

I SPY WITH MY MODEL’S EYE: VISUAL SEARCH AS A BEHAVIOURAL TEST FOR MLLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) achieve strong performance on vision-language tasks, yet their visual processing is opaque. Most black-box evaluations measure task accuracy, but reveal little about underlying mechanisms. Drawing on cognitive psychology, we adapt classic *visual search* paradigms—originally developed to study human perception—to test whether MLLMs exhibit the “pop-out” effect, where salient visual features are detected independently of distractor set size. Using controlled experiments targeting colour, size and lighting features, we find that advanced MLLMs exhibit human-like pop-out effects in colour or size-based disjunctive (single feature) search, as well as capacity limits for conjunctive (multiple feature) search. We also find evidence to suggest that MLLMs, like humans, incorporate natural scene priors such as lighting direction into object representations. We reinforce our findings using targeted fine-tuning and mechanistic interpretability analyses. Our work shows how visual search can serve as a cognitively grounded diagnostic tool for evaluating perceptual capabilities in MLLMs.

1 INTRODUCTION

Despite their impressive performance, Multimodal Large Language Models (MLLMs) remain opaque in how they internally process and represent visual information. This is partly due to the lack of information disclosed by frontier model developers, and partly due to evaluation approaches: traditional evaluation benchmarks focus on accuracy or alignment with human outputs, but reveal little about the intermediate representations or cognitive-like processes these models may develop or employ. To gain deeper insight, we draw inspiration from investigations in cognitive science, particularly visual search paradigms, which have been used to probe the internal structure of human vision systems. Our goal is to develop a diagnostic toolkit for probing how these models respond to controlled visual challenges. Primarily targetting Marr’s computational level of analysis (Marr, 2000), this approach offers a way to investigate how MLLMs internally represent and prioritise information. We then provide a brief investigation into how these effects are represented in internal activations of MLLMs.

Experimental psychologists seeking to understand the ‘black box’ of the human visual attention system have long used controlled search tasks to probe attentional bottlenecks and feature integration (Treisman & Gelade, 1980; Wolfe, 1994). These reveal regularities and rules underlying visual cognition, while remaining agnostic as to the neural architecture producing them. Here, we take the same experimental approach to interrogate ‘black box’ MLLMs to identify rules governing their visual search behaviour. Our goal is not to optimise model performance, but to uncover regularities and divergences in how these models process visual information—particularly in relation to classic cognitive phenomena like pop-out effects and distractor interference. This approach enables a structured comparison between human and model behaviour, offering insight into the kinds of representations and inductive biases these systems may have developed. Concretely, we present a systematic investigation of visual search behaviour in MLLMs, varying structural components of visual scenes such as the variety of features and the size of distractor sets. In doing so, we provide a targeted investigation of fundamental perceptual capabilities in state-of-the-art MLLMs under zero-shot settings.

2 BACKGROUND

2.1 MULTIMODAL LARGE LANGUAGE MODELS

MLLMs extend the capabilities of language models by incorporating visual inputs, enabling them to perform tasks such as image captioning, visual question answering, and multimodal reasoning. While earlier vision-language models (VLMs) like CLIP (Radford et al., 2021) and BLIP (Li et al., 2022) focused on learning modality-consistent embeddings for retrieval or generation, MLLMs are typically built by integrating visual encodings into autoregressive language models. Notable examples include Flamingo (Alayrac et al., 2022), which augments a frozen language model with cross-attention to a visual encoder; GPT-4o (OpenAI et al., 2024); and LLaVA (Liu et al., 2023), an open-source model that injects visual embeddings into a fine-tuned Vicuna model. These models achieve strong performance across standard image-oriented benchmarks such as VQAv2 (Goyal et al., 2017), OK-VQA (Marino et al., 2019), and MMMU (Yue et al., 2024) often without task-specific tuning. MLLMs differ in how they fuse vision and language, and their internal processing remains largely opaque. Many models are proprietary with limited documentation, and even open architectures offer limited insight into how visual inputs are handled within the model. This means current evaluations overwhelmingly rely on aggregated end-task accuracy, which provides little insight into the internal structure of model reasoning. Given the increasing desire to deploy MLLMs in real-world applications that require strong perceptual capabilities (e.g., medical image analysis (Liu et al., 2025), control of physical robots (Luo et al., 2025)), understanding *how* models achieve these benchmark scores, and in what ways they might fail, is critical. Here, we retain output-based evaluation but reframe it through the lens of controlled experimentation. Just as careful experimentation has provided insight into the internal structures and representations of human cognition, we will use cognitive science inspired visual search paradigms to study MLLM perception.

2.2 VISUAL SEARCH

Visual search tasks have long been used by psychologists to study human attention and perception. Participants are shown a scene containing a target object and several distractors, and must quickly judge whether the target is present, or identify its location. What makes these tasks powerful is not their difficulty, but their structure—the systematic variation of parameters such as set size, feature complexity, and target presence enables the exposure of latent properties of the underlying cognitive system (Wolfe, 1998; 2020). A classic distinction is between *disjunctive* search, where a target is identifiable from the distractors along a single dimension (e.g., colour—a red square among blue squares, or shape—a square among circles), and *conjunctive* search, where the target can only be identified by a unique combination of features (e.g., colour and shape—a red square among red and blue circles and blue squares). The compositional nature of human visual representations means that disjunctive search typically yields fast, parallel detection—so-called “pop-out” effects—as a pre-attentive feed-forward pass through the early visual system is sufficient to detect the single distinguishing feature (e.g., red among green). However, as conjunctive search requires the attentional binding of two or more primitive features to distinguish the target from distractors (e.g., ‘red’ and ‘square’), reaction times are slower and increase linearly with the number of distractors as each item requires individual inspection.

The distinction between conjunctive and disjunctive search has been shown to be highly reliable across humans and other animals (Orlowski et al., 2015; Reichenthal et al., 2019). Such behavioural phenomena can be used to understand the nature of the processes the visual system is performing when allocating attention and to predict behaviour in novel scenarios (for example, one can predict that a human will quickly identify and locate a coffee stain on a white carpet, but not on a colourful patterned rug). Using careful manipulation of inputs and observing changes in behaviour, the underlying structure of informational processes can be inferred. This makes visual search tasks particularly suitable for interrogating black-box models like MLLMs where the internal representations may be opaque. Structured behavioural experiments can reveal the presence, predictability and human-likeness of search strategies. Beyond theoretical interest, visual search has clear applied value: the same principles are used to optimise display layouts in cars Smith et al. (2015), and to design and assess professional search in domains such as airport security (Biggs et al., 2018; Mitroff et al., 2018). This makes these paradigms natural probes for MLLMs as well, allowing us to reason about which perceptual and attentional mechanisms they possess and how their failures might translate to visually

demanding deployments. Similarly, understanding visual search in MLLMs may be valuable for determining how to best present tasks to multimodal systems to maximise the likelihood of success. For example, road signs are designed to be salient to human drivers, but it’s becoming increasingly important to identify whether the same features would be salient to self-driving cars.

3 VISUAL SEARCH EXPERIMENTS

We adapted three visual search experiments for MLLMs: **Circle Sizes**, **2 Among 5** and **Light Priors**. Each targets a specific visual feature known to induce the “pop-out” effect in humans: size (Samiei & Clark, 2022), colour (Wolfe et al., 2010; Wolfe & Horowitz, 2004) and light source direction (Adams, 2007), respectively. Each experiment is described in detail in Sections 3.1–3.3. We include two target localisation variants evaluating levels of precision:

- **Cells:** The image is divided into a 2×2 grid, and the model must identify the grid cell containing the target (always present). Accuracy is the proportion of trials in which the correct cell is identified.
- **Coordinates:** The model must return the coordinates of the target. Performance is evaluated using Euclidean distance from the chosen point to the centre of the target.

We evaluate a selection of MLLMs, comprising both closed source and open source models. Specifically, we evaluate GPT-4o (OpenAI et al., 2024), Claude Sonnet 3.5 (Anthropic, 2024), Llama 3.2 90B (Meta AI, 2024). Full details of the model are provided in Appendix C. A selection of other models are also evaluated in Appendix B. We also compare against a human baseline ($N = 90$). Humans are generally highly accurate in visual search tasks, and processing differences between experimental conditions are usually identified using response times. However, by limiting stimulus presentation time (e.g., to 1500ms) we can compare humans and LLMs using accuracy scores alone (see McElree & Carrasco (1999)). Further details of the human portion of the experiment are provided in Appendix I.

3.1 EXPERIMENT 1: CIRCLE SIZES

Our first experiment was designed to investigate whether MLLMs, like humans, show ‘pop-out’ effects in disjunctive search tasks – that is, tasks in which the target can be easily distinguished from distractors by a single change in a simple visual feature. Object size is one such feature that can be detected pre-attentively in human subjects (Proulx, 2010). Figure 1 illustrates how we systematically manipulated the radius of a target circle amongst matched distractors using three experimental conditions: Small (22.5 pixels), Medium (25 pixels) and Large (30 pixels). The number of distractors varied from 0 to 49 across trials, and all circles were randomly placed on a white 400×400 pixel background and rendered without overlap. Target and distractor circles were presented in the same colour, which was randomly sampled from red, green and blue. If MLLMs process size as a visual primitive, we would expect them to locate the target circle in the Large condition with high accuracy and precision, irrespective of distractor numbers. We would also expect that this ‘pop-out’ effect would be attenuated when size differences are less salient in Small and Medium conditions.

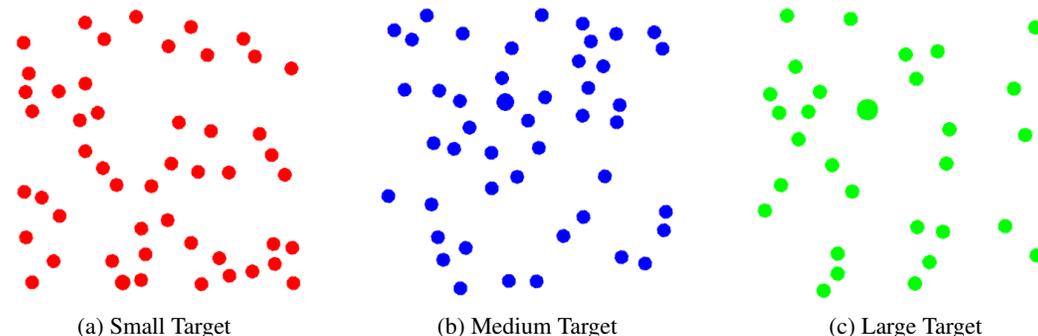


Figure 1: Circle Sizes task. Examples from the three experimental conditions. The target is always the single circle that is larger than the rest. The colour of all circles is always the same and is sampled from red, blue and green.

Table 1: Regression slopes and correlations for the effect of distractor number on accuracy within conditions across all three experiments

Exp.	Model	Condition	Mean Acc	Regression		Correlation	
				Slope	95% CI	r	p
CS	claude	Small	0.302	-0.017	[-0.020, -0.014]	-0.111	< 0.001
	claude	Medium	0.425	-0.023	[-0.026, -0.021]	-0.166	< 0.001
	claude	Large	0.600	-0.019	[-0.022, -0.016]	-0.136	< 0.001
	gpt-4o	Small	0.425	-0.025	[-0.028, -0.022]	-0.177	< 0.001
	gpt-4o	Medium	0.722	-0.016	[-0.019, -0.013]	-0.102	< 0.001
	gpt-4o	Large	0.832	-0.005	[-0.009, -0.002]	-0.028	0.082
	llama	Small	0.281	-0.004	[-0.007, -0.001]	-0.029	0.057
	llama	Medium	0.341	0.005	[0.002, 0.008]	0.032	0.021
	llama	Large	0.465	0.011	[0.008, 0.014]	0.081	< 0.001
2A5	claude	Disjunctive	0.676	-0.005	[-0.006, -0.004]	-0.065	< 0.001
	claude	Shape	0.587	-0.011	[-0.012, -0.010]	-0.159	< 0.001
	claude	Shape-Col.	0.368	-0.016	[-0.017, -0.015]	-0.219	< 0.001
	gpt-4o	Disjunctive	0.847	0.002	[< 0.001, 0.003]	0.017	0.249
	gpt-4o	Shape	0.555	-0.019	[-0.020, -0.018]	-0.267	< 0.001
	gpt-4o	Shape-Col.	0.409	-0.018	[-0.019, -0.017]	-0.244	< 0.001
	llama	Disjunctive	0.548	0.001	[< 0.001, 0.002]	0.010	1.000
	llama	Shape	0.412	-0.007	[-0.008, -0.006]	-0.100	< 0.001
	llama	Shape-Col.	0.307	-0.008	[-0.009, -0.006]	-0.100	< 0.001
LP	claude	Top	0.330	-0.045	[-0.053, -0.037]	-0.109	< 0.001
	claude	Bottom	0.428	-0.039	[-0.047, -0.032]	-0.102	< 0.001
	claude	Left	0.298	-0.058	[-0.066, -0.050]	-0.138	< 0.001
	claude	Right	0.298	-0.059	[-0.067, -0.050]	-0.139	< 0.001
	gpt-4o	Top	0.545	0.027	[0.020, 0.035]	0.072	< 0.001
	gpt-4o	Bottom	0.729	0.089	[0.080, 0.097]	0.200	< 0.001
	gpt-4o	Left	0.380	-0.021	[-0.028, -0.013]	-0.053	< 0.001
	gpt-4o	Right	0.429	-0.046	[-0.054, -0.039]	-0.120	< 0.001
	llama	Top	0.514	0.045	[0.037, 0.052]	0.117	< 0.001
	llama	Bottom	0.506	0.032	[0.025, 0.040]	0.084	< 0.001
	llama	Left	0.441	0.040	[0.032, 0.047]	0.104	< 0.001
llama	Right	0.382	0.015	[0.007, 0.023]	0.038	0.003	

Note: Correlations are Pearson’s r and p values are Bonferroni corrected for multiple comparisons. For the 2 Among 5 experiment, Shape and Shape-Col. (Shape-Colour) refer to conjunctive search conditions. Experiments are abbreviated to CS (Circle Sizes), 2A5 (2 Among 5) and LP (Light Priors). For Models, ‘claude’ refers to claude-sonnet, and ‘llama’ refers to llama-90B.

To assess model performance on the cells variant of the Circle Sizes task, we conducted regression and correlation analyses (reported in Appendix F and summarised in Table 1). In this context, small absolute slope values and non-significant correlation coefficients indicate relatively uniform performance across set sizes. A pronounced negative correlation or slope suggests the target was increasingly harder to find at higher set sizes, while a combination of high accuracy coupled with set-size independence is indicative of a pop-out effect.

Model performance on this task is illustrated in Figure 2. GPT-4o exhibits a clear pop-out effect: its accuracy for Large targets is high ($M = 83\%$) and remains stable across increasing numbers of distractors ($r = -0.028$, $p = 0.082$). Accuracy for Medium targets also remains high ($M = 72\%$), though it shows a modest decline with set size ($r = -0.102$, $p < .001$). In contrast, performance for Small targets is lower ($M = 43\%$), and declines more markedly as distractor count increases ($r = -0.177$, $p < .001$). Notably, GPT-4o’s response pattern closely mirrors that of human participants, who showed pop-out in the Large condition, near pop-out in the Medium condition, but declining accuracy for higher distractor numbers in the Small condition (see Appendix I). Claude Sonnet also demonstrates sensitivity to target size, with performance systematically decreasing across conditions. However, unlike GPT-4o, it does not exhibit set size independence in any condition ($r_s < -0.110$,

216 $ps < .001$). LLaMA 90B shows slight increases in performance across target size conditions, but
 217 is generally poorer and mostly flat across set-sizes. The difference between models is particularly
 218 evident in the coordinates task variant (see Appendix Figure 8). GPT-4o localises Medium and Large
 219 targets with minimal error; Claude Sonnet maintains low error rates only for Large targets; whereas
 220 Llama 90B exhibits high error rates across all three conditions, with unexpectedly high error for
 221 Large targets due to a high number of invalid coordinate responses (see Appendix D). These findings
 222 suggest that more capable models, particularly GPT-4o, exhibit robust and human-like size-driven
 223 salience effects in disjunctive search.

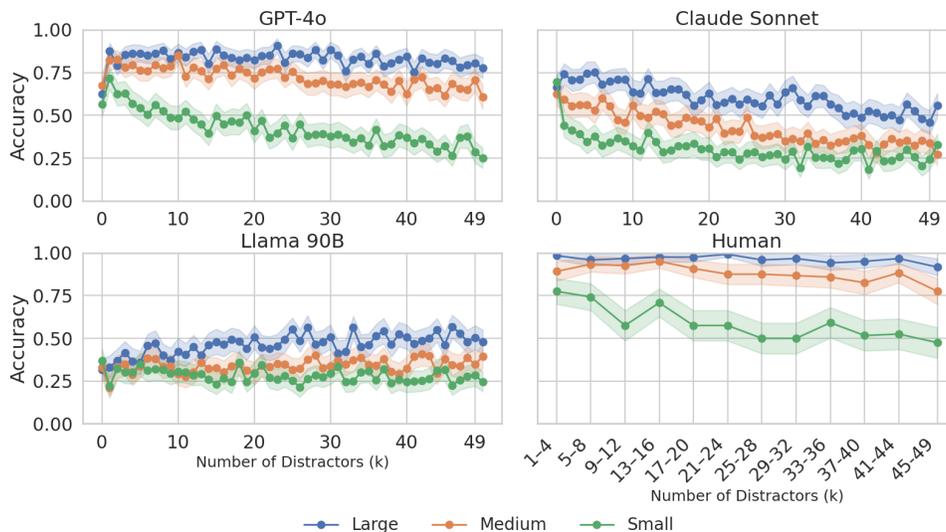


Figure 2: Results for Circle Sizes on Cells.

3.2 EXPERIMENT 2: 2 AMONG 5

249 In Experiment 1, we demonstrated that MLLMs show a pop-out-like pattern of behaviour in disjunctive
 250 search. In Experiment 2, we investigated whether MLLMs demonstrated human-like attentional
 251 limitations for *conjunctive* search by adapting the “2 Among 5” task from Wolfe et al. (2010) – the
 252 goal of which is to locate a target “2” among distractor “5”s (or *vice versa*). These digits are mirror
 253 images of each other, and are more complex than simple shapes (e.g., circles) as they consist of a
 254 combination of five line segments. We manipulated search type across three experimental conditions
 255 (see Figure 3): Disjunctive, Shape Conjunctive and Shape-colour Conjunctive. In the Disjunctive
 256 condition, target digits differ in colour from all other distractors, enabling parallel detection based on
 257 colour alone. In the Shape Conjunctive condition, however, all digits are the same colour, requiring
 258 the representation of both spatial configuration (i.e., the combination of line segments) and chirality
 259 (i.e., a “5” or “2”) to identify the target. Finally, in the Shape-Colour Conjunction condition, the
 260 target must be differentiated by a unique combination of both colour and spatial features (e.g., a “red
 261 2”) amongst distractors that have other shape-colour combinations (e.g., “red 5”, “blue 2” etc.). In
 262 humans, ‘binding’ primitive visual features to form more complex object representations requires
 263 attentional resources (Treisman & Gelade, 1980), resulting in serial search behaviour as potential
 264 targets require individual inspection, leading to search times that increase linearly with distractor
 265 numbers. If MLLMs have similar representational limits, they should show poorer performance in
 Shape Conjunctive and Shape-Colour Conjunctive conditions relative to the Disjunctive condition.

266 As in the previous experiment, targets and distractors were randomly placed on a 400×400 white
 267 background, but were also assigned a random rotation between 0 and 360 degrees. As before, digit
 268 colour was randomly sampled from red, green and blue, though for Disjunctive and Shape-Colour
 269 Conjunctive conditions two colours were sampled without replacement and assigned appropriately to
 target and distractors. The number of distractors present in each trial was also varied from 0 to 99.

