

fully automated pipeline

10K real images from COCO / ImageNet + 10K high-quality AI-generated images
Region-level annotations with bounding boxes highlighting **telltale regions**
Forensic explanations auto-generated via VLM pipeline with consistency checks
Multiple generator sources including FLUX and Stable Diffusion variants
Quality control through automated filtering and human review

10K+10K
Real + AI-Generated

3 regions
Annotated per Image

Query 1: Explanation Generation

You are an expert in classifying AI-generated images and real images. Please carefully observe the image that the user provide, decide whether the image is real or AI-generated. Justify your verdict by listing reasons that support your decision.

Q1 - Response

This image is likely AI-generated due to the dog's unnatural, human-like pose and the seamless integration of text into the scene. Additional evidence includes inconsistent details on the Jenga blocks, overly uniform lighting, and the synthetic appearance of the background elements.

Query 2: Spatial Grounding

The user will show you a (label) image with explanations on why this image looks (label). Please annotate bounding boxes for key regions mentioned in the justification.

Q2 - Response

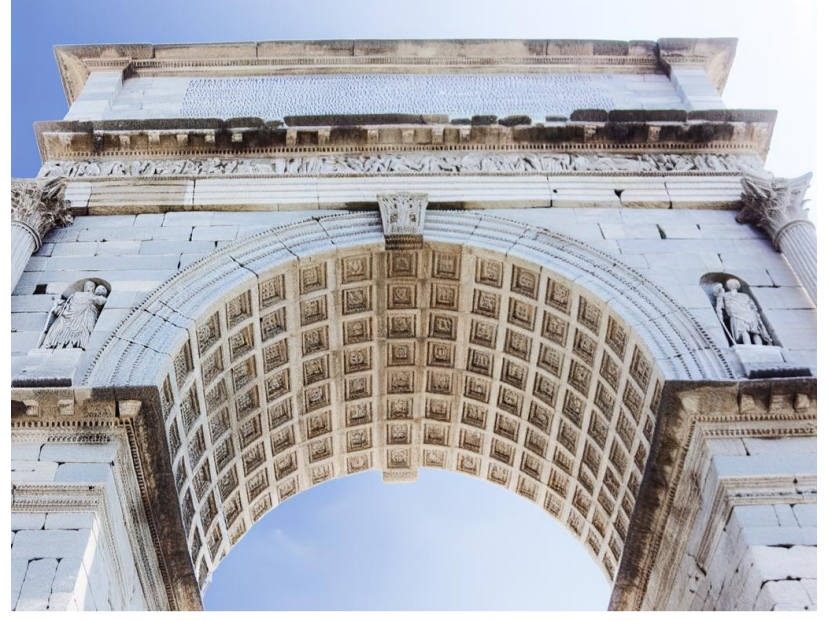
Based on your detailed analysis, here are some key regions...

- * Dog's Paw Interaction with Blocks (270, 398, 485, 630)
- * Jenga Blocks (-)

Roles: S System U User A Assistant

Locate THEN Examine

Please analyze the image and:
1. Provide an *initial verdict* on whether the image is **real** or **AI-generated**.
2. You will need to examine specific regions more closely. Output *bounding boxes* for areas that require *closer inspection*.



Vision Encoder



Image Generators are getting better, **one pass is not enough.**

Just like how humans perceive, we **Let models Look Closer.**

Initial Verdict is weak, but the refinement is strong.

I. Locate

Which specific region(s) in the image, when examined, can tell the authenticity of the image?

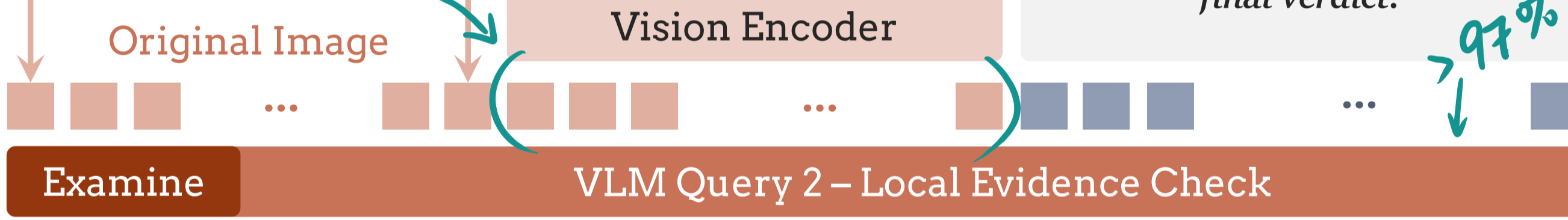
This step typically identifies **suspicious regions** such as facial features, hands, fine text, and complex textures where generators commonly fail. Consequently, correctly rendering these regions makes the image look more authentic.

