
Synthesizing Images on Perceptual Boundaries of ANNs for Uncovering and Manipulating Human Perceptual Variability

Chen Wei^{*1,2} Chi Zhang^{*1} Jiachen Zou¹ Haotian Deng¹ Dietmar Heinke² Quanying Liu¹

Abstract

Human decision-making in cognitive tasks and daily life exhibits considerable variability, shaped by factors such as task difficulty, individual preferences, and personal experiences. Understanding this variability across individuals is essential for uncovering the perceptual and decision-making mechanisms that humans rely on when faced with uncertainty and ambiguity. We propose a systematic Boundary Alignment Manipulation (BAM) framework for studying human perceptual variability through image generation. BAM combines perceptual boundary sampling in ANNs and human behavioral experiments to systematically investigate this phenomenon. Our perceptual boundary sampling algorithm generates stimuli along ANN perceptual boundaries that intrinsically induce significant perceptual variability. The efficacy of these stimuli is empirically validated through large-scale behavioral experiments involving 246 participants across 116,715 trials, culminating in the variMNIST dataset containing 19,943 systematically annotated images. Through personalized model alignment and adversarial generation, we establish a reliable method for simultaneously predicting and manipulating the divergent perceptual decisions of pairs of participants. This work bridges the gap between computational models and human individual difference research, providing new tools for personalized perception analysis. Code and data for this work are publicly available.



Figure 1. Human perceptual variability. For the same set of stimuli, individuals often exhibit varied responses, highlighting differences in their visual perception. For instance, the digit on the left may be interpreted as a “3” by some and as a “5” by others. Similarly, the image on the right might be initially perceived as a cat by some individuals, while others may perceive it as a bird.

1. Introduction

A core goal of cognitive science is to establish models that reflect the relationship between external stimuli and human internal experiences. The development of ANNs has significantly contributed to this goal, particularly through the latent representations of ANNs that have shown a strong correlation with human psychological representations ((Wei et al., 2024a;b; Muttenthaler et al., 2022a; Mahner et al., 2024; Zheng et al., 2019; Hebart et al., 2020; Muttenthaler et al., 2022b)). Current studies focus on the phenomenon that even with identical physical stimuli, individuals may exhibit significant differences in their perceptual experiences. While such perceptual variability has been extensively documented in complex cognitive tasks (e.g., aesthetic or moral judgments), studies on simple visual decision tasks (e.g., handwritten digit classification) often neglect such inter-individual differences. As shown in Figure 1, the stimulus on the left may be recognized as either a “3” or a “5”. This *high perceptual variability* in human perception has been inadequately explored in current research into human visual perception. Inspired by the similarity between ANNs and humans, we hypothesize that the perceptual boundaries of ANNs are also related to inter-individual variability in human perception and that images generated along these boundaries can evoke divergent perceptual experiences across individuals.

Recent advances in using ANNs as perceptual models have revealed novel relationships between stimuli and human experiences. (Veerabadrán et al., 2023; Zhou & Firestone,

^{*}Equal contribution ¹Department of Biomedical Engineering, Southern University of Science and Technology, Shenzhen, China ²Department of Psychology, University of Birmingham, Birmingham, UK. Correspondence to: Quanying Liu <liuqy@sustech.edu.cn>.

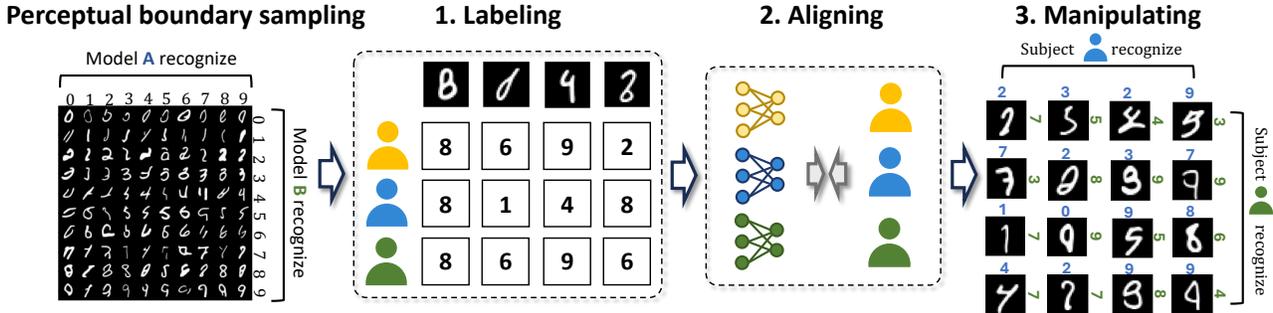


Figure 2. **Overview of BAM.** We sample images from ANN perceptual boundaries using the *perceptual boundary sampling* algorithm for subsequent human evaluations. Our approach consists of three main components: **1. Labeling:** The images generated by perceptual boundary sampling are labeled by human experiments, thus constructing the *variMNIST* dataset. In this process, a single image will be presented to multiple participants; **2. Aligning:** Finetuning models with human behavioral data to align them with human perceptual variability at the group and individual levels, enhancing behavior prediction accuracy and aligning the models with humans; **3. Manipulating:** Employing two individually aligned models, each corresponding to a specific individual, to generate images designed to elicit divergent responses between them, which are then validated through these two human participants.

2019; Elsayed et al., 2018) demonstrated that ANNs not only exhibit vulnerability to small changes but also that these changes can systematically influence human perceptual decisions in controlled environments. Similarly, (Gaziv et al., 2024) leveraged robustified ANNs to identify low-norm perturbation patterns that significantly disrupt human classification behavior. Methodologically, (Feather et al., 2023; 2019) proposed *model metamers* generation—constructing stimuli with identical ANN latent representations but distinct appearances—which revealed systematic discrepancies between model metamers and human perception from the perspective of representational invariance. Additionally, (Golan et al., 2020; 2023) developed *controversial stimuli* that provoke clearly distinct responses among two or more models, further exposing their misalignment with human perception. Building on these frameworks, we extend the concept of controversial stimuli from model-to-model comparisons to human perception, generating stimuli that evoke perceptual divergences among human participants.

We propose a systematic framework BAM (**B**oundary **A**lignment & **M**anipulation framework) for studying *human perceptual variability* through image generation. As illustrated in Figure 2, BAM builds upon the *Perceptual Boundary Sampling* (Sec. 3) algorithm and comprises three interconnected steps: **1. Labeling:** We sample images from ANN perceptual boundaries and construct the *variMNIST* dataset through human behavioral experiments, systematically capturing inter-individual perceptual differences (Sec. 4). **2. Aligning:** By fine-tuning ANN models with human behavioral data, we establish computational models of perceptual variability at both group and individual levels (Sec. 5). **3. Manipulating:** Using individualized models as adversarial generators, we synthesize controversial stimuli that amplify perceptual differences, with experimental

validation of their behavioral manipulation efficacy (Sec. 6).

Our principal contributions are:

- (1) We develop an image sampling method along ANN perceptual boundaries, constructing the *variMNIST* dataset through human behavioral labeling. Experimental results demonstrate successful induction of high perceptual variability.
- (2) We achieve precise alignment between ANN and human perceptual variability using limited behavioral data, proving the computational feasibility of individual differences.
- (3) We design a dual-subject controversial generation framework that synthesizes stimuli inducing targeted divergent decisions. Human experiments confirm these stimuli significantly surpass random baselines in behavioral manipulation.

2. Related Works

Researchers have extensively used synthetic images generated by ANNs to study human perceptual space, uncovering differences between model and human perception while refining generation techniques to enhance their influence on human cognition. For instance, (Golan et al., 2020; 2023) utilized controversial stimuli to highlight classification discrepancies between neural networks. Similarly, (Veerabaran et al., 2023) demonstrated that adversarial perturbations could simultaneously influence ANN classifications and human perceptual choices, revealing shared sensitivities. However, (Gaziv et al., 2024) found that while standard ANN perturbations fail to impact human perception, robustified ANN models can generate low-norm perturbations that significantly disrupt human percepts.

Other studies have approached this problem from different

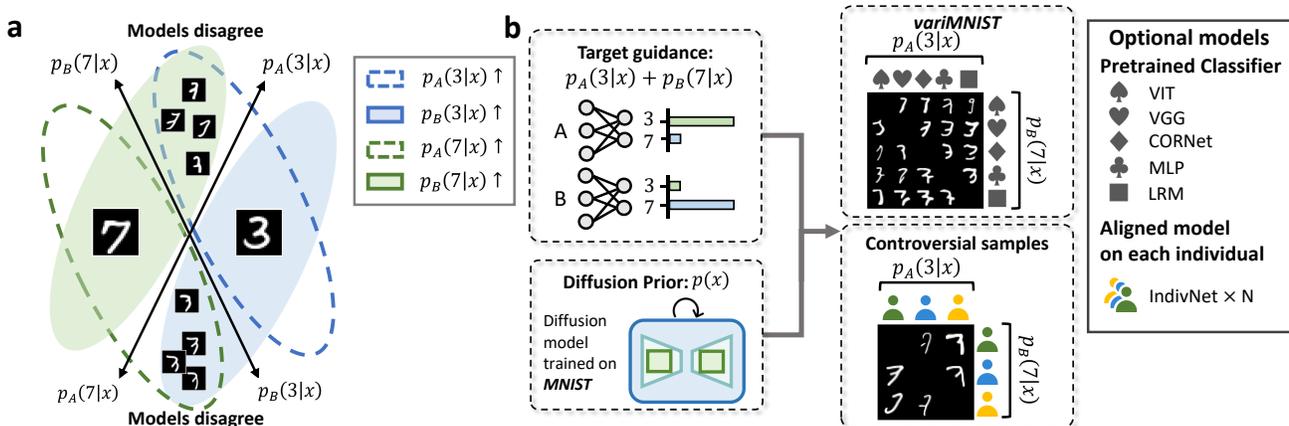


Figure 3. Sampling on perceptual boundaries. (a) The sample space can be partitioned into four distinct regions based on two classification axes. Taking the digit pair (3,7) as an example, our objective is to generate samples that induce disagreement between models A and B, as illustrated in the figure. The upper region consists of stimuli classified as “7” by model B and “3” by model A, whereas the lower region contains stimuli classified as “3” by model A and “7” by model B. The left region includes stimuli that both models classify as “3,” while the right region contains stimuli that both models classify as “7.” (b) Utilizing targeted controversial guidance, we constructed the variMNIST dataset. This approach employs classifier guidance on the diffusion model, directing model A toward “3” and model B toward “7,” thereby constraining the generated samples to lie on perceptual boundaries while preserving the diffusion prior. Following model alignment with human perception, this method was further applied to generate controversial samples designed to modulate human perceptual decisions, as shown in the lower section of the right panel.

angles. For example, works like (Feather et al., 2023; 2019; Nanda et al., 2022; 2023) investigated *model metamers*, revealing fundamental mismatches between model activations and human recognition. Extending beyond perceptual discrepancies, (Fu et al., 2023) introduced DreamSim, a perceptual metric leveraging synthetic data and human experimental data to better reflect human similarity judgments and address shortcomings in conventional perceptual metrics. Building on such synthetic data and behavioral insights, recent efforts have sought to align vision models with human perceptual representations by incorporating human-like conceptual structures, resulting in improved alignment and enhanced performance across diverse tasks (Muttenthaler et al., 2024; Sundaram et al., 2024).