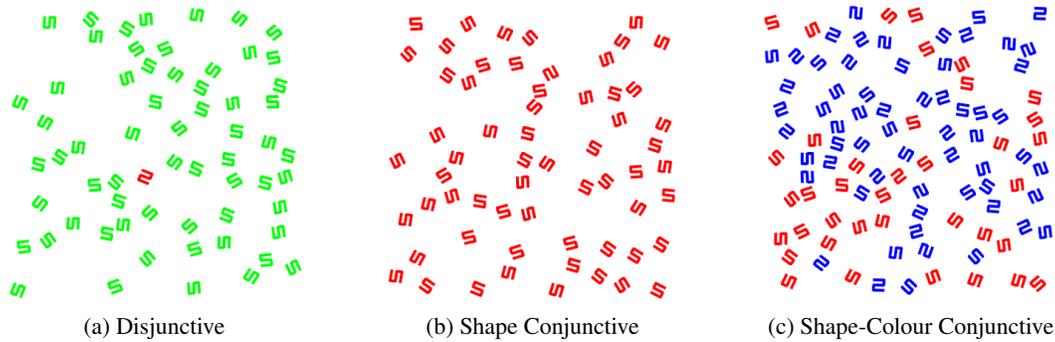


Figure 3: Example stimuli from the 2 Among 5 task, illustrating the three experimental conditions. In all examples, the target “2” is coloured red. (a) **Disjunctive**: The target differs from distractors by colour. (b) **Shape Conjunctive**: All digits share the same colour, requiring shape discrimination. (c) **Shape-Colour Conjunctive**: The target is uniquely defined by both shape and colour, with distractors sharing at most one feature. Target and distractor colours are randomized across trials.

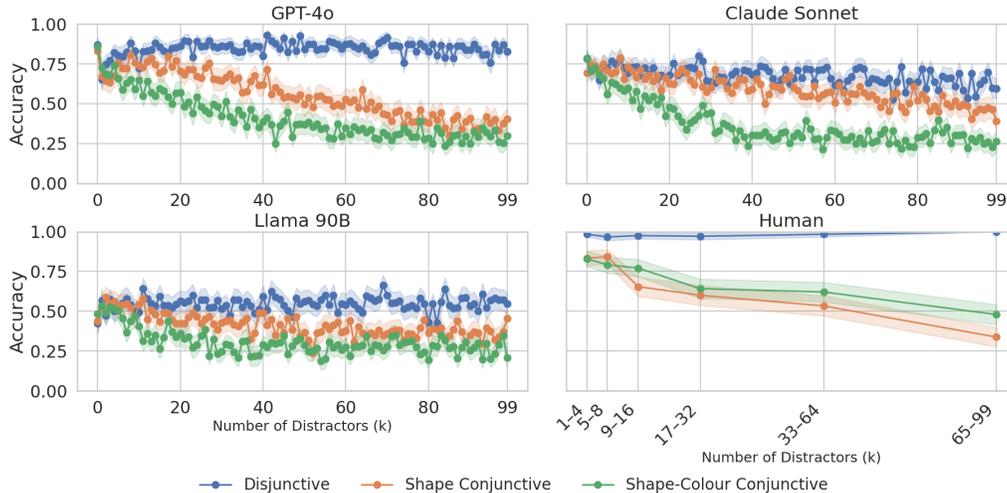


Figure 4: Results for the 2Among5 task — Cells mode. The shaded region denotes the 95% confidence interval.

We tested both “2 among 5” and “5 among 2” cases, in order to control for potential asymmetries or biases in the model’s representations.

Table 1 and Figure 4 show model performance for “2 Among 5” and “5 Among 2” tasks combined in the cells variant. GPT-4o showed high performance in the Disjunctive condition ($M = 85\%$) which was also flat across set sizes ($r = 0.017$, $p = 0.249$), replicating the pop-out effects that we found in Experiment 1 and in our human experiments. GPT-4o also showed human-like limits for conjunctive search, as its performance in Shape Conjunctive ($M = 56\%$, $r = -0.267$, $p < .001$) and Shape-Colour Conjunctive ($M = 41\%$, $r = -0.244$, $p < .001$) conditions declined as the number of distractors increased. Claude Sonnet showed a slight decline in accuracy with increasing set size in the Disjunctive condition ($r = -0.065$, $p < .001$), but was strongest in this condition overall. However, it showed a much smaller performance gap between Disjunctive and Shape Conjunctive, which may be suggestive of representational differences between the models. Although Llama 90B’s performance was generally poor, it did show an advantage in the Disjunctive condition, perhaps suggesting a crude salience heuristic. This pattern was corroborated in the coordinate-based localisation task (Appendix Figure 9), where error rates for GPT-4o and Claude Sonnet show clear differences between

the conditions reflecting highly precise disjunctive search, but increasing error-rates with set-size for conjunctive search. Llama 90B, however, did not distinguish between search tasks, with high error in all three conditions.

3.3 EXPERIMENT 3: LIGHT PRIORS

Thus far, we have demonstrated that some MLLMs, like humans, show set-size *independent* detection of targets when they can be distinguished by a single primitive visual feature (i.e., disjunctive search), but show set-size *dependent* search performance when representational binding of multiple features is required (i.e., conjunctive search). In a third experiment, we investigated whether MLLMs possess more sophisticated representations that incorporate assumptions about how objects appear in the real world. Classic work in cognitive science has found that humans have a ‘light-comes-from-above’ prior due to their experience of the natural world (Enns & Rensink, 1990). This is incorporated into low-level visual perception such that objects lit from the top or bottom can be rapidly detected – and bottom-lit objects, due to their novelty, are particularly salient (Enns & Rensink, 1990).

To test whether MLLMs also incorporate natural scene priors in visual search, we adapt a visual search task from cognitive science (Adams, 2007). In this task, each image contains a number of shaded circles designed to resemble 3D spheres illuminated from a specific direction. The goal is to identify the sphere that is the ‘odd one out’, that is lit from the opposite (180°) direction to distractors, but is otherwise identical. Here, we defined four task variants: Top, Bottom, Left, and Right, corresponding to the direction the target appeared to be lit from. Distractor numbers varied between 0 and 17 per trial. The spheres were rendered in greyscale to avoid introducing colour cues, and were randomly placed at least 20 pixels apart in a medium-toned greyscale circle within the image without overlap. We also included a black border to ensure a consistent maximum distance of spheres from the image centre. Example stimuli are shown in Figure 5.

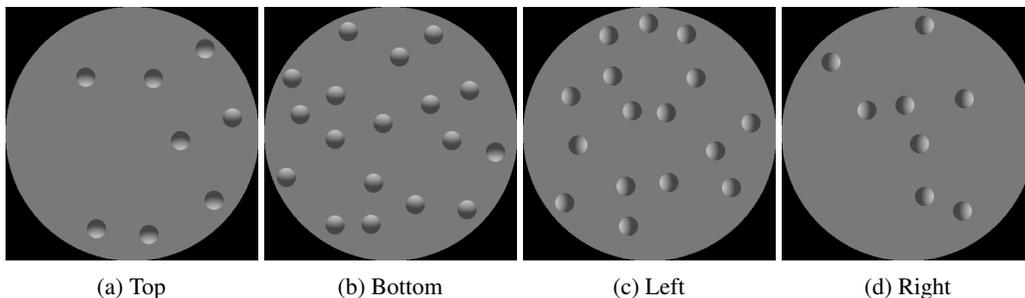


Figure 5: Light Priors Task: Circles are shaded with directional gradients, mimicking the shading on a sphere if lit from different directions. The subject must identify the target lit from a specific direction.

Figure 6 illustrates the results for the Light Priors cells task. Although we don’t see a pattern indicative of pop-out—as models show overall reduced accuracy and flat performance across set-sizes in *all* conditions—we do see clear differences between light-source conditions (see Table 1). For two or more distractors, as at least two distractors are required for an accurate “odd-one-out” decision, the pattern of results from GPT-4o shows remarkable similarity to our human baseline, with performance advantages for vertical (top and bottom-lit) relative to horizontal gradients (left and right-lit). Interestingly, we also see a clear performance advantage for bottom-lit ($M = 73\%$) relative to top-lit ($M = 55\%$) spheres, matching human behaviour. Claude and LLama’s performance on this task was poorer overall (GPT-4o: $M = 52\%$, Claude: $M = 34\%$, Llama: $M = 46\%$), but both showed a vertical-gradient advantage, and Claude (though not Llama) showed the highest accuracy for bottom-lit spheres. In the coordinates task, GPT-4o error rates (Appendix Figure 10) mirror the cell variant almost exactly, Claude Sonnet shows a clear performance advantage for bottom-lit spheres relative to the other lighting directions, whilst Llama 90B demonstrates slightly better performance for vertical-gradient spheres.

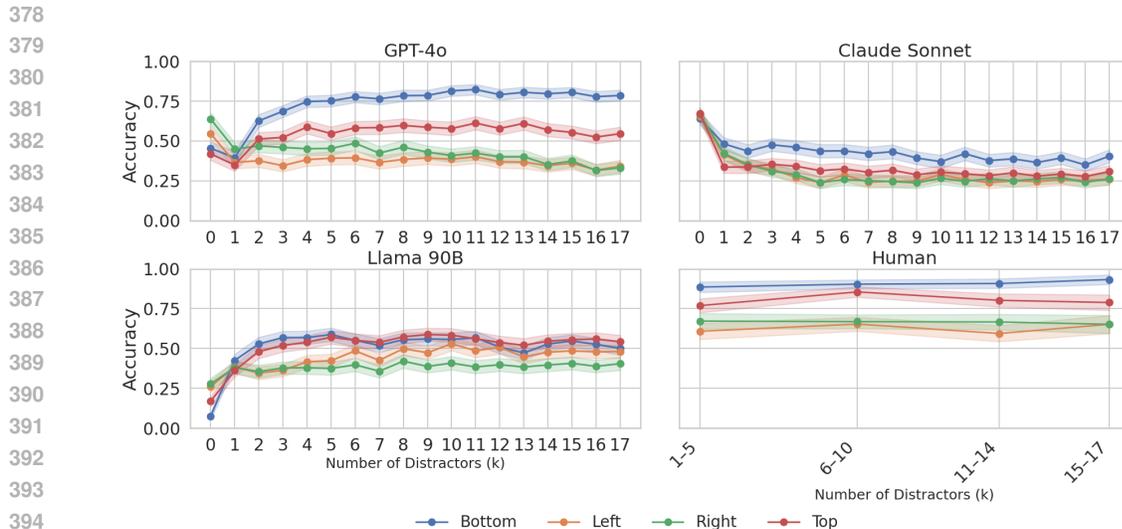


Figure 6: Results for the Light Priors task — Cells mode (three models + human baseline).

3.4 MECHANISTIC INTERPRETABILITY ANALYSIS

As an exploratory extension of our main analyses, we applied techniques from mechanistic interpretability (Lin et al., 2025) to probe the internal structure of MLLM representations. We focus here on Llama 90B, the largest open-weight model we evaluated, and report findings from the 2-Among-5 task, for which it exhibited the most reliable search performance.¹ Drawing on prior work showing that simple visual features (e.g., colour) are represented in early layers of CNNs and transformers (Raghu et al., 2021; Zeiler & Fergus, 2014), we hypothesised that disjunctive search tasks relying on primitive features would primarily engage early network layers, whereas more complex conjunctive search tasks requiring feature binding would recruit deeper layers. Mirroring findings in human vision, where low-level saliency is processed in early visual cortex (Zhaoping & May, 2007) and conjunctive search engages higher visual regions (Chelazzi et al., 1993), we observed a similar early/late division in Llama 90B’s activation patterns.

4 FINE-TUNING

Our experimental results have demonstrated that MLLMs, like humans, show capacity limits in conjunctive search. Human performance on conjunctive search tasks has been shown to improve after training (Czerwinski et al., 1992). Similarly, we tested whether MLLM performance on conjunctive tasks could be improved through fine-tuning. *Supervised Fine-tuning* (SFT) refers to an additional task-specific training step used to improve the performance of a language model in a particular domain. Here, we fine-tuned GPT-4o on examples from the Shape Conjunctive variant of the 2 Among 5 task, paired with ground-truth cell responses. We trained models with datasets generated with a different seed, with sizes 10 and 100 for three epochs, and 1000 for a single epoch.² The training data included only items with 0–49 distractors; test data included the full range up to 99, allowing us to assess generalisation beyond the training distribution.

Figure 7 shows accuracy on the Shape Conjunctive 2 Among 5 task under the *Cells* evaluation. Even minimal fine-tuning (10 examples) yields modest gains, while 100 and 1000 examples produce substantial improvements, though not to the level of pop-out. Notably, these gains extend to out-of-distribution set sizes (50–99 distractors), indicating fine-tuning can provide more generalisable visual search strategies. We probed this further by evaluating transfer to related tasks: Shape-Colour

¹Comprehensive mechanistic interpretability results, including analyses for all three experiments, are provided in Appendix K.

²The 1000-example model was also trained for three epochs, but after the first epoch we observed a collapse in behaviour, with nonsensical outputs. We report only the model after one epoch.

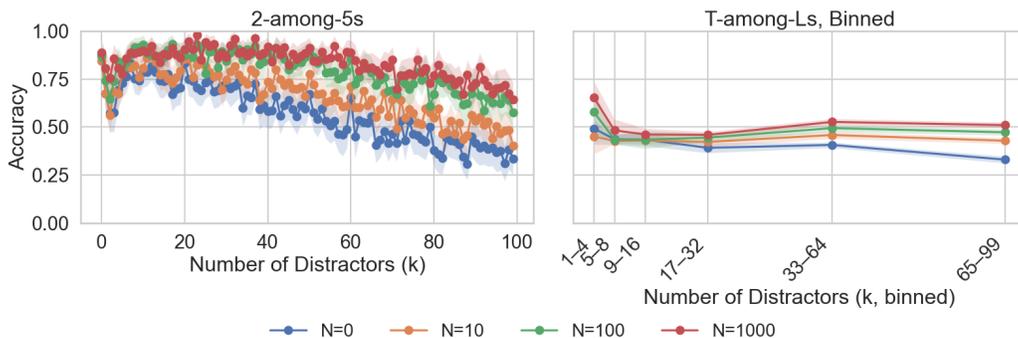


Figure 7: Evaluation of GPT-4o fine-tuned at various dataset sizes for the 2 Among 5 Shape Conjunctive task (left) and the T-among-Ls Shape Conjunctive task (right)

Conjunctive 2 Among 5, Shape Conjunctive “T among L” and the Circle Sizes task (details in Appendix H). As shown in Figure 7, performance improves on the “T among L” task, especially at higher distractor counts, suggesting that the model transfers its search behaviour across shape domains³. No such improvement occurred for Shape-Colour Conjunctive 2 Among 5 or Circle Sizes, indicating transfer is strongest when both training and evaluation conditions target similar feature domains (e.g., shape), and weakest when the downstream task requires integration of additional non-trained feature cues (e.g., shape and colour). The full set of results are plotted in Appendix H.

5 RELATED WORK

We have presented a thorough, systematic investigation of visual-attentional search capabilities in MLLMs, including three distinct search tasks, fine-tuning, and comparisons to human baselines. Our findings align with those of Campbell et al. (2024), who included a brief visual search task within a broader exploration of human-like capacity limits in MLLMs. They reported serial-search-like behaviour in a conjunctive search task and argued that such limitations reflect compositional representations (e.g., representing “blue circle” as “blue” and “circle”, not “BlueCircle”).

However, our work differs from Campbell et al. (2024) in several substantive ways. First, whereas their study focused on target *detection* (present/absent judgements), we examine *localisation* – that is, the ability of models and humans to report *where* a target is, and to do so at multiple levels of spatial precision. Second, we present a more comprehensive evaluation of visual search capabilities. Our experiments span multiple features (e.g., size, shape, colour) across three tasks, and unlike Campbell et al. (2024), we employ matched stimulus conditions to facilitate direct like-for-like comparisons. For example, stimuli in the Shape Conjunctive condition are identical to those in the Disjunctive condition except for the colour manipulation, allowing us to isolate the targeted effect. Third, we directly test whether performance limits in conjunctive search can be mitigated through training. Campbell et al. (2024) speculated that improvements might require the introduction of a serial search mechanism (or an equivalent process for decomposing images) if compositionality was to also be preserved. We tested this hypothesis and found that fine-tuning on a conjunctive search task (2 Among 5) improved performance, even for larger, unseen set sizes, though not to the level of pop-out. Similar training effects have been observed elsewhere (Buschoff et al., 2025), and also in human studies that have shown extensive practice can lead to “unitized” (i.e., non compositional) representations, where test items are perceived more holistically, partially overcoming the binding problem (Czerwinski et al., 1992). Yet, unlike human participants, who typically show limited cross-task transfer (Su et al., 2014; Ding et al., 2023), GPT-4o exhibits mild transfer to a distinct conjunctive search task (T-among-L), suggesting that MLLM performance improvements through fine-tuning are not stimulus-specific (e.g., finding a 5 among 2s), though they may still be limited to a specific feature domain (e.g., shape, not shape and colour). Fourthly, we report a novel result not addressed in Campbell et al. (2024): MLLMs, like humans, incorporate sophisticated priors about physical regularities in the natural world

³Note that these results are binned by number of distractor for clarity. See Appendix H for unbinned.

(e.g., light direction) into their object representations, and these expectations systematically modulate their search performance.

More broadly, visual search can be situated within the field of visual question answering (VQA), where models are evaluated on their ability to reason about visual scenes. MLLMs have been widely applied to this setting, often with architecture-specific optimizations or fine-tuning strategies. Other work has also tried to improve the visual capabilities of MLLMs by augmenting them with specialized mechanisms. V* (Wu & Xie, 2023) incorporates contextual knowledge to guide attention, while the Target and Context-Aware Transformer (Ding et al., 2022) fuses object- and scene-level features for efficient zero-shot visual search. ViSioNS (Travi et al., 2022) introduces scanpath modelling to align image processing with human attentional patterns. For comprehensive overviews of VQA and its extension to multimodal foundation models, we refer the reader to recent surveys (Wu et al., 2017; Kuang et al., 2025). In contrast to most Visual Search work within ML, our goal is not to improve model performance on a specific task. Instead, we treat MLLMs as cognitive systems in their own right, using visual search tasks to examine whether—and *how*—their responses reflect structured visual processing similar to that observed in humans. In this sense our work is more aligned with visual search in cognitive science, where the focus is more on identifying which features test subjects find salient and attention guiding.

6 DISCUSSION AND CONCLUSION

6.1 DISCUSSION

Drawing on classic work in cognitive science, we systematically investigate the visual search capabilities of MLLMs. Our experiments show that the most advanced models (e.g., GPT-4o) closely match human behaviours including: (1) parallel performance – or “pop-out” – for search targets defined by a single primitive feature (i.e., disjunctive search), (2) capacity limits for targets that require feature binding (i.e., conjunctive search), and (3) the incorporation of natural-scene features such as a “light-from-above” prior. These findings not only help us to anticipate future perceptual behaviour of these systems, but also provides insight into the nature of their internal representations, and how these differ between models. For example, Llama 90B exhibited markedly less human-like behaviour, which we attribute to poorer overall perceptual or related auxiliary capabilities rather than differences in architecture, since smaller versions of other models we tested (e.g., GPT-4-Turbo and Claude-Haiku) were also less human-like. Interestingly, MLLMs such as GPT-4o showed evidence of using sophisticated natural scene features, such as lighting direction, to guide visual search, and showed the best performance for objects lit from the most “surprising” direction (below), just like humans. Multimodal large language models are trained on vast, often opaque datasets. While the composition of this training data is unknown—GPT-4o’s system card (OpenAI et al., 2024), for instance, offers only high-level sourcing—much of it likely reflects real-world imagery. Our findings suggest that natural regularities, such as lighting direction, contained in training images has allowed MLLMs to incorporate such features into their object representations—as humans do. Our main findings were also supported by our fine-grained localisation (coordinates) evaluation, measuring spatial precision, and augmented by a fine-tuning experiment that substantially improved conjunctive search performance, though not to the level of parallel pop-out performance. These improvements also generalised to larger, unseen set-sizes and showed mild transfer to a distinct search task targeting the same feature domain (i.e., shape).

