

Dual Thinking and Logical Processing in Human Vision and Multi-modal Large Language Models

Kailas Dayanandan, Nikhil Kumar, Anand Sinha and Brejesh Lall

Abstract—The dual thinking framework considers fast, intuitive, and slower logical processing. The perception of dual thinking in vision requires images where inferences from intuitive and logical processing differ, and the latter is under-explored in current studies. We introduce a novel adversarial dataset to provide evidence for the dual thinking framework in human vision, which also facilitates the study of the qualitative behavior of deep learning models. Our psychophysical studies show the presence of multiple inferences in rapid succession, and analysis of errors shows that the early stopping of visual processing can result in missing relevant information. MLLMs (Multi-modal Large Language Models) and VLMs (Vision Language Models) have made significant progress in correcting errors in intuitive processing in human vision and showed enhanced performance on images requiring logical processing. However, their improvements in logical processing have not kept pace with their advancements in intuitive processing. In contrast, segmentation models exhibit errors similar to those seen in intuitive human processing and lack understanding of sub-structures, as indicated by errors related to sub-components in identified instances. As AI (Artificial Intelligence)-based systems find increasing applications in safety-critical domains like autonomous driving, the integration of logical processing capabilities becomes essential. This not only enhances performance but also addresses the limitations of scaling-based approaches while ensuring robustness and reliability in real-world environments. The code for this paper is available at https://github.com/kailasdayanandan/dual_thinking/

Impact Statement—Reasoning or logical processing is a key aspect of human cognition, yet deep learning models often struggle with tasks beyond pattern recognition. While current methods rely on scaling up data and computation to improve performance, our study focuses on understanding reasoning in human visual processing, which is less resource-intensive. It brings forth some of the limitations of current VLMs, where reasoning can be a more resource-efficient alternative to improve performance with reduced reliance on large datasets and computation. This is significant as AI technologies are progressively being adopted in safety-critical domains as they approach human-level performance, where enhanced reliability and efficiency are essential.

Index Terms—Dual Thinking, Cognitive science, Human intelligence, Image segmentation, Deep Learning, Visual Language Models, Multi-modal Large Language Models, System 1 - System 2, Intuitive and logical processing

Kailas Dayanandan is with Indian Institute of Technology, New Delhi India (e-mail: kailasd@gmail.com).

Nikhil Kumar is with Indraprastha Institute of Information Technology, New Delhi India (e-mail: nikhil21174@iitd.ac.in).

Anand Sinha is with Indian Institute of Technology, New Delhi India (e-mail: anand.sinha85@gmail.com).

Brejesh Lall is with Indian Institute of Technology, New Delhi India (e-mail: brejesh@iitd.ac.in).

Manuscript received 10 February 2025; accepted 29 August 2025. Date of publication XXXX; date of current version XXXX. This article was handled by Associate Editor XXXX upon evaluation of the reviewers' comments. (Corresponding author: Kailas Dayanandan)

I. INTRODUCTION

Dual thinking framework argues that humans have a fast intuitive system and a slow logical system, and has been widely mentioned in recent years [1], [2], [3], [4], [5]. A growing body of electrophysiological research suggests gist formation in initial 150 – 200ms using fast feed-forward processing, followed by slower iterative refinement [1], [6], [7], [8], [2], [9]. The dual thinking framework is not part of studies on visual perception due to the unavailability of methods to study them. Recent advancements in state-of-the-art performance have primarily been driven by scaling models and training on massive datasets using extensive computational resources. However, as the performance gains from scaling begin to saturate [10], the focus is shifting toward enhancing the reasoning abilities of foundation models and large language models.

Gestalt principles have been studied for over a century [11], [12] indicating its presence in human visual processing [13]. Gestalt theory proposes that humans perceive things as a whole rather than as a sum of parts [14], [15], which was a deviation from structuralism that considered analyzing parts as a key to understanding the whole object [11], [16]. Gestalt principles and Structuralism are prominent theories about human visual perception [17], [18]; however, they are insufficient to describe class identification and complex logical processing in human vision. The theories are also independent of findings about initial gist formation in 200ms observed in many electrophysiological studies [1], [7], [3], [6]. There exists a research gap for a dataset that enables the analysis of initial gist formation and logical processing in human vision.

Marr's three-level hypothesis [19], [20], [21] provides a framework for analyzing human visual processing. The analysis of computational aspects or the algorithmic component has been ignored in recent years as end-to-end deep learning models gained prominence [22], [23], [24]. The comparative studies have shown to help improve accuracy, robustness, and generalizability of models by incorporating strategies from human vision [25], [25], [26], [27], [28], [29], [30]. The computational models using deep learning helped enhance the understanding of human vision [31], [7], [32], [33], [34], though, they face criticism for using classification models unlike human vision that localizes objects [35], [36], [17]. Our main contributions in this paper include

(a) We observe that logical processing in the dual thinking framework is essential to define human visual processing completely. Our study shows that logical processing was nascent in segmentation models, and multi-modal LLMs in

recent years have made tremendous progress in resolving many logical errors; however, they still face limitations that require reasoning using minor features like size differences in sub-components in an instance.

(b) Our analysis offers insights into intuitive processing and the logical errors addressed in the later stages of human vision. Our study highlights the role of Gestalt principles in sub-component formation and emphasizes the importance of shape in grouping elements during intuitive processing. Additionally, we observed that segmentation models like YOLACT exhibit characteristics similar to intuitive processing in human vision.

(c) We introduce the Human Confusion Dataset, a diverse dataset to study dual thinking and strategies in human vision. We leverage the dataset to study the behavioral properties of deep learning models.

II. RELATED WORK

Deep learning models are attaining human-level performance in many tasks; however, they lack the robustness and generalization of human vision, which led to a comparison of strategies in deep learning models and human vision [37], [38], [43]. The comparisons are based on external behavior [44], [45], [46], [47] or by comparing internal representations [48], [49], [50], [51]. In the initial part of this section, we examine the dual-thinking framework and explore studies that compare human vision with deep learning models. Subsequently, we discuss existing datasets, highlighting their limitations in capturing the dual-thinking framework in human vision.

A. Dual Thinking

The dual thinking process can be modeled to include an initial feed-forward stage, followed by later-stage generative processing, as present in current studies. For example, Tschantz introduces hybrid predictive coding, which integrates generative and predictive approaches to visual perception [52]. Similarly, studies employing generative classifiers have demonstrated their alignment with human performance in classification tasks [53]. However, these studies are predominantly quantitative, offering limited insights into the underlying reasoning processes in human visual perception, and do not align with current progress in multi-modal LLMs [54].

B. Visual Processing

As deep learning models have increasingly approached human-level performance in classification tasks, researchers have shifted toward qualitative analyses comparing human vision and deep learning models. A significant finding in these studies is the texture bias in deep learning models [37], [38] and shape bias in human vision. The robust and accurate models had more shape bias, which led to techniques to increase shape bias in deep learning models [38], [55], [56], [57], [58], [59]. However, models have shown good performance without shape information (e.g. BagNets [60]), and shape information can be extracted from later layers despite not being used for classification [43]. Subsequent work also shows that image stylization serves as a strong augmentation with

shape bias emerging as a byproduct [61], [56]. While shape may not be essential for discriminatory tasks, these findings highlight the need to explore its significance in tasks such as instance segmentation, which involves object localization and is more similar to human vision. The computational analysis of segmentation task in cognitive science represents segmentation as interval graphs [62], [63], [64], which do not directly extend to image segmentation as instances can be composed of non-adjacent or disconnected segments (e.g., amodal closure in Fig. 2e), requiring separate analysis.

C. Adversarial Datasets

Deep learning models made many recent advances, helping them approach human-level performance on many tasks [65], [66], [53], [45], [67], [68], [69] including on benchmark datasets like ImageNet [70]. This led to the introduction of many datasets that expose different vulnerabilities of deep learning models [71], [72], [73], [74]; including adversarial datasets like Image Net-A, consisting of natural images that were tough for the ResNet-50 [39]. However, they are easier for humans and lack instance-level annotations. With deep learning models making progress on many of these adversarial datasets [75], [76], there is a need for datasets that are challenging for human vision to evaluate the progress and understand the behaviors and strategies in human vision and deep learning models. Stylized Image Net (SIN) enforces conflict between texture and shape by synthetically transferring texture from one class to object of another class [38]; however, it lacks diversity as it can study only texture shape bias and does not have a correct ground truth class. In contrast, images in our dataset have a ground truth class (e.g., Fig. 2f in our dataset is a human hand and not a hand with cat texture, whereas Fig. 1b from SIN is a cat shape with an elephant texture). The images in our datasets have a single stable state compared to bi-stable images. Similarly, another set of datasets focuses only on camouflaged or concealed objects (or figure-ground errors) [77], [78], [79], [80], [81]. These datasets cannot evaluate instance segmentation models trained on benchmark MS-COCO dataset [82] as the class of objects do not match. We do not constrain on any specific cue, which helps us address the gap in research about datasets that help cognitive scientists explore multiple aspects of human visual processing [83], [28], [84]. Our dataset has images with inherent cue conflicts, with the cue for the correct perception contradicting the cues used for intuitive processing. These images can enable the evaluation of the model qualitatively to understand its inductive biases, which can be helpful while designing systems for use in resource-constrained settings [85], [86], [87].

