# Robust Vision

- Image QA performance suggests that VLMs are capable of robust generalist vision

| | Claude 3.5 Sonnet[4] | GPT-4o[5] | Qwen2-VL-72B[6] | Phi 4 Multimodal[7] |
|---|---|---|---|---|
| AI2D[1] | 94.7% | 94.2% | 88.4% | 82.3% |
| ChartQA[2] | 90.8% | 85.7% | 88.3% | 81.4% |
| DocVQA[3] | 95.2% | 92.8% | 96.5% | 93.2% |

[1] Kembhavi et al. [ECCV 2016]
[2] Masry et al. [ACL 2022]
[3] Mathew et al. [WACV 2021]
[4] Anthropic et al. [Anthropic 2024]
[5] Hurst et al. [arXiv 2024]
[6] Wang et al. [arXiv 2024]
[7] Abouelenin et al. [arXiv 2025]

# Strong Chart Performance ≠ Robust Vision

- Despite strong chart performance, these models fail simple perception tests

### VLMs are Blind



"Are these circles touching?"

### HallusionBench



"Are the orange circles the same size?"

| | Claude 3.5 Sonnet | GPT-4o | Qwen2-VL-72B |
|---|---|---|---|
| VLMs are Blind[8] | 74.94% | 48.47% | — |
| HallusionBench[9] | 55.16% | 55.00% | 55.16% |

[8] Rahmanzadehgervi et al. [ACCV 2024]
[9] Guan et al. [CVPR 2024]

# Strong Chart Performance ≠ Robust Vision

- Newer models improve on primitive perception benchmarks
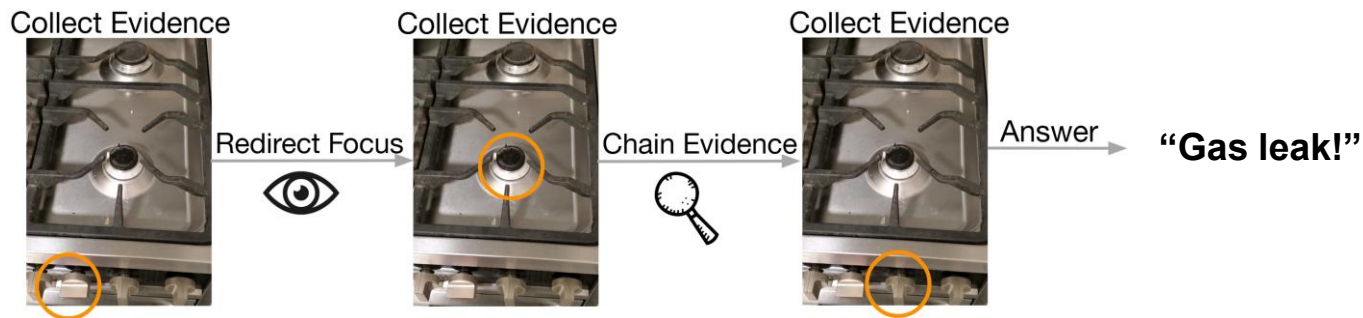
- Do they have robust vision?

| | Claude 3.5 Sonnet | GPT-4o | o4-mini[10] | o3[10] |
|---|---|---|---|---|
| VLMs are Blind[8] | 74.94% | 48.47% | 87.3% | 90.1% |

[8] Rahmanzadehgervi et al. [ACCV 2024]
[10] Ramesh et a. [OpenAI 2025]
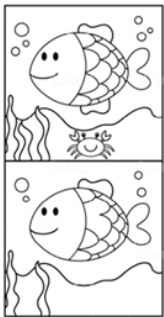
# How do humans process images?

- Humans gather evidence from multiple regions, and use the image itself to redirect focus