6.2 CONCLUSION

Visual search tasks offer a powerful lens on MLLM behaviour, revealing both human-like attentional dynamics and model-specific processing constraints. However, our work has limitations that need to be taken into account. First of all, LLMs are known to be sensitive to the way that prompts are phrased (Alzahrani et al., 2024; Chaudhary et al., 2024) and we explored only a few prompts due to budgetary constraints (see Appendix J). Further, we limited our investigation to three visual features (colour, size, lighting direction), which could be expanded on. Future work could extend this framework to other feature dimensions—such as texture, motion, or occlusion—or investigate how models respond under compositional load or temporal constraints. As multimodal models continue to scale, cognitively informed probes like these may prove important for understanding how they represent, reason about, and act on the visual world.

540 REPRODUCIBILITY STATEMENT

541
542 The implementation of this work will be released upon publication as via a GitHub repository. This
543 will contain all of the code to generate stimuli, run models, and create the results. The raw results
544 themselves will also be made available to allow for additional analysis without running all of the
545 models.

546
547 USE OF LARGE LANGUAGE MODELS

548
549 Large language models were used to assist with grammar and sentence formulation in select places.
550 They were also used to identify additional research literature. LLMs were also used to assist with
551 coding.

552
553 REFERENCES

- 554
555 Wendy J. Adams. A common light-prior for visual search, shape, and reflectance judgments.
556 *Journal of Vision*, 7(11):11, August 2007. ISSN 1534-7362. doi: 10.1167/7.11.11. URL [http://](http://jov.arvojournals.org/article.aspx?doi=10.1167/7.11.11)
557 jov.arvojournals.org/article.aspx?doi=10.1167/7.11.11.
- 558
559 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
560 Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan
561 Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian
562 Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo
563 Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a Visual Lan-
564 guage Model for Few-Shot Learning, 2022. URL <https://arxiv.org/abs/2204.14198>.
Version Number: 2.
- 565
566 Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie,
567 Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M Saiful
568 Bari, and Haidar Khan. When benchmarks are targets: Revealing the sensitivity of large language
569 model leaderboards, 2024. URL <https://arxiv.org/abs/2402.01781>.
- 570
571 Anthropic. Claude 3.5 sonnet model card addendum, 2024. URL [https://www-cdn.](https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf)
572 [anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_](https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf)
Card_Claude_3_Addendum.pdf.
- 573
574 Adam T. Biggs, Michelle R. Kramer, and Stephen R. Mitroff. Using cognitive psychology re-
575 search to inform professional visual search operations. *Journal of Applied Research in Mem-*
576 *ory and Cognition*, 7(2):189–198, 2018. ISSN 2211-3681. doi: [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.jarmac.2018.04.001)
577 [jarmac.2018.04.001](https://www.sciencedirect.com/science/article/pii/S2211368118300044). URL [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/S2211368118300044)
[pii/S2211368118300044](https://www.sciencedirect.com/science/article/pii/S2211368118300044).
- 578
579 Luca M Schulze Buschoff, Konstantinos Voudouris, Elif Akata, Matthias Bethge, Joshua B Tenen-
580 baum, and Eric Schulz. Testing the limits of fine-tuning to improve reasoning in vision language
581 models. *arXiv preprint arXiv:2502.15678*, 2025.
- 582
583 Declan Iain Campbell, Sunayana Rane, Tyler Giallanza, C. Nicolò De Sabbata, Kia Ghods, Amogh
584 Joshi, Alexander Ku, Steven M Frankland, Thomas L. Griffiths, Jonathan D. Cohen, and Tay-
585 lor Whittington Webb. Understanding the limits of vision language models through the lens of
586 the binding problem. In *The Thirty-eighth Annual Conference on Neural Information Processing*
Systems, 2024. URL <https://openreview.net/forum?id=Q5RYn6jagC>.
- 587
588 Manav Chaudhary, Harshit Gupta, Savita Bhat, and Vasudeva Varma. Towards understanding the
589 robustness of llm-based evaluations under perturbations, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2412.09269)
[abs/2412.09269](https://arxiv.org/abs/2412.09269).
- 590
591 Leonardo Chelazzi, Earl K Miller, John Duncan, and Robert Desimone. A neural basis for visual
592 search in inferior temporal cortex. *Nature*, 363(6427):345–347, 1993.
- 593
Mary Czerwinski, Nancy Lightfoot, and Richard M Shiffrin. Automatization and training in visual
search. *The American journal of psychology*, pp. 271–315, 1992.

- 594 Yulong Ding, Tingni Li, and Zhe Qu. Is a new feature learned behind a newly efficient color-
595 orientation conjunction search? *Psychonomic Bulletin & Review*, 30(1):250–260, 2023.
- 596
- 597 Zhiwei Ding, Xuezhe Ren, Erwan David, Melissa Vo, Gabriel Kreiman, and Mengmi Zhang. Efficient
598 zero-shot visual search via target and context-aware transformer, 2022. URL <https://arxiv.org/abs/2211.13470>.
- 599
- 600 James T Enns and Ronald A Rensink. Influence of scene-based properties on visual search. *Science*,
601 247(4943):721–723, 1990.
- 602
- 603 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in
604 vqa matter: Elevating the role of image understanding in visual question answering, 2017. URL
605 <https://arxiv.org/abs/1612.00837>.
- 606
- 607 Jiayi Kuang, Ying Shen, Jingyou Xie, Haohao Luo, Zhe Xu, Ronghao Li, Yinghui Li, Xianfeng
608 Cheng, Xika Lin, and Yu Han. Natural language understanding and inference with mllm in visual
609 question answering: A survey. *ACM Comput. Surv.*, 57(8), March 2025. ISSN 0360-0300. doi:
610 10.1145/3711680. URL <https://doi.org/10.1145/3711680>.
- 611 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping Language-Image
612 Pre-training for Unified Vision-Language Understanding and Generation, 2022. URL <https://arxiv.org/abs/2201.12086>. Version Number: 2.
- 613
- 614 Zihao Lin, Samyadeep Basu, Mohammad Beigi, Varun Manjunatha, Ryan A Rossi, Zichao Wang,
615 Yufan Zhou, Sriram Balasubramanian, Arman Zarei, Keivan Rezaei, et al. A survey on mechanistic
616 interpretability for multi-modal foundation models. *arXiv preprint arXiv:2502.17516*, 2025.
- 617
- 618 Fenglin Liu, Zheng Li, Qingyu Yin, Jinfa Huang, Jiebo Luo, Anshul Thakur, Kim Branson, Patrick
619 Schwab, Bing Yin, Xian Wu, Yefeng Zheng, and David A. Clifton. A multimodal multidomain
620 multilingual medical foundation model for zero shot clinical diagnosis. *npj Digital Medicine*,
621 8(1):86, February 2025. ISSN 2398-6352. doi: 10.1038/s41746-024-01339-7. URL <https://www.nature.com/articles/s41746-024-01339-7>.
- 622
- 623 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning, 2023. URL
624 <https://arxiv.org/abs/2304.08485>. Version Number: 2.
- 625
- 626 Gen Luo, Ganlin Yang, Ziyang Gong, Guanzhou Chen, Haonan Duan, Erfei Cui, Ronglei Tong, Zhi
627 Hou, Tianyi Zhang, Zhe Chen, Shenglong Ye, Lewei Lu, Jingbo Wang, Wenhai Wang, Jifeng
628 Dai, Yu Qiao, Rongrong Ji, and Xizhou Zhu. Visual Embodied Brain: Let Multimodal Large
629 Language Models See, Think, and Control in Spaces, May 2025. URL <http://arxiv.org/abs/2506.00123>. arXiv:2506.00123 [cs].
- 630
- 631 Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A Visual
632 Question Answering Benchmark Requiring External Knowledge. In *2019 IEEE/CVF Conference*
633 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 3190–3199, Long Beach, CA, USA,
634 June 2019. IEEE. ISBN 978-1-72813-293-8. doi: 10.1109/CVPR.2019.00331. URL <https://ieeexplore.ieee.org/document/8953725/>.
- 635
- 636 David Marr. *Vision: a computational investigation into the human representation and processing of*
637 *visual information*. Freeman, New York, 14. print edition, 2000. ISBN 978-0-7167-1567-2.
- 638
- 639 Brian McElree and Marisa Carrasco. The temporal dynamics of visual search: evidence for parallel
640 processing in feature and conjunction searches. *Journal of Experimental Psychology: Human*
641 *Perception and Performance*, 25(6):1517, 1999.
- 642
- 643 Meta AI. Llama 3.2-90b vision model card. [https://huggingface.co/meta-llama/](https://huggingface.co/meta-llama/Llama-3.2-90B-Vision)
644 [Llama-3.2-90B-Vision](https://huggingface.co/meta-llama/Llama-3.2-90B-Vision), 2024. Accessed: 2025-05-13.
- 645
- 646 Stephen R. Mitroff, Justin M. Ericson, and Benjamin Sharpe. Predicting Airport Screening Officers’
647 Visual Search Competency With a Rapid Assessment. *Human Factors*, 60(2):201–211, 2018. doi:
10.1177/0018720817743886. URL <https://doi.org/10.1177/0018720817743886>.
_eprint: <https://doi.org/10.1177/0018720817743886>.

648 OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan
649 Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-
650 Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex
651 Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau,
652 Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian,
653 Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein,
654 Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrew
655 Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia,
656 Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben
657 Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake
658 Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon
659 Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo
660 Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li,
661 Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,
662 Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim,
663 Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley
664 Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler,
665 Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki,
666 Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay,
667 Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric
668 Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani,
669 Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh,
670 Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang
671 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik
672 Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung,
673 Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu,
674 Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon,
675 Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie
676 Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe,
677 Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi
678 Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers,
679 Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan
680 Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh
681 Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn
682 Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra
683 Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe,
684 Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman,
685 Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng,
686 Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Fevrier, Lu Zhang, Lukas Kondraciuk,
687 Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine
688 Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin
689 Zhang, Marwan Aljubei, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank
690 Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna
691 Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle
692 Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles
693 Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho
694 Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine,
695 Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige,
696 Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko,
697 Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick
698 Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan,
699 Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal,
700 Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo
701 Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob
Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory
Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi
Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara
Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu
Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer

- 702 Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal
703 Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas
704 Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao
705 Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan,
706 Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie
707 Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra,
708 Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang,
709 Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024.
710 URL <https://arxiv.org/abs/2410.21276>.
- 711 Julius Orłowski, Christian Beissel, Friederike Rohn, Yair Adato, Hermann Wagner, and Ohad
712 Ben-Shahar. Visual pop-out in barn owls: Human-like behavior in the avian brain. *Journal of*
713 *Vision*, 15(14):4, October 2015. ISSN 1534-7362. doi: 10.1167/15.14.4. URL [http://jov.](http://jov.arvojournals.org/article.aspx?doi=10.1167/15.14.4)
714 [arvojournals.org/article.aspx?doi=10.1167/15.14.4](http://jov.arvojournals.org/article.aspx?doi=10.1167/15.14.4).
- 715 Michael J. Proulx. Size Matters: Large Objects Capture Attention in Visual Search. *PLoS ONE*,
716 5(12):e15293, December 2010. ISSN 1932-6203. doi: 10.1371/journal.pone.0015293. URL
717 <https://dx.plos.org/10.1371/journal.pone.0015293>.
- 718 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
719 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
720 Learning Transferable Visual Models From Natural Language Supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
721 Version Number: 1.
- 722 Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy.
723 Do vision transformers see like convolutional neural networks? *Advances in neural information*
724 *processing systems*, 34:12116–12128, 2021.
- 725 Adam Reichenenthal, Mor Ben-Tov, Ohad Ben-Shahar, and Ronen Segev. What pops out for you pops
726 out for fish: Four common visual features. *Journal of Vision*, 19(1):1, January 2019. ISSN 1534-
727 7362. doi: 10.1167/19.1.1. URL [http://jov.arvojournals.org/article.aspx?](http://jov.arvojournals.org/article.aspx?doi=10.1167/19.1.1)
728 [doi=10.1167/19.1.1](http://jov.arvojournals.org/article.aspx?doi=10.1167/19.1.1).
- 729 Manoosh Samiei and James J. Clark. Target features affect visual search, a study of eye fixations,
730 2022. URL <https://arxiv.org/abs/2209.13771>.
- 731 Missie Smith, Jillian Streeter, Gary Burnett, and Joseph L Gabbard. Visual search tasks: the effects
732 of head-up displays on driving and task performance. In *Proceedings of the 7th international*
733 *conference on Automotive User Interfaces and Interactive Vehicular Applications*, pp. 80–87, 2015.
- 734 Yuling Su, Yunpeng Lai, Wanyi Huang, Wei Tan, Zhe Qu, and Yulong Ding. Short-term perceptual
735 learning in visual conjunction search. *Journal of Experimental Psychology: Human Perception*
736 *and Performance*, 40(4):1415, 2014.
- 737 Fermín Travi, Gonzalo Ruarte, Gaston Bujia, and Juan Esteban Kamienkowski. Visions:
738 Visual search in natural scenes benchmark. In S. Koyejo, S. Mohamed, A. Agar-
739 wal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Pro-*
740 *cessing Systems*, volume 35, pp. 11987–12000. Curran Associates, Inc., 2022. URL
741 [https://proceedings.neurips.cc/paper_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/4eal4e6090343523ddcd5d3ca449695f-Paper-Datasets_and_Benchmarks.pdf)
742 [4eal4e6090343523ddcd5d3ca449695f-Paper-Datasets_and_Benchmarks.](https://proceedings.neurips.cc/paper_files/paper/2022/file/4eal4e6090343523ddcd5d3ca449695f-Paper-Datasets_and_Benchmarks.pdf)
743 [pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/4eal4e6090343523ddcd5d3ca449695f-Paper-Datasets_and_Benchmarks.pdf).
- 744 Anne M. Treisman and Garry Gelade. A feature-integration theory of attention. *Cognitive Psychology*,
745 12(1):97–136, January 1980. ISSN 00100285. doi: 10.1016/0010-0285(80)90005-5. URL
746 <https://linkinghub.elsevier.com/retrieve/pii/0010028580900055>.
- 747 Jeremy M. Wolfe. Guided Search 2.0 A revised model of visual search. *Psychonomic Bulletin &*
748 *Review*, 1(2):202–238, June 1994. ISSN 1069-9384, 1531-5320. doi: 10.3758/BF03200774. URL
749 <http://link.springer.com/10.3758/BF03200774>.
- 750 Jeremy M. Wolfe. What Can 1 Million Trials Tell Us About Visual Search? *Psychological Science*, 9
751 (1):33–39, January 1998. ISSN 0956-7976, 1467-9280. doi: 10.1111/1467-9280.00006. URL
752 <https://journals.sagepub.com/doi/10.1111/1467-9280.00006>.

- 756 Jeremy M. Wolfe. Visual Search: How Do We Find What We Are Looking For? *Annual Review of Vision Science*, 6(1):539–562, September 2020. ISSN 2374-4642, 2374-4650.
757 doi: 10.1146/annurev-vision-091718-015048. URL [https://www.annualreviews.org/](https://www.annualreviews.org/doi/10.1146/annurev-vision-091718-015048)
758 doi/10.1146/annurev-vision-091718-015048.
759
- 760 Jeremy M. Wolfe and Todd S. Horowitz. What attributes guide the deployment of visual attention
761 and how do they do it? *Nature Reviews Neuroscience*, 5(6):495–501, June 2004. ISSN 1471-
762 003X, 1471-0048. doi: 10.1038/nrn1411. URL [https://www.nature.com/articles/](https://www.nature.com/articles/nrn1411)
763 nrn1411.
764
- 765 Jeremy M. Wolfe, Evan M. Palmer, and Todd S. Horowitz. Reaction time distributions constrain
766 models of visual search. *Vision Research*, 50(14):1304–1311, June 2010. ISSN 00426989. doi:
767 10.1016/j.visres.2009.11.002. URL [https://linkinghub.elsevier.com/retrieve/](https://linkinghub.elsevier.com/retrieve/pii/S0042698909005021)
768 pii/S0042698909005021.
- 769 Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms,
770 2023. URL <https://arxiv.org/abs/2312.14135>.
771
- 772 Qi Wu, Damien Teney, Peng Wang, Chunhua Shen, Anthony Dick, and Anton van den Hengel. Visual
773 question answering: A survey of methods and datasets. *Computer Vision and Image Understanding*,
774 163:21–40, 2017. ISSN 1077-3142. doi: <https://doi.org/10.1016/j.cviu.2017.05.001>. URL <https://www.sciencedirect.com/science/article/pii/S1077314217300772>. Lan-
775 guage in Vision.
776
- 777 Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu
778 Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin,
779 Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.
780 MMMU: A Massive Multi-Discipline Multimodal Understanding and Reasoning Benchmark
781 for Expert AGI. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
782 (CVPR), pp. 9556–9567, Seattle, WA, USA, June 2024. IEEE. ISBN 9798350353006. doi:
783 10.1109/CVPR52733.2024.00913. URL [https://ieeexplore.ieee.org/document/](https://ieeexplore.ieee.org/document/10656299/)
784 10656299/.
- 785 Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
786 *European conference on computer vision*, pp. 818–833. Springer, 2014.
787
- 788 Li Zhaoping and Keith A May. Psychophysical tests of the hypothesis of a bottom-up saliency map
789 in primary visual cortex. *PLoS Computational Biology*, 3(4):e62, 2007.
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A FINE-GRAINED LOCALISATION RESULTS

Here we present the results for the fine-grained localisation ("coordinates") evaluation. The addition of coordinates (instead of just the Cells variant) allows us to capture different levels of visual understanding. *Cells* requires only a coarse-grained ability to locate the object, while *Coordinates* demands fine-grained spatial precision. Although classical visual search studies typically focus on detection, we argue that localisation is especially relevant in the context of MLLMs, which are intended to act on or reason over visual scenes—tasks that depend on accurate visual perception. In both *Cells* and *Coordinates* conditions, correct localisation is defined by the centre of the target object. For *Cells*, this determines the ground-truth grid cell label; for *Coordinates*, it serves as the reference point for Euclidean distance evaluation.

We note that Claude Sonnet frequently declined to provide coordinates in this task. In these cases, it instead responded that no target could be identified. We treat such refusals as maximum reasonable localisation error trials ($\sqrt{400^2 + 400^2} \approx 566$ pixels). Similarly, Llama 90B would often report coordinates outside of the 400×400 range. Here we score models normally according to Euclidean distance. We examine the phenomena of invalid responses in more detail in Appendix D.

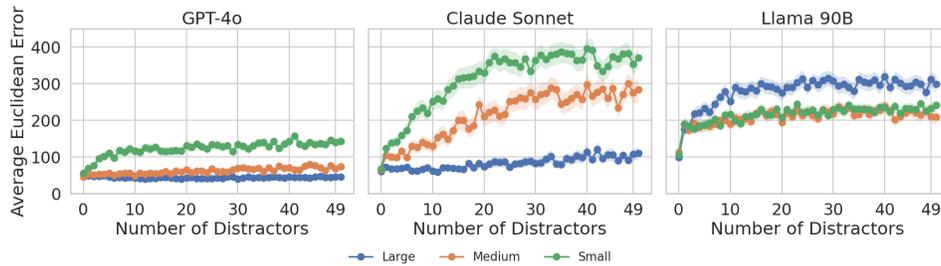


Figure 8: Results for Circle Sizes on Coordinates. The shaded region denotes the 95% confidence interval.

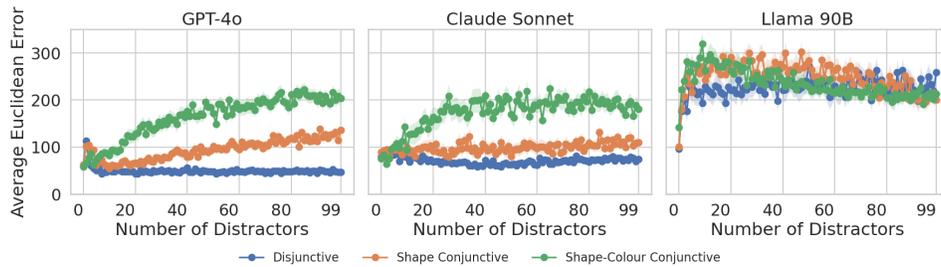


Figure 9: Results for the 2Among5 task — Coordinates mode. The shaded region denotes the 95% confidence interval.