II. Examine

After programmatically cropping the regions, models should now become more evident, effectively polishing their initial verdict.

Local evidence check will utilize the regional close-ups to detect subtle artifacts, or determine that the details are evident that the artwork is indeed realistic.

During training, we used SFT + pGRPO with IoU, BLEU, Accuracy, Formatting as rewards. We only take final verdict's accuracy into account.



LTE Refined

from **Real to Fake**

- The fan hanging above has a wrong geometry
- Title mismatch; the towel should not reflect
- Extra claw; missing eye

from **Fake to Real**

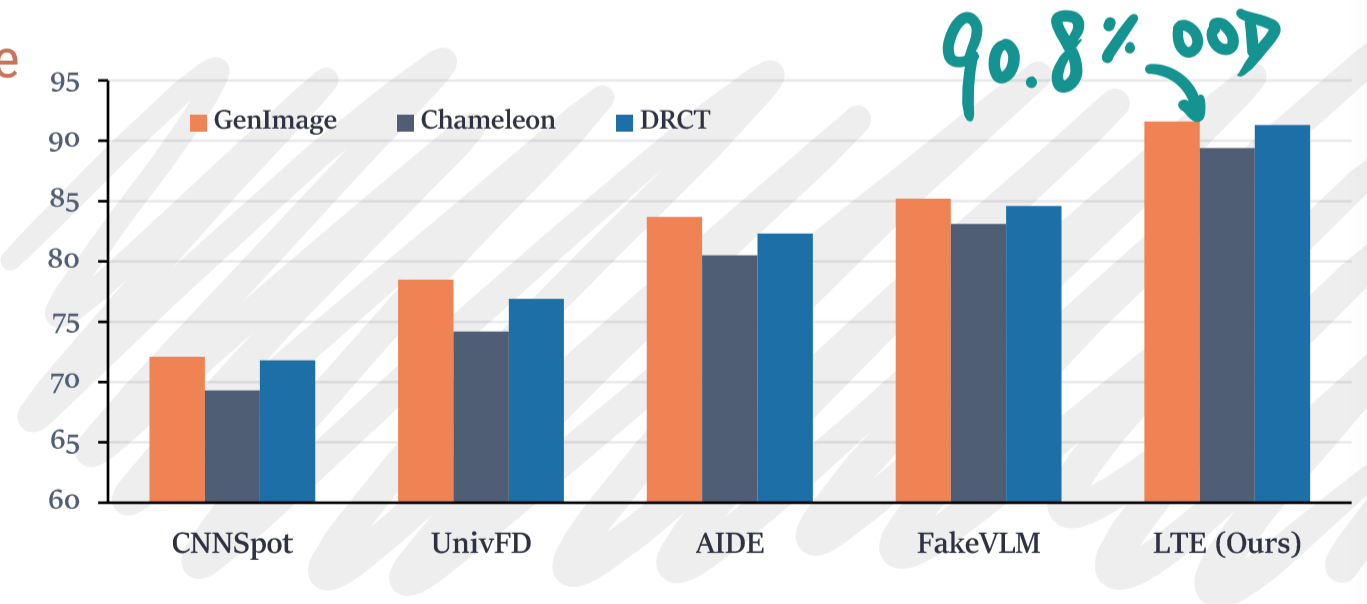
- Clear blueberry contour, natural camera blur
- The CAD diagram is clean with identifiable text

LTE Unchanged

- Branding consistent, correct geometry

LTE In Action

- Missing petals replaced by leaves
- The gates are too far away from each other
- Peanuts not cracked and poorly-textured
- Precise metallic keyrings; consistent shadows
- Symmetrical structure with intricate details
- Natural fog and scenery



In an era where generative AI booms, We tackle OoD with explainability, And took it further via the LTE paradigm.

LOCATE-THEN-EXAMINE

METHOD