To study the variability of human perception, it is essential that generated images significantly influence human cognition. Given that we sample from the perceptual boundaries of ANNs, which often contain high noise levels, better methods are needed to ensure that the generated images appear natural. Recently, the fields of adversarial examples and counterfactual explanations in machine learning have adopted effective techniques to help deal with this problem, such as (Jeanneret et al., 2023), (Wei et al., 2024b), (Chen et al., 2023), (Jeanneret et al., 2022), (Vaeth et al., 2023), and (Atakan Bedel & Çukur, 2023). These studies use diffusion models with training-free guidance (Yu et al., 2023; Ma et al., 2023; Yang et al., 2024) as regularizers to introduce prior distributions, thereby enhancing the naturalness of generated images and their impact on human perception.

3. Generating Images by Sampling from the Perceptual Boundary of ANNs

The image perturbations that significantly affect ANN perception also influence human perception (Gaziv et al., 2024; Veerabadrán et al., 2023; Wei et al., 2024a; Muttenthaler et al., 2022a), suggesting that ANNs and humans may share similar perceptual boundaries. Based on this, we hypothesize that samples on these boundaries (which exhibit high perceptual variability for ANNs) may also lead to ambiguous perception in humans, resulting in different internal experiences for the same stimuli. You can find the schematic diagram in Figure 3.

We adopted two guidance strategies: uncertainty guidance and controversial guidance. Uncertainty guidance aims to generate images that lie near the perceptual boundaries of classifiers. Its loss function is defined as:

$$\mathcal{L} = H(p_1(y|x), q_1(y))$$

where $H(p, q)$ is the cross-entropy function that measures the discrepancy between the predicted distribution $p(y)$ and the target distribution $q(y)$. The target distribution ensures equal probabilities for two categories (e.g., “3” and “7”), resulting in high-uncertainty images. Controversial guidance generates images that cause conflicting predictions between two classifiers. Its loss function is defined as:

$$\mathcal{L} = H(p_1(y|x), q_1(y)) + H(p_2(y|x), q_2(y)),$$

where $p_1(y|x)$ and $p_2(y|x)$ are the predicted probability distributions of classifiers 1 and 2, and $q_1(y)$ and $q_2(y)$ are their

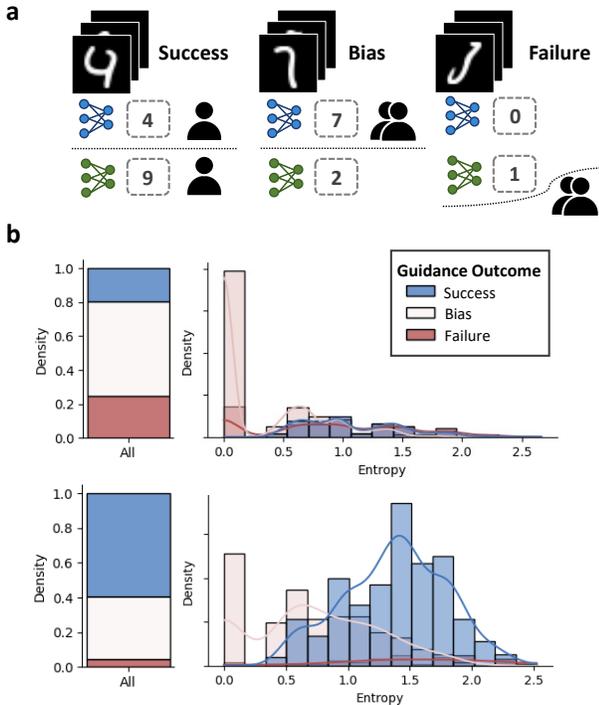


Figure 4. **Controversial guidance influence human perception.** (a) Examples of three types of *guidance outcome*: *success*, *bias*, and *failure*. (b) We present the proportion and entropy distribution of our generated datasets based on handwritten digits and natural images. The upper section displays the results for handwritten digits, while the lower section corresponds to natural images. As observed, eliciting human perceptual variability is more challenging for handwritten digits, as human observers tend to exhibit high agreement on such a straightforward classification task.

corresponding target distributions. The target distributions ensure that classifier 1 predicts one category (e.g., “3”) with high confidence, while classifier 2 predicts another category (e.g., “7”) with high confidence, generating controversial images. Figure A.1 illustrates the guidance methods.

Details of additional analyses and comparisons of guidance methods can be found in Appendix A.2. Previous studies have shown that when using generated images to investigate models and human perception (e.g., (Golan et al., 2020; Gaziv et al., 2024; Veerabadran et al., 2023; Feather et al., 2023)), a common issue is the lack of naturalness in the generated images. This often makes the images difficult for participants to recognize, thereby weakening their impact on human cognition (see Figure A.2). Recent research has demonstrated that diffusion models, when used as regularizers, can introduce prior information and help generate more natural images (Jeanneret et al., 2023; Wei et al., 2024b; Chen et al., 2023; Jeanneret et al., 2022; Vaeth et al., 2023; Atakan Bedel & Çukur, 2023). Building on these findings,

we employ a classifier-guided diffusion model for image generation. This method produces images that are closer to the true distribution of handwritten digits, thereby significantly enhancing their impact on human perception (see Appendix A.1).

4. Collecting Human Perceptual Variability by Recognition Experiment

4.1. Digit Recognition Experiment

We used the image dataset generated through uncertainty or controversial guidance and digit judgment surrogate guidance (section B.1.3) as experimental samples to measure human behavior in a digit recognition task. An illustration of the experiment procedure can be found at Figure 2. For each test image, participants were asked, “What number is this image?” with responses restricted to one of the digits from 0 to 9. We collected the probability distributions of human responses and calculated the average response time and entropy distribution for all test images (Figure A.6). The experiment collected behavioral data from 400 participants, each completing 500 trials, resulting in a total of 200,000 trials across 20,000 stimuli. During data preprocessing, 154 participants were excluded based on Sentinel trials, leaving data from 246 participants (116,715 trials and 19,943 valid stimuli). Using this cleaned dataset, we constructed a high perceptual variability dataset, variMNIST, which serves as a foundation for subsequent analysis and modeling.

4.2. Quantitative Analysis of variMNIST

Evaluation metrics. To comprehensively evaluate the guiding effectiveness of the generation method, we define three types of *guidance outcome*, as illustrated in Figure 4a: *success*, *bias*, and *failure*. For the guidance targets o_1 and o_2 , let p_1 and p_2 represent the probabilities of participants choosing o_1 and o_2 , respectively. A result is considered *success* if $p_1 + p_2 \geq 80\%$ and $\min(p_1, p_2) \geq 10\%$, indicating the generated stimuli guide participants to make a balanced choice between the two targets. A result is labeled as *bias* if $p_1 + p_2 \geq 80\%$ but $\min(p_1, p_2) < 10\%$, indicating a strong bias toward one target. A result is classified as *failure* if $p_1 + p_2 < 80\%$, meaning the stimuli fail to guide participants effectively. These definitions allow us to evaluate and compare the performance of different guidance strategies and classifiers.

ANN variability can arouse human variability. To evaluate whether the images generated by sampling on the perceptual boundaries of ANNs can arouse human perceptual variability, we first calculated the entropy of participants’ choice probabilities in the digit recognition experiment. As shown in Figure. A.6 (bottom left), the entropy values for

more than half of the generated images were significantly greater than zero, indicating substantial variability in human choices. This suggests that the generated images successfully elicited human perceptual variability. Furthermore, as illustrated in Figure 4b, the average sum of success rate and bias rate across all generated images was close to 80%. This indicates that, in the majority of cases, human choices aligned with either both or one of the guidance targets. This demonstrates that the generation method effectively guided human digit recognition behavior. To further analyze the effects of different models and different guidance methods on the outcome, we compared the outcome results, shown in Figure A.15.

5. Predicting Human Perceptual Variability

5.1. Model Fine-tuning for Human Alignment

To align the base models pretrained on MNIST with the performance of both group and individual levels, we adopted a mixed fine-tuning approach with an 80:20 split for training and validation. For individual-level datasets (variMNIST-i), which are subsets of variMNIST corresponding to specific individuals, the validation set was designed to avoid overlap with the group validation set. For group-level fine-tuning, we combined the MNIST and variMNIST datasets in a 1:1 ratio, ensuring performance on MNIST while fine-tuning for perceptual variability. For individual-level fine-tuning, we mixed variMNIST-i, variMNIST, and MNIST datasets in a 2:1:1 ratio, ensuring the models performed effectively on individual-specific, group, and original datasets. For each individual, a customized individual model was trained. For the whole procedure, the group model (GroupNet) was finetuned from the base model (BaseNet), and the individual model (IndivNet) was finetuned from the group model. See Appendix C.2 for more details.

5.2. Alignment Analysis on Validation Datasets

Fine-tuning improves both group-level and individual-level prediction performance. As shown in Figure 5a, BaseNet, GroupNet, and IndivNet achieve nearly identical prediction accuracy on the MNIST dataset, indicating no significant loss of baseline performance after fine-tuning. On the variMNIST dataset, both GroupNet and IndivNet outperform BaseNet by $\sim 20\%$. Furthermore, IndivNet achieves an additional $\sim 5\%$ accuracy improvement over GroupNet on the variMNIST-i dataset, demonstrating its superior adaptability to individual differences. After individual fine-tuning, accuracy improved for 241 participants, while only 5 participants experienced a slight decrease, highlighting the effectiveness of individual fine-tuning in adapting to unique participant behavior and capturing human perceptual variability more accurately.

Different classifiers exhibit inconsistent performance.

Figure 5b and A.8 compares the fine-tuning performance of five classifiers. On the MNIST dataset, group fine-tuning improved the prediction accuracy of VIT (Dosovitskiy et al., 2021) and VGG (Simonyan & Zisserman, 2015), while CORNet (Kubilius et al., 2019) and MLP (Rumelhart et al., 1986) remained unchanged, and LRM (Seber & Lee, 2003) showed a significant decrease in accuracy. On the variMNIST dataset, all classifiers exhibited improvements, with VIT and MLP achieving the largest gains and LRM the smallest. Individual fine-tuning further improved all classifiers, with VIT and MLP showing the greatest adaptability to fine-tuning, while LRM demonstrated weaker generalization ability. These results highlight that both group-level and individual-level fine-tuning can significantly enhance classifier performance, but the degree of improvement depends on the classifier architecture.

Human variability can be predicted by models.

To evaluate the alignment between model and human perceptual variability, we analyzed the correlation between model and human entropy, as shown in Figure 5c and A.16. Taking VGG as an example, group fine-tuning increased the Spearman rank correlation between model and human entropy from $\rho = 0.08$ to $\rho = 0.74$. This significant improvement indicates that fine-tuning enables the model to better capture human uncertainty, aligning model predictions more closely with human perceptual behavior.

Performance of behavior prediction across images with varying entropy levels.

Image entropy reflects task difficulty, with higher entropy indicating more challenging samples. To examine the impact of entropy levels on prediction accuracy, we analyzed model performance across varying entropy levels, as shown in Figure 5d and A.17. Both GroupNet and IndivNet outperform BaseNet across all entropy levels, demonstrating that fine-tuning enhances prediction accuracy regardless of task difficulty. Notably, IndivNet’s performance gains over GroupNet are most pronounced for high-entropy images, suggesting that individual fine-tuning primarily improves prediction accuracy for difficult samples. These findings highlight the ability of fine-tuned models to better handle challenging stimuli, capturing subtle variations in human perceptual behavior more effectively.

Extended validation on ImageNet.