D. Visual Reasoning Datasets

Visual question-answering datasets introduced earlier [88], [89], [90], [91], [92] are comparatively easy for current large language models. MMMU is a multi-discipline, multi-modal benchmark dataset for visual reasoning with questions on science, technology, and others [41] and requires specific knowledge of these fields. Similarly, the MATH-Vision dataset has 3,040 mathematical problems with visual contexts sourced

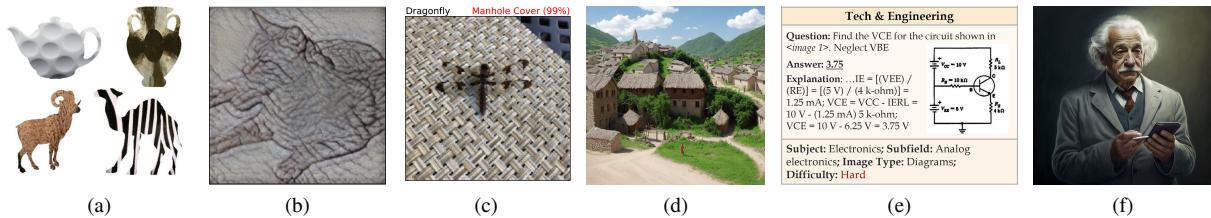


Fig. 1: (a) Global Shape [37] (b) Shape vs Texture from SIN dataset [38] (c) Images in ImageNet-A confuse deep learning models but not human vision [39] (d) An example from synthetically generated IllusionBench dataset [40] (e) An example from MMMU dataset for MLLMs [41] (f) Example from WHOOPS dataset showing Einstein using a phone that was not available during his time [42].

from real math competitions [93]. While the images in these datasets have visual information that is simple to understand, answering the questions requires complex analysis and reasoning on the visual information gathered from these images (Fig.1e). In contrast, gathering visual information from images in our dataset requires reasoning to select the correct solution from different possible solutions. As deep learning models increasing applications find use in real-world settings like autonomous driving, there is a research gap for datasets that can help analyze and understand model limitations.

E. Multi-modal Behavioral Datasets

Deep learning models may display unanticipated biases, prompting the creation of datasets designed to uncover their behavioral patterns [71], [72], [73], [74], [39]. Stylized ImageNet (SIN) [38] and IllusionBench [40] investigate shape bias by introducing conflicts between texture and shape through synthetically generated illusory images. Similarly, several datasets have been developed to specifically study camouflaged or concealed objects [77], [78], [79], [80], [81]. WHOOPS dataset [42] comprises synthetic images depicting unusual scenarios, requiring models to recognize anomalous elements. However, these images often necessitate additional contextual knowledge beyond generic object recognition (Fig.1f). Multimodal Visual Patterns (MMVP) benchmark exploits differences in CLIP and DINOv2 embeddings to identify seemingly simple queries that lead to failure in multimodal large language models (MLLMs) [94]. While several studies have investigated specific capabilities—such as fine-grained classification [95], [96], spatial arrangement [96], state of object [97], camera angle [97], color [98], rotation [98], [97], reflection [98], size transformations [98], relative positioning [97], background [97], binding problem in feed-forward processing [99] and counting of objects [96], [98], [97]—these investigations primarily assess specific capabilities rather than their impact on improving visual perception or their role in logical processing in visual perception.

While existing studies and datasets assess various capabilities and reveal specific vulnerabilities, they do not enable a direct comparison of processing in human vision and deep learning models. Such a comparison is crucial for understanding the factors contributing to human perception's generalizability and robustness. As deep learning models increasingly achieve human-level performance on quantitative metrics, there is a

need for a diverse and challenging dataset that can effectively benchmark these models and provide insights into their inductive biases.

III. DUAL THINKING FRAMEWORK

In this section, we discuss the computational advantages of using a dual-thinking framework for visual processing and the significance of existing theories on human vision within this framework.

Theorem : *If we have multiple possible groupings involving a common region, with different prior probabilities, the lowest average number of iterations to reach correct grouping or segmentation is achieved by completely assigning the component to highest probabilities sequentially.*

Let $p_1, p_2, p_3 \dots p_n$ be the probability in the order of most probable grouping, such that $p_1 > p_2$ and $p_2 > p_3$ and so on and there are n possible groupings. The most probable grouping will be correct in the first iteration for p_1 fraction of the time, and the next highest probable grouping will take two iterations and be correct p_2 fraction of the time, and so on. The average length is given by

$$L_{avg} = 1p_1 + 2p_2 + 3p_3 \dots + np_n = \sum_{i=1}^n ip_i \quad (1)$$

We can prove the theorem by contradiction by showing that deviation from the condition in theorem increases the average length. Let us consider a case where a combination with probabilities p_x and p_y (such that $p_x > p_y$) are assigned position j and i , where $j > i$ to prove the theorem by contradiction. In this case the higher probability combination is tried out at later stage. The difference in average length in this case compared to correct assignment would be $(jp_x + ip_y) - (ip_x + jp_y) = p_x(j - i) - p_y(j - i) = (p_x - p_y)(j - i)$. Since $p_x > p_y$ and $j > i$, this will result in a positive quantity. As deviating from the condition stated in the theorem causes an increase in the average length, the best possibility is the condition stated in the theorem.

We can also note that the $\sum_{i=1}^n p_i = 1$ and when $p_1 \rightarrow 1$ other p_i tends to 0 for $i \neq 1$.

$$\lim_{p_1 \rightarrow 1} L_{avg} \rightarrow 1 \quad (2)$$

The intuitive processing in human vision is quite accurate; hence, we usually do not perceive the first inference. However,

multiple inferences are perceivable in images in our dataset as the first inference is wrong and differs from the final inference. For example, in Fig.2a referenced later, interpreting the shadow as a skirt worn by the second lady is the human intuitive perception, and ignoring the shadow and considering only the second lady is the correct perception.

a) *Sub-component Formation:* Let us consider a group of pixels $\mathbf{I} = \{i_{x,y}\}$ where $i_{x,y}$ represent pixel at (x, y) position. Instance segmentation creates disjoint and nonempty subsets of pixels (s_m) that correspond to object instances or the background and satisfy the conditions that each pixel can belong to only one instance, and the segments must collectively cover the entire image (eq. 3).

$$S = \{s_0, s_1, \dots, s_T\} \quad \text{where} \quad s_m \cap s_n = \emptyset \quad \text{and} \quad \bigcup_{m=0}^T s_m = I \quad (3)$$

$$c_j = \{i_{x,y}\} \quad \text{where} \quad \bigcup_j c_j \cup s_0 = I \quad \text{and} \quad c_j \cap c_k = \emptyset \quad (4)$$

In segmentation, the complexity or hardness is informally considered as the size of the search space [63], [62], [100], and the number of possible instances or the subsets possible for the set \mathbf{I} is $O(2^K)$, where K is the number of pixels in the image. A common approach is to reduce the number of elements [101], [64]. We can observe that figure-ground extraction can reduce a significant number of pixels by ignoring the background (s_0). The formation of sub-components (C) using Gestalt principles (e.g., proximity and similarity) or edges can reduce possible combinations from $O(2^K)$ to $O(2^{card(C)})$. This sub-component formation is also amenable for iterative processing as rectification of incorrect grouping in pixel-based processing can result in redundant intermediate steps observed in segmentation models that use iterative processing [102], [103], [104], and prevent partial sub-component assignments (Fig.4e in the results section) different types of boundary errors [105].

b) *Instance and Class Identification:* The shape can help evaluate the possibility of the groupings in this search space to form an instance (e.g. in Fig.4c, with more probable shapes that correspond to common postures can be prioritized for faster computation in intuitive processing (similar to Theorem 1), and can be evaluated till $\bigcup_j c_j \cup s_0 = I$ is satisfied (e.g. Fig.4c). The space of all possible groupings of components can be searched optimally to identify instances using shape as it prioritizes sub-groups that fit the different object shapes or postures.

c) *Logical Processing:* Evaluating all possible combinations is computationally infeasible (eq.1). Instead, combinations can be examined iteratively until no errors are found in the generated instances. The absence of errors in identified instances can serve as an early stopping criterion, improving energy efficiency in human vision. Processing time can be further optimized by using insights from previous evaluations to reorder subsequent assessments based on updated proba-

bilities or by limiting the search to predefined sub-groups or sub-regions (eq.5).

$$L_{avg} = 1p_1 + 2p'_2 + 3p'_3 \dots + np'_n = p_1 + \sum_{i=2}^n ip'_i \quad (5)$$

Early stopping criteria can enhance the energy efficiency of human vision; however, they may also allow specific errors, undetected by logical processing, to persist. Recent research on intuitive processing (or System 1) and logical processing (or deliberate System 2 analysis) emphasizes that users consciously examine each region [106] to achieve complete image coverage in critical domains like medical imaging [107]. Our study examines errors corrected through logical processing and those that persist in the results section.

IV. METHODS

In the first part, we introduce Human confusion dataset consisting of images adversarial for human vision to study dual thinking framework and deep learning models. We then propose a method to automatically evaluate model correctness and similarity to intuitive processing in human vision.

A. Human confusion dataset

Grouping elements with similar properties, including brightness, color, contrast, texture, etc, is called similarity (e.g., Fig.2c). Continuity is important in human vision and is concerned with grouping elements that form smooth contours (e.g. Fig.2c), and proximity refers to closer elements likely to be grouped [11] (e.g. Fig.2b). Human vision focuses only on salient regions of the image (or foreground) for analysis, and an example of a figure-ground error is shown in Fig.2a and 2d. Amodal closure denotes the completion of occluded elements [108], [109], and errors mainly occur when we perceive parts of different instances as a single instance (e.g., Fig.2e). The errors can also be due to a combination of many properties, for example, similarity, proximity, and continuity in Fig.2c and similarity and proximity in Fig.2b. A comprehensive review of Gestalt principles is provided in [11], [16]; and examples in supplementary data.

1) *Data Collection:* We crowd-source the initial data using Amazon Turk, where participants collect 150 confusing images by searching for "confusing photos", "images that need to be looked twice to understand," and similar terms. The initial three participants collected generic images, and the following four participants collected confusing images related to common objects. Crowd-sourcing with the initial seven Amazon Mechanical Turk workers predominantly yielded popular, ambiguous images, with substantial overlap and a significant proportion of synthetic content. This outcome is consistent with the power-law distribution governing online content popularity [110]. To address this limitation and expand the dataset, we engaged three additional participants to collect confusing images that were not part of the crowd-sourced set. In total, ten participants contributed to the image collection, including seven from Amazon Turk and three for dataset expansion.