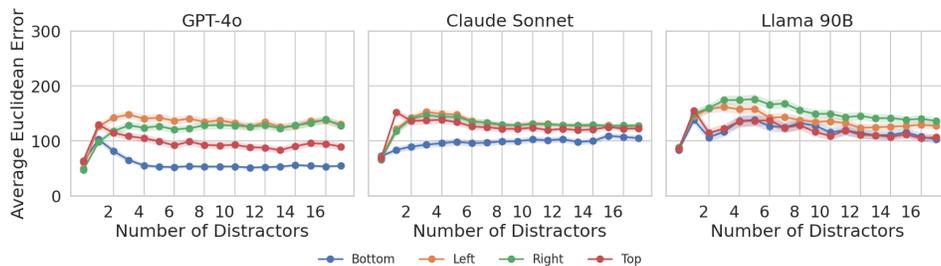


Figure 10: Results for the Light Priors task — Coordinates mode. The shaded region denotes the 95% confidence interval.

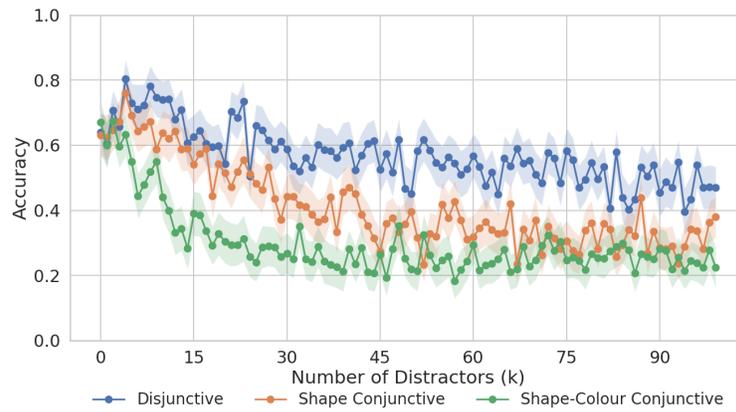
B RESULTS FOR ADDITIONAL MODELS

Due to space constraints it was not possible to provide figures for each model tested on each task within the main body of the paper. In this section we provide the results for other tested models and compare the results to humans and other models. We principally compare against smaller models (or earlier models in the same families) as those in the main paper; intending to enable a comparison of how visual search mechanisms change with scale. Therefore we include GPT-4-Turbo, Claude-Haiku (3.5), and Llama (3.2) 11B. With the exception of GPT-4-Turbo, these models perform significantly worse than their larger versions, and quickly fall to chance level performance, even in disjunctive cases, or other variants where pop-out effects were present. We additionally provide results for models from the Qwen family on the cells variant of the task (Qwen 2.5 7B VL Instruct and Qwen 2.5 32B VL Instruct). These two models also both perform worse than the models evaluated in the main body of the paper (likely due to model size).

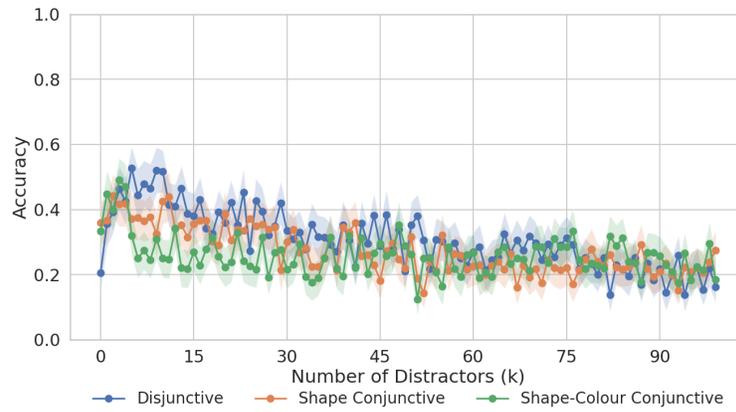
Two Among Five additional Results Figures 11 and 13 detail the results for our smaller models. For Claude Haiku and Llama 11B in the Cells evaluation there is little difference between their performance and effectively random chance. GPT-4-Turbo on the other hand performs better, in similar pattern to the MLLMs in the main paper and humans, but with reduced performance. In the Coordinates evaluation, all three perform only marginally better in the disjunctive setting, though the distinction can become clearer at a higher set-size. Figure 12 provide the results for the two Qwen models. Qwen 7B performs particularly poorly and is largely indistinguishable from guessing. Qwen 32B performs marginally better, notably so in the disjunctive case.

Light Priors Figures 14 and 16 detail the results for GPT-4-Turbo, Claude-Haiku, and Llama 11B in the Light Priors task, for both the Cells and Coordinates Evaluation. Similar to the Two Among Five results, these less capable models perform much worse than the models in the main body of the paper. Notably, Claude-Haiku, manages to robustly attain worse-than-chance performance in the Cells evaluation. Upon inspection it became evident that Haiku was selecting invalid options such as “Cell (2,4)” and “Cell (2,3)” frequently. It is unclear what caused this specifically for the Light Priors task and Claude-Haiku, when other models and tasks were largely unaffected. Figure 15 presents the results for our Qwen models. Once again, Qwen 7B performs extremely poorly. Qwen 32B also performs poorly here but for low distractor numbers does perform marginally better on the bottom / top conditions when compared to left/right.

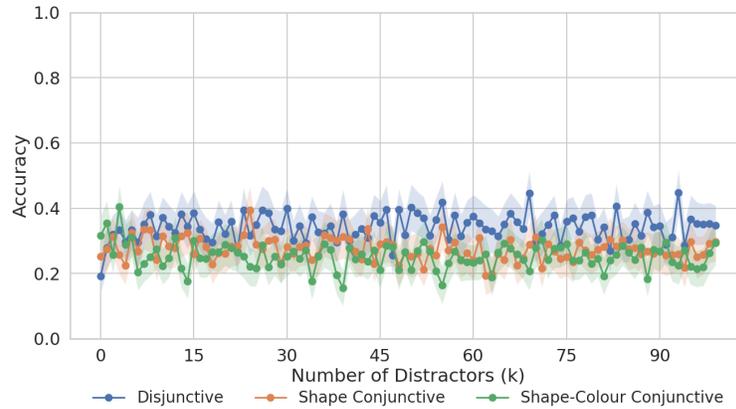
Circle Sizes Figures 18 and 19 detail the results for our smaller models on the Circle Sizes task, again for both Cells and Coordinates evaluations. Once again, we see substantially worse performance compared to the larger models in the main paper, with Llama 11B in particular not performing better than chance, and Haiku and GPT-4-turbo only performing marginally better than chance for the Large variant of the task. In terms of coordinates set up, GPT-4-Turbo and Claude-Haiku again perform slightly better in the Large variant. The results for Qwen on the coordinates set up are given in Figure 18. Here both models perform poorly. Yet, again, Qwen 32B performs marginally better in the easiest condition “Large”.



(a) GPT-4-Turbo

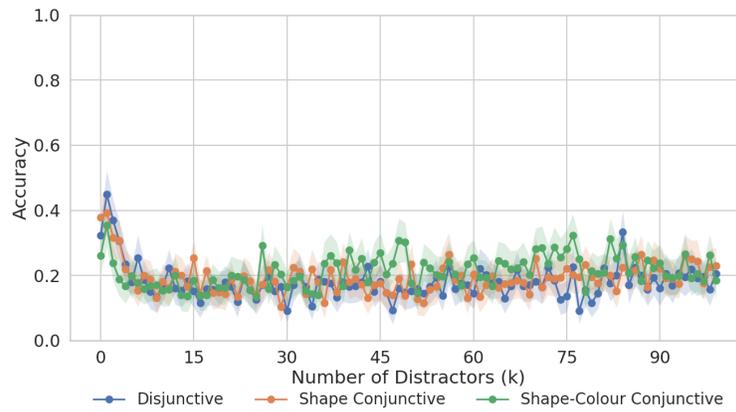


(b) Claude Haiku

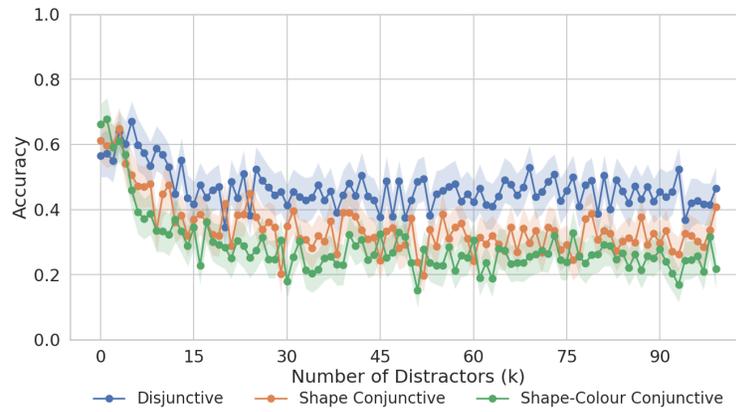


(c) Llama 11B

Figure 11: Results for the 2Among5 task using Cells modes for our three smaller or earlier models. The shaded region denotes the 95% confidence interval.



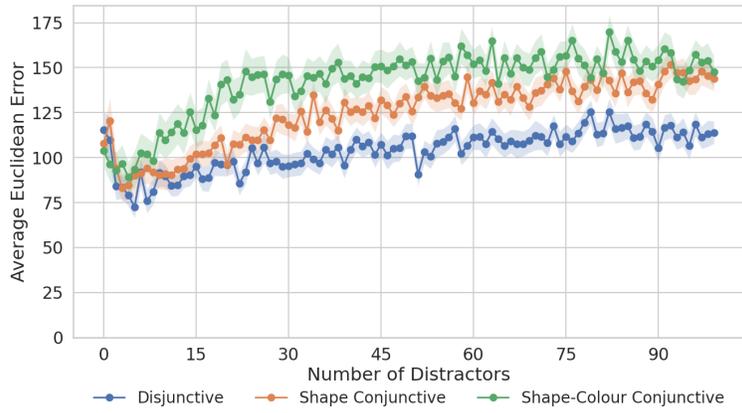
(a) Qwen 7B



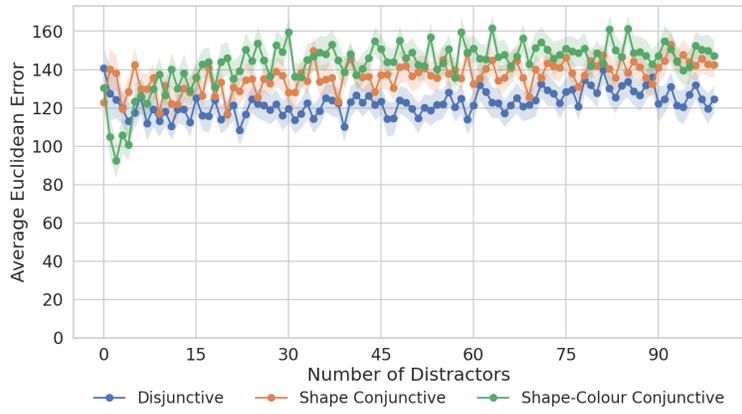
(b) Qwen 32B

Figure 12: Results for the 2Among5 task using Cells modes for Qwen 7B and Qwen 32B. The shaded region denotes the 95% confidence interval.

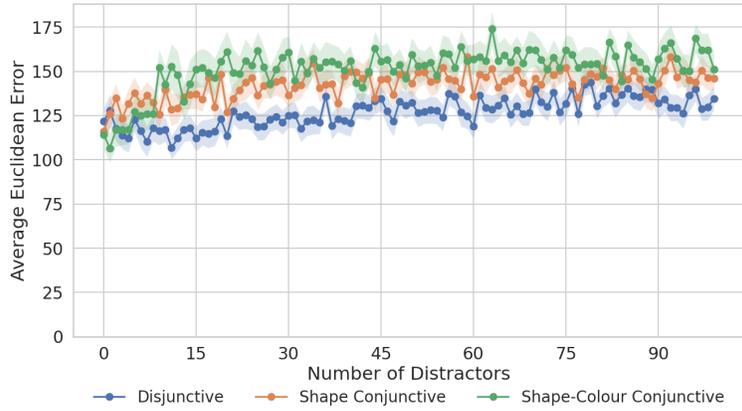
1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079



(a) GPT-4-Turbo



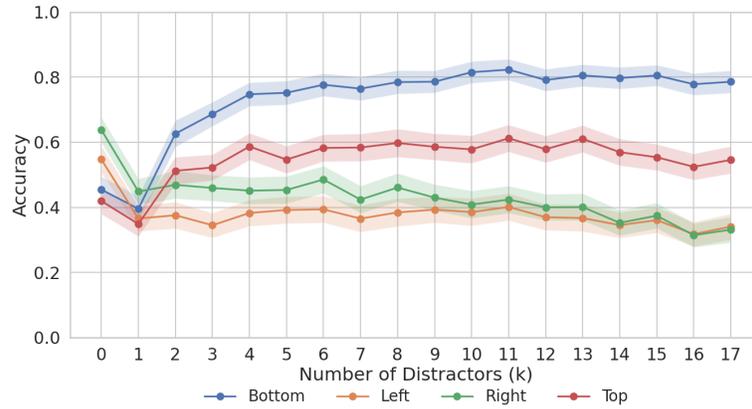
(b) Claude Haiku



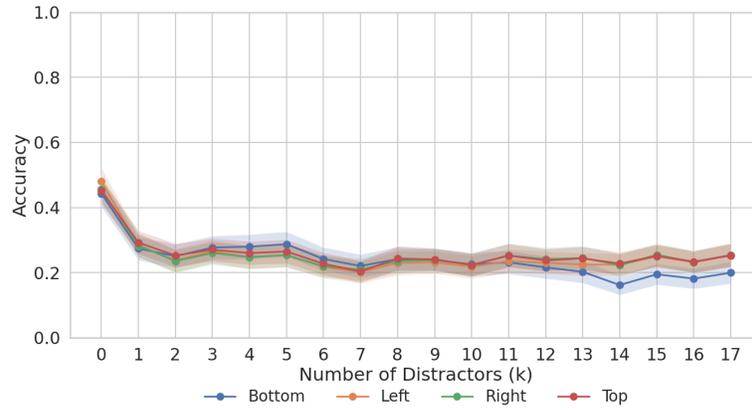
(c) Llama 11B

Figure 13: Results for the 2Among5 task using Coordinates modes for our three smaller or earlier models. The shaded region denotes the 95% confidence interval.

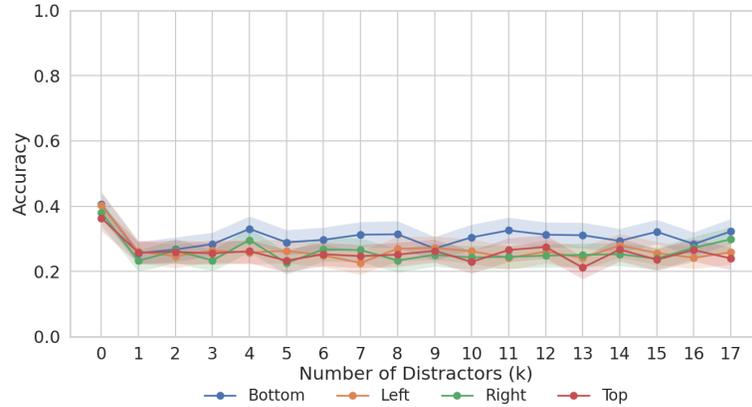
1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133



(a) GPT-4-Turbo

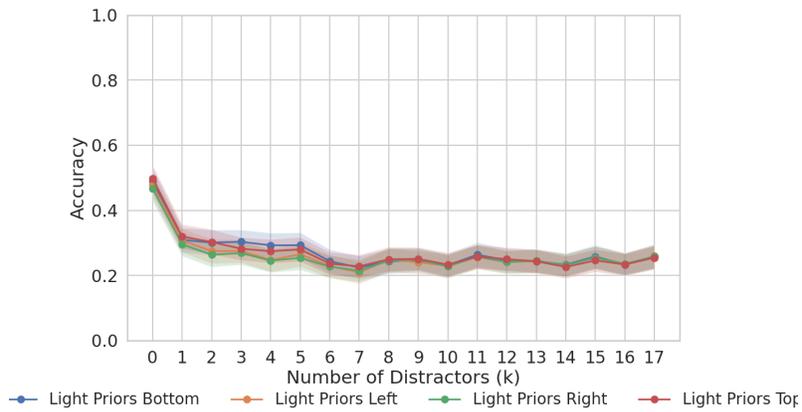


(b) Claude Haiku

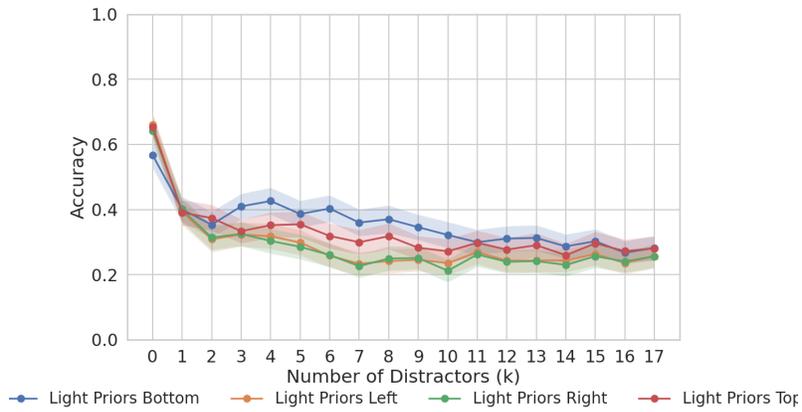


(c) Llama 11B

Figure 14: Results for the Light Priors task using Cells modes for our three smaller or earlier models. The shaded region denotes the 95% confidence interval.



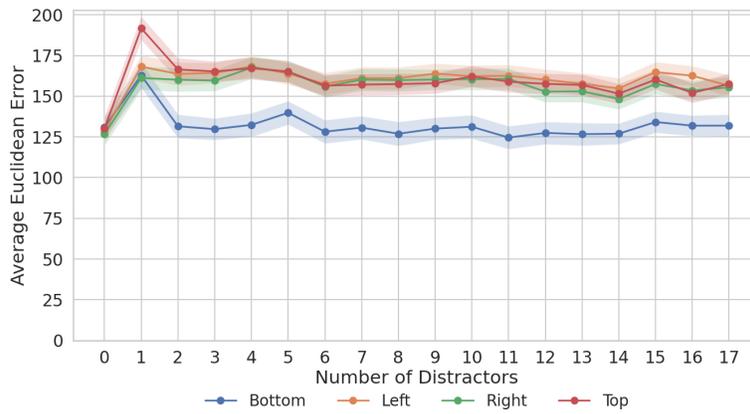
(a) Qwen 7B



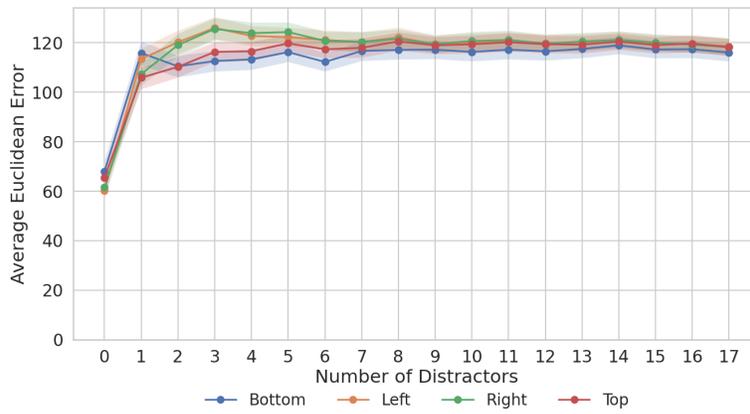
(b) Qwen 32B

Figure 15: Results for the Light Priors task using Cells modes for Qwen 7B and Qwen 32B. The shaded region denotes the 95% confidence interval.