To demonstrate the generalizability of this method, we conducted the same experiments on natural images, instead of only on digits. Following the generation methods outlined in (Gaziv et al., 2024; Wei et al., 2024a;b), we generated and evaluated data based on the ImageNet dataset, with detailed method description provided in the Appendix B.4. It can be observed that these results are consistent with those based on the digits. Examples and further analysis of the results are provided

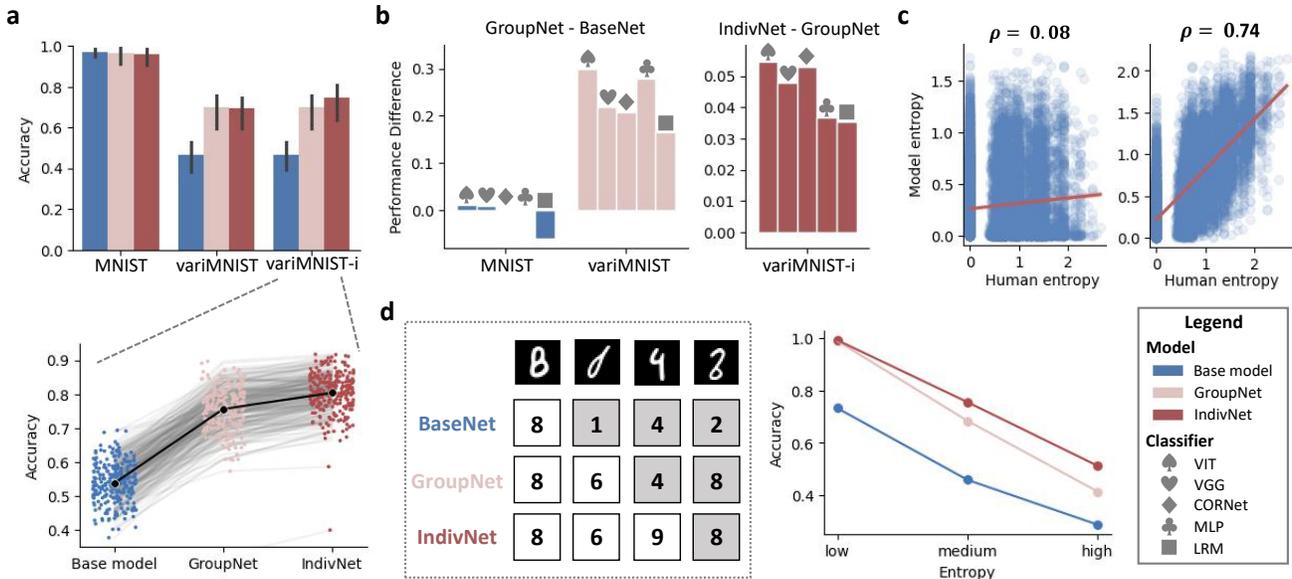


Figure 5. **Human alignment results.** (a) Accuracy of BaseNet, GroupNet, and IndivNet on MNIST, variMNIST, and variMNIST-i. All models performed similarly on MNIST. On variMNIST, GroupNet and IndivNet improved accuracy by $\sim 20\%$ over BaseNet, with IndivNet outperforming GroupNet by $\sim 5\%$ on variMNIST-i. Accuracy improved for 241 participants and decreased for 5 after individual fine-tuning. (b) Fine-tuning results for five classifiers. On MNIST, group fine-tuning improved VIT and VGG, while others remained unchanged or declined. On variMNIST, all classifiers improved, with VIT and MLP showing the largest gains and LRM the smallest. Individual fine-tuning further improved all classifiers with the same trend. (c) For VGG, Spearman rank correlation between model and human entropy increased from $\rho = 0.08$ to $\rho = 0.74$ after group fine-tuning. (d) Performance of BaseNet, GroupNet, and IndivNet of varying entropy levels. The choices from selected subject for the example images are 8, 6, 9, 6, with increasing entropy levels. Here, the gray background indicates that the model’s choice is inconsistent with the subject. GroupNet and IndivNet improved over BaseNet on all entropy levels, while IndivNet’s gains over GroupNet were focused on high-entropy images.

in Appendix B.4.

6. Manipulating Human Perceptual Variability

6.1. Experimental Paradigm

Building on variMNIST and alignment experiments, we designed a paradigm to test whether individually fine-tuned models can amplify perceptual differences and guide decision-making (Figure 6a). This experiment evaluates the ability of targeted stimuli to reveal individual variability and achieve precise manipulation of perceptual outcomes, highlighting the potential of personalized modeling in understanding human perception. For the *first round* of experiments, we initially selected around 500 balanced samples from the variMNIST dataset as stimuli. After collecting behavioral data from pairs of participants, we fine-tuned their individual models using the method described in Section 5.1. Controversial stimuli were then generated using the updated models, aiming to elicit distinct choices between the two participants, with each choosing their respective guidance targets.

In the *second round* of experiments, these controversial stimuli were presented to participants in pairs, with each pair

completing trials designed to test whether the fine-tuned models could effectively guide their decisions in opposite directions. The goal was to evaluate whether the generated stimuli amplified perceptual differences and aligned participants’ responses with their respective guidance targets. For each subject pair, approximately 180 controversial samples were generated, ensuring the sample distribution remained as balanced as possible. A total of 18 participants were recruited for in-lab experiments, grouped into six sets of three participants each. Within each group, participants were paired in all possible combinations, resulting in three pairs per group and 18 pairs overall. Each participant completed 500 trials in the first round and approximately 360 trials (180 per pair, across two pairs) in the second round.

6.2. Manipulating Results

Evaluation metrics. To analyze the effects of individual manipulation, we employed two key metrics. The first metric, referred to as the **guidance outcome** (Figure 6b), was adapted from Section 4.2. It categorizes outcomes for two participants, s_1 and s_2 , with respective guidance targets o_1 and o_2 , and choices c_1 and c_2 . A result is labeled as *success* if both participants’ choices fall within their respec-

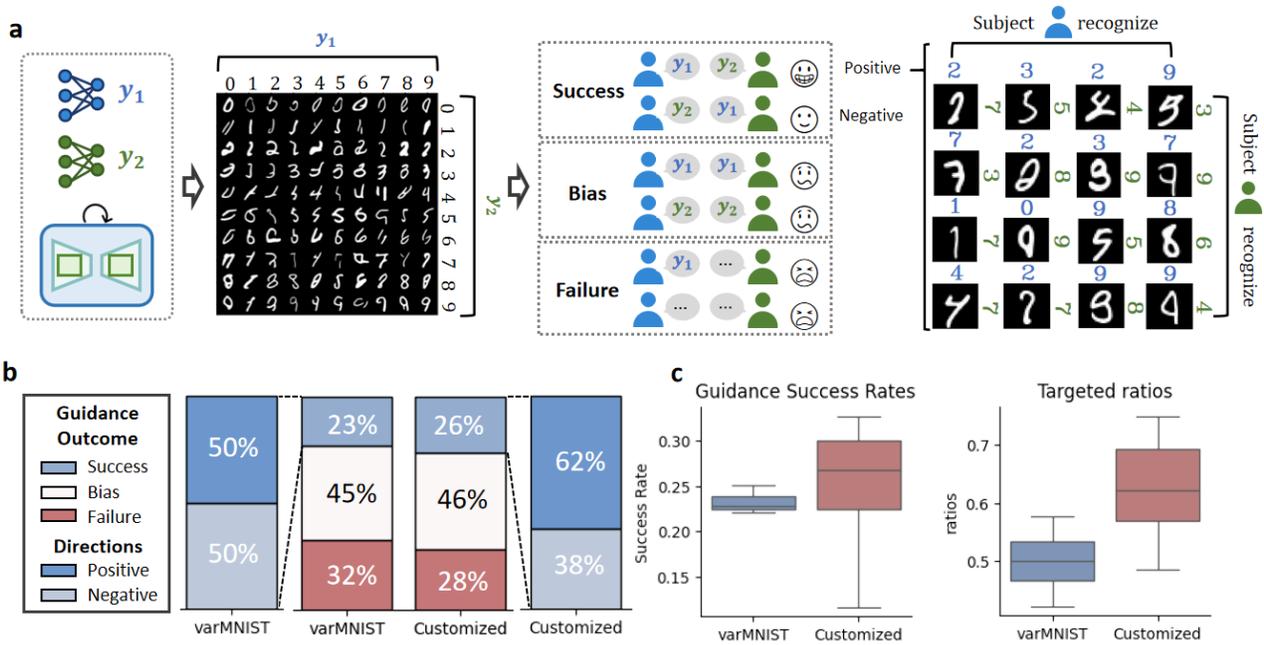


Figure 6. Manipulation analysis. (a) We first utilize behavioral data from the experiments to fine-tune the models, resulting in aligned models, as depicted in the leftmost section with green and blue regions. Subsequently, controversial guidance is applied, directing one model toward output y_1 and the other toward y_2 . The generated images are then presented to participants in a follow-up experiment. Based on the behavioral responses, the images are categorized into three groups: *Success*, *Bias*, and *Failure*. To further analyze the effects of the manipulation, the *Success* category is subdivided into *Positive* and *Negative* cases. Representative examples of *Positive* stimuli are displayed in the rightmost section. (b) The middle two bars show the guidance outcomes for variMNIST and the individually customized dataset, which is generated under the guidance of finetuned individual models to better evoke variability and bias participants toward certain directions, with the individually customized dataset achieving a higher success rate. The left and right bars further analyze the successful samples, where the dark blue indicates the participant’s choices aligned with the guidance direction, and the light blue indicates the opposite. Compared to variMNIST, IndividNets also improves the directionality of guiding perceptual changes. (c) The left panel shows the guidance success rates for the first-round stimuli and the second-round stimuli generated by the finetuned models, with an improvement of $\sim 3\%$ ($p < 0.001$). The right panel shows the *targeted ratios* (i.e., the proportion of participant choices aligned with the guidance direction) for these two groups of stimuli, with an increase of $\sim 12\%$ ($p < 0.001$).

tive guidance targets and are distinct, i.e., $c_1, c_2 \in \{o_1, o_2\}$ and $c_1 \neq c_2$. If both choices are biased toward the same target, such as $c_1 = c_2 = o_1$ or o_2 , it is categorized as *bias*. Finally, if at least one choice is outside the targets ($c_1, c_2 \notin \{o_1, o_2\}$), the outcome is labeled as *failure*. The second metric, called the **targeted ratio** (Figure 6c), quantifies the directionality of successful guidance. Within successful trials, participant choices are classified as either *positive*, where $c_1 = o_1$ and $c_2 = o_2$, meaning both choices align with their respective targets, or *negative*, where $c_1 = o_2$ and $c_2 = o_1$, indicating swapped choices. The targeted ratio is defined as the proportion of positive trials among all success trials, providing a measure of the effectiveness of directional guidance. We present examples of stimuli demonstrating various guidance outcomes and directions in Figure A.18.

Improvement in guidance outcome. We first analyzed the improvements in the guidance outcome achieved through individual manipulation. As shown in Figure. 6, A.27, A.28, A.29, compared to variMNIST, the success rate in the individually customized dataset increased by 3%, the bias rate increased by 1%, and the failure rate decreased by 4%. Considering that each participant completed only around 200 samples in the experiment, compared to 20,000 samples in variMNIST, this represents a very small sample size. Therefore, these results indicate that even with a limited sample size, we successfully captured the perceptual differences among individual participants. These findings validate the feasibility and effectiveness of individual manipulation using small, customized datasets, demonstrating that precise modeling and manipulation of human perceptual behavior can be achieved even at low cost.

