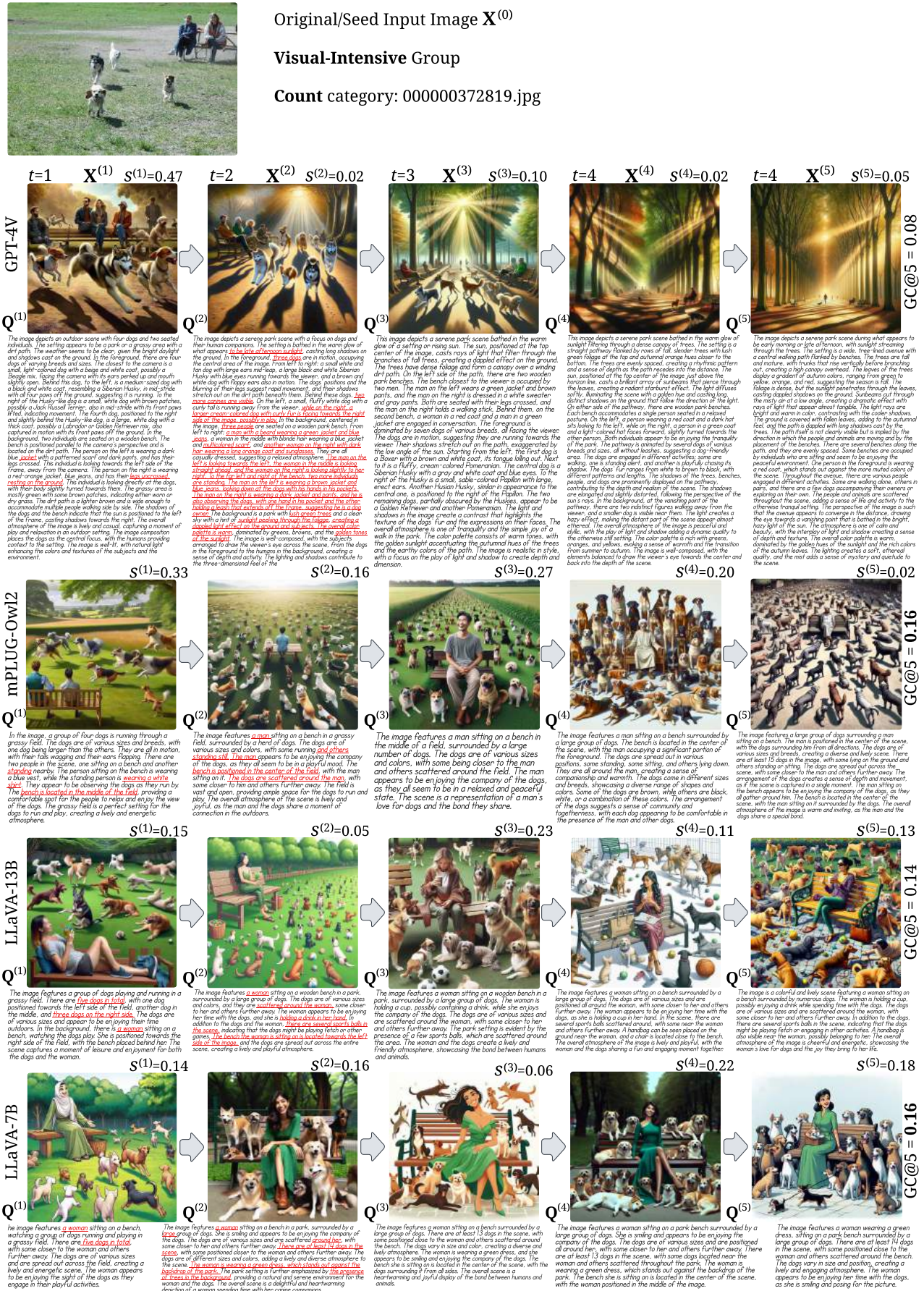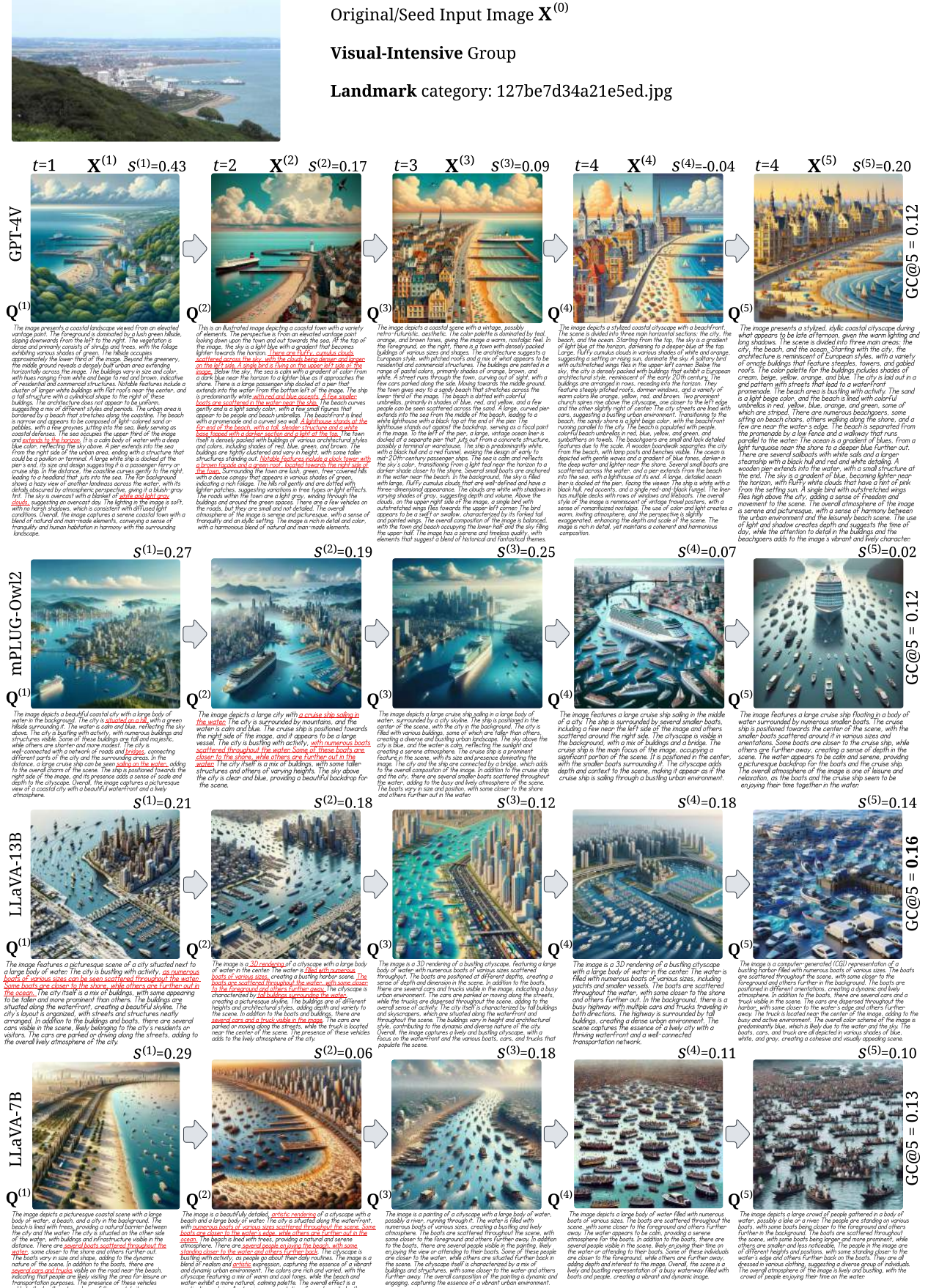
Figure 4: Illustration of a 5-iteration GenCeption procedure run on commonsense_reasoning images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.