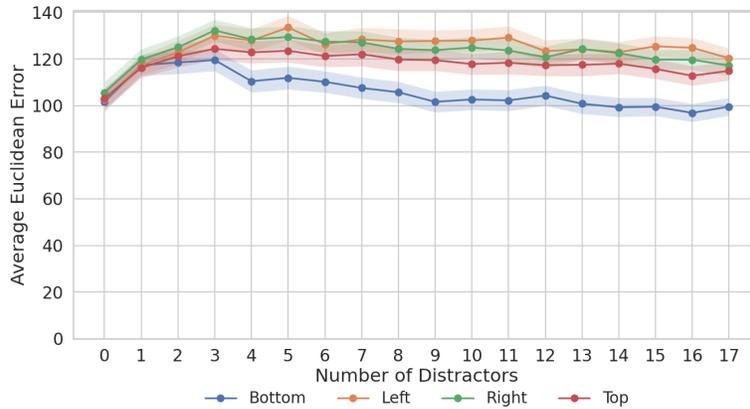
1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241



(a) GPT-4-Turbo

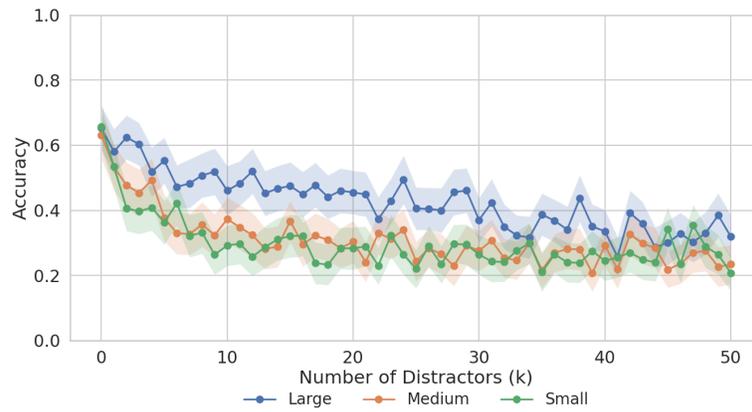


(b) Claude Haiku

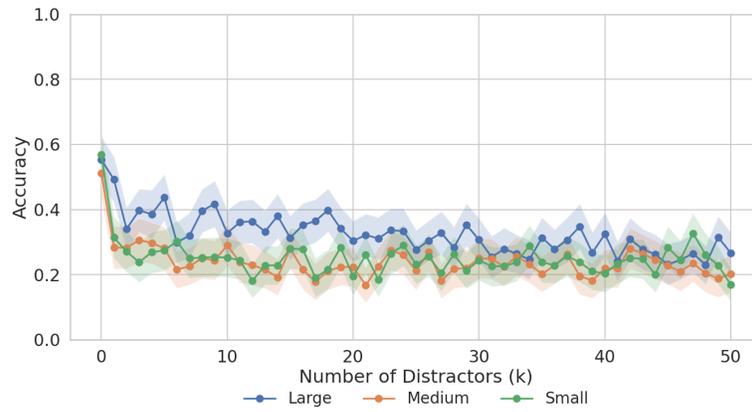


(c) Llama 11B

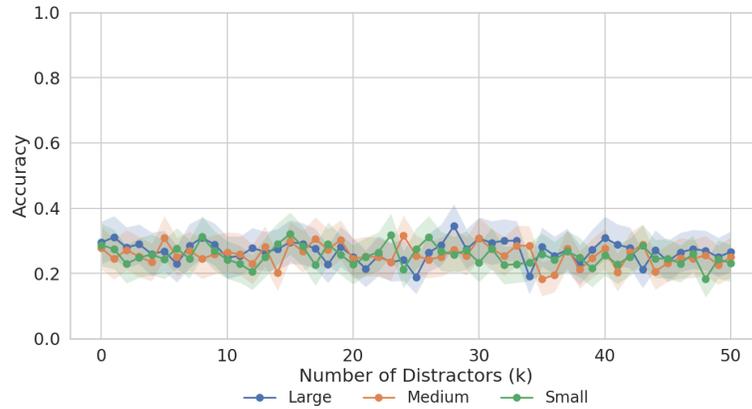
Figure 16: Results for the Light Priors task using Coordinates modes for our three smaller or earlier models. The shaded region denotes the 95% confidence interval.



(a) GPT-4-Turbo

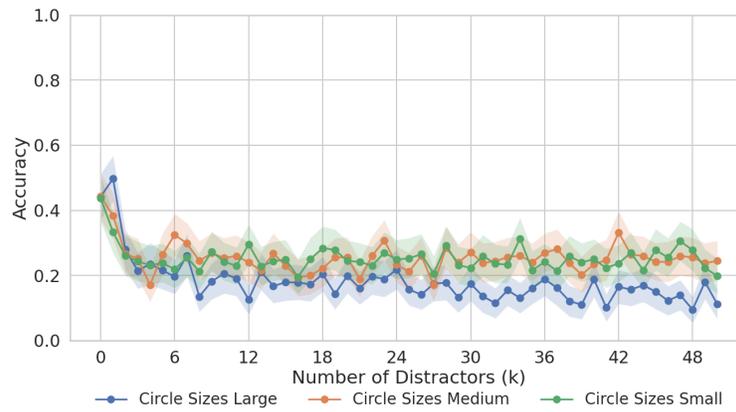


(b) Claude Haiku

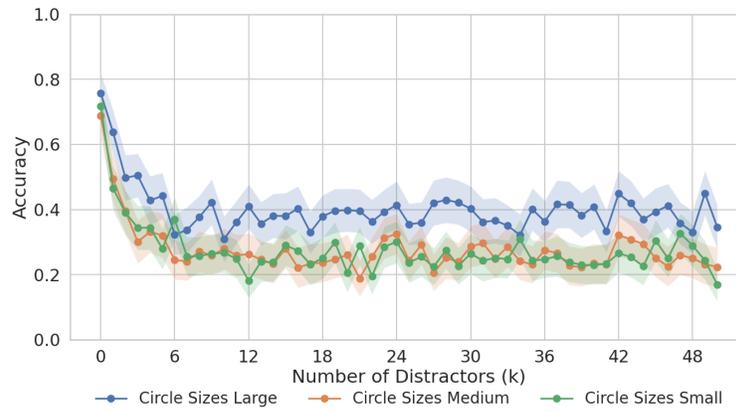


(c) Llama 11B

Figure 17: Results for the Circle Sizes task using Cells modes for our three smaller or earlier models. The shaded region denotes the 95% confidence interval.



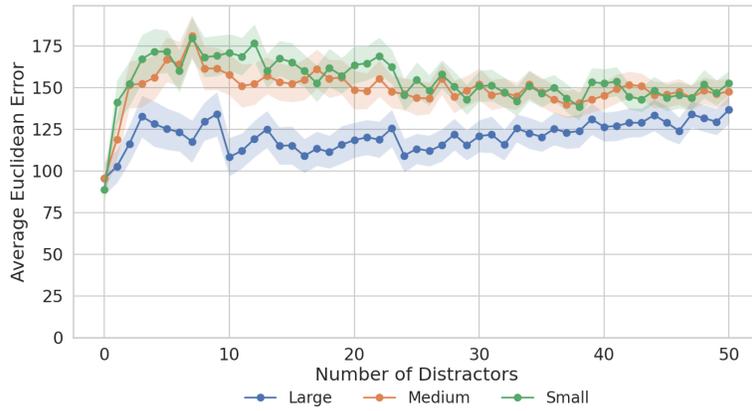
(a) Qwen 7B



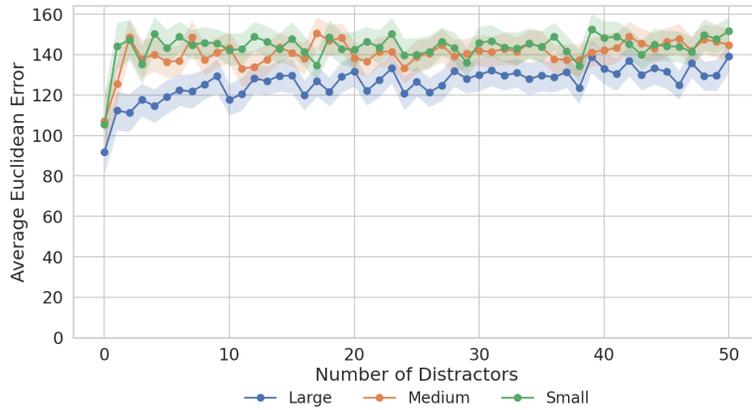
(b) Qwen 32B

Figure 18: Results for the Circle Sizes task using Cells modes for Qwen 7B and Qwen 32B. The shaded region denotes the 95% confidence interval.

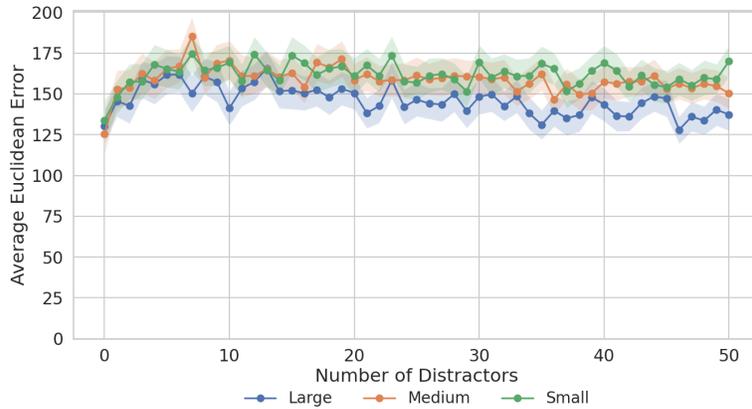
1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403



(a) GPT-4-Turbo



(b) Claude Haiku



(c) Llama 11B

Figure 19: Results for the Circle Sizes task using Coordinates modes for our three smaller or earlier models. The shaded region denotes the 95% confidence interval.

C MODEL DETAILS

In this section we provide more detail on the specific models used, hyper-parameters and overall settings. Most hyper-parameters were left at their default settings in order to be the most “out-of-the-box” representation of each model. We did set for each model a temperature of 0.0 to ensure as deterministic a response as possible. The specific versions of models used are described in Table 2.

When running experiments for the Llama models we utilised a pre-existing implementation via Hugging Face Transformers ⁴. Models were instantiated with `MllamaForConditionalGeneration` and paired with their corresponding `AutoProcessor`. All inference ran on GPU using automatic device mapping and `bfloat16` precision. Crucially, we did not enable sampling (`do_sample=False` by default), so generation was greedy — i.e., the model always selected the highest-probability token at each step. Thus, the model’s output was deterministic and equivalent to using temperature 0. Qwen models were run via an API on FireworksAI ⁵, again with 0 temperature, and full FP16 bit precision.

Table 2: Language models and their specific versions used in our experiments.

Model	Version Used
GPT-4o	gpt-4o-2024-08-06
GPT-4-Turbo	gpt-4-turbo-2024-04-09
Claude Sonnet	claude-3-5-sonnet-20241022)
Claude Haiku	claude-3-5-haiku-20241022
Llama 90B	Meta LLaMA 3.2 90B
Llama 11B	Meta LLaMa 3.2 11B
Qwen 7B	Qwen 2.5 7B VI Instruct
Qwen 32B	Qwen 2.5 32B VI Instruct

D INVALID RESPONSES

Despite efforts to ensure models responded appropriately, the stochastic nature of MLLMs led to invalid or nonsensical results. For the Cells evaluations, these were infrequent, and detailed in Tables 5, 6, 7, 8, 9, 10, 11, 12, and 13. To summarise, in 2 Among 5, All models would occasionally respond with invalid Cells (such as Cell (2,3)). The biggest offender was Claude-Haiku, however, GPT-4o and Claude-Sonnet would also do this between 3 and 5% of the time in some variations. For Light Priors, Claude-Haiku was again the biggest culprit, in some cases providing invalid responses over 78% of the time. Again, GPT-4o would do this over 8% of the time in the Top and Bottom variations—interesting because it performed better in these variants overall. For Circle Sizes, GPT-4o responded invalidly for 2.78% of instances for the small variation, while Llama 90B did the same for over 1% of cases in all variations. Again, Haiku frequently provided invalid responses (between 4.49 and 13.18 %). For the Qwen models, invalid rates were general low, with the exception of Qwen 7B for the 2Among5 task, where they ranged from 15.83% to 46.02%).

In the Coordinates evaluations, we can distinguish between models refusing to provide coordinates, and providing implausible coordinates outside of the image size.

In 2 Among 5, we see models failing to provide coordinates as outlined in Table 3. Again, Sonnet is the largest contributor, rising to over 8% invalid in the Conjunctive variant. Other models are all less than 1%.

Llama 90B is the only model to provide coordinates outside of the expected range (outside the 400x400 pixel range). Strangely, it does this more frequently for the easier disjunctive variant (42.16%) than the more difficult Shape Conjunctive (39.34%) or Shape-Colour Conjunctive (18.2%) variations.

⁴<https://huggingface.co/docs/transformers/index>

⁵<https://fireworks.ai/>

1458 Table 3: Rate of failure to provide coordinates (%) by model and logical condition in the 2 Among 5
 1459 task.

1461 Model	1462 Disjunctive	1463 Shape Conjunctive	1464 Shape-Colour Conjunctive
1465 claude-haiku	0.00	0.00	0.00
1466 claude-sonnet	0.17	3.21	8.34
1467 gpt-4-turbo	0.00	0.00	0.00
1468 gpt-4o	0.03	0.39	0.91
1469 llama11B	0.02	0.05	0.18
1470 llama90B	0.00	0.00	0.01

1471 For the Light Priors task, we detail failure to provide coordinate rates in Table 4. Models generally
 1472 provided coordinates in this task, with no model ever approaching even a 1% failure rate. Llama 90B
 1473 continues to provide coordinates that do not make sense given the question. Returning coordinates
 1474 outside of the range mostly frequently in Bottom (9.58%) and Top (6.86%), but also present in Right
 1475 (0.93% and Left (1.02%). This is odd, as generally Llama 90B performs much better on the Vertical
 1476 (Top and Bottom) stimuli, but when mistakes are made, it’s frequently because of this coordinate
 1477 error.

1478 Table 4: Rate of failure to provide coordinates (%) by model and lighting direction in the Light Priors
 1479 Task.

1480 Model	1481 Bottom	1482 Left	1483 Right	1484 Top
1485 claude-haiku	0.00	0.00	0.00	0.00
1486 claude-sonnet	0.00	0.00	0.00	0.00
1487 gpt-4-turbo	0.00	0.00	0.00	0.00
1488 gpt-4o	0.09	0.02	0.01	0.16
1489 llama11B	0.00	0.00	0.00	0.00
1490 llama90B	0.00	0.00	0.00	0.00

1491 In the Circle Sizes task, Claude would fail to provide a coordinate pair in 2.7% (Large), 23.42%
 1492 (Medium), and 39.2% (Small). GPT-4o failed to provide coordinates in 0.24% of Small instances.
 1493 All other models always provided coordinates.

1494 Llama 90B would again often provide coordinates outside of the expected range. Again, more often
 1495 in the easier Large variation (52.7% of the time) than the Medium (10.56%) or Small variations
 1496 (1.6%). It is unclear why Llama 90B responded like this. No other model ever provided coordinates
 1497 outside of the 400x400 range in this task.

1498 For the finetuned models described in Section 4 and Appendix H, the number of invalid responses
 1499 decayed as finetuning went on, eventually reaching 0 invalid responses at n=100 and n=1000.

1500 E SPATIAL BIAS

1501 In this section we investigate the spatial bias and preferences exhibited by models in our experiments.
 1502 For each experiment in the Cells evaluation we break down performance by cell category (i.e., which
 1503 quadrant of the screen the target was in). We present Precision, Recall, and the Selection proportion
 1504 (i.e., in what proportion of trials did the model select this cell?).

1505 Overall we uncover that different models have different tendencies to pick particular cells more
 1506 frequently than others, even though the set of tasks had approximately an equal number of trials
 1507 where the target was in each quadrant. However, these preferences or biases do not seem to carry
 1508 over from task to task.

1509 We also observed small elements of spatial bias in the human results (See tables in I.3. This makes
 1510 more sense as the humans were exposed the stimulus for a short period of time and consistent search
 1511 strategies (e.g., top-down) would lead to apparent spatial bias if there wasn’t sufficient time to find

the target. However, MLLMs were not subject to these timing constraints and are simply required to transform the input image-prompt pair into an answer.

E.1 2 AMONG 5

The spatial bias results for GPT-4o, Claude-Sonnet, and Llama 90B are presented in Table 5. From the selection proportions it is apparent that GPT-4o has a preference for selecting the top-right corner, and to some extent, the bottom right quadrant. The top-right preference is most clear in the hardest task, Shape-Colour Conjunctive, where GPT-4o selects it 47.9% of the time. Conversely, Claude-sonnet clearly prefers the bottom-left quadrant, and in the two harder task variants (Shape Conjunctive and Shape-Colour Conjunctive) it selects it over 50% of the time. Finally, Llama 90B also exhibits a preference for the bottom-right quadrant. The spatial bias results for the additional models are provided in Table 6. Here all models are heavily biased towards the bottom right quadrant. This likely explains these models lower performance. It is unclear where this preference has come from. The Qwen models are presented in Table 7. Qwen 7B presents high rates of responding with invalid cells, and a bias towards the bottom of the screen when correct. Qwen 32B is more likely to give a valid answer but is extremely biased to the bottom right cell.

Table 5: Model classification performance and predicted proportions on the 2Among5 task GPT-4o, Claude-Sonnet, and Llama 90B. “Sel” denotes the proportion of predictions (in %) assigned to each label.

(a) GPT-4o

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.89	0.78	22.1	0.94	0.30	8.0	0.50	0.23	11.6
Top Right	0.76	0.98	32.9	0.58	0.70	30.8	0.36	0.68	47.9
Bottom Left	0.90	0.82	23.0	0.71	0.48	16.8	0.48	0.37	19.0
Bottom Right	0.92	0.81	21.2	0.47	0.75	38.6	0.49	0.36	17.9
Invalid			0.80			5.80			3.54

(b) Claude-Sonnet

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.96	0.26	6.8	0.93	0.21	5.7	0.54	0.17	7.7
Top Right	0.61	0.75	31.2	0.59	0.58	24.7	0.47	0.28	15.1
Bottom Left	0.58	0.94	40.5	0.46	0.94	50.6	0.32	0.74	58.1
Bottom Right	0.85	0.76	21.5	0.84	0.63	18.4	0.46	0.28	15.1
Invalid			0.02			0.53			3.90

(c) LLaMA 90B

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.71	0.48	17.3	0.79	0.28	9.0	0.33	0.46	34.8
Top Right	0.48	0.21	11.0	0.42	0.09	5.5	0.35	0.04	3.1
Bottom Left	0.65	0.64	24.9	0.52	0.47	22.7	0.29	0.25	22.2
Bottom Right	0.45	0.88	46.8	0.35	0.82	56.9	0.30	0.48	39.4
Invalid			0.07			5.83			0.56

Table 6: Classification performance and predicted proportions on the 2Among5 Task for GPT-4-Turbo, Claude-Haiku, and Llama 11B. “Sel” indicates the proportion of model predictions assigned to each label (in %).