Improvement in guiding directionality. We further evaluated the guiding directionality in successful trials. As shown in Figure. 6, A.27, A.28, A.29., compared to variMNIST, the target ratio of IndivNets improved by 12%, indicating a significant enhancement in the directional guidance achieved with individually customized datasets. This result suggests that individual fine-tuning not only improves the model’s guiding capability but also enables more precise directional guidance, leading participants to make choices aligned with the intended targets. This finding further validates the effectiveness of the individual manipulation experiment, demonstrating that small, customized datasets can achieve more efficient and precise human behavior manipulation.

7. Discussion

Human decision-making in cognitive tasks and daily life exhibits considerable variability, shaped by factors such as task difficulty, individual preferences, and personal experiences. Understanding this variability across individuals is essential for uncovering the perceptual and decision-making mechanisms that humans rely on when faced with uncertainty and ambiguity. We present a computational framework that combines perceptual boundary sampling with behavioral manipulation to systematically investigate this phenomenon. By generating stimuli along ANN perceptual boundaries through large-scale behavioral experiments, we construct the variMNIST dataset demonstrating significantly enhanced perceptual variability. Through subject-specific fine-tuning of ANN models using behavioral data, we develop predictive models that capture individual perceptual patterns with high fidelity. We implement adversarial generation strategies to synthesize stimuli that systematically shift perceptual boundaries between individuals, enabling targeted manipulation of perceptual judgments. This work bridges the gap between computational models and human individual difference research, providing new tools for personalized perception analysis.

From the perspectives of cognitive science and neuroscience, our method significantly enhances the utility and flexibility of generated images in the study of human perception. Unlike the methods employed by (Golan et al., 2020) and (Feather et al., 2023), which reveal the disparities between model and human perception by generating images that strongly affect ANNs while having minimal impact on human cognition, our method is capable of influencing both models and human perception simultaneously. This dual impact allows for a nuanced counterfactual examination of the subtle differences in perceptual variability between the two. In contrast to the approaches taken by (Veerabadran et al., 2023) and (Gaziv et al., 2024), which focus on improving ANNs to produce images that can influence human percep-

tion, our use of diffusion models with prior distributions allows for broader applicability across various ANN models and perturbation methods. This also expands the range of image sampling, enabling sampling from high-noise areas like perceptual boundaries. Moreover, the incorporation of prior distributions ensures that our generated images more closely resemble natural images, enhancing their effectiveness in influencing human perception. With this significant improvement in the usability and flexibility of generated images, we successfully explored individual differences in human perception and opened the door for personalized manipulation, increasing the efficiency and scope of human perception studies.

From the perspective of computer science methodology, we have made significant improvements upon existing methods, opening new avenues for fields of AI for science and AI-human alignment. Drawing on the controversial stimuli from (Golan et al., 2020) and adversarial perturbations from (Veerabadran et al., 2023), we integrated these concepts with diffusion model priors to create two new loss-guiding methods: controversial guidance and uncertainty guidance. This enhancement increases the naturalness of the generated images and their influence on human perception. Additionally, inspired by works such as (Jeanneret et al., 2023), (Wei et al., 2024a), (Wei et al., 2024b), (Chen et al., 2023), (Jeanneret et al., 2022), (Vaeth et al., 2023), (Atakan Bedel & Çukur, 2023), we introduced counterfactual methodologies into the study of human perceptual variability, allowing us to explore this relatively under-researched area in greater depth. Our experiments demonstrate that the variMNIST dataset we generated significantly evokes human perceptual variability, providing a novel approach for aligning AI and human by harmonizing their perceptual variabilities. Furthermore, variMNIST can reveal individual differences among humans, enabling the generation of customized images that reflect these differences through ANNs aligned with individual participants.

Despite our progress in exploring human perceptual variability, several limitations remain. Our datasets, generated by sampling along ANN perceptual boundaries, cannot fully capture human variability, especially influences like culture, as some ANNs are trained on data from specific groups. To address this, we plan to include participants from diverse cultural backgrounds for a more comprehensive understanding. Furthermore, the dataset is focused on object recognition tasks, while effective for evoking perceptual variability, limits the exploration of broader visual phenomena. Expanding beyond object recognition to tasks like similarity judgments, emotion recognition, visual attention, and scene memory could offer deeper insights. However, exploring such complex tasks remains challenging given the limited number of trials available in individual behavioral experiments.

In terms of aligning AI with humans, although ANNs finetuned with individual behavioral data showed a notable improvement in predicting perceptual variability, there remains a significant gap when compared to their performance in standard classification tasks. This indicates that perceptual variability is a promising but underexplored method for AI-human alignment, with ample room for improvement. To address this, we propose incorporating optimal experimental design (Rainforth et al., 2024; Foster et al., 2019; 2021) into human experiments, using ANNs finetuned with individual behavioral data to generate customized images that maximize individual variability. These new behavioral data could then be fed back into the training of ANNs, dramatically improving AI-human alignment with fewer experimental trials. This approach would significantly increase the efficiency of human behavior data collection, reduce the cost of AI-human alignment, and accelerate the advancement of both cognitive science and artificial intelligence.

Code and Data Availability

Please see our website <https://eaterminator.github.io/BAM/> for a description of our work. Code and data for this paper are released at <https://github.com/ncclab-sustech/HumanPerceptualVariability>.

Ethics

In this study, All participants signed an informed consent form before the experiments. This study was approved by the local ethics committee of Southern University of Science and Technology (Resolution Number: 20240326).

Acknowledgements

This work was funded in part by the National Natural Science Foundation of China (62472206), Shenzhen Science and Technology Innovation Committee (2022410129, KJZD20230923115221044), Guangdong Provincial Key Laboratory of Advanced Biomaterials (2022B1212010003), and the Center for Computational Science and Engineering at Southern University of Science and Technology.

Impact Statement

This work advances the integration of machine learning methods within cognitive science, opening new avenues for modeling human perception. Our algorithms and datasets are intended solely for scholarly research in cognitive science. We emphasize that they should **not** be repurposed for deceptive or potential misleading manipulative applications under any circumstances.

References

- Atakan Bedel, H. and Çukur, T. Dreamr: Diffusion-driven counterfactual explanation for functional mri. *arXiv e-prints*, pp. arXiv-2307, 2023.
- Bansal, A., Chu, H.-M., Schwarzschild, A., Sengupta, S., Goldblum, M., Geiping, J., and Goldstein, T. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.
- Chen, X., Gao, X., Zhao, J., Ye, K., and Xu, C.-Z. Advdiffuser: Natural adversarial example synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4562–4572, 2023.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houslyby, N. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2021.
- Elsayed, G., Shankar, S., Cheung, B., Papernot, N., Kurakin, A., Goodfellow, I., and Sohl-Dickstein, J. Adversarial examples that fool both computer vision and time-limited humans. *Advances in neural information processing systems*, 31, 2018.
- Engstrom, L., Ilyas, A., S., H., S., S., and T., D. Robustness (python library). <https://github.com/MadryLab/robustness>, 2019. Accessed: 2025-01-25.
- Feather, J., Durango, A., Gonzalez, R., and McDermott, J. Metamers of neural networks reveal divergence from human perceptual systems. *Advances in Neural Information Processing Systems*, 32, 2019.
- Feather, J., Leclerc, G., Madry, A., and McDermott, J. H. Model metamers reveal divergent invariances between biological and artificial neural networks. *Nature Neuroscience*, 26(11):2017–2034, 2023.
- Foster, A., Jankowiak, M., Bingham, E., Horsfall, P., Teh, Y. W., Rainforth, T., and Goodman, N. Variational bayesian optimal experimental design. *Advances in Neural Information Processing Systems*, 32, 2019.
- Foster, A., Ivanova, D. R., Malik, I., and Rainforth, T. Deep adaptive design: Amortizing sequential bayesian experimental design. In *International conference on machine learning*, pp. 3384–3395. PMLR, 2021.
- Fu, S., Tamir, N., Sundaram, S., Chai, L., Zhang, R., Dekel, T., and Isola, P. Dreamsim: Learning new dimensions

- of human visual similarity using synthetic data. *arXiv preprint arXiv:2306.09344*, 2023.
- Gaziv, G., Lee, M., and DiCarlo, J. J. Strong and precise modulation of human percepts via robustified anns. *Advances in Neural Information Processing Systems*, 36, 2024.
- Golan, T., Raju, P. C., and Kriegeskorte, N. Controversial stimuli: Pitting neural networks against each other as models of human cognition. *Proceedings of the National Academy of Sciences*, 117(47):29330–29337, 2020.
- Golan, T., Siegelman, M., Kriegeskorte, N., and Baldassano, C. Testing the limits of natural language models for predicting human language judgements. *Nature Machine Intelligence*, 5(9):952–964, 2023.
- Hebart, M. N., Zheng, C. Y., Pereira, F., and Baker, C. I. Revealing the multidimensional mental representations of natural objects underlying human similarity judgements. *Nature human behaviour*, 4(11):1173–1185, 2020.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jeanneret, G., Simon, L., and Jurie, F. Diffusion models for counterfactual explanations. In *Proceedings of the Asian Conference on Computer Vision*, pp. 858–876, 2022.
- Jeanneret, G., Simon, L., and Jurie, F. Adversarial counterfactual visual explanations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16425–16435, 2023.
- Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35: 26565–26577, 2022.
- Kubilius, J., Schrimpf, M., Kar, K., Rajalingham, R., Hong, H., Majaj, N., Issa, E., Bashivan, P., Prescott-Roy, J., Schmidt, K., Nayebi, A., Bear, D., Yamins, D. L., and DiCarlo, J. J. Brain-like object recognition with high-performing shallow recurrent anns. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/7813d1590d28a7dd372ad54b5d29d033-Paper.pdf.
- Liang, Y., He, J., Li, G., Li, P., Klimovskiy, A., Carolan, N., Sun, J., Pont-Tuset, J., Young, S., Yang, F., et al. Rich human feedback for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19401–19411, 2024.
- Ma, J., Hu, T., Wang, W., and Sun, J. Elucidating the design space of classifier-guided diffusion generation. *arXiv preprint arXiv:2310.11311*, 2023.
- Mahner, F. P., Muttenthaler, L., Güçlü, U., and Hebart, M. N. Dimensions underlying the representational alignment of deep neural networks with humans. *arXiv preprint arXiv:2406.19087*, 2024.
- Muttenthaler, L., Dippel, J., Linhardt, L., Vandermeulen, R. A., and Kornblith, S. Human alignment of neural network representations. *arXiv preprint arXiv:2211.01201*, 2022a.
- Muttenthaler, L., Zheng, C. Y., McClure, P., Vandermeulen, R. A., Hebart, M. N., and Pereira, F. Vice: Variational interpretable concept embeddings. *Advances in Neural Information Processing Systems*, 35:33661–33675, 2022b.
- Muttenthaler, L., Greff, K., Born, F., Spitzer, B., Kornblith, S., Mozer, M. C., Müller, K.-R., Unterthiner, T., and Lampinen, A. K. Aligning machine and human visual representations across abstraction levels. *arXiv preprint arXiv:2409.06509*, 2024.
- Nanda, V., Speicher, T., Kolling, C., Dickerson, J. P., Gum-madi, K., and Weller, A. Measuring representational robustness of neural networks through shared invariances. In *International Conference on Machine Learning*, pp. 16368–16382. PMLR, 2022.
- Nanda, V., Majumdar, A., Kolling, C., Dickerson, J. P., Gum-madi, K. P., Love, B. C., and Weller, A. Do invariances in deep neural networks align with human perception? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9277–9285, 2023.
- Peebles, W. and Xie, S. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- Rainforth, T., Foster, A., Ivanova, D. R., and Bickford Smith, F. Modern bayesian experimental design. *Statistical Science*, 39(1):100–114, 2024.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- Seber, G. A. and Lee, A. J. Linear regression analysis. *Wiley-Interscience*, 2003.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2015.

Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.

Sundaram, S., Fu, S., Muttenthaler, L., Tamir, N. Y., Chai, L., Kornblith, S., Darrell, T., and Isola, P. When does perceptual alignment benefit vision representations? *arXiv preprint arXiv:2410.10817*, 2024.

Vaeth, P., Fruehwald, A. M., Paassen, B., and Gregorova, M. Diffusion-based visual counterfactual explanations—towards systematic quantitative evaluation. *arXiv preprint arXiv:2308.06100*, 2023.

Veerabadrán, V., Goldman, J., Shankar, S., Cheung, B., Papernot, N., Kurakin, A., Goodfellow, I., Shlens, J., Sohl-Dickstein, J., Mozer, M. C., et al. Subtle adversarial image manipulations influence both human and machine perception. *Nature Communications*, 14(1):4933, 2023.

Wei, C., Zou, J., Heinke, D., and Liu, Q. Cocog: Controllable visual stimuli generation based on human concept representations. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, 2024a.

Wei, C., Zou, J., Heinke, D., and Liu, Q. Cocog-2: Controllable generation of visual stimuli for understanding human concept representation. *arXiv preprint arXiv:2407.14949*, 2024b.

Yang, L., Ding, S., Cai, Y., Yu, J., Wang, J., and Shi, Y. Guidance with spherical gaussian constraint for conditional diffusion. *arXiv preprint arXiv:2402.03201*, 2024.

Yu, J., Wang, Y., Zhao, C., Ghanem, B., and Zhang, J. Freedom: Training-free energy-guided conditional diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23174–23184, 2023.

Zheng, C. Y., Pereira, F., Baker, C. I., and Hebart, M. N. Revealing interpretable object representations from human behavior. *arXiv preprint arXiv:1901.02915*, 2019.

Zhou, Z. and Firestone, C. Humans can decipher adversarial images. *Nature communications*, 10(1):1334, 2019.

A. Details of perceptual boundary sampling

A.1. Classifier guidance diffusion model

A.1.1. DIFFUSION MODELS.

Diffusion models (Song et al., 2020; Karras et al., 2022) consist of two main phases: forward and reverse. The forward phase transforms an image into Gaussian noise over

time $t \in [0, T]$, while the reverse phase reconstructs the image from noise by reversing this process. At any time t , the state x_t is defined as:

$$x_t = a_t x_0 + b_t \epsilon_t, \quad (1)$$

where $a_t = \sqrt{\alpha_t}$, $b_t = \sqrt{1 - \alpha_t}$, α_t increases with t , and $\epsilon_t \sim \mathcal{N}(0, I)$. A neural network is trained to predict the added noise:

$$\min_{\theta} \mathbb{E}_{x_t, \epsilon_t} [\|\epsilon_{\theta}(x_t, t) - \epsilon_t\|_2^2], \quad (2)$$

where the loss depends on the noise and the probability distribution $p_t(x_t)$. The reverse process follows an ordinary differential equation (ODE):

$$\frac{dx_t}{dt} = f(t)x_t - \frac{g^2(t)}{2} \nabla_x \log p_t(x_t), \quad (3)$$

with $f(t) = -\frac{d \log a_t}{dt}$ and $g^2(t) = \frac{db_t^2}{dt} - 2 \frac{d \log \sqrt{\alpha_t}}{dt} b_t^2$. This ODE enables the reconstruction of the image by reversing the noise-adding process.

The specific steps for both phases are determined by the sampling algorithm. We use the DDPM algorithm (Ho et al., 2020), where the forward and reverse steps are represented as:

$$x_t = DDPM^+(x_{t-1}) \quad \text{and} \quad x_{t-1} = DDPM^-(x_t).$$

A.1.2. CLASSIFIER GUIDANCE

Classifier guidance is also known as Training-free guidance. Using a diffusion model and the conditional information y , we define the conditional probability of the generative process as:

$$p(x_t|y) = \frac{p(y|x_t)p(x_t)}{p(y)}$$

where x_t is the generated stimuli at time step t .

The gradient of this probability is calculated as follows:

$$\nabla_{x_t} \log p_t(x_t|y) = \nabla_{x_t} \log p_t(x_t) + \nabla_{x_t} \log p_t(y|x_t)$$

In the training-free approach, we utilize a network f_{ϕ} and define a loss function $\ell(f_{\phi}(x_t), y)$ for conditional generation. Thus, we obtain:

$$\nabla_{x_t} \log p_t(y|x_t) = \nabla_{x_t} \ell(f_{\phi}(x_t), y)$$

In the reverse sampling process, we introduce a correction step:

$$x_{t-1} = DDPM^-(x_t) - \gamma \nabla_{x_t} \ell(f_{\phi}(x_t), y)$$

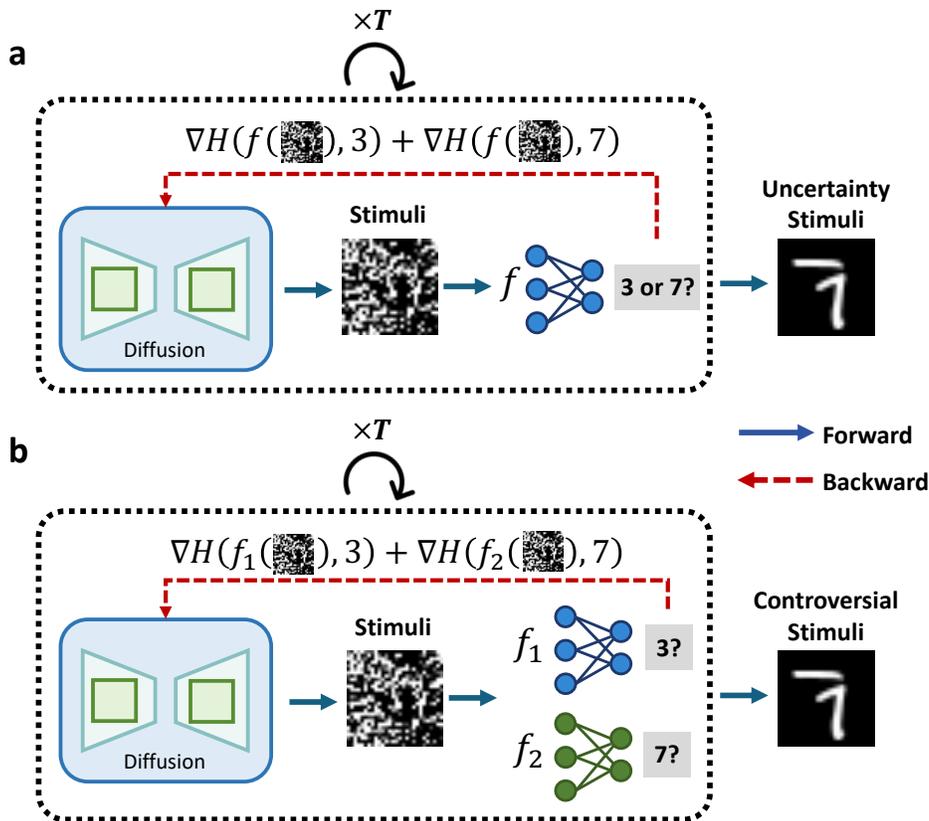


Figure A.1. **Guidance methods.** (a) The uncertainty guidance. It utilizes the classification uncertainty of the generated images from model f to guide the diffusion model in generating stimuli toward specific directions. Model f ensures the image is a digit. (b) The controversial guidance. It employs the classification differences between the generated images from model f_1 and model f_2 to guide the diffusion model in generating stimuli toward specific directions.

Therefore we can generate certain stimuli by designing the loss function ℓ . To obtain stimuli that can disrupt human perception, we explored two potentially suitable approaches: uncertainty sampling and controversial sampling. High uncertainty sampling aims to generate stimuli that challenge the model’s judgment, while controversial sampling seeks to produce stimuli that maximize the difference in probability distributions between two models.

A.2. Details of guidance algorithms

A.2.1. DETAILS OF TARGETED GUIDANCE

In targeted guidance, we specify directions for the guidance. For instance, at position (3, 5), uncertainty guidance directs towards both categories 3 and 5. For controversial guidance, classifier 1 is directed towards category 3, while classifier 2 is directed towards category 5. Similarly, at position (5, 3), uncertainty guidance directs towards both categories 5 and 3. For controversial guidance directs classifier 1 towards category 5 and classifier 2 towards category 3. To ensure balanced targeted guidance, we generate samples in

multiples of 100 for each targeted guidance. This ensures that all stimuli corresponding to positions from 0 to 9 × 0 to 9 (i.e., covering all guidance directions) are included, thereby maximizing sampling uniformity. In the generation the guidance scale is set to 0.1, resampling steps is set to 5, and the inference steps is set to 50.

When generating stimuli using this guidance strategy, we ensure that each term in the loss function is effectively utilized. While this approach guarantees category-balanced sampling during the generation process, the final retained stimuli may not necessarily exhibit category balance. For the stimuli intended for human digit recognition experiments, we apply additional filtering to the generated images. Specifically, for uncertainty sampling, we require that the top two p-values exceed 0.4 and the digit surrogate score is above 0.5. For controversial sampling, we ensure that the classification outputs of both classifiers correspond to the intended guidance direction, with the highest p-value exceeding 0.9 and the digit surrogate score above 0.5. A detailed analysis of the filtered dataset derived from uniform sampling was performed, and the distribution of category counts is presented in Figure

A.13. By comparing this distribution with the cognitive data shown in Figure A.14, a correlation can be observed.

A.2.2. THE ROLE OF DIFFUSION PRIOR

In previous studies that employed generated images to investigate model and human perception ((Golan et al., 2020; Gaziv et al., 2024; Veerabadran et al., 2023; Feather et al., 2023)), a common issue was that the generated images lacked sufficient naturalness and failed to significantly influence human perception. This issue is particularly crucial within the context of our research objectives. Using previous methods often resulted in images that were unrecognizable to human participants, leading to nearly random classification results (see Figure. A.2). Recent advances in adversarial examples and counterfactual explanations in machine learning ((Jeanneret et al., 2023; Wei et al., 2024b; Chen et al., 2023; Jeanneret et al., 2022; Vaeth et al., 2023; Atakan Bedel & Çukur, 2023)) have addressed this issue by employing diffusion models as regularizers to introduce prior information. This technique allows for the generation of natural images capable of influencing human perception.

Inspired by these advances, we utilize a classifier-free diffusion model as the core of image generation process. By sampling noise from the target dataset (MNIST) distribution and feeding it into the diffusion model for denoising, we effectively incorporate prior information. This approach enhances the naturalness of the generated images, making them more reflective of the real distribution of handwritten digits and thereby increasing their impact on human perception, as shown in Figure. A.2.

A.2.3. EDITING EXISTING DATASETS VS. GENERATING DATA FROM SCRATCH

The process of adding MSE loss to the loss function can be seen as editing existing datasets. The MSE loss is used to constrain the pixel space of the stimuli. Without the MSE loss, the model is more likely to sample from distributions of stimuli that are very similar in pixel space within a certain class. We aim to enforce a constraint in the pixel space that encourages the stimuli to be closer to the original distribution of randomly sampled samples from the MNIST dataset. This approach is intended to enhance the diversity of the generated stimuli. We define α as the pixel-level restraint scale.