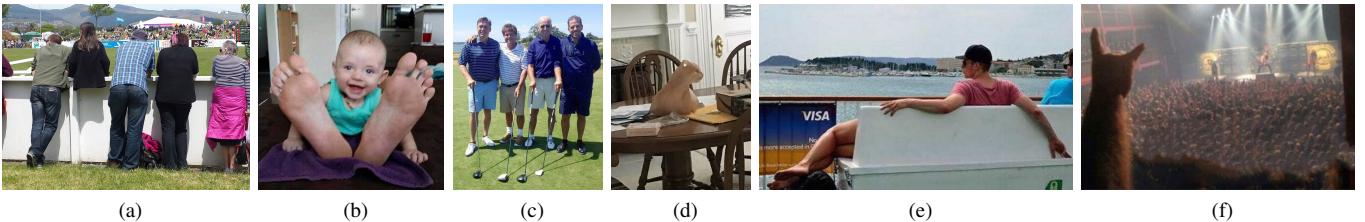


Fig. 2: (a) Figure ground with a part of background grouped with the second lady (b) Proximity and similarity error where the legs are grouped with the child. (c) Similarity, continuity and proximity error occurring together, the shade on hand makes it similar to adjacent shirt and forms a continuous smooth curve. (d) Figure ground error where cover is wrongly identified as cat (e) Amodal Closure and Similarity error where components from two instances that are apart grouped together as a single instance. (f) An example of image that are not analyzed using segmentation models as the confusion is about class of identified instance

2) *Data Preparation*: We remove duplicate and synthetic images from the collected dataset and filter images where the instance segmentation task cannot identify the confusion or the confusion pertains to objects that do not belong to classes in the MS-COCO dataset. The final dataset consists of 1000 images without any NSFW content in correct perception with confusing instances belonging to classes in MS-COCO dataset [82]. Using open-clip embeddings, we identified over 350 images matching entries in the LAION-5B dataset [111] (NeurIPS-2022), making them available for research. Annotations include details about image presence in existing datasets and those available in the open domain. We have shared the URL's of all 1000 images in the dataset.

We annotate each image on (a) Gestalt principles (figure-ground, proximity, continuity, similarity, and amodal closure) except modal closure to analyze qualitative behavior, (b) logical errors (size and count difference) to evaluate vulnerabilities, (c) instances with global shape of one instance but containing instances of different class to evaluate shape bias, and (d) images in wild that can help understand out of distribution behavior. The annotation also includes the operations and masks to check the correctness and similarity to the intuitive perception.

B. Psychophysical Experiments

We conducted experiments with 100 participants recruited through the external agency InFolks, comprising 67 males and 33 females. The participants had an average age of 26.89 years (median: 27), with ages ranging from 23 to 33. Each participant was shown a random subset of 100 or 200 images from our dataset, ensuring that the entire dataset was covered across all participants. During early trials, we collected feedback on perceived confusion in the images to verify that it aligned with the confusion patterns represented in our masks. The first two questions—whether participants “*found anything confusing in the image*”—helps us confirm that the images contained ambiguous regions, while the follow-up question—whether they “*could immediately recognize the correct interpretation*”—allowed us to assess for multiple inferences and dual thinking. The questions in our study included: (a) Did you find anything confusing? (b) Did you

immediately recognize it? (c) Would you have identified the confusion if the image had not been labeled as confusing? (d) For participants who initially found no confusion, we showed them the confusion masks and then asked if they found that confusing. These questions also allowed us to investigate early stopping in human visual processing.

C. Model Evaluation

Multi-modal LLM's currently do not support instance segmentation and respond only in text. To evaluate LLM's, we frame questions that allow us to assess whether the model's response aligns with intuitive or correct perception. The questions had answers in Yes/No format or numeric format and prompts were used to instruct the LLM's to respond accordingly (“*Answer only Yes or No*” or “*Answer only in Numeric*”). The model outputs were then compared against annotated correct answers and intuitive perceptions. To ensure fine-grained evaluation, we separately assess confusing regions in segmentation models to prevent them from being overwhelmed by overall segmentation accuracy [105]. In our analysis of images in the dataset, we observe that the errors in human perception mainly fall into the following categories (a) an object is missed out completely, (b) a part of an instance or background is wrongly grouped with another instance (c) two instances group to form a single instance (d) a part of the background identifies as an object or (e) an instance is misclassified to be of a different class. We identify the confusing instance using the Intersection over Union (IoU) score with the instance mask and the model outputs to identify confusing objects. We use separate masks to identify human perception and correct perception to make evaluation independent of the size of the confusing region. The algorithm to evaluate the different errors is shown in Alg.1 and more details are present in the supplementary data.

D. Experimental Setting

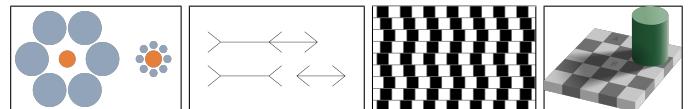
In our study, we consider twelve multi-modal LLM's, GPT-4o mini [112], GPT-4o [113] and reasoning model O3 from OpenAI and open source model LIAMA 3.2 11B and LIAMA 3.2 90B [114] from Meta, Deepseek-VL from Deepseek, Qwen2.5-VL-7B-Instruct and Qwen2.5-VL-72B-Instruct [115]

Algorithm 1 Basic automatic evaluation algorithm

```

Input: modelOutput, instanceMask, confusionMask, includeConfusion, cThreshold
Output: correctStatus
1: matches ← []
2: for each instance ∈ modelOutput do
3:   score ← IoU(instance, instanceMask)
4:   if (score > match_threshold) then
5:     matches.insert((instance,score))
6:   end if
7: end for
8: matches.sort(key=score,reverse=True)
9: result = False
10: for each match, score ∈ matches do
11:   overlap ← IoU((match*confusionMask,
12:   confusionMask)
13:   if (overlap > cThreshold) and includeConfusion then
14:     result = True
15:   end if
16:   if (overlap < cThreshold) and not includeConfusion
17:     then
18:       result = True
19: end for
return result

```



(a) Optical illusions are usually shortcomings of human vision (Ebbinghaus, Müller-Lyer, Café Wall, Checker Shadow Illusion).



(b) In dual thinking the false perception in intuitive processing is easily identified and immediately corrected.



(c) Shape can group components even without specific regions like eyes in intuitive perception.



(d) Shape is an important factor in grouping sub-components.

Fig. 3: Multiple inferences indicates dual thinking framework.

V. RESULTS

In this section, we provide evidence for the dual thinking framework in human vision and discuss different strategies used in these stages. We then analyze deep learning models using the method described earlier.

A. Human Visual Processing

In this section, we first examine the evidence for dual thinking and aspects of intuitive and logical processing using our Human Confusion dataset. We then present the results of psychophysical experiments and demonstrate the usefulness of our dataset to study dual thinking in human vision.

a) Dual thinking framework: We could perceive multiple inferences, including the wrong first inference followed by a correct inference in our dataset, indicating the presence of a dual-thinking framework. The initial inference is corrected almost immediately (e.g., Fig.3d,3c), which is in contrast to optical illusions used to demonstrate the shortcomings of human vision, where the decision does not usually change (Fig.3a). In some images, the logical processing does not identify the errors; however, the users can recognize them retrospectively, which validates a common assumption in computational models of human vision about processing till a confidence threshold is attained [52]. For example, the camouflaged bird in Fig.4d is not easily identifiable, but as these are not limitations of human vision, the users can identify the bird when mentioning the location.

from Alibaba, Gemini-2.0-Flash [116] from Google, Sonnet-4 from Anthropic [117], Molmo-7B-D [118] from Allen AI and Pixtral-12B [119] from Mistral AI. We use the API from the respective companies for models from OpenAI, Anthropic and Google, we use API from Hyperbolic for Qwen2.5-VL-7B-Instruct, Qwen2.5-VL-72B-Instruct and Pixtral-12B. We had setup LlAMA and Deepseek-VL (version 2) locally and Molmo-7B-D on Google Colab for our experiments. While O3, GPT-4o mini, GPT-4o, Sonnet-4 and Gemini-2.0-Flash were closed, LlAMA, Qwen2.5-VL, Pixtral-12B and Molmo-7B-D were open weight models.

We use segmentation models from MMDetection framework [120], trained on the MS-COCO dataset [82], and consider model outputs above the confidence threshold of 0.7. We evaluate models on 983 images except for 17 images with confusion related to the labels (e.g., Fig.2f). We select the thresholds for the automatic evaluation by comparing human annotations on a randomly selected subset of 250 images from our dataset on seven models. The threshold for recognizing instance is 73%, 25% overlap of confusion region not part of an instance, and 70% for being part of an instance. The automatic detection algorithm matched the human annotation decisions in 94.94% of cases. We randomly selected 1,000 images from the MS-COCO dataset and generated questions similar in style to those in our Human Confusion Dataset. We then evaluated three models—4o-mini, Qwen-7B, and sonnet-4—on these questions. The models achieved higher accuracy on the MS-COCO data (93.8%, 92.9%, and 88.9%, respectively) compared to their overall accuracy on our dataset (70%, 65%, and 62%).

Instance and Class Identification Shape as a method for grouping components to form instances is observable in grouping errors due to common posture (e.g. Fig.4c), and the examples where two components from two instances with the shape of an instance (e.g., Fig.3c). The priority for shape over color for grouping components in intuitive processing is observable in amodal closure errors, where sub-components of a single class with widely different colors are grouped as they have the shape of a single instance. The shape also had priority over length in grouping components that are widely separated. Our dataset provides evidence for the importance of shape in instance identification compared to existing studies that focus on shape for classification. We can observe the formation of sub-components, which is evident in errors that occur in groups of related pixels. Gestalt principles and edge-based processing can help in sub-component formation; however, the adversarial examples only indicate sub-component formation but cannot distinguish between these approaches.

b) Psychophysical Trials: Our analysis of psychophysical trials showed that participants found on average, 96% of the images confusing. In 51.69% images the confusion was immediately recognizable, indicating the errors in intuitive processing that are corrected immediately and can be considered examples of dual thinking. Participants also noted that, for many images, they might not have noted the confusing regions without prior mention as confusing images before the experiment. This suggests early stopping in human vision, where processing stops after the perception deems the information sufficient.