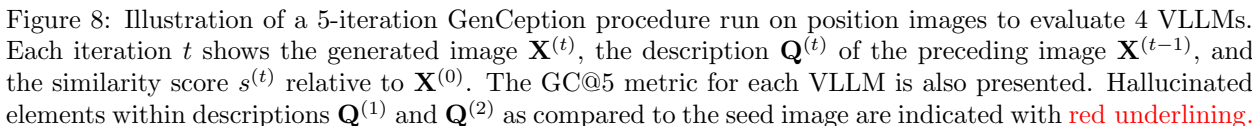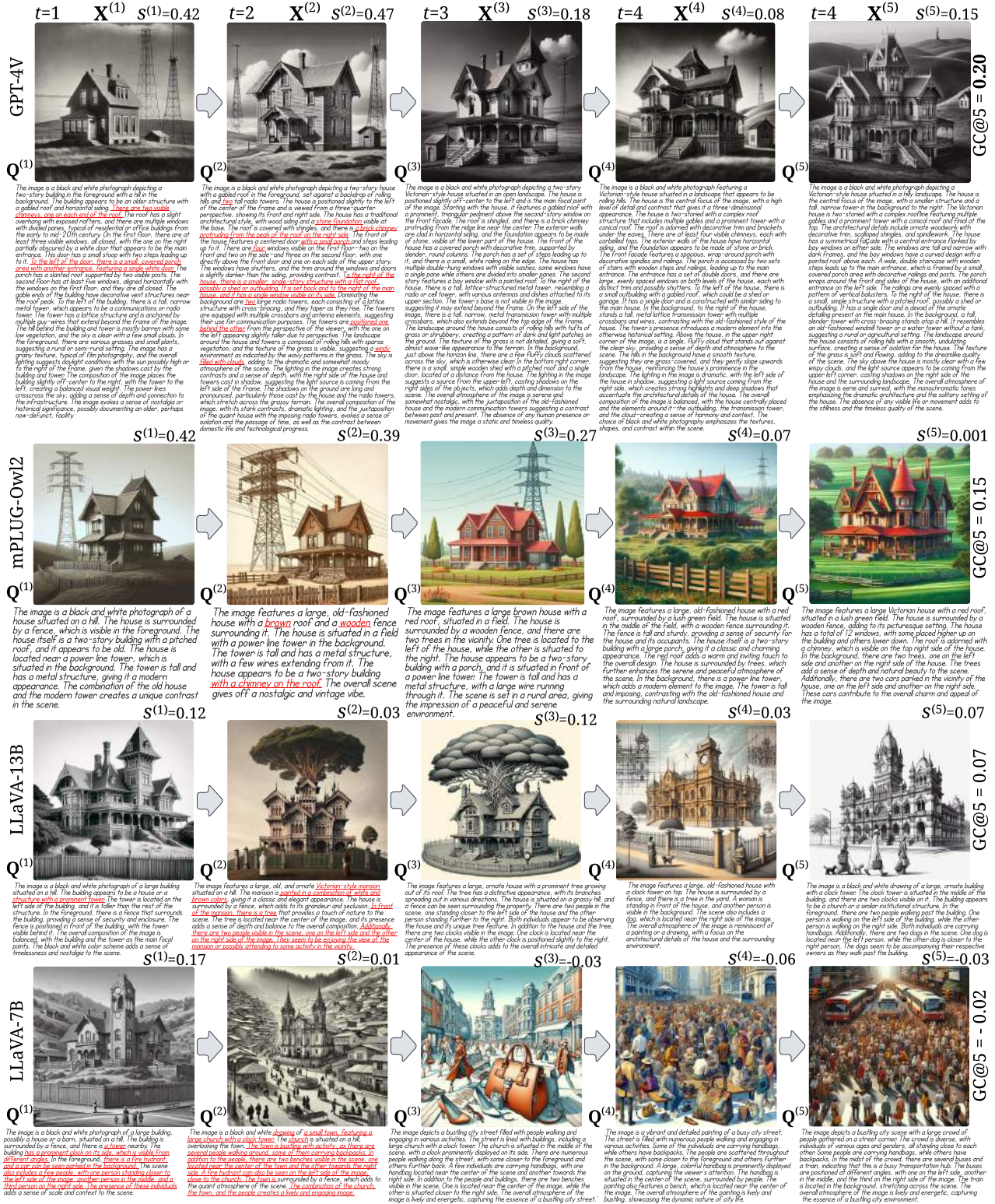Figure 5: Illustration of a 5-iteration GenCeption procedure run on count images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with <u>red underlining</u>.