For uncertainty guidance with MSE constraint, we have:

$$\ell(f(x_t), y) = H(y|x_t) + \alpha \|x_t - x_{ref}\|^2,$$

where H represents the entropy, α represents the strength of the MSE loss. For controversial guidance with MSE constraint, we have:

$$\ell(f(x_t), y) = D_{KL}(p(y_1|x_t), p(y_2|x_t)) + \alpha \|x_t - x_{ref}\|^2,$$

where D_{KL} represents the KL divergence between two distributions. We conducted experiments on five models. For uncertainty sampling, we generated stimuli for each model, resulting in five groups of stimuli. In the controversial sampling experiments, we pitted the models against each other in pairs, creating ten groups of stimuli. However, this approach can be perceived as manipulating one class into another, which is similar to our goal of sampling along the perceptual boundaries of ANNs, but not exactly the same. In generation, the guidance scale is set to 0.1, resampling steps is set to 5, the inference steps is set to 50, and α is set to 50.

A.2.4. TARGETED GUIDANCE VS. UNTARGETED GUIDANCE

Untargeted guidance focuses solely on increasing the variability of the generated images, without considering the overall distribution of the images. We adopted an untargeted guidance method to generate stimuli for the digit recognition experiment. To sample at the perceptual boundary of the model, we drew on previous research and proposed two guidance methods: uncertainty guidance and controversial guidance. Uncertainty guidance ensures that the generated images are as close as possible to the model’s perceptual boundary by maximizing the entropy of the classification probability distribution of a single ANN model for the generated images, thereby obtaining images with high perceptual variability for the model. For uncertainty guidance, this can be represented as:

$$\ell(f(x_t), y) = H(y|x_t),$$

Where H is the entropy, y is the output probability of the neural network. Controversial guidance, on the other hand, utilizes two different ANN models and generates images that maximize the KL divergence between their classification probability distributions, thereby maximizing perceptual differences between the models. For controversial guidance, this can be represented as:

$$\ell(f(x_t), y) = D_{KL}(p(y_1|x_t), p(y_2|x_t)),$$

Where D_{KL} is the KL divergence, $p(y_1|x_t)$ is the output probability of the first neural network, $p(y_2|x_t)$ is the output probability of the second neural network. In generation the guidance scale is set to 0.1, resampling steps is set to 5, and the inference steps is set to 50. Targeted and untargeted guidances are compared in Figure.A.11. Losses with and without MSE are also compared in Figure. A.10.

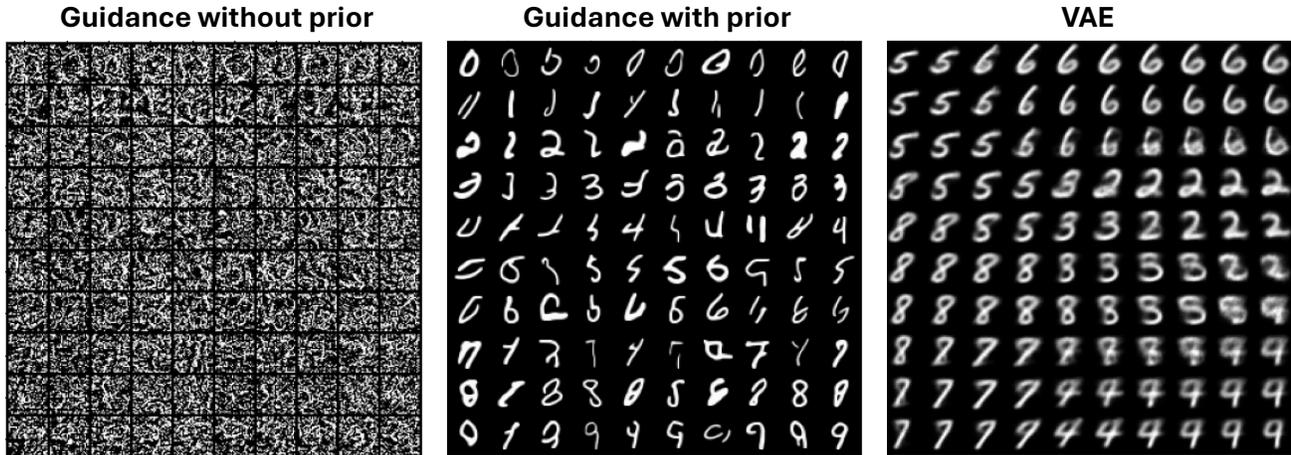


Figure A.2. **Comparison of our method with other approaches.** Images generated by a diffusion model without prior distribution exhibit severe noise. Images produced by a Variational Autoencoder (VAE) show minimal differences and are generally blurry. Our method (with prior), however, yields images that are not only clear and noise-free but also exhibit substantial diversity.

B. Details of collecting human perceptual variability

B.1. Model configuration and training

B.1.1. DIFFUSION MODEL

Configuration of DiT. The Diffusion Transformer (DiT) ((Peebles & Xie, 2023)) is a Transformer-based diffusion model tailored for generative tasks. In our configuration, the model processes 28×28 grayscale images using a patch size of 2×2 , resulting in patch embeddings transformed into sequences of hidden size 128, with 1 input channel and 10 output classes. The architecture includes 4 Transformer layers with 8 attention heads per layer and an MLP ratio of 4.0.

DiT incorporates Patch Embedding, Timestep Embedding, and Label Embedding modules. These embeddings are combined with fixed sinusoidal positional encodings to provide spatial and temporal context. AdaLN (Adaptive Layer Normalization) layers condition the model on timestep and label embeddings, with zero-initialized manipulation for training stability.

The model outputs spatial predictions through a final linear layer followed by an unpatching operation, restoring the input image dimensions. Classifier-free guidance is supported by computing conditional and unconditional outputs, enabling control over generated samples.

Training of Diffusion Model. For prior diffusion model, we use the MNIST dataset as the training dataset. The dataset consists of grayscale images of size 28×28 , which are directly used without further resizing. The training

process is conducted using a single GPU (NVIDIA GeForce RTX 4090) with the Adam optimizer.

The data is loaded into the training pipeline using a PyTorch DataLoader with a batch size of 128, and the number of worker threads for data loading is set to 128. The model is trained for 150 epochs, with a learning rate of $1e - 4$ and an unconditional training rate of 0.1, and the weight decay is not applied. Dropout is applied to the class embedding with a probability of 0.1, while the model does not learn the variance (sigma).

B.1.2. CLASSIFIERS

The classifier models were trained on the MNIST dataset using 28×28 grayscale images, normalized with the ‘ToTensor’ transformation. Training and testing sets were loaded with a batch size of 100, and the models were implemented with 5 different configurations (ViT, VGG, MLP, CORNET, LRM, see Table. A.1) to map input images to 10 output classes. Training was performed on an NVIDIA GPU using the AdamW optimizer ($lr = 1 \times 10^{-3}$) for 16 epochs, and CrossEntropyLoss function was used to compute the classification loss.

B.1.3. DIGIT JUDGMENT SURROGATE

Digit judgment experiment. We used the synthetic images as experimental stimuli to measure human behavior in a digit judgment task. The purpose of this experiment was to collect human judgments on whether a given image qualifies as a digit, thereby establishing a human criterion for handwritten digit. For each image, participants were asked the question, “Is this image a digit?” with responses limited to “True” or “False.” More experimental details can

Table A.1. Configurations and MNIST Accuracy of Classifiers

Model Name	Model Type	MNIST Accuracy (%)
ViT	Vision Transformer	97.2
VGG	Small VGG	98.2
CORNet	CORnet-Z	98.9
MLP	Multi-Layer Perceptron	98.3
LRM	Logistic Regression Model	92.7

be found in Appendix B.2.1.

The experiment collected behavioral data from 400 participants, comprising 200,000 trials and 20,000 stimuli. During the data cleaning process, 124 participants were excluded based on Sentinel trials, leaving data from 276 participants (138,000 trials and 19,878 valid stimuli). This dataset provides a robust foundation for analyzing perceptual standards for handwritten digits.

Training of digit judgment surrogate. For the training of the digit judgment surrogate model, we constructed a dataset based on the results of the human digit judgment experiment. Specifically, for any given image, the frequency of participants responding "True" was taken as the probability of the image being judged as a digit. These images and their corresponding probabilities were then used to train the digit judgment surrogate. The dataset was split into a training set and a test set in a ratio of 8:2.

The surrogate model is based on the SmallVGG architecture, with a final output layer designed for regression tasks. The model was trained using the AdamW optimizer with a learning rate of 0.001, and the mean squared error (MSE) was used as the loss function. The training process lasted for 8 epochs with a batch size of 128. After each epoch, the validation loss was monitored, and the model with the best validation performance was saved for further evaluation.

Performance of digit judgment surrogate. To ensure the validity of the digit judgment surrogate’s predictions, we computed the correlation between the predicted scores and human scores. For any given image, the human score was defined as the frequency of participants responding "True," indicating the image is a digit, while the predicted score was the probability assigned by the model classifying the image as a digit.

As shown in Figure. A.3a, the Spearman rank correlation coefficient between the predicted scores and human scores is 0.8035. This indicates that the model’s digit judgment is highly consistent with human. Additionally, Figure. A.3b presents image examples corresponding to different predicted scores. For scores of 0.10, 0.25, 0.50, 0.75, and 0.90, eight samples were randomly selected for each score. The

examples reveal that as the predicted score increases, the images progressively resemble digits more closely. These results demonstrate that the digit judgment surrogate effectively simulates human digit judgment behavior.

Guiding generative process by digit judgment surrogate.

For any given image, we use the frequency of participants responding "True" as the probability of the image being a digit. The initial image dataset, along with the corresponding probabilities, was used to train a digit judgment surrogate. As the previous works of image generation by human preferences ((Liang et al., 2024; Bansal et al., 2023)), this surrogate, functioning as a image quality predictor, was then employed to guide the image generation process (see Appendix B.1.3). The guidance formula can be expressed as:

$$\mathcal{L}_{\text{surr}} = \mathcal{L} + \max((1 - f_{\text{surr}}(x))^2, 0.5)$$

In this formula, $\mathcal{L}_{\text{surr}}$ represents the total loss. $f_{\text{surr}}(x)$ represents the probability give by the digit judge model. The probability of the digit judge is combined to the formula to ensure the generated image is considered as a digit by humans. The max function is used so that when the score is above a certain threshold, the gradient of the digit judge will not effect generation.

B.2. Online human behaviroal measurement

B.2.1. DIGIT JUDGMENT

We use the initial synthetic dataset as experimental stimuli to measure human behavior in the judgment task. The purpose of the experiment is to collect human judgments on whether any given image test is a digit, in order to filter out images that do not meet the standards of handwritten digits.