B. Segmentation Models

Our study observes that human vision uses a top-down approach similar to Gestalt theory, focusing on overall structure and shape, compared to deep learning models that focus more on details, similar to the structuralist theory of human vision. In many cases, the models generate outputs with both correct perception and perception similar to intuitive processing, as they do not learn the constraints in eq.3, a trivial logic in human visual processing. In general, the texture bias in deep learning models helped segmentation models perform slightly better on figure-ground errors (24 – 40%) (e.g., Fig.4a, 4b and 6a), and the smaller receptive field helps segmentation models avoid amodal closure errors in our adversarial dataset (23 – 66%) (Fig.6a). The dataset and the proposed approach can help evaluate changes in the qualitative behavior of models with changes in model architecture, training procedure, or dataset and augmentation approaches. Despite better performance in categories like amodal closure and figure-ground errors, all the segmentation models predicted wrongly on 47.81% of images, and swin-transformer was most accurate at 34%.

YOLACT model [122], developed for real-time processing, utilizes an architecture that generates and groups sub-components. YOLACT demonstrated a shape bias of 17, significantly higher than the next highest model, GroIE [123], which exhibited a shape bias of only 5.5 on images in which the global shape aligned with one class but comprised two

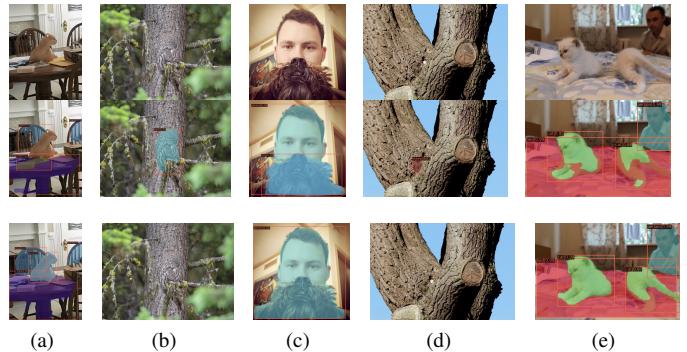


Fig. 4: Top row has the input image, second row contains response of Instaboost with ResNet-101 backbone [121] and third row contain outputs of Swin-Transformer (a) Shape focus in human vision does not text for texture and identifies as cat which is a cover (b,c) Texture focus help models identify the camouflaged bird and the dog (d) The texture focus identifies the bird but assigns wrong class as sheep (e) Amodal closure and the second row shows absence of sub-component formation causing errors on boundaries.

distinct instances from different classes (e.g., Fig. 3d). In these instances, intuitive processing predominantly relies on shape, whereas accurate perception requires attention to texture. YOLACT's pronounced shape bias suggests that, like human vision, segmentation models inherently learn to use shape cues for grouping sub-components. Additionally, YOLACT exhibited a higher ratio of human intuitive references to correct inferences overall and across Gestalt properties, except for amodal closure errors. Future research can further investigate segmentation models as computational frameworks for modeling intuitive processing in human vision and address a limitation of using classification models, which lack object localization capabilities.

C. Logical Processing

In this section, we examine logical processing in human vision, which helps detect and correct errors that arise in instances formed during intuitive processing. We explore how human vision assesses the validity, integrity, and coherence of instances perceived in intuitive processing. Using thematic analysis, we identify common logical errors in perception, analyze the mechanisms of logical processing in vision, and compare the ability of segmentation and vision-language models to interpret such images. The mistakes in intuitive processing in human vision are corrected during subsequent processing, whereas segmentation models face a significant limitation in their inability to identify and rectify logical errors. The segmentation models lack sub-component structure and depend on probabilistic grouping of components, leading to errors related to count mismatches, size differences, and failure to recognize impossible patterns.

Sub-component Differences : Human vision can distinguish instances where objects of different sizes appear together with the shape of a single object in the second stage of processing. Using overall shape to group components can make intuitive



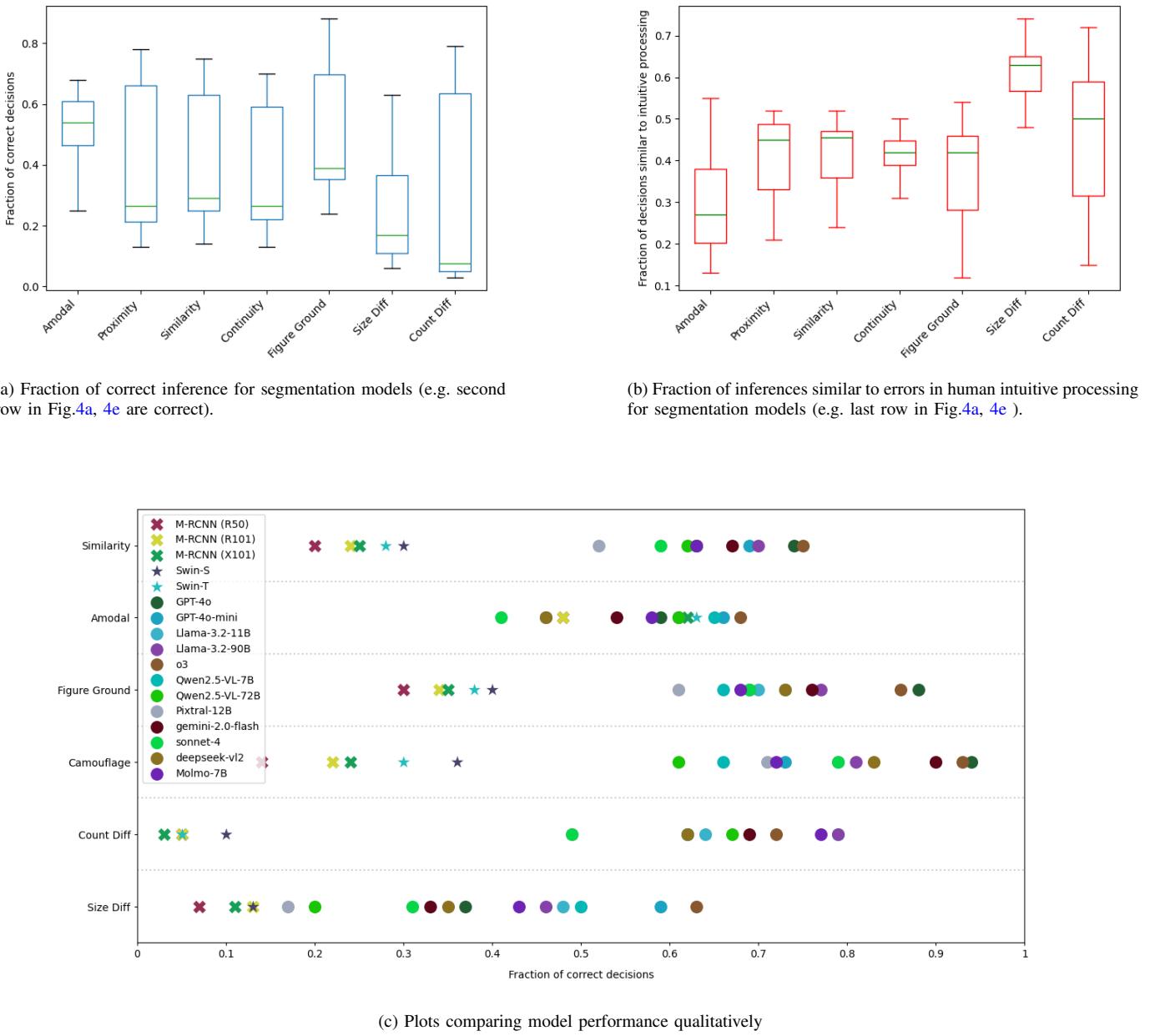
Fig. 5: Input image, segmentation model (swin-t) output and response of multi-modal LLMs. Segmentation models shows poor performance on logical errors, where as VLMs have made progress. Size differences are difficult for VLMs still to understand.

processing prone to these errors. Multi-modal LLMs have made progress but are still in the lower ranges (20 – 63%) on images with size difference in components (Fig.5a, 6c). Segmentation models had lower accuracy (6 – 17%) and made errors similar to human intuitive processing (56 – 81%), as they are unable to recognize the relative size of components. Similarly, multi-modal LLMs also learns amodal closure as an emergent property and tend to merge disconnected components by shape, hence performing poorly on adversarial examples (Fig.6c). They lack finer reasoning of sub-components as in human vision (e.g., wrongly merged components of different genders are identified from clothing and accessories), and such examples requiring finer analysis may be rare in the dataset, making it difficult to learn during training [124].

Count Mismatch : Segmentation models made the same mistake as intuitive processing in human vision in 56 – 62% of images (Fig.6b). The components grouped using the most probable shape tend to group nearby components without sub-component structure, and most models had poor accuracy (3 – 10%). For example, in Fig.5d, the common shape is used to form the first instance, and the remaining components form the

next instance in human vision. In many cases, we also observe all the components to be grouped in intuitive processing (e.g., Fig.5b). The models seem to rely on a probabilistic grouping of the components without substructures and identify them as a single instance. Multi-modal LLMs made progress at (49 – 79%), where the errors in count difference occur primarily with sub-components higher than possible for the class and are learnable [125].

Impossible Patterns : Human vision can identify objects in poses that are not feasible (e.g., Fig.5c,5d). In the case of segmentation models, this capability is missing as it relies on probabilistic grouping. Human vision also assesses the validity of instances formed by grouping sub-components during intuitive processing, facilitating the identification and correction of errors when a sub-component merges with the background—a limitation inherent in human intuitive processing (e.g., a girl's skirt blending with the background would render an instance consisting solely of the upper body invalid). In these instances, the human visual system can infer the missing part by re-evaluating nearby probable regions to validate the object. Additionally, segmentation models split sub-components and



partially assign them to two instances by splitting on proximity rather than assigning a sub-component entirely to one instance. The formation of sub-components in human vision can help avoid these errors.

