

low-level vision tasks that are trivial to humans.

# Vision Language Models are Blind

Pooyan Rahmanzadehgervi\*, Logan Bolton\*, Mohammad Reza Taesiri, Anh Totti Nguyen



Results

Project, code, data



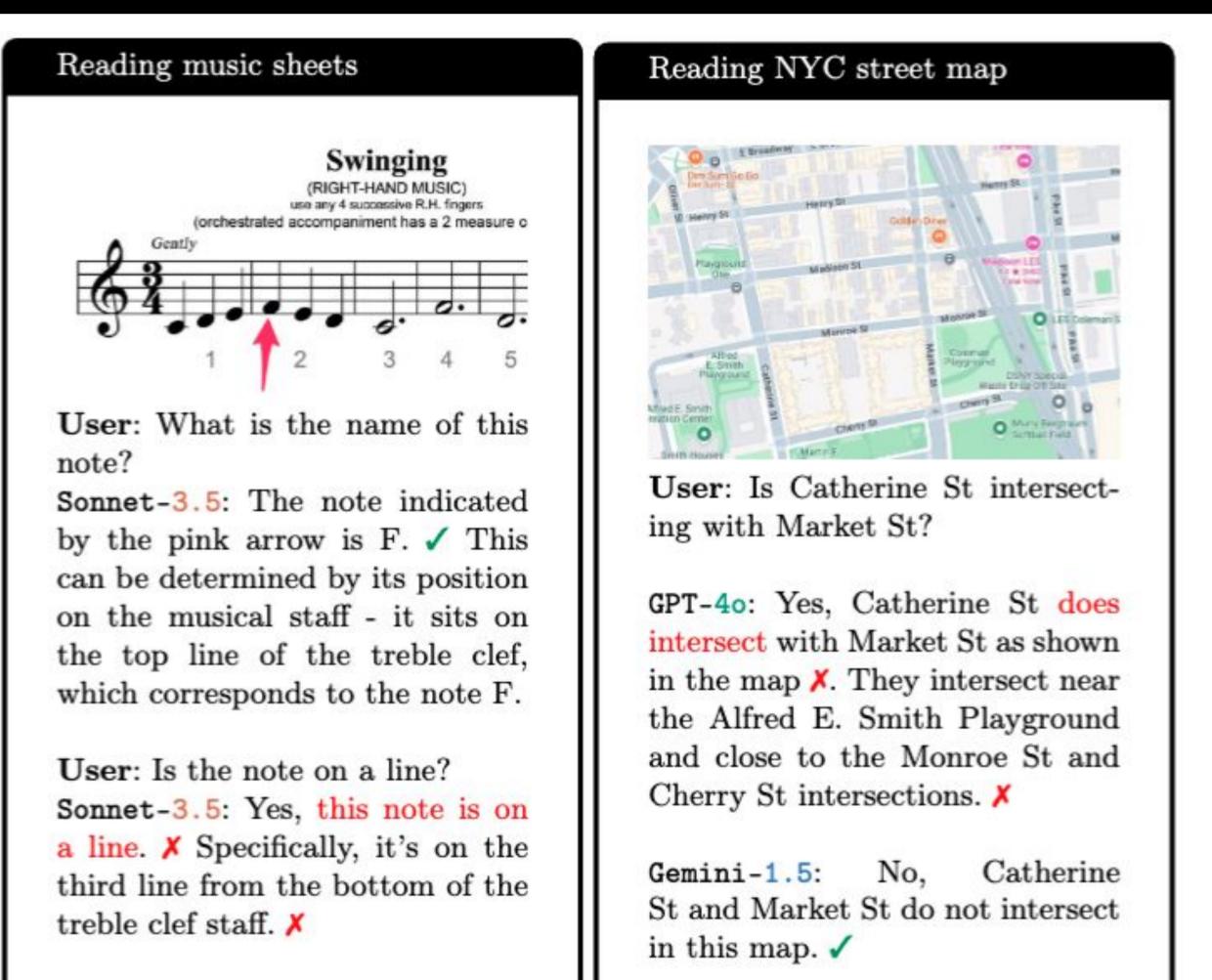
Summary

Large language models with vision capabilities still struggle with

### Motivation

- Gemini can solve 42.9% of the questions in MMMU benchmark without seeing the input image. Are we on the right way... Chen et al. 2024
- Most VQA benchmarks do not exclusively test vision capabilities.
- Text-only LLMs can reach > 80% of SOTA on DocVQA, TextVQA, ChartQA, AI2D (images are serialized). Analyzing...Hedge et al. 2023