Figure 6: Illustration of a 5-iteration GenCeption procedure run on existence images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.
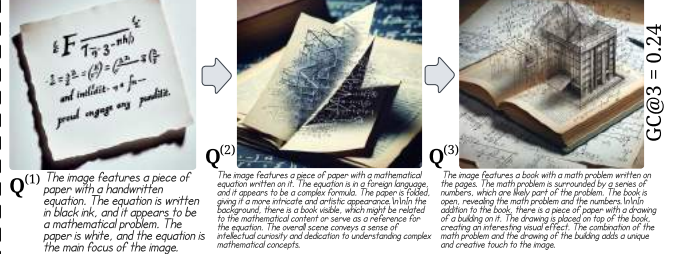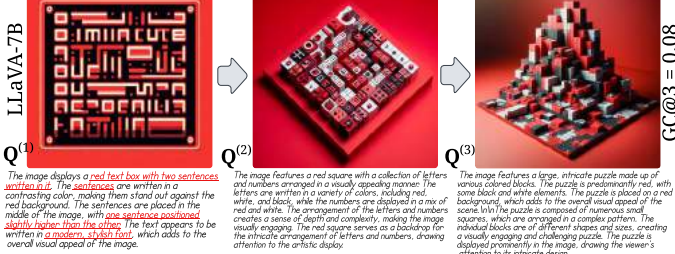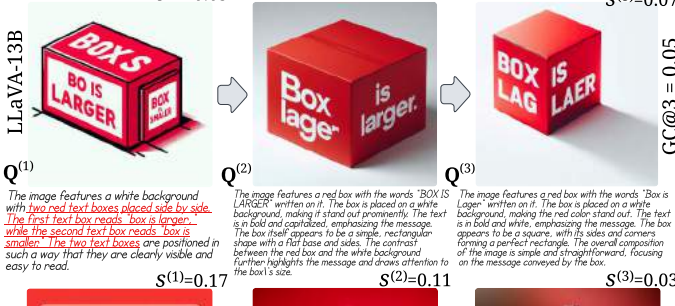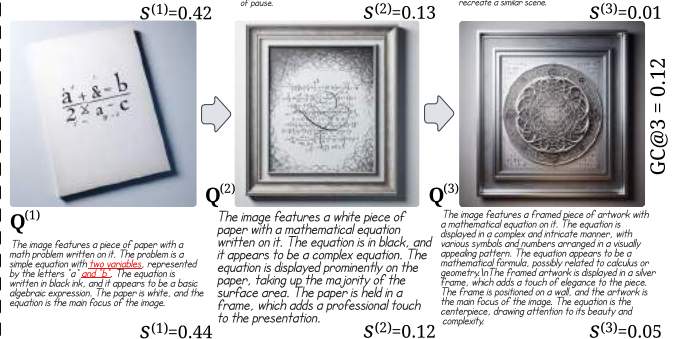
Figure 7: Illustration of a 5-iteration GenCeption procedure run on landmark images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.