Task paradigm Before the formal experiment, participants will first complete a pre-experiment. Each round of the pre-experiment consists of two stages. (1) Selection Stage: A test image appears at the center of the screen, with two buttons labeled "True" and "False" displayed below it. Participants are required to judge whether the image represents a number. (2) Feedback Stage: After making their choice, participants will receive feedback below the image

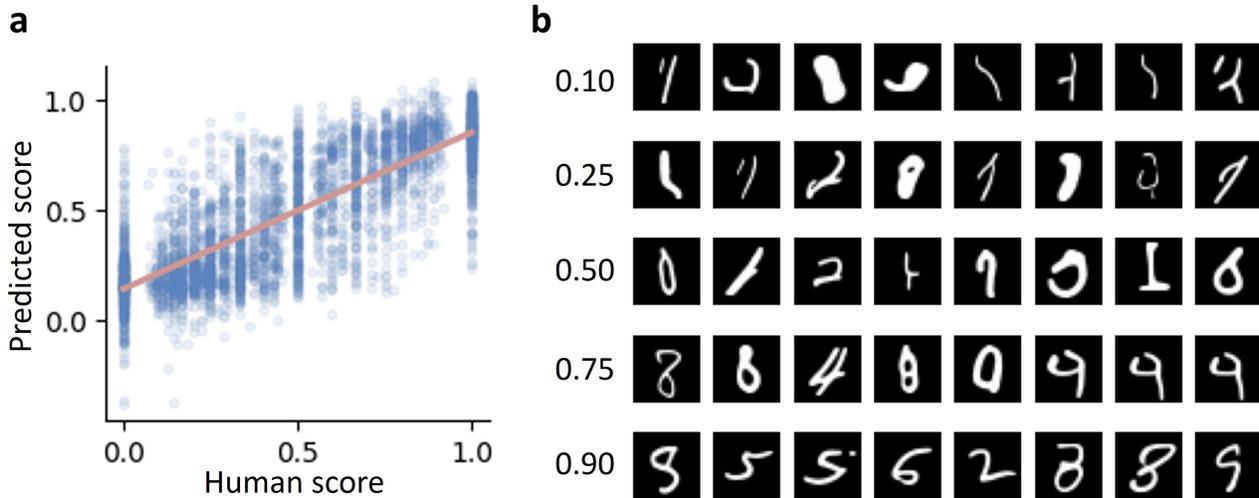


Figure A.3. Performance of digit judgment surrogate. (a) The predicted scores and human scores show a strong correlation. For any given image, the human score is defined as the frequency of participants answering "True" for the image being a digit, while the predicted score is the model's probability of classifying the image as a digit. The Spearman rank correlation coefficient between the two scores is 0.8035. (b) Examples of images with different scores. For predicted scores of 0.10, 0.25, 0.50, 0.75, and 0.90, 8 samples are randomly displayed for each score. As the score increases, the images increasingly resemble digits.

indicating whether it is a number. The pre-experiment includes a total of 10 rounds, after which participants will proceed to the formal experiment. In formal experiment, participants performed multiple rounds of a choice task (see Figure. A.4). Each trial consisted of two phases: (1) Fixation Phase: A black cross was displayed at the center of the screen for 300 ms to direct participants' attention to the center. (2) Selection Phase: A test image appears at the center of the screen, with two buttons labeled "True" and "False" displayed below it. The positions of the buttons were fixed and remained unchanged throughout the trials. Participants were asked to judge whether the image represents a figure by selecting the corresponding button with the mouse or pressing the key on the keyboard (A represents True and D represents False). There was no time limit for responding. Each session of formal experiment comprised 500 trials, divided into two types: (1) Sentinel trials ($n = 10$), in which participants are shown a set of 10 pre-selected MNIST images, i.e., the correct response should be True. We screened participants based on their accuracy in the sentinel trials to ensure high-quality responses. (2) Random Trials ($n=490$), where images were randomly selected from the dataset, excluding the fixed images. The two trial types were presented in a random alternating order. No feedback was provided after participants made their selection. The experiment was programmed using JSPsych, with stimuli presented via the JSPsych-Psychophysics component.

Human data collection The experiment got ethics approval from the local University. We recruited participants

($N=400$) and collected data through the NAODAO platform. Prior to the experiment, participants read an informed consent form detailing any potential risks associated with participation. Participants were allowed to withdraw from the experiment at any time. No personal identification information was collected. We only included data from participants with sentinel trial accuracy greater than 70%, resulting in data from 276 participants and 135240 trials involved in the following analyses.

B.2.2. DIGIT RECOGNITION EXPERIMENT

We used the filtered synthetic dataset as experimental stimuli to measure human behavior in a digit recognition task. The goal of the experiment was to collect the probability distribution of human choices for any given test image. In this task, participants were presented with ten possible choices, represented by the digits 0 to 9.

Task paradigm Participants performed multiple rounds of a category comparison task. Each trial consisted of two phases (see Figure. A.5): (1) Fixation Phase: A black cross was displayed at the center of the screen for 300 ms to direct participants' attention to the center. (2) Selection Phase: A test image appeared at the center of the screen, accompanied by ten labeled buttons below it, with labels ranging from 0 to 9. The positions of the buttons were fixed and remained unchanged throughout the trials. Participants were asked to identify the digit in the image by selecting the corresponding button with the mouse or pressing the number key on the

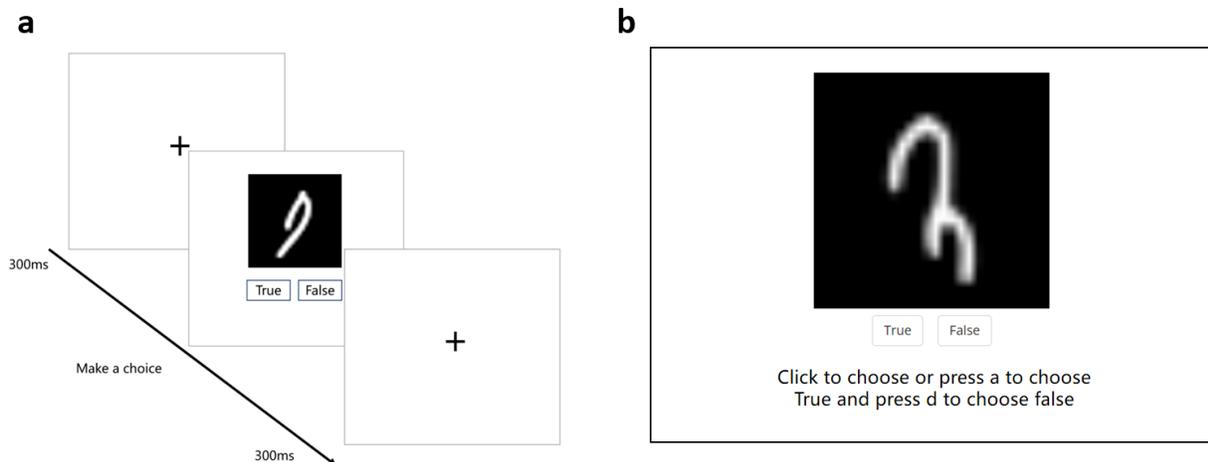


Figure A.4. Human digit judgment experiment procedure. In each trial, participants first observe a fixation cross (“+”) for 300 milliseconds. Following the fixation, a stimulus image is presented along with 2 clickable buttons labeled “True” and “False”. Participants are instructed to judge whether the image represents a digit and either click the corresponding button or pressing the key on the keyboard (A represents True and D represents False). The images shown to participants are generated by our model. After each selection, no feedback is provided, and the next trial begins immediately. Each participant first performed 10 rounds of pre-experiments with feedback, followed by 500 formal trials without feedback, including 10 sentinel trials and 490 random trials.

keyboard. There was no time limit for responding.

Each session comprised 500 trials, where images were randomly selected from the dataset. No feedback was provided after participants made their selection. The experiment was programmed using JSPsych, with stimuli presented via the JSPsych-Psychophysics component.

Human data collection The experiment got ethics approval from the local University. The experiment collected behavioral data from 400 participants through the NAO-DAO platform, comprising 200,000 trials and 20,000 stimuli. Prior to the experiment, participants read an informed consent form detailing any potential risks associated with participation. Participants were allowed to withdraw from the experiment at any time. No personal identification information was collected. During data preprocessing, 154 participants were excluded based on Sentinel trials (accuracy < 0.7), leaving data from 246 participants (123,000 trials and 19,952 valid stimuli). Table A.2 and Table A.3 shows the stimuli distribution across guidance strategies and classifier. Using this cleaned dataset, we constructed a high perceptual variability dataset, variMNIST, which serves as a foundation for subsequent analysis and modeling.

B.3. additional dataset details

B.3.1. EVALUATION METRICS

Judgment distribution. As shown in Figure. A.6 (top left), we evaluated the distribution of human judgments across the ten digit classes (0–9). The results indicate that the probabilities are relatively uniform, with all categories exhibiting values close to 0.1. Notably, digits 0, 6, and 9 were judged with slightly higher probabilities (around 0.15) compared to other digits, while digits 1 through 5 demonstrated lower probabilities (around 0.06).

RT and entropy. We further examined the relationship between response time (RT) and entropy to gain insights into the cognitive process underlying human judgments. RTs were predominantly distributed between 500 and 1500 ms, following a long-tail distribution, indicating that most decisions were made quickly, with a few requiring significantly more time (Figure. A.6 (top right)). The entropy of human judgments primarily concentrated near 0, reflecting high confidence in about half of the trials. Values between 0.5 and 2 also appeared, indicating uncertainty or ambiguity (Figure. A.6 (bottom left)). A positive correlation (Spearman rank correlation coefficient = 0.55) was observed between entropy and RT, suggesting that higher uncertainty in judgment often corresponds to longer decision times (Figure. A.6 (bottom right)).

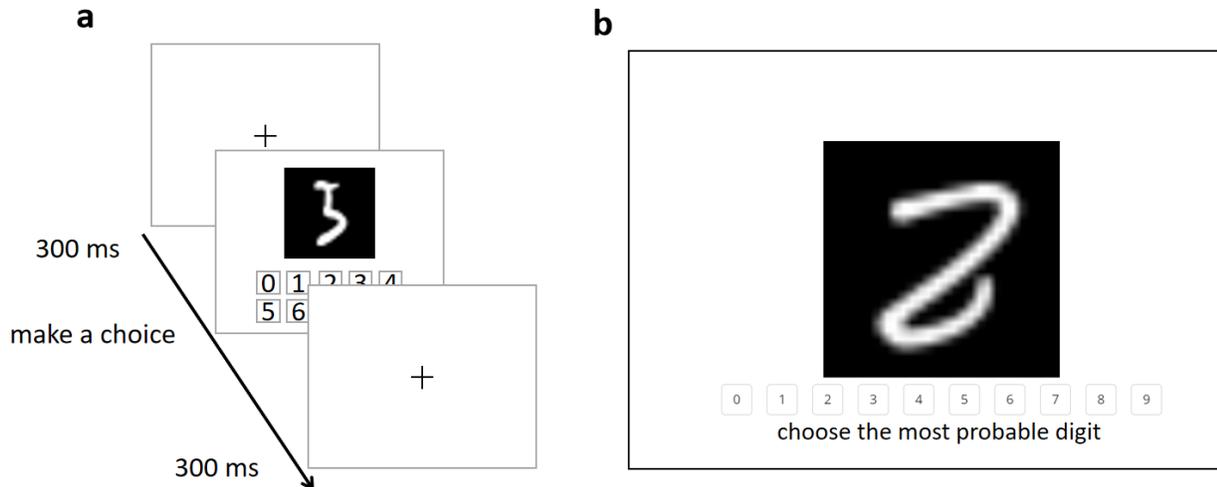


Figure A.5. Human digit recognition experiment procedure. In each trial, participants first observe a fixation cross (“+”) for 300 milliseconds. Following the fixation, a stimulus image is presented along with 10 clickable buttons representing the digits 0 to 9. Participants are instructed to identify the most likely digit represented by the image and either click the corresponding button or press the corresponding number on the keyboard. The images shown to participants are generated by our model. After each selection, no feedback is provided, and the next trial begins immediately. Each participant completes a total of 500 trials, consisting of 10 sentinel trials and 490 random trials.

Classifier configurations influence the guidance outcome.