### How would VLMs perform on tasks that exclusively test vision?



## **Findings**

#### VLMs cannot reliably:

- 1. Tell if two circles are touching
- 2. Count the number of times two lines intersect
- 3. Follow paths from one point to another
- 4. Reliably count how many rows and columns are in a table

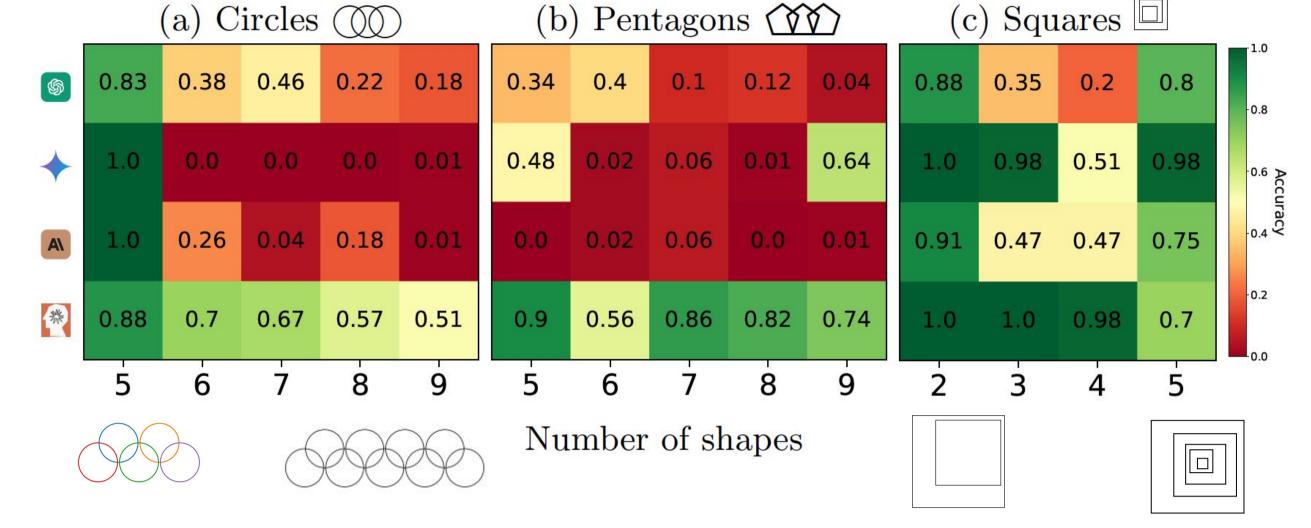
accuracy

5. Tell which letter is being circled

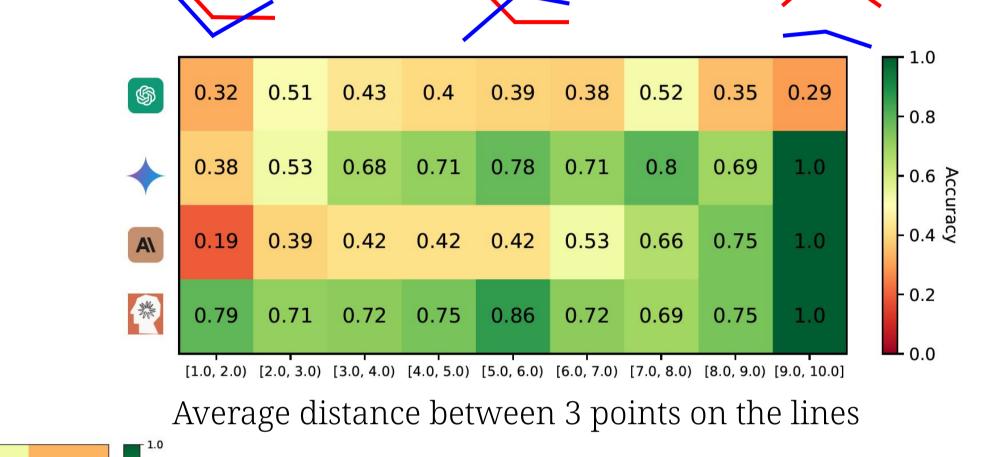
VLMs can accurately perform these same tasks when the shape primitives are distanced away.

