

HIDDEN IN PLAIN SIGHT: Evaluating Abstract Shape Recognition in Vision-Language Models



UNIVERSITY OF OXFORD



UNIVERSITY OF ILLINOIS
URBANA-CHAMPAIGN

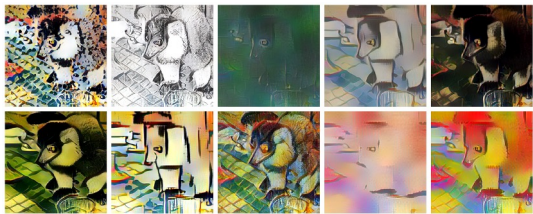


Arshia Hemmat¹, Adam Davies^{*2}, Tom A. Lamb^{*1}, Jianhao Yuan^{*1}, Philip Torr¹, Ashkan Khakzar¹, Francesco Pinto¹

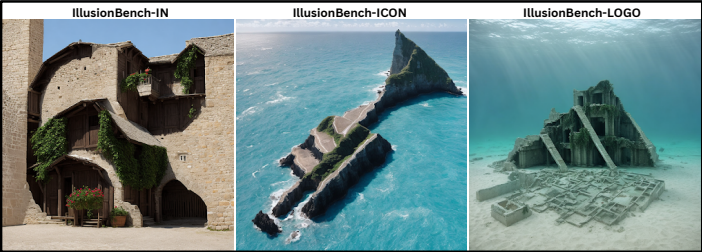
^{*}Equal Contribution. ¹University of Oxford. ²University of Illinois at Urbana-Champaign

Problem

- Do image classifiers rely on shape or texture?
- Limitations of existing benchmarks:
 - Lack of coherent, naturalistic, complex visual scenes.
 - Missing shape information, poor fine-grained details.

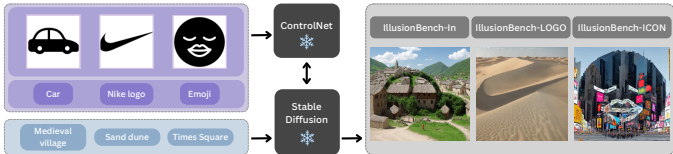


Can SOTA vision-language models recognise these shapes while ignoring scene/texture?



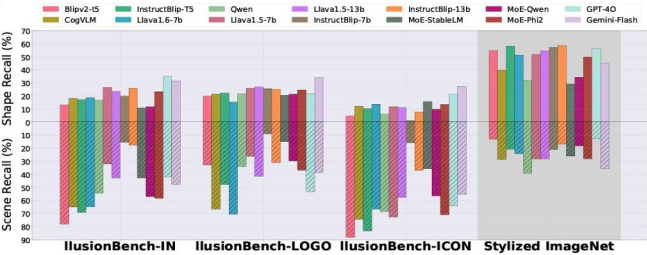
Approach

To address these issues, we introduce IllusionBench:



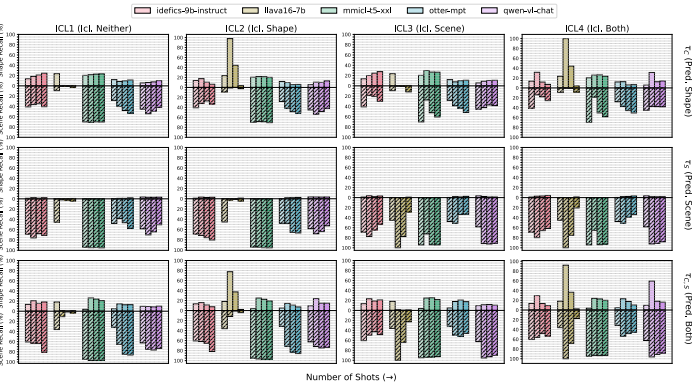
- Use ControlNet to condition Stable Diffusion on...
 - Conditioning images (binary masks of target shape).
 - Scene description (e.g., "medieval village" or "sand dune").

Zero-shot Shape Recognition



- SOTA VLMs biased towards scene/texture.
- Shape recognition gap between open and closed-source models.

Few-shot Shape Recognition



- ICL does not consistently solve the problem.
- VLMs still biased towards scene/texture.

Zero-shot? ❌
Few-shot? ❌
Fine-tuned? 🤔

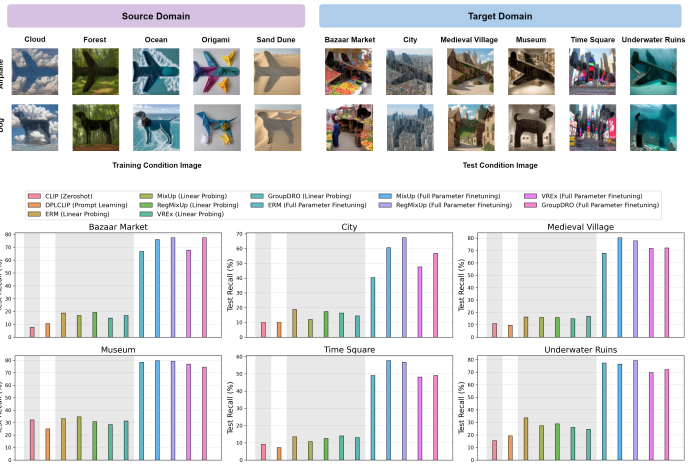


PAPER



WEBSITE

Domain Generalisation



- CLIP models fail zero-shot or via probing.
- Can be fine-tuned to learn domain-generalisable features.