Multi-modal LLMs have made progress and show ability to understand substructures and the relationships between sub-components. When queried without constraints limiting responses to Yes/No, these models can identify and localize substructures, likely learning their relationships from textual descriptions rather than explicit annotations. In contrast, segmentation models lack sub-component structure and rely on grouping sub-components without logical underpinnings. Their poor performance with errors similar to intuitive processing indicates the absence of the second stage of processing

that helps detect and recover from these errors (More examples are present in the supplementary data).

In general, models should be correct, robust, and generalizable and not limited by human performance or similarity in strategies to human vision. Deep learning models have been criticized for shortcut learning by biasing on texture [26], while our dataset shows that human vision is also prone to errors due to shape bias; however, logical processing helps human vision recover from many of these errors. Deep learning models have outperformed human vision in many tasks that require the analysis of finer patterns [126], [127], [128], [129], and performance on these tasks can be improved with scaling [130] (for example, the biggest improvements have been made in figure ground, camouflage and similarity

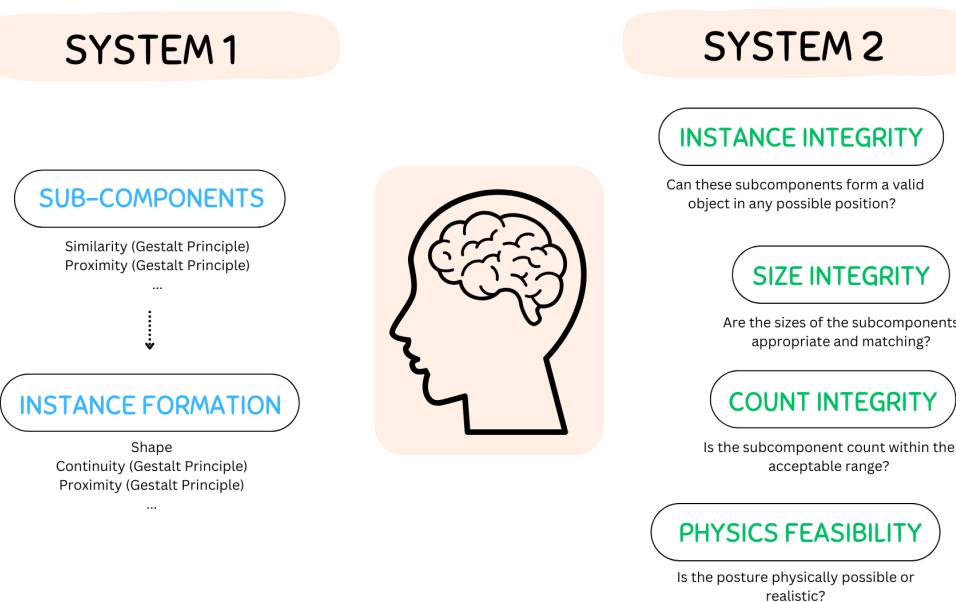


Fig. 7: Human visual processing uses a dual thinking framework which comprises of a intuitive (or System 1 / fast) and logical (or System 2 / deliberate / slow). Intuitive inference uses Gestalt principles (like similarity and proximity) to quickly group sub-components and form instances. Logical processing validates those instances using logical processing before final perception.

which require identifying finer patterns Fig.6c). The coarse to fine processing and early stopping can enable efficient processing in human vision. From an energy standpoint, the human brain is the most expensive organ [7]; hence, visual processing must operate within energy and size constraints.

Our results empirically support a dual-process account of human vision, wherein rapid intuitive judgments are followed by logical refinement. Psychophysical experiments validate this two-stage process, highlighting humans' ability to recognize and correct perceptual errors. Segmentation models capture only the initial intuitive stage, lacking the capacity for logical correction. Multi-modal LLMs have made notable progress in reasoning and sub-component interpretation, outperforming segmentation models. However, they remain sensitive to fine-grained analysis such as size and count mismatches. In summary, our findings underscore logical reasoning as a key differentiator in human vision, and integrating structured reasoning or hybrid processing strategies could improve model robustness in real-world, ambiguous scenarios.

VI. DISCUSSION

Initial approaches to reasoning in LLM involve planning the solution approach and guiding the answer generation through multiple steps or 'thoughts' (e.g. chain of thought). This approach was further improved by generating many answers and choosing the final response by voting (e.g. Self-Consistency with Chain of Thought [131]). However, this method does not necessarily select the most suitable answer through reasoning and encounters performance saturation by scaling [132]. In contrast, human visual processing makes an initial inference and then verifies for logical errors. In case of error, the solution search is constrained to the region with error to select

the appropriate solutions for that specific region, which can also be more efficient. This is conceptually similar to the tree of thought method of reasoning [133], which enables backtracking in case of error in a multi-step approach (more similar to eq.5 rather than eq.1). However, unlike the step-by-step evaluations of the tree of thought method to generate the solution [134], human visual processing appears to create a complete initial inference before verifying for errors.

In this study, we investigate the role of reasoning in low-complexity visual perception to enhance robustness and address rare errors often overlooked in supervised learning. A key aspect of efficient perception is determining when reasoning is necessary. Human vision achieves this by assessing the integrity of perceived objects using basic factual knowledge. Incorporating verification steps, such as follow-up questioning, to enforce factual consistency (e.g., Fig.8) can help correct count-related errors in LLMs, which possess relevant knowledge but fail to apply it during inference. Large reasoning models (LRMs) underperform on simple tasks compared to LLMs, but outperform them on moderately complex ones [135]. While complex tasks (e.g., Fig.1e) demand more advanced reasoning, it is crucial to determine when such reasoning is needed. Human vision avoids the overthinking problem by choosing the first (or most probable) correct answer without errors and stopping early [135]. Additionally, human visual reasoning is influenced by cognitive biases. For instance, subtle cues such as clothing style, fit, or skin texture may be associated with gender in human perception, even though such features may not be present in representations learned by deep models. These differences underscore the distinct nature of reasoning in human vision compared to that in large language models (LLMs). Bridging this gap may

involve aligning model behavior more closely with human reasoning using methods such as Reinforcement Learning from Human Feedback (RLHF), which can help models better handle ambiguous or visually subtle scenarios.



How many legs of the woman are visible in this picture? Answer with a Numeric value

4

How many legs can a woman have?
Answer with a numeric value

2

so how many legs of the woman are visible in this picture? Answer with a numeric value

2

Fig. 8: Incorporating reasoning improves factual consistency in large language models. Example showing the LLaMA 3.2 11B Vision-Instruct model using internal knowledge to self-correct its initial response.

Reasoning and evaluation play a crucial role in identifying and addressing rare errors, often termed the "Curse of Rarity"—anomalies that arise infrequently in real-world

("in-the-wild") deployments yet remain critical to resolve in safety-sensitive systems [136]. Traditional machine learning approaches, which rely on pattern recognition from training data, may struggle with such cases due to entropy [137] or the scarcity of these patterns in training datasets, making it difficult for models to handle these scenarios effectively. This challenge is particularly significant as autonomous systems are increasingly deployed in safety-critical domains such as autonomous driving. Our dataset provides a more accurate representation of real-world scenarios than specialized datasets for testing model capabilities without requiring domain-specific information or external tool integration. For instance, solving questions in the MMMU dataset for multimodal LLMs [41] (Fig. 1e) often necessitates a structured, step-by-step reasoning approach and specialized domain knowledge. In contrast, our dataset primarily relies on immediate perception, which aligns more closely with real-world visual cognition.

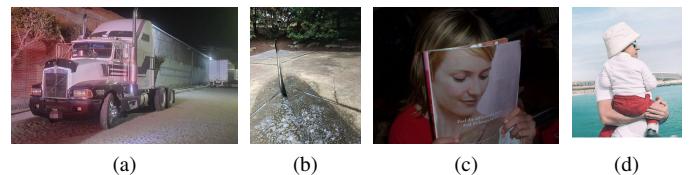


Fig. 9: Examples from the human confusion dataset that can be useful for specific domains. (a) Automatic parking assist can ignore available spaces. (b) Driving assist can fail to recognize presence of cars and can cause collision. (c) Facial spoofing can compromise facial recognition systems using aligned photos or masks. (d) Person recognition systems can fail to identify instances correctly.

Large language models (LLMs) are computationally demanding and energy-intensive, making them impractical for many real-time or resource-constrained environments. In such cases, smaller, task-specific models are more suitable—such as those used in autonomous driving, parking assistance, facial recognition, robotics, retail surveillance, and drone operations. These systems are more likely to encounter visually ambiguous or partially obscured objects in real-world settings, increasing the risk of misinterpretation. For instance, autonomous vehicles may misidentify static objects or other vehicles as drivable paths (e.g. Fig. 9b) or drivable areas for parking vehicles (e.g. Fig. 9a). Security and surveillance systems may fail to detect individuals due to facial spoofing techniques, such as printed photos, background blending, clothing-based obfuscation, or challenging lighting. In high-stakes settings like airports, facial recognition systems can be tricked by fake or printed faces (e.g. Fig. 9c), resulting in false identifications or undetected access breaches. These examples highlight the need for robust and reliable models capable of handling real-world ambiguity, particularly where deploying large models is infeasible. Systems could also be designed to flag potentially confusing inputs for human review or more intensive processing, adding a layer of safety and effective design of human-in-the-loop systems. Our dataset provides a valuable source for assessing such edge cases, which can be

extended and tailored to these specific safety-critical domains to support the development and evaluation of models.

Human vision employs logical reasoning to detect errors similar to "tree of thought" strategy that helps in efficient computation, compared to approaches, such as chain-of-thought with self-consistency that face limitations in selecting optimal solutions. A key aspect of human reasoning is the integration of factual knowledge to validate initial percepts—a process that current LLMs, despite possessing such knowledge, often fail to invoke without explicit prompting. Bridging this gap may require structured reasoning frameworks, reinforcement learning from human feedback (RLHF), or hybrid architectures that can perform logical processing for more robust and generalizable AI systems.