Figure 8: Illustration of a 5-iteration GenCeption procedure run on position images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.

Figure 9: Illustration of a 5-iteration GenCeption procedure run on poster images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with red underlining.

Figure 10: Illustration of a 5-iteration GenCeption procedure run on scene images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@5 metric for each VLLM is also presented. Hallucinated elements within descriptions $\mathbf{Q}^{(1)}$ and $\mathbf{Q}^{(2)}$ as compared to the seed image are indicated with <u>red underlining.</u>

Figure 11: Illustration of a 3-iteration GenCeption procedure run on code_reasoing (left) and numerical_calculation (right) images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@3 metric for each VLLM is also presented. Hallucinated parts within descriptions $\mathbf{Q}^{(1)}$ as compared to the seed image are indicated with red underlining.

Figure 12: Illustration of a 3-iteration GenCeption procedure run on text_translation (left) and OCR (right) images to evaluate 4 VLLMs. Each iteration $t$ shows the generated image $\mathbf{X}^{(t)}$, the description $\mathbf{Q}^{(t)}$ of the preceding image $\mathbf{X}^{(t-1)}$, and the similarity score $s^{(t)}$ relative to $\mathbf{X}^{(0)}$. The GC@3 metric for each VLLM is also presented. Hallucinated parts within descriptions $\mathbf{Q}^{(1)}$ as compared to the seed image are indicated with red underlining.