- We call this **nonlocal visual reasoning**

# Nonlocal Visual Reasoning Skills

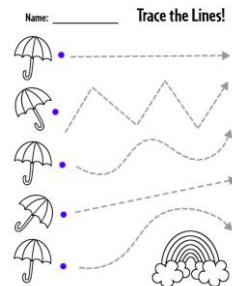- We define tasks that test three skills

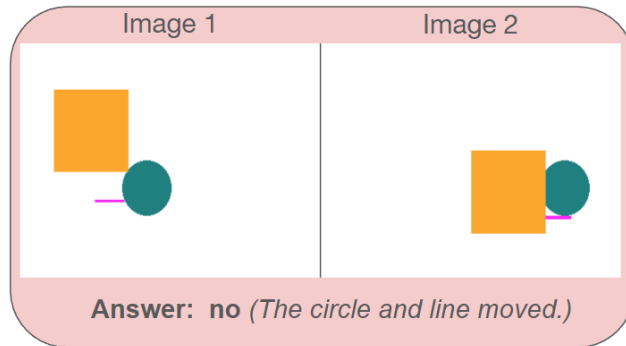Comparative Perception      Saccadic Search      Smooth Visual Search
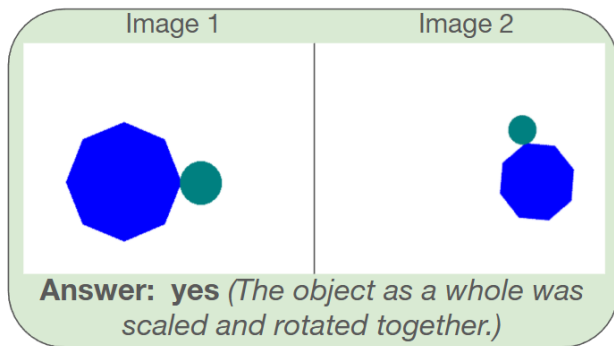


- Our tasks are designed to minimize necessary background knowledge and require the tested skill

# Comparative Perception

- The ability to compare two similar objects

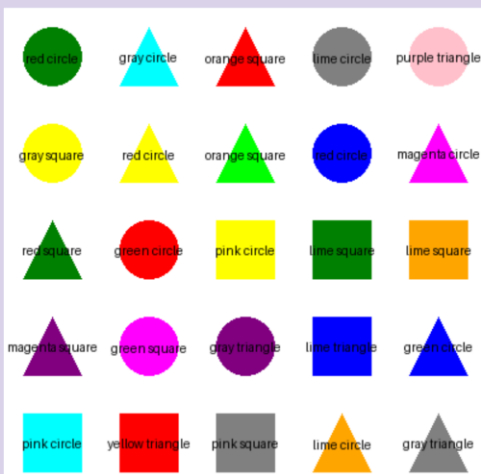- We evaluate three variants: Standard, Unconnected, Pixel-Perfect



**Question:** Does the object in Image 1 appear in Image 2? Answer no if it has been corrupted.

Image 1 | Image 2
**Answer: yes** *(The object as a whole was scaled and rotated together.)*

Image 1 | Image 2
**Answer: no** *(The circle and line moved.)*

# Saccadic Search

- The ability to make discrete, evidence-driven jumps across an image
- We evaluate 2, 3, and 4 jumps

# Smooth Visual Search

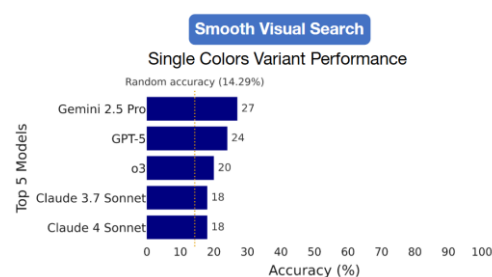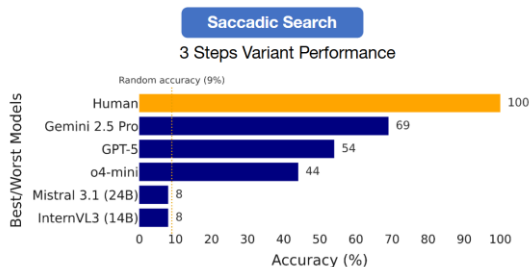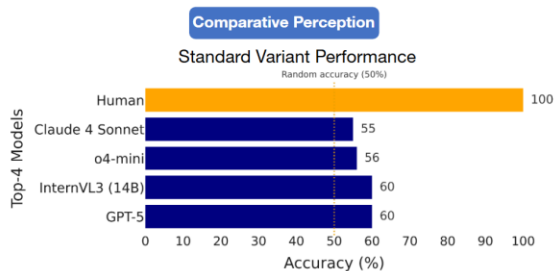- The ability to trace a line or a contour
- We evaluate over the wire colors: Single Color, Standard, Unique Colors



**Question:** "Which component does the wire from port 6 on the breadboard connect to?"

**Answer:** C1 (It crosses the blue wire and then terminates on component C1}

# Main Results

- Models fail catastrophically when comparing connected objects
- VLMs struggles to make discrete steps
- Cannot trace lines without color cues or other heuristics

# Conclusion

- Current models achieve strong benchmark scores despite lacking the human-like skills typically used to solve those tasks

- Visual reasoning is not robust even in frontier VLMs

- To the extent VLMs can visually reason, they perform it inconsistently and fail in non-intuitive ways

**Paper**

**Website**