Limitations : Our dataset is adversarial, requiring models to rely on cues that are not commonly encountered. As a result, strong performance on this dataset may not directly translate to improved metrics on standard benchmark datasets. While our study highlights the shortcomings of deep learning models compared to human visual processing, the solutions to these challenges may differ due to the distinct constraints within which each system operates. An optimal approach for deep learning models may not necessarily align with the mechanisms of human vision. For instance, human vision prioritizes rapid initial responses, which can be evolutionarily advantageous—for example, immediately recognizing a potential threat, such as a dangerous animal, and preparing to act. Our dataset primarily includes common object categories to ensure broad compatibility with general-purpose pretrained models; it may not fully represent challenges involving rare or domain-specific objects. Developing robust deep learning systems for safety-critical and resource-constrained environments may require domain-specific extensions of the dataset, incorporating adversarial examples that reflect the unique challenges and failure modes relevant to the target application domain.

VII. CONCLUSION

Our adversarial dataset provides insights into the dual-thinking framework in human visual processing, and computational analysis shows the advantages of human visual processing. The analysis of model behavior shows the limitations of segmentation models in performing logical processing in the second stage of human vision and the tremendous progress made by multi-modal LLMs. However, LLMs still lack fine-grained reasoning abilities, including recognizing size differences. Our study shows that segmentation models can replicate characteristics of human vision, indicating their potential to serve as computational models for human vision. Additionally, our dataset can help cognitive scientists evaluate new theories on human vision and validate assumptions about human visual processing.

ACKNOWLEDGMENT

The authors thank Dr. Sumeet Agarwal for the initial guidance and discussions, and Shubham Jain and Abhilash Kankokaran for their insightful discussions and contributions. The authors also acknowledge the contributions of Gayathri,

Shafi, Suneer, Nikhil VM, Tanushka, and others for their support in data collection and annotation. The authors acknowledge the use of Grammarly and ChatGPT for text correction.

REFERENCES

- [1] R. VanRullen, "The power of the feed-forward sweep," *Advances in Cognitive Psychology*, vol. 3, no. 1-2, p. 167, 2007.
- [2] R. S. van Bergen and N. Kriegeskorte, "Going in circles is the way forward: the role of recurrence in visual inference," *Current Opinion in Neurobiology*, vol. 65, pp. 176–193, 2020.
- [3] Y. Mohsenzadeh, S. Qin, R. M. Cichy, and D. Pantazis, "Ultra-rapid serial visual presentation reveals dynamics of feedforward and feedback processes in the ventral visual pathway," *Elife*, vol. 7, p. e36329, 2018.
- [4] K. Daniel, "Thinking, fast and slow," 2017.
- [5] C. Chen, K. Hammernik, C. Ouyang, C. Qin, W. Bai, and D. Rueckert, "Cooperative training and latent space data augmentation for robust medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 149–159, Springer, 2021.
- [6] T. Grootswagers, A. K. Robinson, and T. A. Carlson, "The representational dynamics of visual objects in rapid serial visual processing streams," *NeuroImage*, vol. 188, pp. 668–679, 2019.
- [7] G. Kreiman and T. Serre, "Beyond the feedforward sweep: feedback computations in the visual cortex," *Annals of the New York Academy of Sciences*, vol. 1464, no. 1, pp. 222–241, 2020.
- [8] S. Thorpe, D. Fize, and C. Marlot, "Speed of processing in the human visual system," *nature*, vol. 381, no. 6582, pp. 520–522, 1996.
- [9] H. Tang, M. Schrimpf, W. Lotter, C. Moerman, A. Paredes, J. Ortega Caro, W. Hardesty, D. Cox, and G. Kreiman, "Recurrent computations for visual pattern completion," *Proceedings of the National Academy of Sciences*, vol. 115, no. 35, pp. 8835–8840, 2018.
- [10] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.
- [11] J. Wagemans, J. H. Elder, M. Kubovy, S. E. Palmer, M. A. Peterson, M. Singh, and R. von der Heydt, "A century of gestalt psychology in visual perception: I. perceptual grouping and figure-ground organization," *Psychological bulletin*, vol. 138, no. 6, p. 1172, 2012.
- [12] E. Van Geert and J. Wagemans, "Prägnanz in visual perception," *Psychonomic Bulletin & Review*, pp. 1–27, 2023.
- [13] J. R. Pomerantz and M. C. Portillo, "Grouping and emergent features in vision: toward a theory of basic gestalts," *Journal of Experimental Psychology: Human Perception and Performance*, vol. 37, no. 5, p. 1331, 2011.
- [14] M. Wertheimer, "Experimentelle studien über das sehen von bewegung," *Zeitschrift für psychologie*, vol. 61, 1912.
- [15] J. R. Pomerantz, L. C. Sager, and R. J. Stoever, "Perception of wholes and of their component parts: some configural superiority effects," *Journal of Experimental Psychology: Human Perception and Performance*, vol. 3, no. 3, p. 422, 1977.
- [16] J. Wagemans, J. Feldman, S. Gepshtain, R. Kimchi, J. R. Pomerantz, P. A. Van der Helm, and C. Van Leeuwen, "A century of gestalt psychology in visual perception: II. conceptual and theoretical foundations," *Psychological bulletin*, vol. 138, no. 6, p. 1218, 2012.
- [17] V. Biscione and J. S. Bowers, "Mixed evidence for gestalt grouping in deep neural networks," *Computational Brain & Behavior*, vol. 6, no. 3, pp. 438–456, 2023.
- [18] B. Kim, E. Reif, M. Wattenberg, S. Bengio, and M. C. Mozer, "Neural networks trained on natural scenes exhibit gestalt closure," *Computational Brain & Behavior*, vol. 4, no. 3, pp. 251–263, 2021.
- [19] B. L. Anderson, "Can computational goals inform theories of vision?," *Topics in Cognitive Science*, vol. 7, no. 2, pp. 274–286, 2015.
- [20] B. C. Love, "The algorithmic level is the bridge between computation and brain," *Topics in cognitive science*, vol. 7, no. 2, pp. 230–242, 2015.
- [21] D. Peebles and R. P. Cooper, "Thirty years after marr's vision: Levels of analysis in cognitive science," 2015.
- [22] C. T. Zahn, "Graph-theoretical methods for detecting and describing gestalt clusters," *IEEE Transactions on computers*, vol. 100, no. 1, pp. 68–86, 1971.
- [23] G. Papari and N. Petkov, "Adaptive pseudo dilation for gestalt edge grouping and contour detection," *IEEE transactions on image processing*, vol. 17, no. 10, pp. 1950–1962, 2008.