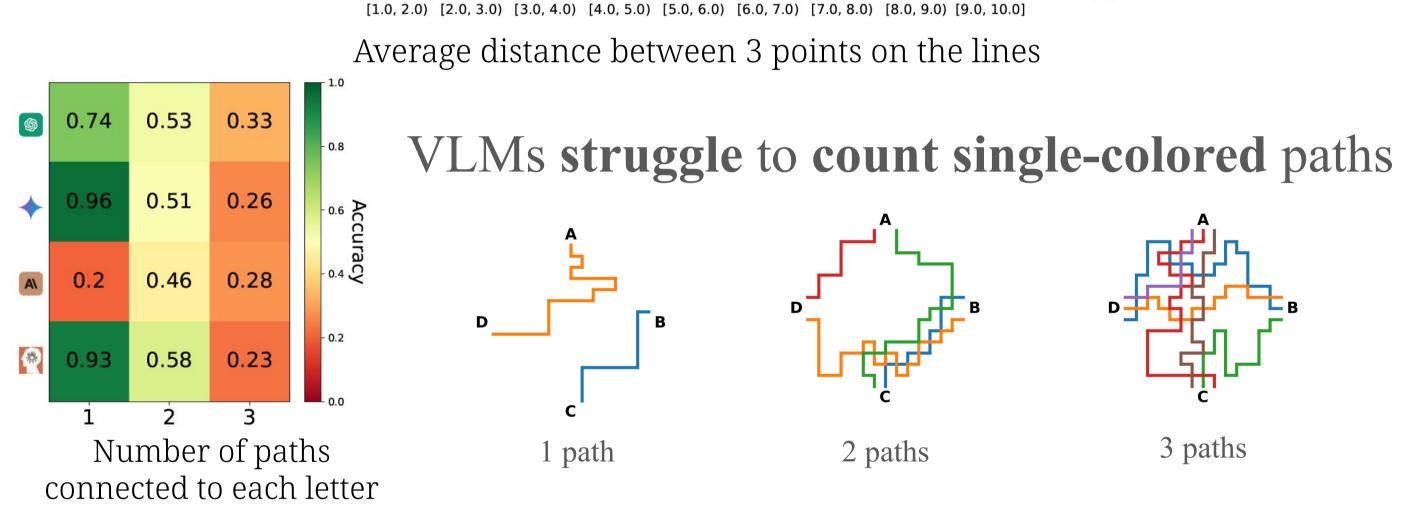
Model	$\triangleright$		0	$\bigcirc$	(W)			-4	Task mean
Random	33.33	50.00	5.77	20.00	20.00	25.00	4.55	33.33	24.00
◎ GPT-4o	41.61	75.91	74.23	41.25	20.21	55.83	39.58	53.19	50.23
♦ Gemini 1.5 Pro	66.94	93.62	83.29	20.25	24.17	87.08	39.39	53.13	58.48
Sonnet 3	43.41	86.46	72.06	29.79	1.87	65.00	36.17	31.11	45.73
Sonnet 3.5	75.36	90.82	87.88	66.46	77.71	92.08	74.26	58.19	77.84
Mean	56.83	86.70	79.37	39.44	30.99	75.00	47.35	48.91	58.07