(a) GPT-4-turbo

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.93	0.20	5.3	0.88	0.11	3.1	0.72	0.05	1.7
Top Right	0.60	0.63	26.9	0.65	0.21	8.1	0.42	0.17	9.9
Bottom Left	0.75	0.54	18.3	0.56	0.47	20.7	0.32	0.41	32.6
Bottom Right	0.46	0.93	48.0	0.32	0.88	66.7	0.27	0.58	52.8
Invalid			1.49			1.43			3.02

(b) Claude-Haiku

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	1.00	0.07	1.7	1.00	0.03	0.7	0.93	0.01	0.3
Top Right	0.46	0.09	4.8	0.50	0.03	1.6	0.34	0.01	0.7
Bottom Left	0.68	0.25	9.4	0.68	0.18	6.6	0.50	0.05	2.5
Bottom Right	0.29	0.85	69.9	0.26	0.88	81.8	0.25	0.99	95.9
Invalid			14.29			9.33			0.57

(c) LLaMA 11B

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.69	0.19	7.0	0.32	0.11	9.0	0.29	0.05	4.5
Top Right	0.27	0.11	10.4	0.28	0.20	18.5	0.26	0.09	8.6
Bottom Left	0.33	0.35	27.1	0.27	0.26	24.1	0.26	0.35	34.1
Bottom Right	0.32	0.73	55.6	0.27	0.53	48.3	0.25	0.53	52.0
Invalid			0.01			0.03			0.81

Table 7: Classification performance and predicted proportions on the 2Among5 Task for Qwen 7B and Qwen 32B. “Sel” indicates the proportion of model predictions assigned to each label (in %).

(a) Qwen 7B

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.65	0.23	8.8	0.65	0.11	4.3	0.47	0.08	4.2
Top Right	0.01	0.00	1.6	0.03	0.00	0.6	0.17	0.03	4.0
Bottom Left	0.44	0.41	23.6	0.41	0.35	21.6	0.29	0.10	9.0
Bottom Right	0.09	0.08	19.9	0.19	0.31	39.9	0.24	0.65	66.9
Invalid			46.02			33.60			15.83

(b) Qwen 32B

Label	Disjunctive			Shape Conjunctive			Shape-Colour Conjunctive		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.92	0.29	8.0	0.94	0.11	3.0	0.42	0.10	5.9
Top Right	0.75	0.31	10.3	0.72	0.15	5.5	0.34	0.16	11.4
Bottom Left	0.72	0.34	11.9	0.57	0.17	7.4	0.32	0.18	14.2
Bottom Right	0.34	0.95	67.6	0.28	0.97	83.2	0.26	0.73	68.0
Invalid			2.13			0.98			0.54

E.2 LIGHT PRIORS

Within the Light Priors task, we again see preferences. GPT-4o, will routinely prefer a specific quadrant. For example, in Left and Right it selects the top-right cell over 50% of the time. On the other hand, Claude-Sonnet exhibits a strong bias towards selecting the bottom left cell in the Left, Right, and Top variants. Similarly, Llama90B seems to prefer the bottom right quadrant in all task variants. The smaller models all seem to lean towards selecting the Bottom Right quadrant. Once again, as in the 2 Among 5 task, the Qwen models (Table 10) display heavy bottom right bias.

Table 8: Model classification performance and predicted proportions on the **Light Priors Task** for GPT-4o, Claude-Sonnet, Llama 90B. Direction corresponds to the direction the target appears to be lit from. “Sel” indicates the percentage of model predictions assigned to each label.

(a) GPT-4o

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.82	0.55	17.3	0.46	0.09	4.8	0.71	0.14	5.2	0.78	0.25	8.1
Top Right	0.64	0.86	33.2	0.45	0.38	20.7	0.43	0.37	21.0	0.54	0.63	29.0
Bottom Left	0.78	0.77	24.6	0.45	0.35	19.6	0.55	0.46	20.8	0.62	0.58	23.6
Bottom Right	0.74	0.74	24.8	0.32	0.71	54.7	0.35	0.76	53.0	0.46	0.73	39.3
Invalid			0.13			0.11			0.06			0.04

(b) Claude-Sonnet

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.67	0.25	9.6	0.46	0.11	6.2	0.51	0.12	6.2	0.54	0.16	7.5
Top Right	0.47	0.34	17.9	0.35	0.20	13.9	0.34	0.19	13.9	0.36	0.23	15.6
Bottom Left	0.33	0.80	60.8	0.28	0.74	67.2	0.27	0.75	68.4	0.28	0.76	67.5
Bottom Right	0.69	0.32	11.5	0.28	0.14	12.8	0.28	0.13	11.5	0.45	0.17	9.3
Invalid			0.05			0.00			0.01			0.01

(c) LLaMA-90B

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.59	0.70	30.6	0.61	0.47	19.5	0.56	0.33	14.9	0.61	0.73	30.7
Top Right	0.65	0.20	7.4	0.46	0.21	11.4	0.39	0.19	12.4	0.71	0.19	6.6
Bottom Left	0.49	0.41	21.2	0.41	0.47	28.8	0.36	0.40	27.6	0.49	0.44	22.5
Bottom Right	0.49	0.70	35.4	0.41	0.62	36.9	0.36	0.61	41.9	0.48	0.69	35.4
Invalid			5.43			3.36			3.24			4.88

Table 9: Classification performance and predicted proportions on the Light Priors Task for GPT-4-Turbo, Claude-Haiku, and Llama 11B. “Sel” indicates the percentage of predictions assigned to each cell.

(a) GPT-4 Turbo

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.35	0.34	25.1	0.28	0.42	38.1	0.28	0.29	26.5	0.26	0.32	31.6
Top Right	0.42	0.05	2.7	0.42	0.01	0.6	0.22	0.01	1.3	0.35	0.01	0.9
Bottom Left	0.97	0.01	0.4	0.00	0.00	0.0	0.25	0.00	0.0	0.00	0.00	0.0
Bottom Right	0.28	0.80	71.9	0.26	0.63	61.3	0.25	0.72	72.2	0.25	0.69	67.4
Invalid			0.00			0.00			0.00			0.00

(b) Claude-Haiku

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.70	0.01	0.3	0.33	0.00	0.0	1.00	0.00	0.0	0.75	0.00	0.1
Top Right	0.43	0.03	1.5	0.42	0.03	1.8	0.41	0.03	1.7	0.42	0.03	1.6
Bottom Left	0.68	0.07	2.7	0.60	0.03	1.3	0.61	0.03	1.2	0.61	0.04	1.6
Bottom Right	0.24	0.88	88.5	0.25	0.96	94.0	0.25	0.98	95.1	0.25	0.98	94.9
Invalid			7.01			2.90			2.03			1.78

(c) LLaMA-11B

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.68	0.09	3.4	0.36	0.06	4.3	0.36	0.07	4.8	0.45	0.07	3.8
Top Right	0.17	0.04	5.1	0.23	0.10	10.7	0.22	0.10	11.0	0.17	0.06	8.5
Bottom Left	0.29	0.36	30.9	0.26	0.37	35.5	0.25	0.34	34.2	0.25	0.35	34.4
Bottom Right	0.30	0.74	60.2	0.27	0.53	49.3	0.27	0.54	49.9	0.26	0.57	53.0
Invalid			0.29			0.17			0.19			0.28

Table 10: Classification performance and predicted proportions on the Light Priors Task for Qwen 7B and Qwen 32B “Sel” indicates the percentage of predictions assigned to each cell.

(a) Qwen 7B

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.63	0.08	3.1	0.50	0.04	2.0	0.50	0.04	2.0	0.56	0.05	2.4
Top Right	0.09	0.00	0.5	0.17	0.00	0.2	0.24	0.00	0.2	0.23	0.00	0.3
Bottom Left	0.47	0.07	3.6	0.46	0.06	3.2	0.45	0.05	2.8	0.52	0.08	4.0
Bottom Right	0.26	0.96	92.0	0.25	0.97	94.0	0.25	0.96	94.6	0.26	0.96	92.6
Invalid			0.71			0.50			0.28			0.69

(b) Qwen 32B

Label	Bottom			Left			Right			Top		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.62	0.12	4.9	0.53	0.08	4.0	0.51	0.08	4.0	0.53	0.10	4.8
Top Right	0.66	0.05	1.8	0.52	0.02	1.1	0.50	0.02	1.1	0.63	0.02	0.8
Bottom Left	0.41	0.41	25.0	0.31	0.28	22.7	0.30	0.28	22.8	0.41	0.33	20.6
Bottom Right	0.32	0.86	66.8	0.28	0.80	71.8	0.27	0.79	71.8	0.29	0.88	73.7
Invalid			1.47			0.43			0.29			0.14

E.3 CIRCLE SIZES

Finally, in the Circle Sizes task, we again see spatial biases from models. Claude-sonnet heavily prefers the bottom left across all size variations. Meanwhile Llama 90B prefers the top-left. For both of these models the strength of preference increases with task difficulty, indicating a sort of “uncertainty response” with the preference. On the other hand GPT-4o eventually heavily prefers the bottom-right, but this is only present in the hardest Small variant. Both GPT-4-Turbo and Claude-Haiku heavily favour the bottom right quadrant, while Llama 11B is split between the bottom right and bottom left. Similarly, the Qwen models also display high levels of right-bottom bias (13).

Table 11: Model classification performance and predicted proportions on the **Circle Sizes Task** for GPT-4o, Claude-Sonnet, and Llama 90B. “Sel” indicates the percentage of model predictions assigned to each label.

(a) GPT-4o

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.93	0.74	20.0	0.92	0.51	13.8	0.73	0.15	5.2
Top Right	0.72	0.97	33.3	0.68	0.73	26.8	0.58	0.20	8.6
Bottom Left	0.90	0.85	23.5	0.75	0.82	27.7	0.46	0.58	31.4
Bottom Right	0.91	0.76	21.2	0.69	0.83	29.6	0.37	0.76	51.9
Invalid			2.04			2.14			2.78

(b) Claude-Sonnet

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.81	0.27	8.3	0.48	0.29	15.1	0.33	0.28	21.2
Top Right	0.57	0.67	28.9	0.50	0.27	13.4	0.36	0.13	9.0
Bottom Left	0.50	0.93	46.8	0.34	0.79	59.0	0.27	0.70	63.4
Bottom Right	0.84	0.53	15.9	0.68	0.34	12.4	0.42	0.10	6.3
Invalid			0.03			0.04			0.00

(c) LLaMA 90B

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.51	0.75	36.6	0.35	0.78	56.4	0.28	0.75	66.3
Top Right	0.66	0.12	4.6	0.54	0.14	6.3	0.35	0.08	5.7
Bottom Left	0.35	0.42	30.0	0.23	0.09	9.7	0.28	0.07	6.4
Bottom Right	0.52	0.56	27.4	0.33	0.35	26.4	0.28	0.22	20.4
Invalid			1.35			1.21			1.22

E.4 FINETUNING

In Table 14 we present the preference results within the Shape Conjunctive variation for GPT-4o after various levels of finetuning, corresponding to the fine-tuned models described in Section 4 and Appendix H. Here we see that that increased levels of finetuning can reduce the extent of spatial bias inherent in the model. At larger levels ($n=100$, $n=1000$) the cell preferences have almost completely disappeared.

Table 12: Classification performance and predicted proportions on the Circle Sizes Task for GPT-4-turbo, Claude-Haiku, and LLaMA 11B. “Sel” indicates the percentage of predictions assigned to each cell.

(a) GPT-4-turbo

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.70	0.25	8.9	0.44	0.16	9.0	0.35	0.12	8.6
Top Right	0.43	0.32	18.1	0.37	0.18	11.7	0.35	0.16	11.2
Bottom Left	0.52	0.33	15.9	0.31	0.18	14.7	0.33	0.21	15.5
Bottom Right	0.36	0.81	57.0	0.28	0.74	64.7	0.27	0.70	64.6
Invalid			0.08			0.02			0.01

(b) Claude-Haiku

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.98	0.02	0.6	0.93	0.01	0.3	0.81	0.01	0.3
Top Right	0.41	0.01	0.8	0.30	0.00	0.3	0.54	0.01	0.3
Bottom Left	0.73	0.38	13.1	0.57	0.05	2.1	0.47	0.03	1.3
Bottom Right	0.30	0.87	72.3	0.25	0.91	89.3	0.26	0.95	93.6
Invalid			13.18			8.09			4.49

(c) LLaMA 11B

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.35	0.15	10.8	0.28	0.15	13.7	0.29	0.15	13.0
Top Right	0.24	0.08	8.7	0.26	0.10	9.4	0.24	0.09	9.7
Bottom Left	0.26	0.42	40.9	0.26	0.40	40.0	0.25	0.40	40.5
Bottom Right	0.26	0.41	39.4	0.25	0.37	36.9	0.25	0.37	36.7
Invalid			0.15			0.06			0.12

Table 13: Classification performance and predicted proportions on the Circle Sizes Task for Qwen 7B and Qwen 32. “Sel” indicates the percentage of predictions assigned to each cell.

(a) Qwen 7B

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.47	0.37	19.7	0.30	0.29	24.0	0.27	0.27	25.7
Top Right	0.10	0.04	9.0	0.21	0.12	14.4	0.25	0.14	14.2
Bottom Left	0.37	0.06	4.0	0.57	0.01	0.5	0.65	0.01	0.5
Bottom Right	0.16	0.26	39.9	0.26	0.60	57.7	0.25	0.58	57.5
Invalid			27.38			3.34			2.16

(b) Qwen 32B

Label	Large			Medium			Small		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.84	0.15	4.3	0.56	0.04	1.7	0.52	0.04	1.9
Top Right	0.43	0.08	4.5	0.57	0.01	0.4	0.74	0.01	0.4
Bottom Left	0.65	0.38	14.7	0.43	0.09	5.4	0.38	0.07	4.8
Bottom Right	0.33	0.98	76.4	0.26	0.97	92.4	0.26	0.97	93.0
Invalid			0.12			0.03			0.01

Table 14: GPT-4o’s classification performance and predicted proportions across levels of finetuning set-sizes. (n). “Sel” indicates the percentage of predictions assigned to each label.

Label	$n = 0$			$n = 10$			$n = 100$			$n = 1000$		
	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel	Prec	Rec	Sel
Top Left	0.91	0.30	8.1	0.84	0.53	15.7	0.84	0.77	22.7	0.78	0.89	28.4
Top Right	0.55	0.74	33.8	0.65	0.74	28.2	0.73	0.86	29.4	0.89	0.77	21.8
Bottom Left	0.72	0.53	18.8	0.64	0.77	30.2	0.81	0.83	25.9	0.84	0.85	25.7
Bottom Right	0.51	0.72	35.2	0.61	0.63	25.6	0.81	0.72	22.0	0.85	0.82	24.1
Invalid			4.07			0.32			0.00			0.00

F REGRESSION AND CORRELATION ANALYSES

We also analysed model performance in the Cells variant of the three visual search experiments using logistic regression models. For each experiment, we fitted a generalised linear model with binomial error distribution:

$$\text{logit}(P(\text{Accuracy})) = \beta_0 + \beta_1 \text{NumDistractors}_c + \beta_2 \text{Condition} + \beta_3 (\text{NumDistractors}_c \times \text{Condition}).$$

where the number of distractors was centred around its mean with each model. We also included the $\text{NumDistractors} \times \text{Condition}$ interaction to investigate whether set-size effects differed across experimental conditions. This analysis is summarised for large models in Table 15 and for small models in Table 16. We also computed marginal slopes for each condition, which estimates the rate of change in log-odds of accuracy per unit increase in the number of distractors. All slopes were estimated on the logit scale with 95% asymptotic confidence intervals. These are summarised in Table 1 (large models) and Table 17 (smaller models) alongside the mean accuracy of each model per condition.

A series of Pearson’s correlations were also conducted to analyse whether model performance was associated with increasing distractor numbers within task conditions. A pronounced negative correlation suggests the target was increasingly harder to find at higher set sizes, and is indicative of serial search in human participants. A lack of a correlation suggests performance was flat across

1998 set sizes. A combination of high performance coupled with set-size independence is indicative of a
 1999 pop-out effect. The results of this analysis are displayed in Table 1 and Table 17.

2000

2001 F.1 CIRCLE SIZES

2002

2003 Focusing on the larger models (Table 15), accuracy decreased with Distractor Number in the Small-
 2004 target (reference) condition, whereas accuracy was higher overall in Medium and Large conditions.
 2005 The Distractor Number \times Condition interaction effects suggest that Medium and Large targets
 2006 attenuated the negative effect of distractor number in gpt-4o and llama90B, but not in claude-sonnet.
 2007 This pattern is reflected in Table 1: for gpt-4o, Distractor Number slopes flatten and mean accuracy
 2008 increases for larger targets. In contrast, although claude-sonnet shows higher mean accuracy for larger
 2009 targets, its slope estimates remain similar for Small (-0.017) and Large (-0.019) targets, suggesting
 2010 limited modulation by target size.

2011

2012 F.2 TWO AMONG FIVE

2013

2014 In this experiment, the disjunctive condition served as the reference. As expected, it showed positive
 2015 intercepts and very small Distractor Number coefficients (e.g., all < 0.01), consistent with accurate set-
 2016 size independent performance. In contrast, both the Shape Conjunctive and Shape-Colour Conjunctive
 2017 conditions showed negative effects, indicating reduced accuracy relative to the disjunctive baseline.
 2018 These conditions also produced negative Distractor Number \times Condition interactions, reflecting
 2019 a stronger accuracy decline as distractor number increased. As shown in Table 1, these steeper
 2020 Distractor Number slopes in conjunctive conditions are accompanied by corresponding decreases in
 mean accuracy.

2021

2022 F.3 LIGHT PRIORS

2023

2024 For the Light Priors experiment, we observed positive effects of the Bottom condition relative to
 2025 the Top (reference) condition (though not for llama-90B) alongside consistently negative effects
 2026 of the Left and Right conditions relative to Top across all three larger models. This aligns with
 2027 improved performance for vertical gradient spheres, particularly those lit from below. For gpt-4o and
 2028 claude-sonnet, the Left and Right conditions also showed negative Distractor Number \times Condition
 2029 interactions, indicating a stronger set-size-related decline in accuracy. In gpt-4o, the Bottom condition
 2030 additionally attenuated the Distractor Number effect. Overall, however, mean accuracy was low
 2031 and Distractor Number slopes were generally shallow across all conditions in this experiment (see
 Table 1).

2032

2033 G EXPERIMENT REPLICATION AND COMPUTING RESOURCES

2034

2035 All of the experiments in our paper are designed to be replicable and upon publication will be
 2036 accompanied with a code repository ⁶ which contains detailed guidance on which scripts to run and
 2037 with what arguments in the README file. The code contains everything needed to generate the exact
 2038 images we used, including any random seeds we set. The images were generated with seed 42.

2039

2040 The code contains scripts to generate and submit batches to either Anthropic or OpenAI, or to submit
 2041 models to a local SLURM cluster (though users will need to provide their own appropriately set up
 2042 SLURM scripts. Finally, the code also contains our scripts for generating figures used in the paper.
 2043 Upon publication we also intend to release our instance-level responses from each model to each
 question.

2044

2045 The biggest obstacle to replication is cost or access to compute. The experiments using the GPT and
 2046 Claude Models rely on paid API access and in total cost multiple hundreds of dollars to run all of
 2047 these experiments. Nevertheless, this reduces compute burden on the user. For both providers the
 2048 API batching option was used to reduce cost, but this makes it difficult to anticipate the total compute
 used or time required.

2049

2050 To run the Llama models we used an internal HPC cluster with exclusive access to nodes consisting
 2051 of 4 Nvidia A100s. Each Llama experiment took around 16-18 hours, but was split into two batches

2051

⁶Code link removed while under review.