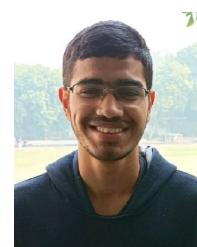
- [24] I.-C. Shen and W.-H. Cheng, "Gestalt rule feature points," *IEEE Transactions on Multimedia*, vol. 17, no. 4, pp. 526–537, 2015.
- [25] L. Muttenthaler, L. Linhardt, J. Dippel, R. A. Vandermeulen, K. Hermann, A. Lampinen, and S. Kornblith, "Improving neural network representations using human similarity judgments," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [26] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann, "Shortcut learning in deep neural networks," *Nature Machine Intelligence*, vol. 2, no. 11, pp. 665–673, 2020.
- [27] J. Dapello, T. Marques, M. Schrimpf, F. Geiger, D. Cox, and J. J. DiCarlo, "Simulating a primary visual cortex at the front of cnns improves robustness to image perturbations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 13073–13087, 2020.
- [28] D. Linsley, D. Shiebler, S. Eberhardt, and T. Serre, "Learning what and where to attend," in *International Conference on Learning Representations*, 2019.
- [29] D. Linsley, A. K. Ashok, L. N. Govindarajan, R. Liu, and T. Serre, "Stable and expressive recurrent vision models," *Advances in neural information processing systems*, 2020.
- [30] T. Fel, I. F. Rodriguez Rodriguez, D. Linsley, and T. Serre, "Harmonizing the object recognition strategies of deep neural networks with humans," *Advances in neural information processing systems*, vol. 35, pp. 9432–9446, 2022.
- [31] L. Vogelsang, S. Gilad-Gutnick, E. Ehrenberg, A. Yonas, S. Diamond, R. Held, and P. Sinha, "Potential downside of high initial visual acuity," *Proceedings of the National Academy of Sciences*, vol. 115, no. 44, pp. 11333–11338, 2018.
- [32] A. Saxe, S. Nelli, and C. Summerfield, "If deep learning is the answer, what is the question?," *Nature Reviews Neuroscience*, vol. 22, no. 1, pp. 55–67, 2021.
- [33] J. Kim, D. Linsley, K. Thakkar, and T. Serre, "Disentangling neural mechanisms for perceptual grouping," in *International Conference on Learning Representations*, 2019.
- [34] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman, "Building machines that learn and think like people," *Behavioral and brain sciences*, vol. 40, p. e253, 2017.
- [35] J. S. Bowers, G. Malhotra, M. Dujmović, M. L. Montero, C. Tsvetkov, V. Biscione, G. Puebla, F. Adolphi, J. E. Hummel, R. F. Heaton, et al., "Deep problems with neural network models of human vision," *Behavioral and Brain Sciences*, pp. 1–74, 2022.
- [36] T. Serre, "Deep learning: the good, the bad, and the ugly," *Annual review of vision science*, vol. 5, no. 1, pp. 399–426, 2019.
- [37] N. Baker, H. Lu, G. Erlikhman, and P. J. Kellman, "Deep convolutional networks do not classify based on global object shape," *PLoS computational biology*, vol. 14, no. 12, p. e1006613, 2018.
- [38] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel, "Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness," in *International Conference on Learning Representations*, 2018.
- [39] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song, "Natural adversarial examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15262–15271, 2021.
- [40] A. Hemmat, A. Davies, T. Lamb, J. Yuan, P. Torr, A. Khakzar, and F. Pinto, "Hidden in plain sight: Evaluating abstract shape recognition in vision-language models," *Advances in Neural Information Processing Systems*, vol. 37, pp. 88527–88556, 2025.
- [41] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun, et al., "Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- [42] N. Bitton-Guetta, Y. Bitton, J. Hessel, L. Schmidt, Y. Elovici, G. Stanovsky, and R. Schwartz, "Breaking common sense: Whoops! a vision-and-language benchmark of synthetic and compositional images," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2616–2627, 2023.
- [43] K. L. Hermann, T. Chen, and S. Kornblith, "The origins and prevalence of texture bias in convolutional neural networks," *Advances in neural information processing systems*, 2020.
- [44] R. Geirhos, K. Meding, and F. A. Wichmann, "Beyond accuracy: quantifying trial-by-trial behaviour of cnns and humans by measuring error consistency," *Advances in Neural Information Processing Systems*, vol. 33, pp. 13890–13902, 2020.
- [45] R. Geirhos, K. Narayananappa, B. Mitzkus, T. Thieringer, M. Bethge, F. A. Wichmann, and W. Brendel, "Partial success in closing the gap between human and machine vision," *Advances in Neural Information Processing Systems*, vol. 34, pp. 23885–23899, 2021.
- [46] R. Geirhos, C. Medina Temme, J. Rauber, H. Schütt, M. Bethge, and F. Wichmann, "Generalisation in humans and deep neural networks," in *Thirty-second Annual Conference on Neural Information Processing Systems 2018 (NeurIPS 2018)*, pp. 7549–7561, Curran, 2019.
- [47] S. Tuli, I. Dasgupta, E. Grant, and T. Griffiths, "Are convolutional neural networks or transformers more like human vision?," in *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 43, 2021.
- [48] C. F. Cadieu, H. Hong, D. L. Yamins, N. Pinto, D. Ardila, E. A. Solomon, N. J. Majaj, and J. J. DiCarlo, "Deep neural networks rival the representation of primate it cortex for core visual object recognition," *PLoS Comput Biol*, vol. 10, no. 12, p. e1003963, 2014.
- [49] G. Jacob, R. Pramod, H. Katti, and S. Arun, "Qualitative similarities and differences in visual object representations between brains and deep networks," *Nature communications*, vol. 12, no. 1, pp. 1–14, 2021.
- [50] J. Mehrer, C. J. Spoerer, N. Kriegeskorte, and T. C. Kietzmann, "Individual differences among deep neural network models," *Nature communications*, vol. 11, no. 1, p. 5725, 2020.
- [51] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy, "Do vision transformers see like convolutional neural networks?," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [52] A. Tschantz, B. Millidge, A. K. Seth, and C. L. Buckley, "Hybrid predictive coding: Inferring, fast and slow," *PLOS Computational Biology*, vol. 19, no. 8, p. e1011280, 2023.
- [53] P. Jaini, K. Clark, and R. Geirhos, "Intriguing properties of generative classifiers," 2024.
- [54] C. Snell, J. Lee, K. Xu, and A. Kumar, "Scaling lilm test-time compute optimally can be more effective than scaling model parameters," *arXiv preprint arXiv:2408.03314*, 2024.
- [55] Z. Xu, D. Liu, J. Yang, C. Raffel, and M. Niethammer, "Robust and generalizable visual representation learning via random convolutions," 2021.
- [56] Y. Li, Q. Yu, M. Tan, J. Mei, P. Tang, W. Shen, A. Yuille, and C. Xie, "Shape-texture debiased neural network training," in *International Conference on Learning Representations*, 2021.
- [57] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial examples are not bugs, they are features," *Advances in neural information processing systems*, vol. 32, 2019.
- [58] A. Subramanian, E. Sizikova, N. Majaj, and D. Pelli, "Spatial-frequency channels, shape bias, and adversarial robustness," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [59] M. Sun, Z. Li, C. Xiao, H. Qiu, B. Kailkhura, M. Liu, and B. Li, "Can shape structure features improve model robustness under diverse adversarial settings?," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7526–7535, 2021.
- [60] W. Brendel and M. Bethge, "Approximating cnns with bag-of-local-features models works surprisingly well on imagenet," in *International Conference on Learning Representations*, 2018.
- [61] C. K. Mummadi, R. Subramaniam, R. Hutmacher, J. Vitay, V. Fischer, and J. H. Metzen, "Does enhanced shape bias improve neural network robustness to common corruptions?," 2021.
- [62] F. Adolphi, T. Wareham, and I. van Rooij, "Computational complexity of segmentation," in *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 44, 2022.
- [63] F. Adolphi, T. Wareham, and I. van Rooij, "A computational complexity perspective on segmentation as a cognitive subcomputation," *Topics in Cognitive Science*, 2022.
- [64] K. J. Friston, N. Sajid, D. R. Quiroga-Martinez, T. Parr, C. J. Price, and E. Holmes, "Active listening," *Hearing research*, vol. 399, p. 107998, 2021.
- [65] V. Boutin and F. Thomas, "Diffusion models as artists: Are we closing the gap between humans and machines?," in *International Conference on Machine Learning*, 2023.
- [66] V. Boutin, L. Singhal, X. Thomas, and T. Serre, "Diversity vs. recognizability: Human-like generalization in one-shot generative models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 20933–20946, 2022.
- [67] A. J. O'Toole and C. D. Castillo, "Face recognition by humans and machines: three fundamental advances from deep learning," *Annual Review of Vision Science*, vol. 7, pp. 543–570, 2021.
- [68] C. J. Parde, V. E. Strehle, V. Banerjee, Y. Hu, J. G. Cavazos, C. D. Castillo, and A. J. O'Toole, "Twin identification over viewpoint change: A deep convolutional neural network surpasses humans," *ACM Transactions on Applied Perception*, vol. 20, no. 3, pp. 1–15, 2023.

- [69] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- [70] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [71] A. Abbas and S. Deny, "Progress and limitations of deep networks to recognize objects in unusual poses," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, pp. 160–168, 2023.
- [72] M. A. Alcorn, Q. Li, Z. Gong, C. Wang, L. Mai, W.-S. Ku, and A. Nguyen, "Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4845–4854, 2019.
- [73] H. Hosseini and R. Poovendran, "Semantic adversarial examples," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1614–1619, 2018.
- [74] Y. Dong, Q.-A. Fu, X. Yang, T. Pang, H. Su, Z. Xiao, and J. Zhu, "Benchmarking adversarial robustness on image classification," in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 321–331, 2020.
- [75] M. Lee and D. Kim, "Robust evaluation of diffusion-based adversarial purification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 134–144, 2023.
- [76] Y. Li, B. Xie, S. Guo, Y. Yang, and B. Xiao, "A survey of robustness and safety of 2d and 3d deep learning models against adversarial attacks," *ACM Computing Surveys*, vol. 56, no. 6, pp. 1–37, 2024.
- [77] D.-P. Fan, G.-P. Ji, G. Sun, M.-M. Cheng, J. Shen, and L. Shao, "Camouflaged object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2777–2787, 2020.
- [78] T.-N. Le, Y. Cao, T.-C. Nguyen, M.-Q. Le, K.-D. Nguyen, T.-T. Do, M.-T. Tran, and T. V. Nguyen, "Camouflaged instance segmentation in-the-wild: Dataset, method, and benchmark suite," *IEEE Transactions on Image Processing*, vol. 31, pp. 287–300, 2021.
- [79] D.-P. Fan, G.-P. Ji, M.-M. Cheng, and L. Shao, "Concealed object detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 10, pp. 6024–6042, 2021.
- [80] Y. Lv, J. Zhang, Y. Dai, A. Li, B. Liu, N. Barnes, and D.-P. Fan, "Simultaneously localize, segment and rank the camouflaged objects," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11591–11601, 2021.
- [81] P. Skurowski, H. Abdulameer, J. Błaszczyk, T. Depta, A. Kornacki, and P. Koziel, "Animal camouflage analysis: Chameleon database," *Unpublished manuscript*, vol. 2, no. 6, p. 7, 2018.
- [82] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*, pp. 740–755, Springer, 2014.
- [83] T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to predict where humans look," in *2009 IEEE 12th international conference on computer vision*, pp. 2106–2113, IEEE, 2009.
- [84] D. B. Walther, B. Choi, E. Caddigan, D. M. Beck, and L. Fei-Fei, "Simple line drawings suffice for functional mri decoding of natural scene categories," *Proceedings of the National Academy of Sciences*, vol. 108, no. 23, pp. 9661–9666, 2011.
- [85] G. Bachmann, S. Anagnostidis, and T. Hofmann, "Scaling mlps: A tale of inductive bias," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [86] Y. Xu, Q. ZHANG, J. Zhang, and D. Tao, "Vitae: Vision transformer advanced by exploring intrinsic inductive bias," in *Advances in Neural Information Processing Systems*, 2021.
- [87] K. Dayanandan and B. Lall, "Enabling multi-modal conversational interface for clinical imaging," in *Extended abstracts of the 2020 CHI conference on human factors in computing systems*, 2024.
- [88] D. A. Hudson and C. D. Manning, "Gqa: A new dataset for real-world visual reasoning and compositional question answering," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6700–6709, 2019.
- [89] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the v in vqa matter: Elevating the role of image understanding in visual question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6904–6913, 2017.
- [90] K. Marino, M. Rastegari, A. Farhadi, and R. Mottaghi, "Ok-vqa: A visual question answering benchmark requiring external knowledge," in *Proceedings of the IEEE/cvpr conference on computer vision and pattern recognition*, pp. 3195–3204, 2019.
- [91] A. F. Biten, R. Tito, A. Mafra, L. Gomez, M. Rusinol, E. Valveny, C. Jawahar, and D. Karatzas, "Scene text visual question answering," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4291–4301, 2019.
- [92] D. Teney, E. Abbasnejad, and A. van den Hengel, "Unshuffling data for improved generalization in visual question answering," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1417–1427, 2021.
- [93] K. Wang, J. Pan, W. Shi, Z. Lu, M. Zhan, and H. Li, "Measuring multimodal mathematical reasoning with math-vision dataset," *arXiv preprint arXiv:2402.14804*, 2024.
- [94] S. Tong, Z. Liu, Y. Zhai, Y. Ma, Y. LeCun, and S. Xie, "Eyes wide shut? exploring the visual shortcomings of multimodal llms," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9568–9578, 2024.
- [95] J. Kim and H. Ji, "Finer: Investigating and enhancing fine-grained visual concept recognition in large vision language models," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 6187–6207, 2024.
- [96] S. Chandhok, W.-C. Fan, and L. Sigal, "Response wide shut: Surprising observations in basic vision language model capabilities," *arXiv preprint arXiv:2408.06721*, 2024.
- [97] M. Gaur, D. Singh, and M. Tapaswi, "Detect, describe, discriminate: Moving beyond vqa for mllm evaluation," *CoRR*, 2024.
- [98] E. Yiu, M. Qraitem, C. Wong, A. N. Majhi, Y. Bai, S. Ginosar, A. Gopnik, and K. Saenko, "Kiva: Kid-inspired visual analogies for testing large multimodal models," *CoRR*, 2024.
- [99] D. Campbell, S. Rane, T. Giallanza, C. N. De Sabbata, K. Ghods, A. Joshi, A. Ku, S. Frankland, T. Griffiths, J. D. Cohen, et al., "Understanding the limits of vision language models through the lens of the binding problem," *Advances in Neural Information Processing Systems*, vol. 37, pp. 113436–113460, 2025.
- [100] N. T. Franklin, K. A. Norman, C. Ranganath, J. M. Zacks, and S. J. Gershman, "Structured event memory: A neuro-symbolic model of event cognition," *Psychological Review*, vol. 127, no. 3, p. 327, 2020.
- [101] M. R. Brent, "Speech segmentation and word discovery: A computational perspective," *Trends in Cognitive Sciences*, vol. 3, no. 8, pp. 294–301, 1999.
- [102] D. Kailas and B. Lall, "Spurious equilibrium in segmentation models and recurrent processing in human vision," in *Workshop on Spurious Correlation and Shortcut Learning: Foundations and Solutions, ICLR*, 2025.
- [103] W. Wang, K. Yu, J. Hugonot, P. Fua, and M. Salzmann, "Recurrent u-net for resource-constrained segmentation," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2142–2151, 2019.
- [104] N. K. Tomar, D. Jha, M. A. Riegler, H. D. Johansen, D. Johansen, J. Rittscher, P. Halvorsen, and S. Ali, "Fanet: A feedback attention network for improved biomedical image segmentation," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [105] B. Cheng, R. Girshick, P. Dollár, A. C. Berg, and A. Kirillov, "Boundary iou: Improving object-centric image segmentation evaluation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15334–15342, 2021.
- [106] Z. Buçinca, M. B. Malaya, and K. Z. Gajos, "To trust or to think: cognitive forcing functions can reduce overreliance on ai in ai-assisted decision-making," *Proceedings of the ACM on Human-computer Interaction*, vol. 5, no. CSCW1, pp. 1–21, 2021.
- [107] W. B. Gefter, B. A. Post, and H. Hatabu, "Commonly missed findings on chest radiographs: Causes and consequences," *Chest*, 2022.
- [108] B. Nanay, "The importance of amodal completion in everyday perception," *i-Perception*, vol. 9, no. 4, 2018.
- [109] M. M. Murray, D. M. Foxe, D. C. Javitt, and J. J. Foxe, "Setting boundaries: brain dynamics of modal and amodal illusory shape completion in humans," *Journal of Neuroscience*, vol. 24, no. 31, pp. 6898–6903, 2004.
- [110] L. A. Adamic and B. A. Huberman, "Power-law distribution of the world wide web," *science*, vol. 287, no. 5461, pp. 2115–2115, 2000.
- [111] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman, et al., "Laion-5b: An open large-scale dataset for training next generation image-text models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 25278–25294, 2022.
- [112] OpenAI, "Gpt-4o mini: advancing cost-efficient intelligence," 2024.