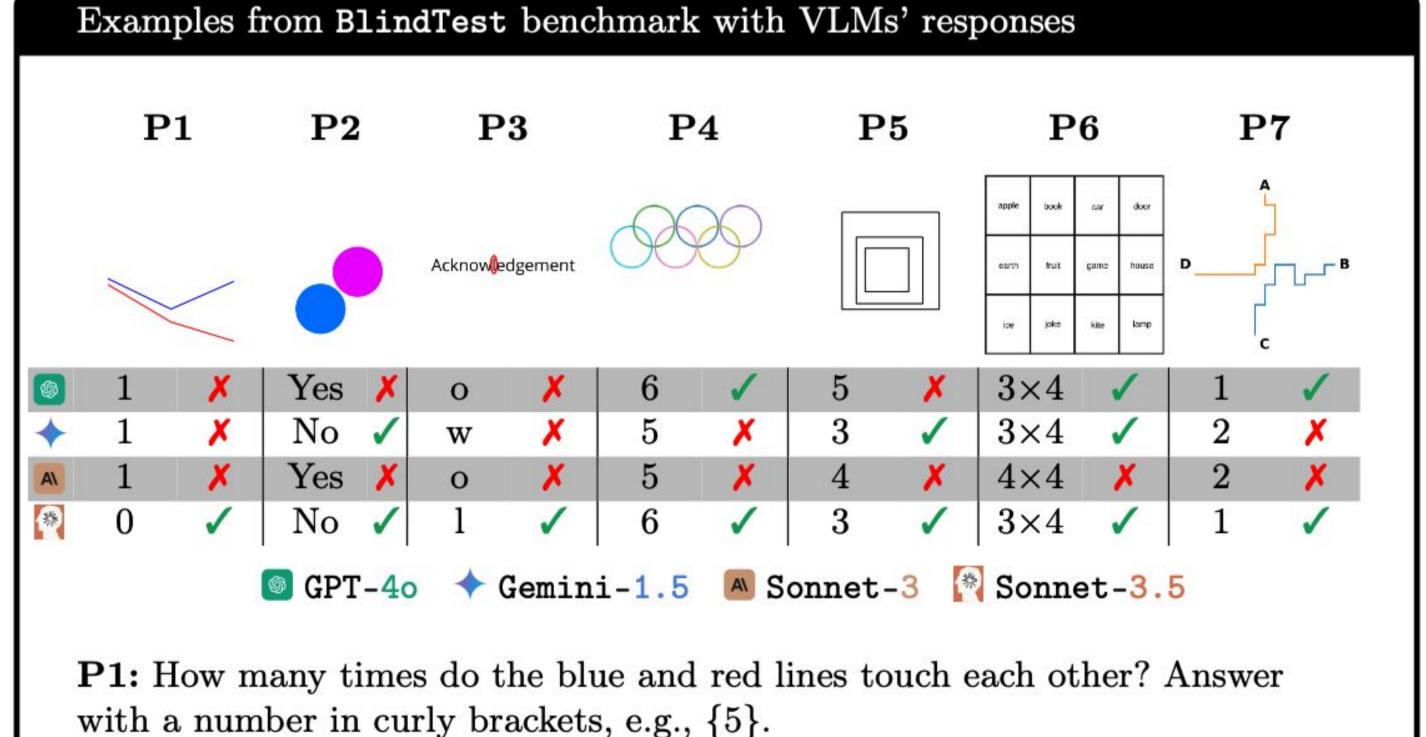
# VLMs struggle to count overlapped and nested shapes