Table 15: Regression coefficients across all experiments (Larger Models)

Experiment	Parameter	claude-sonnet	gpt-4o	llama90B
Circle Sizes	Intercept	-0.851***	-0.314***	-0.938***
	Ndist	-0.017***	-0.025***	-0.004**
	Cond:Medium	0.540***	1.279***	0.279***
	Cond:Large	1.264***	1.912***	0.798***
	Ndist×Medium	-0.007***	0.009***	0.009***
	Ndist×Large	-0.003	0.020***	0.016***
2 Among 5	Intercept	0.738***	1.709***	0.191***
	Ndist	-0.005***	0.002*	0.001
	Cond:Shape Conjunctive	-0.376***	-1.466***	-0.547***
	Cond:Shape-Colour Conjunctive	-1.309***	-2.102***	-1.017***
	Ndist×Shape Conjunctive	-0.006***	-0.021***	-0.008***
	Ndist×Shape-Colour Conjunctive	-0.011***	-0.019***	-0.008***
Light Priors	Intercept	-0.716***	0.183***	0.057**
	Ndist	-0.045***	0.027***	0.045***
	Cond:Bottom	0.421***	0.854***	-0.033
	Cond:Left	-0.162***	-0.676***	-0.296***
	Cond:Right	-0.158***	-0.472***	-0.540***
	Ndist×Bottom	0.005	0.061***	-0.013*
	Ndist×Left	-0.013*	-0.048***	-0.005
	Ndist×Right	-0.014*	-0.074***	-0.030***

Note: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$. Ndist = num distractors (centered); Cond: = condition level; × = interaction.

across two nodes. In total, for Llama 90B, we used approximately 450 hours of compute on these A100s.

H FINETUNING DETAILS

Our finetuning was conducted using OpenAI’s supervised finetuning API⁷ for GPT-4o. Models were trained on images from the Shape Conjunctive 2 Among 5s task generated with seed 1745313698. The evaluation seeds were: 1745332147 for Shape Conjunctive 2 Among 5, 1745566567 for Shape Conjunctive T Among 5s, 1746005099 for Disjunctive T Among 5s, and 1746104336 for Shape-Colour Conjunctive 2 Among 5, and 1746104336 for Circle Sizes.

The Ts Among 5s task is functionally the same as the 2 Among 5 task except the targets are representations of Ts and the distractors 5s. An example instance is given in Figure 20. While similar, the specific targets and distractors are visually distinct from the 2s and 5s, yet the same broad skill is required (identifying a target based on spatial features alone, not colours).

Figure 21 shows plots for all out-of-distribution evaluations we performed on the fine-tuned models. We see mild, but significant, improved performance on the Shape Conjunctive task with T among 5s.

I HUMAN BASELINES

I.1 PARTICIPANTS

Human participants for the three visual search tasks (Circle Sizes, Two Among Five and Light Priors) were recruited online using *Prolific* (www.prolific.co). Following ethics approval from our institution’s Research Ethics Committee. We recruited 30 participants for each task (Total N = 90, Age $M = 35$, 45 female, 45 male). Participants were pre-screened to ensure they were in the age range 18-60, were fluent in reading English, and self-reported normal or corrected-to-normal colour

⁷<https://platform.openai.com/docs/guides/fine-tuning>

Table 16: Regression coefficients across all experiments (Smaller Models)

Exp.	Parameter	claude-haiku	gpt-4-turbo	llama11B	Qwen7B	Qwen32B
CS	Intercept	-1.098***	-0.869***	-1.075***	-1.089***	-0.980***
	Ndist	-0.006***	-0.015***	-0.003	-0.003	-0.012***
	Cond:Medium	-0.063	0.054	0.011	0.013	0.002
	Cond:Large	0.355***	0.564***	0.070*	-0.454***	0.573***
	Ndist×Medium	-0.001	-0.004	0.000	-0.000	0.002
	Ndist×Large	-0.009***	-0.007***	0.001	-0.017***	0.004*
2A5	Intercept	-0.836***	0.285***	-0.654***	-1.515***	-0.136***
	Ndist	-0.012***	-0.010***	0.001*	-0.001*	-0.003***
	Cond:Shape	-0.155***	-0.654***	-0.317***	0.084**	-0.505***
	Cond:Shape-Col.	-0.215***	-1.154***	-0.428***	0.206***	-0.781***
	Ndist×Shape	0.002**	-0.006***	-0.002**	0.001	-0.004***
	Ndist×Shape-Col.	0.008***	-0.002*	-0.002**	0.005***	-0.006***
LP	Intercept	-1.060***	-1.078***	-1.056***	-0.996***	-0.720***
	Ndist	-0.027***	-0.010*	-0.013**	-0.038***	-0.054***
	Cond:Bottom	-0.092**	0.224***	0.235***	0.014	0.120***
	Cond:Left	-0.033	0.061	0.033	-0.037	-0.171***
	Cond:Right	-0.017	-0.004	0.018	-0.048	-0.192***
	Ndist×Bottom	-0.023***	-0.016**	0.012	0.001	0.003
	Ndist×Left	-0.002	0.001	-0.001	0.008	-0.006
	Ndist×Right	0.003	-0.002	0.005	0.011	-0.007

Note: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$. Ndist = num distractors (centered); Cond: = condition level; × = interaction. For the 2 Among 5 experiment, Shape and Shape-Col. (Shape-Colour) refer to conjunctive search conditions. Experiments are abbreviated to CS (Circle Sizes), 2A5 (2 Among 5) and LP (Light Priors).

vision. To maintain data quality, we replaced participants whose mean accuracy fell below chance level (25%, $N = 0$). Participants were compensated at the standard rate for their time and participation (£1.25, approximately £7.50/hour).

I.2 STIMULI AND PROCEDURE

Stimuli for the human baseline experiments were drawn from subsets of the AI test-sets. To construct the human stimulus sample, we used a stratified random sampling procedure within each experimental condition. Specifically, four stimuli were randomly selected from different distractor size ranges (or bins; see Table 18). Within each bin, target locations were evenly distributed, with each of the four Cells represented once. Colour combinations in the Two Among Five and Circle Size tasks were balanced within conditions. Prior to sampling, we excluded all stimuli with targets located within the coordinate range of 170–230 on either axis to prevent targets from appearing along the borders between Cells.

2160
 2161
 2162
 2163
 2164
 2165
 2166
 2167
 2168
 2169
 2170
 2171
 2172
 2173
 2174
 2175
 2176
 2177
 2178
 2179
 2180
 2181
 2182
 2183
 2184
 2185
 2186
 2187
 2188
 2189
 2190
 2191
 2192
 2193
 2194
 2195
 2196
 2197
 2198
 2199
 2200
 2201
 2202
 2203
 2204
 2205
 2206
 2207
 2208
 2209
 2210
 2211
 2212
 2213

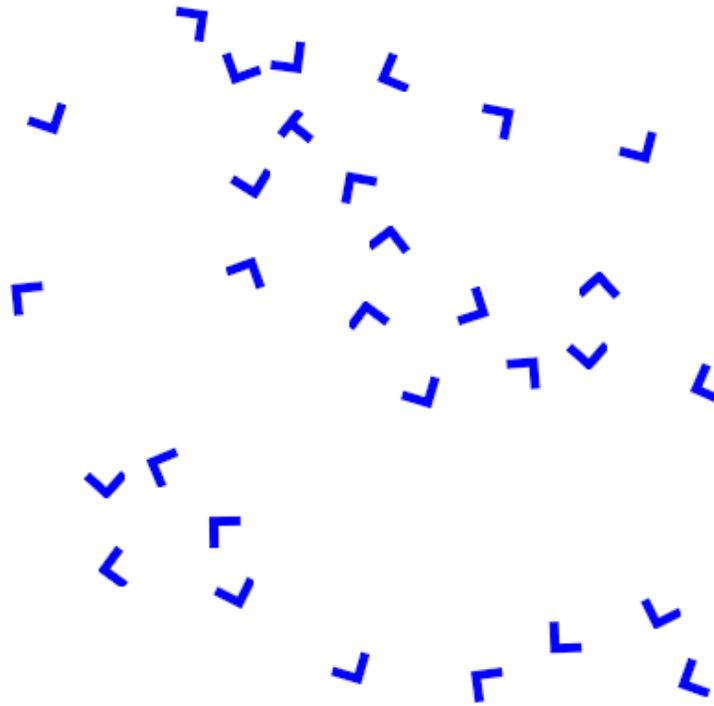


Figure 20: An example of the T Among Ls task.

Table 18: Distractor bin sizes

Circle Sizes	Two Among Five	Light Priors
1–4	1–4	2–5
5–8	5–8	6–9
9–12	9–16	10–13
13–16	17–32	14–17
17–20	33–64	
21–24	65–99	
25–28		
29–32		
33–36		
37–40		
41–44		
45–49		

I.2.1 CIRCLE SIZES

In the Circle Sizes task, we implemented a 3 (Condition: Small, Medium, Large) \times 12 (Distractor Number: twelve bins) factorial design, with four trials per bin, resulting in 144 trials overall.

The experiment was conducted online using *Gorilla* (<https://gorilla.sc/>). After obtaining informed consent and checking their use of a QWERTY keyboard, participants were informed they would need to locate the largest circle amongst distractors, responding with the corresponding Cell using the keys Q (top-left), P (top-right), A (bottom-left) or L (bottom-right). Participants were provided with two examples followed by eight practice trials with feedback to familiarise with the response keys. Participants were also informed that when in doubt as to which cell the target was in, they should respond with the Cell that they believe the target was “most in”, and if the image disappeared before they located the target, they should respond with their best guess.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

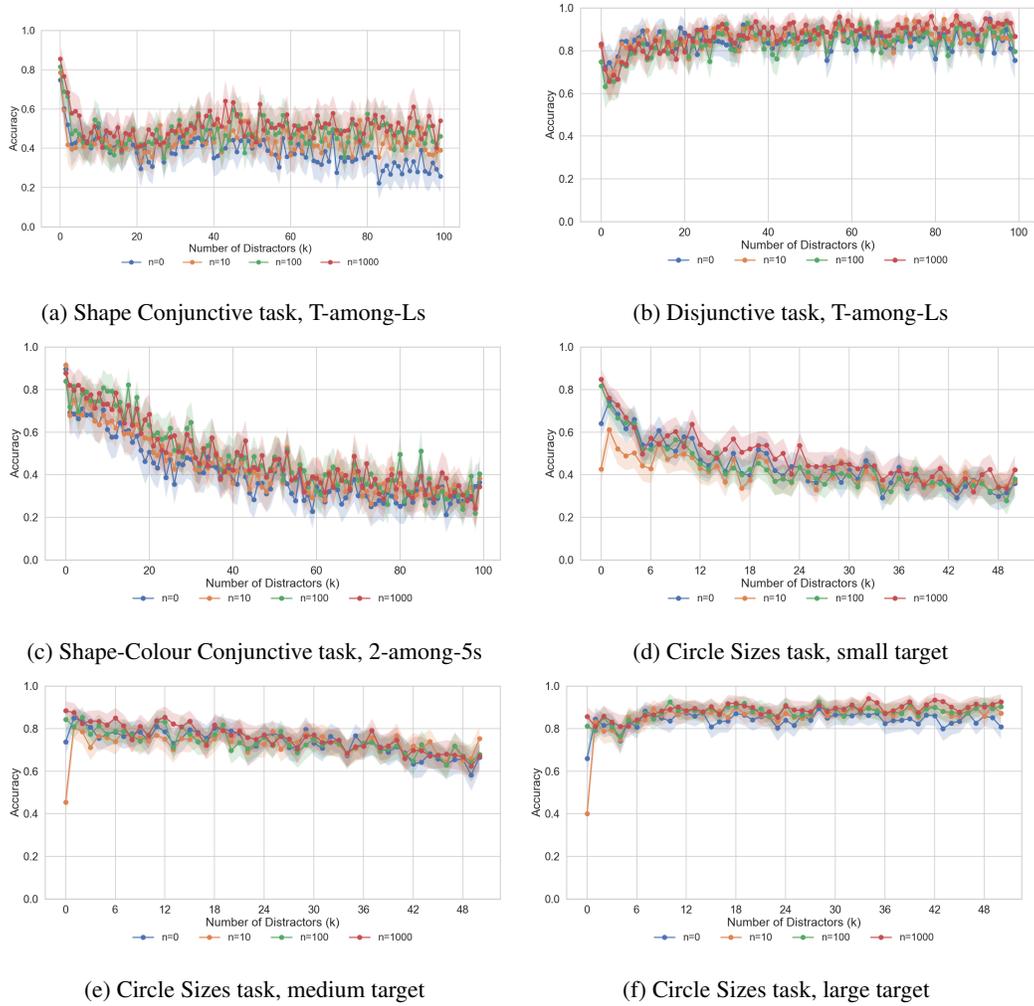
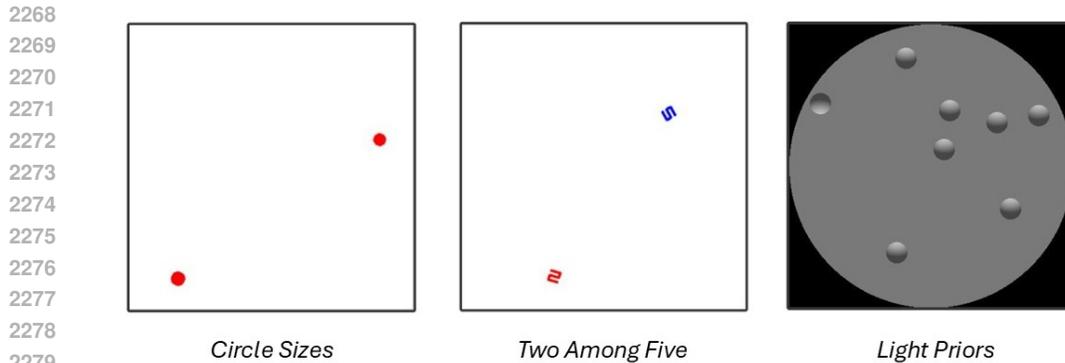


Figure 21: Fine-tuned GPT-4o evaluation



2280
2281 Figure 22: Screenshots from experimental trials presented to online participants in Two Among Five
2282 (Left; "Find the red 2"), Light Priors (Middle; "Find the odd one out") and Circle Sizes
2283 (Right; "Find the largest circle").

2284
2285 Experimental trials were preceded with a central fixation cross (for 500ms). The 400×400 pixel
2286 image stimulus then appeared centrally at full resolution, accompanied by the prompt "Find the
2287 largest circle" (see Figure 22). The image stimulus remained on screen for 1500 milliseconds before
2288 being replaced by a white image mask.
2289

2290 I.2.2 TWO AMONG FIVE

2291 For the Two Among Five task, we employed a 3 (Condition: Disjunctive, Shape Conjunctive, Shape-
2292 Colour Conjunctive) \times 2 (Version: 2 Among 5, 5 Among 2) \times 6 (Distractor Number: six bins)
2293 factorial design, again with four trials per bin, resulting in a total of 144 trials. The experimental
2294 procedure was the same as in the Circle Sizes task, except the prompt said "Find the *colour number*"
2295 — with *colour* and *number* replaced with the correct target features (e.g., "Find the red 2"), and
2296 the stimulus appeared on screen for 3000ms to maintain similar difficulty levels on this task which
2297 included trials with a higher number of distractors (e.g., 99).
2298

2299 I.2.3 LIGHT PRIORS

2300 Finally, in the Light Priors task, we used a 2 (Condition: Vertical Gradient, Horizontal Gradient) \times
2301 2 (Version: Original [Top-lit, Left-lit], Reversed [Bottom-lit, Right-lit]) \times 4 (Distractor Number:
2302 Four bins) factorial design, with twelve trials per bin, yielding 192 trials in total. The experiment
2303 procedure was the same as the Two Among Five task, except participants were told the image would
2304 contain spheres lit from different directions, and were instructed to "Find the odd one out" (e.g., a
2305 bottom-lit sphere amongst top-lit spheres; see Figure 22). The Light Priors task uses fewer distractors
2306 (here 2–17) to prevent the spheres becoming tightly clustered and the target easily identified through
2307 contrast with proximal distractors (Adams, 2007). We also did not include single-distractor trials as
2308 at least two distractors are required to identify an "odd one out". Given these reduced set sizes, we
2309 decreased the image presentation time to 1500 milliseconds to maintain difficulty levels.
2310

2311 I.3 RESULTS

2312 I.3.1 CIRCLE SIZES TASK

2313
2314 Table 19 displays the results for human performance in the Circle Sizes task, which we analysed
2315 using a 3 (Condition: Small, Medium, Large) \times 12 (Distractor Number: twelve bins) \times 4 (Cell:
2316 top-left, top-right, bottom-left, bottom-right) repeated measures ANOVA. Performance in the Large
2317 target condition was close to ceiling ($M = 96\%$), and accuracy was higher than in the Medium
2318 condition ($M = 88\%$, $t = 8.42$, $p < .001$), and accuracy in both of these conditions was higher than
2319 the Small condition ($M = 59\%$, $t_s > 18.9$, $p_s < .001$). Overall performance declined with increasing
2320 Distractor Number, but a Distractor Number \times Condition interaction indicates that this effect was not
2321 uniform across all size conditions. For Large targets, all comparisons between distractor bin pairs
were non-significant ($t_s < 2.8$, $p_s > .06$), indicating a clear pop-out effect. In the Medium condition,

2322 only the difference between bins 13–16 and 45–49 were significant ($t = 4.76$, $ps = .004$, all other ts
 2323 < 3.5 , $ps > .05$), suggesting a large degree of set size independence. Performance declined steadily
 2324 with set size in the Small condition (e.g., bin 1–4 vs 45–49: $t = 5.11$, $p < .001$).

2325 We also found a main effect of Cell, and a marginal Cell \times Condition interaction. While participants
 2326 showed similarly high accuracy across Cells in the Large condition (all $Ms > 95\%$), differences were
 2327 found in the Small and Medium condition. Overall, accuracy was lowest for the bottom-right cell (M
 2328 $= 78\%$), perhaps reflecting a serial search strategy ending in the bottom-right quadrant of the screen.
 2329

2330 I.3.2 TWO AMONG FIVE TASK

2331 To analyse the performance of human participants in the Two Among Five task, we conducted a 3
 2332 (Condition: Disjunctive, Shape Conjunctive, Shape-Colour Conjunctive) \times 2 (Version: 2 Among 5, 5
 2333 Among 2) \times 6 (Distractor Number: six bins) \times 4 (Cell: top-left, top-right, bottom-left, bottom-right)
 2334 repeated measures Analysis of Variance (ANOVA), summarised in Table 20.
 2335

2336 Participants' performance in the Disjunctive condition was at ceiling ($M = 98\%$) and overall higher
 2337 than in the Shape Conjunctive ($M = 64\%$, $t = 26.14$, $p < .001$), and Shape-Colour Conjunctive
 2338 conditions ($M = 69\%$, $t = 22.88$, $p < .001$)⁸. Shape-Colour Conjunctive performance was also
 2339 higher than Shape Conjunctive ($t = 3.22$, $p = .004$). Participant performance overall decreased with
 2340 Distractor Number (see Figure 3), but unlike other conditions, differences between Distractor Number
 2341 bins within the Disjunctive condition were all non-significant (all $ts < 2.9$, $ps > .06$), suggesting
 2342 pop-out performance in that condition. In the Shape Conjunctive (e.g., bin 1–4 vs 65–99: $t = 12.96$,
 2343 $p < .001$) and Shape-Colour Conjunctive (e.g., bin 1–4 vs 65–99: $t = 8.44$, $p < .001$) conditions,
 2344 accuracy declined across set size, indicating serial search.

2345 Participant performance was similar across Versions, but like the previous experiment, performance
 2346 differed depending on which Cell the target number was located. This effect interacted with Condition,
 2347 and post-hoc tests indicated performance between Cells were not significantly different in the
 2348 Disjunctive or Shape Conjunctive conditions (all $ts < 2.3$, $ps > .15$). In the Conjunctive condition,
 2349 participants showed reduced accuracy in the bottom-right cell ($M = 54\%$) relative to the others (all ts
 2350 > 4.9 , $ps < .001$), likely reflecting a serial search strategy ending in the bottom-right quadrant of the
 2351 screen.