- [113] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, *et al.*, "Gpt-4o system card," *arXiv preprint arXiv:2410.21276*, 2024.
- [114] Meta, "Llama 3.2: Revolutionizing edge ai and vision with open, customizable models," 2024.
- [115] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang, *et al.*, "Qwen2. 5-vl technical report," *arXiv preprint arXiv:2502.13923*, 2025.
- [116] G. DeepMind, "Introducing gemini 2.0: our new ai model for the agentic era," 2024.
- [117] Anthropic, "Introducing claude 4," 2025.
- [118] M. Deitke, C. Clark, S. Lee, R. Tripathi, Y. Yang, J. S. Park, M. Salehi, N. Muennighoff, K. Lo, L. Soldaini, *et al.*, "Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models," *CoRR*, 2024.
- [119] P. Agrawal, S. Antoniak, E. B. Hanna, B. Bout, D. Chaplot, J. Chudnovsky, D. Costa, B. De Monicault, S. Garg, T. Gervet, *et al.*, "Pixtral 12b," *arXiv preprint arXiv:2410.07073*, 2024.
- [120] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, *et al.*, "Mmdetection: Open mmlab detection toolbox and benchmark," *arXiv preprint arXiv:1906.07155*, 2019.
- [121] H.-S. Fang, J. Sun, R. Wang, M. Gou, Y.-L. Li, and C. Lu, "Instaboost: Boosting instance segmentation via probability map guided copy-pasting," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 682–691, 2019.
- [122] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "Yolact: Real-time instance segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9157–9166, 2019.
- [123] L. Rossi, A. Karimi, and A. Prati, "A novel region of interest extraction layer for instance segmentation," in *2020 25th international conference on pattern recognition (ICPR)*, pp. 2203–2209, IEEE, 2021.
- [124] E. Francazi, M. Baity-Jesi, and A. Lucchi, "A theoretical analysis of the learning dynamics under class imbalance," in *International Conference on Machine Learning*, pp. 10285–10322, PMLR, 2023.
- [125] T. Gupta and A. Kembhavi, "Visual programming: Compositional visual reasoning without training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14953–14962, 2023.
- [126] P. Shourie, V. Anand, and S. Gupta, "An efficient cnn framework for radiologist level pneumonia detection using chest x-ray images," in *2023 3rd International Conference on Intelligent Technologies (CONIT)*, pp. 1–6, IEEE, 2023.
- [127] A. Y. Hannun, P. Rajpurkar, M. Haghpanahi, G. H. Tison, C. Bourn, M. P. Turakhia, and A. Y. Ng, "Cardiologist-level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network," *Nature medicine*, vol. 25, no. 1, pp. 65–69, 2019.
- [128] P. Ruamviboonsuk, J. Krause, P. Chotcomwongse, R. Sayres, R. Raman, K. Widner, B. J. Campana, S. Phene, K. Hemarat, M. Tadarati, *et al.*, "Deep learning versus human graders for classifying diabetic retinopathy severity in a nationwide screening program," *NPJ digital medicine*, vol. 2, no. 1, p. 25, 2019.
- [129] J. B. Stephanen, A. N. Olesen, M. Olsen, A. Ambati, E. B. Leary, H. E. Moore, O. Carrillo, L. Lin, F. Han, H. Yan, *et al.*, "Neural network analysis of sleep stages enables efficient diagnosis of narcolepsy," *Nature communications*, vol. 9, no. 1, p. 5229, 2018.
- [130] M. I. U. Haque, A. K. Dubey, I. Danciu, A. C. Justice, O. S. Ovchinnikova, and J. D. Hinkle, "Effect of image resolution on automated classification of chest x-rays," *Journal of Medical Imaging*, vol. 10, no. 4, pp. 044503–044503, 2023.
- [131] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," in *The Eleventh International Conference on Learning Representations*, 2023.
- [132] B. Brown, J. Juravsky, R. Ehrlich, R. Clark, Q. V. Le, C. Ré, and A. Mirhoseini, "Large language monkeys: Scaling inference compute with repeated sampling," *arXiv preprint arXiv:2407.21787*, 2024.
- [133] J. Long, "Large language model guided tree-of-thought," *arXiv preprint arXiv:2305.08291*, 2023.
- [134] X. L. Li, V. Shrivastava, S. Li, T. Hashimoto, and P. Liang, "Benchmarking and improving generator-validator consistency of language models," in *The Twelfth International Conference on Learning Representations*, 2024.
- [135] P. Shojaei*†, I. Mirzadeh*, K. Alizadeh, M. Horton, S. Bengio, and M. Farajtabar, "The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity," 2025.



Kailas D is currently pursuing PhD at IIT Delhi, and has completed M.Tech from IIT Delhi in Communication and Radar and B.Tech in Electronics and Telecommunications from SCT College of Engineering. He joined Samsung India in February 2003 and worked for almost a decade in the Wireless Terminal Division. He also worked as Principal Applications Engineer at Oracle for three and half years and have authored two patents related to interpreting human cognition for automatic application generation.



Nikhil Kumar is currently pursuing his B.Tech in Electronics and Communication Engineering at Indraprastha Institute of Information Technology Delhi. His research interests include Multimedia system development, FPGA-based computing, hardware acceleration and embedded systems development.



Anand Sinha completed M.Tech from Indian Institute of Technology Delhi in Computer Technology and B.Tech from Haldia Institute of Technology. He currently is General Manager at Samsung, where he has worked since July 2013. Prior to that he worked at Cognizant Technology Solutions and Wipro Technologies. His areas of interest include Linux system programming, multimedia drivers, audio certification and automotive infotainment.



Brejesh Lall completed PhD in 1997 from IIT Delhi in the area of Multirate Signal Processing. He joined Hughes Software Systems in September 1997 and worked there for nearly 8 years in the Signal Processing group. He returned to his alma mater and joined IIT Delhi as a faculty member in 2005. He is actively working on image processing, signal processing, and communication. These include object representation, tracking and classification, odometry, depth-map generation, and presentation and rendering. He has also served as an expert in numerous

government and private agencies in aspects related to signal processing. He is the former head of Bharti School of Telecom Technology and Management, and the coordinator of two centres of excellence, viz. Airtel IIT Delhi Centre of Excellence in Telecommunications and Ericsson IIT Delhi 5G Center of Excellence, and incharge of an IoT laboratory set up in collaboration with Samsung.