#### VLMs cannot reliably count line intersections



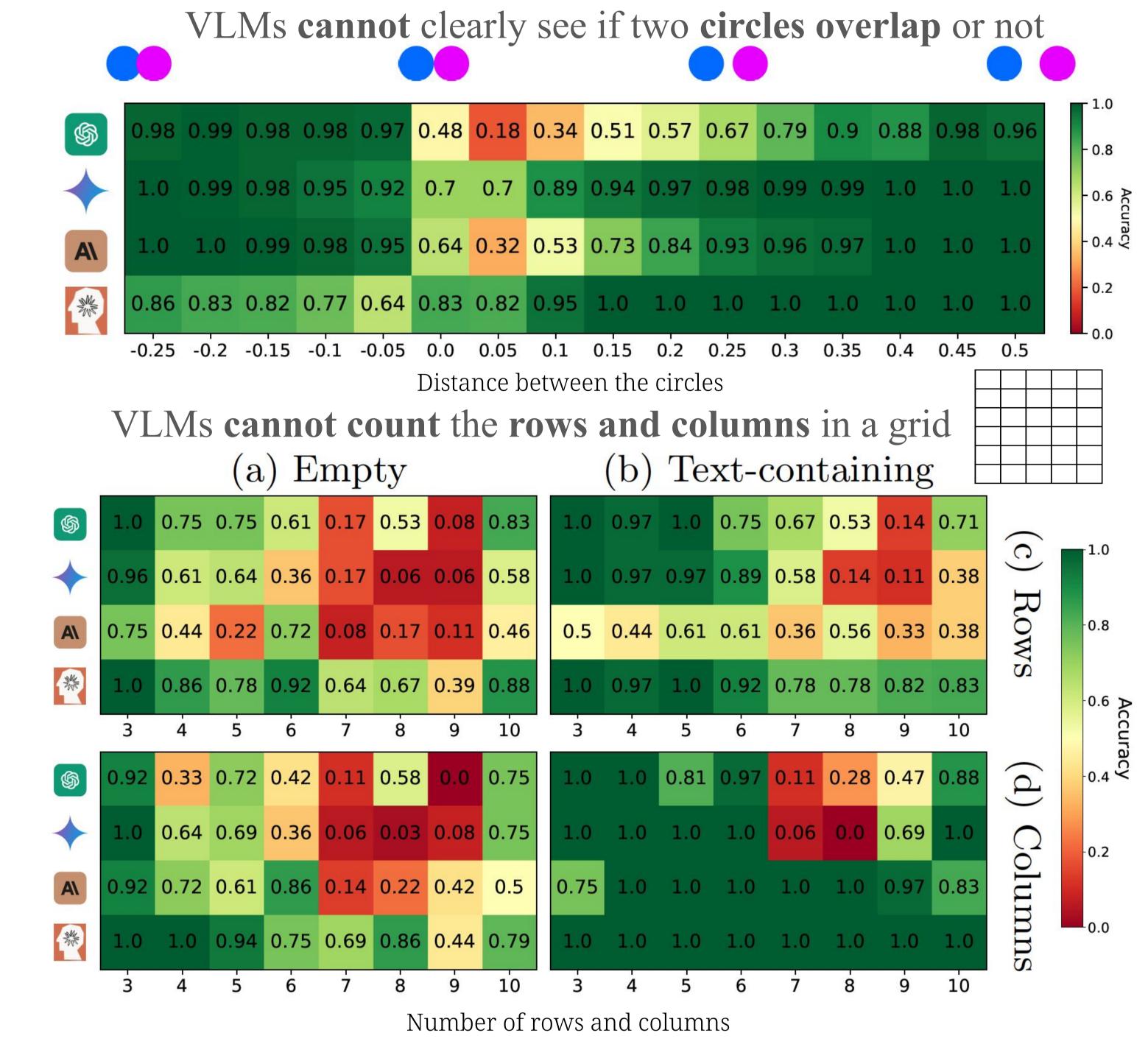




with a number in curly brackets, e.g., {5}.

- **P2:** Are the two circles overlapping? Answer with Yes/No.
- **P3:** Which character is being highlighted with a red oval? Please provide your answer in curly brackets, e.g. {a}
- **P4:** How many circles are in the image? Answer with only the number in numerical format.
- P5: How many squares are in the image? Please answer with a number in curly brackets e.g., {10}.
- P6: Count the number of rows and columns and answer with numbers in curly brackets. For example, rows={5} columns={6}.
- **P7:** How many single-color paths go from A to D? Answer with a number in curly brackets e.g. {3}.

# vlmsareblind.github.io



# When can VLMs understand low-level details?

VLMs can get near 100% accuracy on the tasks when the shapes are more separated.