2352 I.3.3 LIGHT PRIORS TASK

2353 For the Light Priors task, we conducted a 2 (Condition: Vertical Gradient, Horizontal Gradient) \times
 2354 2 (Version: Original, Reversed) \times 4 (Distractor Number: Four bins) \times 4 (Cell: top-left, top-right,
 2355 bottom-left, bottom-right) repeated measures ANOVA, the results of which are displayed in Table 21.
 2356 An effect of Condition showed humans performed better in the vertical ($M = 86\%$) relative to the
 2357 horizontal gradient condition overall ($M = 65\%$), though performance in both conditions was flat
 2358 across across set sizes. We also found an effect of Version, and a marginal Condition \times Version
 2359 interaction. Post-hoc tests indicate that participants were slightly better at detecting right-lit (M
 2360 $= 67\%$) compared to left-lit ($M = 63\%$) lit spheres ($t = 2.30$, $p = 0.021$), but substantially better
 2361 at detecting bottom-lit ($M = 91\%$) compared to top-lit ($M = 81\%$) spheres ($t = 8.23$, $p = <.001$).
 2362 Improved performance for vertical gradients, and for the novel 'bottom-lit' variant in particular, is
 2363 consistent with findings from classic studies (Enns & Rensink, 1990).
 2364
 2365
 2366
 2367
 2368
 2369
 2370
 2371
 2372
 2373
 2374
 2375

⁸All post-hoc p-values are adjusted for multiple comparisons using the Bonferroni correction

2376 J PROMPTING DETAILS

2377

2378 In Tables 22, 23, and 24 we detail the prompts that were used throughout our experiments with the
2379 MLLMs.

2380

2381

2382

2383

2384

2385

2386

2387

2388

2389

2390

2391

2392

2393

2394

2395

2396

2397

2398

2399

2400

2401

2402

2403

2404

2405

2406

2407

2408

2409

2410

2411

2412

2413

2414

2415

2416

2417

2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2429

In Tables 22, 23, and 24 we detail the prompts that were used throughout our experiments with the MLLMs.

We aimed to keep these as similar to each other as feasible, but were forced to adapt a few prompts to either provide sufficient information to the model, or to force models to respond. For instance, the Shape-Colour Conjunctive search in the Two Among Five task requires additional information to solve compared to the other variations (both the colour and shape need to be provided instead of just shape). On the other hand, when requesting some models for coordinates, some models (Claude Sonnet and Llama90B in particular) would rarely provide coordinates unless you offered them the chance to explain their uncertainty.

Table 22: Prompts used for the 2Among5 tasks. Bolded terms indicate that the prompt includes the appropriate target, distractor, or colour for the image instance.

Experiment	Prompt
2Among5 Shape Conjunctive Search (Cells)	The image is divided into a 2x2 grid. Each element of the grid is referred to as a cell. In the presented image there are a number of objects. Almost all of the objects are the number distractor written as a numeral. There is a single target in the image, similarly represented by a numeral. In which cell is the target in? In the case where the target overlaps multiple cells, please provide the cell where the centre of the target is located. Please only respond with 'Cell (i,j)' where (i,j) corresponds to the ith row and jth column of the grid. The top left cell is Cell (1,1). Do not reply with anything else.
2Among5 Shape-Colour Conjunctive (Cells)	The image is divided into a 2x2 grid. Each element of the grid is referred to as a cell. In the presented image there are a number of objects. There are '2's and '5's written as numerals. In which cell is the colour 'shape'? In the case where the colour 'target' overlaps multiple cells, please provide the cell where the centre of the shape is located. Please only respond with 'Cell (i,j)' where (i,j) corresponds to the ith row and jth column of the grid. The top left cell is Cell (1,1). If you are uncertain you may optionally add a note explaining that but please start your response with 'Cell (i,j)'.
2Among5 Shape Conjunctive Search (Coordinates)	The presented image is 400x400 pixels large, and the origin (0,0) is in the top left of the image. In the presented image there are a number of objects. Almost all of the objects are the number distractor written as a numeral. There is a single target in the image, similarly represented by a numeral. What are the coordinates of the centre of the target ? Please give your best estimate. If you are uncertain please guess but optionally add a description to note this. However, for ease of processing please begin your response with a set of coordinates using round brackets.
2Among5 Shape-Colour Conjunctive Search (Coordinates)	The presented image is 400x400 pixels large, and the origin (0,0) is in the top left of the image. In the presented image there are a number of objects. There are '2's and '5's written as numerals. What are the coordinates of the colour 'target' ? Please give your best estimate. If you are uncertain please guess but optionally add a description to note this. However, for ease of processing please begin your response with a set of coordinates using round brackets.

K MECHANISTIC INTERPRETABILITY

K.1 EXPERIMENTAL DETAILS

Tasks and stimuli. We evaluate on the 3 main tasks (Circle Sizes, 2Among5, and Light Priors, breaking down into the constituent datasets (e.g., CircleSizes Small, Medium, and Large are evaluated separately).

For each experiment we sample 300 images uniformly across the number of distractors (field size n), and store per-image metadata (e.g., `num_distractors`, `size`, `quadrant`).

Model and activations. Unless stated, we use Llama-3.2-90B-Vision-Instruct (abbrev. llama90B). During generation we hook three representative language layers (early, middle, late) and record: residual stream, attention output, MLP output, and layer-norm outputs. We also expose minimal vision encoder residuals and projector I/O. Activations are captured at the last prompt token. This yields per-category dicts keyed by `layer_<idx>` with feature shape $[1, 1, d]$.

Probes. We train linear and small non-linear readouts. The linear probe is a single-layer logistic model with L1 penalty (trained in PyTorch with BCE-with-logits and Adam). We use an 80/20 stratified split, grid-search $\lambda_{\ell_1} \in \{10^{-3}, 5 \times 10^{-4}, 10^{-4}, 5 \times 10^{-5}, 10^{-5}\}$, batch size 64, and 8 epochs. We also report two MLP probes: `mlp_small` (256 hidden units) and `mlp_tiny` (128→32). Optional stability selection bootstraps estimate feature-selection frequency. We report validation accuracy per (category, layer).

Hardware. Initial activation extraction and probe training were run on 4×A100 (CUDA enabled).

K.2 EXPERIMENTAL RESULTS

We report results from a series of mechanistic probing experiments across tasks designed to test visual search phenomena in LLaMA-90B. For context, Llama 90B showed limited performance on Circle sizes and Light Priors, only really demonstrating clear differences between conditions on the 2 Among 5 task. Therefore we focus our analysis here on the 2 Among 5 task (Figure 23a). The disjunctive condition starts relatively high for 2 Among 5 and remains relatively flat – indicating that early parts of the network are used to locate the target. On the other hand, Shape-Colour conjunctive begins much lower, and increases almost linearly as the network becomes deeper – indicating that later parts of the network are used to locate the target. We have taken this as evidence, in this task at least, of Llama 90B recruiting different layers to solve disjunctive and conjunctive tasks. We present the Mechanistic probing results for the other two tasks for completion.

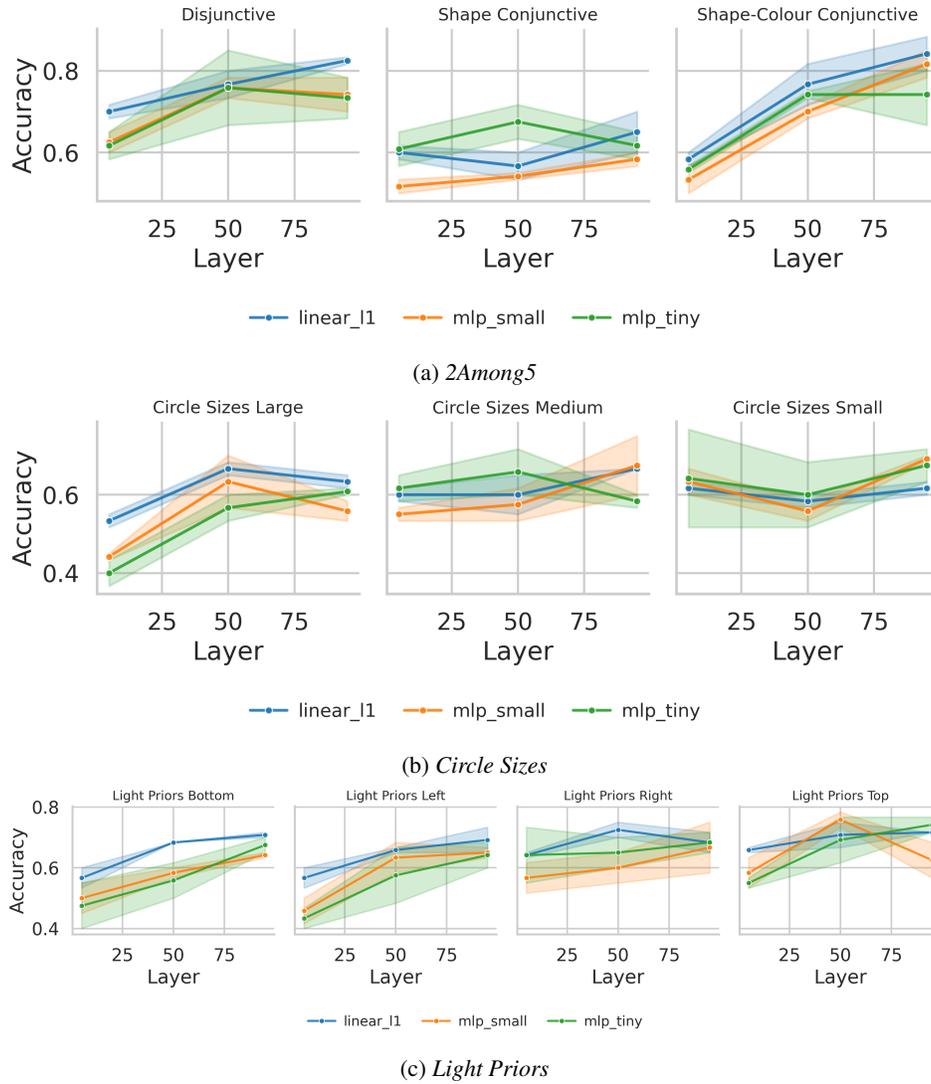


Figure 23: Layer-wise probe accuracy across tasks: *2Among5*, *Circle Sizes*, and *Light Priors*. Each plot contains three data points, at layers 5, 50, and 95.

2538 Table 17: Regression slopes and correlations for the effect of distractor number on accuracy within
 2539 conditions across all three experiments (smaller models)
 2540
 2541

Exp.	Model	Condition	Mean Acc	Regression		Correlation	
				Slope	95% CI	r	p
CS	claude-haiku	Small	0.250	-0.006	[-0.009, -0.003]	-0.039	0.003
	claude-haiku	Medium	0.239	-0.008	[-0.011, -0.004]	-0.047	< 0.001
	claude-haiku	Large	0.324	-0.015	[-0.018, -0.012]	-0.101	< 0.001
	gpt-4-turbo	Small	0.297	-0.015	[-0.018, -0.012]	-0.100	< 0.001
	gpt-4-turbo	Medium	0.310	-0.019	[-0.022, -0.016]	-0.126	< 0.001
	gpt-4-turbo	Large	0.426	-0.022	[-0.025, -0.020]	-0.159	< 0.001
	llama11B	Small	0.254	-0.003	[-0.006, < 0.001]	-0.018	1.000
	llama11B	Medium	0.257	-0.003	[-0.006, 0.001]	-0.016	1.000
	llama11B	Large	0.268	-0.001	[-0.004, 0.002]	-0.009	1.000
	Qwen32B	Small	0.274	-0.012	[-0.015, -0.009]	-0.076	< 0.001
	Qwen32B	Medium	0.274	-0.010	[-0.013, -0.007]	-0.066	< 0.001
	Qwen32B	Large	0.400	-0.007	[-0.010, -0.004]	-0.052	< 0.001
	Qwen7B	Small	0.252	-0.003	[-0.006, < 0.001]	-0.018	1.000
	Qwen7B	Medium	0.254	-0.003	[-0.006, < 0.001]	-0.020	1.000
Qwen7B	Large	0.180	-0.020	[-0.023, -0.016]	-0.110	< 0.001	
2A5	claude-haiku	Disjunctive	0.307	-0.012	[-0.013, -0.011]	-0.156	< 0.001
	claude-haiku	Shape	0.274	-0.010	[-0.011, -0.009]	-0.124	< 0.001
	claude-haiku	Shape-Col.	0.260	-0.004	[-0.005, -0.003]	-0.048	< 0.001
	gpt-4-turbo	Disjunctive	0.570	-0.010	[-0.011, -0.009]	-0.135	< 0.001
	gpt-4-turbo	Shape	0.412	-0.015	[-0.016, -0.014]	-0.215	< 0.001
	gpt-4-turbo	Shape-Col.	0.300	-0.011	[-0.012, -0.010]	-0.148	< 0.001
	llama11B	Disjunctive	0.342	0.001	[< 0.001, 0.002]	0.015	0.934
	llama11B	Shape	0.275	-0.001	[-0.002, < 0.001]	-0.017	0.353
	llama11B	Shape-Col.	0.253	-0.001	[-0.002, < 0.001]	-0.014	1.000
	Qwen32B	Disjunctive	0.466	-0.003	[-0.004, -0.002]	-0.049	< 0.001
	Qwen32B	Shape	0.346	-0.007	[-0.008, -0.006]	-0.098	< 0.001
	Qwen32B	Shape-Col.	0.289	-0.009	[-0.010, -0.008]	-0.118	< 0.001
	Qwen7B	Disjunctive	0.180	-0.001	[-0.003, < 0.001]	-0.015	0.821
	Qwen7B	Shape	0.193	< 0.001	[-0.001, 0.001]	-0.001	1.000
Qwen7B	Shape-Col.	0.213	0.003	[0.002, 0.005]	0.041	< 0.001	
LP	claude-haiku	Top	0.258	-0.027	[-0.036, -0.019]	-0.063	< 0.001
	claude-haiku	Bottom	0.243	-0.050	[-0.059, -0.041]	-0.112	< 0.001
	claude-haiku	Left	0.252	-0.029	[-0.038, -0.021]	-0.067	< 0.001
	claude-haiku	Right	0.255	-0.024	[-0.033, -0.016]	-0.056	< 0.001
	gpt-4-turbo	Top	0.254	-0.010	[-0.019, -0.002]	-0.024	0.554
	gpt-4-turbo	Bottom	0.299	-0.027	[-0.035, -0.019]	-0.064	< 0.001
	gpt-4-turbo	Left	0.266	-0.009	[-0.018, -0.001]	-0.022	0.994
	gpt-4-turbo	Right	0.254	-0.013	[-0.021, -0.004]	-0.029	0.117
	llama11B	Top	0.258	-0.013	[-0.021, -0.004]	-0.029	0.117
	llama11B	Bottom	0.305	-0.001	[-0.009, 0.007]	-0.002	1.000
	llama11B	Left	0.265	-0.014	[-0.022, -0.006]	-0.032	0.038
	llama11B	Right	0.262	-0.007	[-0.016, 0.001]	-0.017	1.000
	Qwen32B	Top	0.330	-0.054	[-0.062, -0.046]	-0.132	< 0.001
	Qwen32B	Bottom	0.357	-0.051	[-0.059, -0.043]	-0.127	< 0.001
Qwen32B	Left	0.295	-0.060	[-0.069, -0.052]	-0.142	< 0.001	
Qwen32B	Right	0.291	-0.061	[-0.069, -0.052]	-0.142	< 0.001	
Qwen7B	Top	0.272	-0.038	[-0.047, -0.030]	-0.089	< 0.001	
Qwen7B	Bottom	0.274	-0.037	[-0.046, -0.029]	-0.087	< 0.001	
Qwen7B	Left	0.264	-0.031	[-0.039, -0.022]	-0.071	< 0.001	
Qwen7B	Right	0.261	-0.027	[-0.036, -0.019]	-0.062	< 0.001	

2590 *Note: Correlations are Pearson's r and p values are Bonferroni corrected for multiple comparisons.*
 2591 *For the 2 Among 5 experiment, Shape and Shape-Col. (Shape-Colour) refer to conjunctive search conditions. Experiments are abbreviated to CS (Circle Sizes), 2A5 (2 Among 5) and LP (Light Priors).*

2592

2593

Table 19: ANOVA results for human performance in the Circle Sizes task.

2594

2595

2596

2597

2598

2599

2600

2601

2602

2603

2604

2605

2606

2607

Table 20: ANOVA results for human performance in the Two Among Five task.

2608

2609

2610

2611

2612

2613

2614

2615

2616

2617

2618

2619

2620

2621

2622

2623

2624

2625

2626

2627

Table 21: ANOVA results for human performance in the Light Priors task.

2628

2629

2630

2631

2632

2633

2634

2635

2636

2637

2638

2639

2640

2641

2642

2643

2644

2645

Effect	DF_n	DF_d	F	p	η_p^2
Condition	1.00	29.00	91.528	<.001	0.759
Inversion	1.00	29.00	18.007	<.001	0.383
Cell	3.00	87.00	1.578	0.200	0.052
Distractor Number	3.00	87.00	1.109	0.350	0.037
Condition \times Inversion	1.00	29.00	3.005	0.094	0.094
Condition \times Cell	3.00	87.00	2.293	0.084	0.073
Inversion \times Cell	3.00	87.00	0.789	0.503	0.027
Condition \times Distractor Number	3.00	87.00	0.742	0.530	0.025
Inversion \times Distractor Number	3.00	87.00	2.245	0.089	0.072
Cell \times Distractor Number	9.00	261.00	2.401	0.012	0.076
Condition \times Inversion \times Cell	3.00	87.00	2.849	0.042	0.089
Condition \times Inversion \times Distractor Number	3.00	87.00	2.466	0.068	0.078
Condition \times Cell \times Distractor Number	9.00	261.00	3.336	0.001	0.103
Inversion \times Cell \times Distractor Number	9.00	261.00	1.376	0.199	0.045
Condition \times Inversion \times Cell \times Distractor Number	9.00	261.00	1.825	0.064	0.059

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

Table 23: Prompts used for the Light Priors tasks. Bolded terms indicate that the prompt includes the appropriate shading directions for the image instance.

Experiment	Prompt
Light Priors (Cells)	The image is divided into a 2x2 grid. Each element of the grid is referred to as a cell. In the presented image there are a number of spheres lit from different directions. Almost all of the spheres are lit from the same direction, but one sphere is lit from the opposite direction. In which cell is this oppositely lit sphere? In the case where the sphere overlaps multiple cells, please provide the cell where the centre of the sphere lit from the opposite direction is located. Please only respond with 'Cell (i,j)' where (i,j) corresponds to the ith row and jth column of the grid. The top left cell is Cell (1,1). If you are uncertain please guess but optionally add a description to note this. However, for ease of processing please begin your response with 'Cell (i,j)'.
Light Priors (Coordinates)	The presented image is 400x400 pixels large, and the origin (0,0) is in the top left of the image. In the presented image there are a number of spheres lit from different directions. Almost all of the spheres are lit from the same direction, but one sphere is lit from the opposite direction. What are the coordinates of the centre of the oppositely lit sphere? If you are uncertain please guess but optionally add a description to note this. However, for ease of processing please begin your response with a set of coordinates using round brackets."

Table 24: Prompts used for the Circle Sizes tasks.

Experiment	Prompt
Circle Sizes (Cells)	The image is divided into a 2x2 grid. Each element of the grid is referred to as a cell. In the presented image there are a number of circles. One of the circles is larger than the rest. In which cell is the larger circle? In the case where the larger circle overlaps multiple cells, please provide the cell where the centre of the larger circle is located. Please only respond with 'Cell (i,j)' where (i,j) corresponds to the ith row and jth column of the grid. The top left cell is Cell (1,1). If you are uncertain you may optionally add a note explaining that but please start your response with 'Cell (i,j)'.
Circle Sizes (Coordinates)	The presented image is 400x400 pixels large, and the origin (0,0) is in the top left of the image. In the presented image there are a number of circles. One of the circles is larger than the others. What are the coordinates of the larger circle? Please give your best estimate. If you are uncertain please guess but optionally add a description to note this. However, for ease of processing please begin your response with a set of coordinates using round brackets.