We evaluated how different classifier configurations affected the guidance outcome under controversial guidance conditions. The overall guidance success was determined by measuring the probability that participants selected digit x when the model guided the judgment toward x . As shown in Figure A.7 (left), the results show that CORNet and VGG achieved the highest success rates, both nearing 0.6, indicating their strong ability to influence human judgments. VIT and MLP followed with moderate success rates of approximately 0.3, while LRM had the lowest success rate at around 0.2, reflecting its weaker guidance capability.

Further analysis compared the guidance outcome differences between classifiers when used as adversarial pairs in controversial guidance (see Figure A.7 (right)). CORNet and VGG consistently outperformed other classifiers, showing significantly higher success rates. In contrast, LRM exhibited the lowest success rates compared to other classifiers. These findings suggest that the choice of classifiers significantly impacts the effectiveness of controversial guidance, with certain architectures like CORNet and VGG being more effective at aligning human responses with their intended guidance.

Guidance targets influence the guidance outcome. We analyzed how different guidance targets influenced the guidance outcome, defined as the proportion of successful stimuli generated for each target pair. As shown in Figure A.14, the results revealed significant variability across guidance

targets. Target pairs such as (1, 7), (1, 2), and (4, 9) demonstrated the highest success rates, each exceeding 0.35. This suggests that these pairs may align better with human perceptual biases or model representations, leading to more effective guidance. Conversely, pairs such as (1, 8), (2, 9), and (7, 8) exhibited the lowest success rates, with values below 0.03, indicating greater difficulty in guiding these pairs. These findings highlight the importance of selecting appropriate guidance targets to maximize the effectiveness of the generated stimuli.

B.3.2. ADDITIONAL VISUALIZATION RESULTS

We generated 900 images using both targeted and untargeted approaches under the guidance of uncertainty and controversial methods. A t-SNE analysis was conducted on the targeted and untargeted methods for both the controversial and uncertainty approaches. To ensure fairness, the t-SNE analysis was performed directly on the raw pixel space for dimensionality reduction. The results are shown in Figure A.12. It can be observed that the distribution is more uniform when targeted guidances are adopted.

B.4. Details of validation on ImageNet

B.4.1. METHOD

Similar to handwritten digits, we constructed a high perceptual variability dataset consisting of nine classes based on the ImageNet dataset. Consistent with (Gaziv et al., 2024),

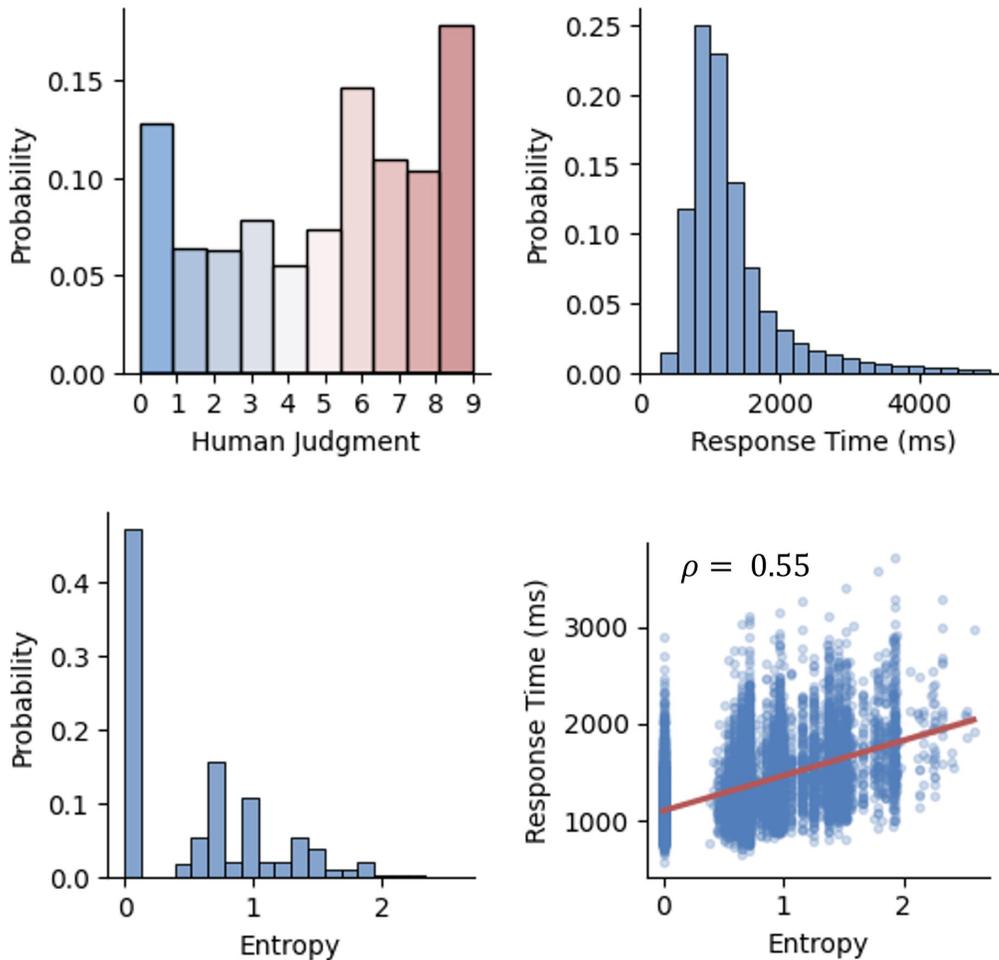


Figure A.6. **Behavioral results of the Digit recognition task.** Top left: In the digit recognition task on the variMNIST dataset, human judgment probabilities are relatively uniform, with values close to 0.1 for each category. Among these, digits 0, 6, and 9 have relatively higher probabilities, while digits 1 to 5 have lower probabilities. Top right: Human response times for the digit recognition task are concentrated between 500 and 1500 ms, showing a long-tail distribution. Bottom left: The entropy of human judgment results is primarily distributed around 0, with additional values observed between 0.5 and 2. Bottom right: Entropy and response time exhibit a positive correlation, with a Spearman rank correlation coefficient of 0.55.

we utilize the restricted ImageNet dataset that contains nine classes. Meanwhile, we refer to the generation schemes presented in (Wei et al., 2024a) and (Wei et al., 2024b). The CLIP latent was first generated with controversial guidance as before, and then fed into the second stage diffusion model to generate images. The classifier models were constructed by adding an additional linear regression layer to the original CLIP model’s image encoder. The finetune process was only conducted on the additional layer, with the original parameters of CLIP unchanged. Following the procedure of collecting human perceptual variability based on the MNIST dataset, the models were first finetuned on the group level, then finetuned on the individual level based on the group level model.

B.4.2. RESULTS

We show some example stimuli in Figure A.20. After generating the stimuli using controversial guidance, an online experiment was conducted and we show the result of our analysis in Figure A.22. Our guidance successfully achieved the goal of sampling on the perceptual boundary. Then group-level finetuning and individual level finetuning are conducted, and the results are shown in Figure A.24. An increase between the group-level finetuned model and the individual-level finetuned model is observed, showing that individual differences also exist in the task of classifying natural images. After individual finetuning, we also tested controversial generation on the individual-finetuned models, with the results shown in Figure A.25. It can be observed

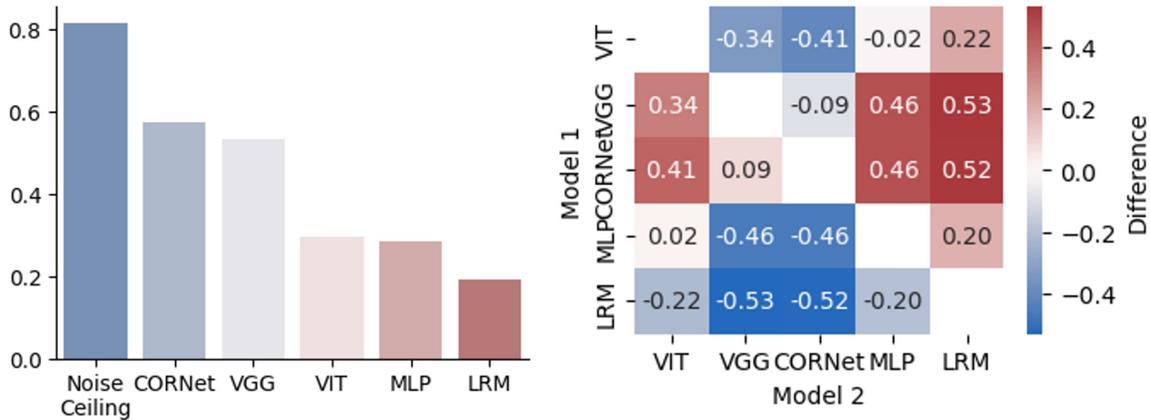


Figure A.7. **Guidance outcome under different classifiers configurations.** Left: Overall controversial guidance outcomes for different classifiers. Under controversial guidance conditions, the success rate is measured as the probability that participants chose digit x when the model guided the judgment to x . CORNet and VGG achieved the highest success rates, nearing 0.6, followed by ViT and MLP with success rates of approximately 0.3. LRM had the lowest success rate at around 0.2. Right: Differences in guidance outcomes among classifiers during controversial guidance (when using two classifiers as adversarial classifiers, the difference in the guidance outcome of one classifier and the other). CORNet and VGG exhibited significantly higher success rates compared to other classifiers, while LRM showed notably lower success rates than the rest.

that images at some positions can cause variability and have the potential to manipulate human perceptual variability.

C. Additional Results of predicting human perceptual variability

C.1. Effects of Fine-Tuning Across Classifiers

Prediction accuracy. As shown in Figure. A.8, on MNIST, group/individual fine-tuning resulted in slight accuracy improvements for ViT and VGG, while CORNet and MLP showed no significant changes. LRM’s accuracy decreased after fine-tuning, indicating limited generalization. On variMNIST, all classifiers exhibited significant accuracy gains after fine-tuning, highlighting the benefits of group and individual fine-tuning for datasets with high perceptual variability.

Model and human entropy. Figure. A.16 highlights the changes in correlation between model-predicted entropy and human behavioral entropy before and after fine-tuning. A positive correlation was observed across all baseline classifiers, indicating that even in the baseline condition, models capture human perceptual variability. Fine-tuning on variMNIST significantly enhanced this correlation, demonstrating improved alignment with human perceptual variability.

Impact of image difficulty. As shown in Figure. A.17, fine-tuned models outperformed baseline models across all entropy levels, confirming the general effectiveness of fine-

tuning. For classifiers other than LRM, individual fine-tuned models achieved greater accuracy improvements on high-entropy images compared to group-tuned models, indicating that individual fine-tuning is particularly effective for challenging stimuli.

Estimation of the effects of finetuning on different datasets. As illustrated in Figure A.26, increasing the volume of training data yields diminishing returns for GroupNet accuracy, whereas IndivNet accuracy continues to improve more substantially. This suggests that the available data are largely sufficient for GroupNet, but remain insufficient for IndivNet. Our dataset was deliberately designed to capture perceptual variability, resulting in many samples that GroupNet struggles to predict reliably. To better understand these effects, accuracy can be decomposed into three sources of uncertainty: group epistemic uncertainty, individual epistemic uncertainty, and aleatoric uncertainty. Here, individual epistemic uncertainty reflects variability arising from individual differences, while aleatoric uncertainty corresponds to intra-subject variability. However, because repeated presentations of the same stimulus introduce trial-to-trial interference, quantifying the precise impact of aleatoric uncertainty is challenging. Therefore, we estimate that individual epistemic uncertainty accounts for at least a

C.2. Model Fine-Tuning

For group-level fine-tuning, The original classifier models were trained on the mixed (ratio = 1:1) MNIST, variMNIST datasets using 28×28 grayscale images, normalized with

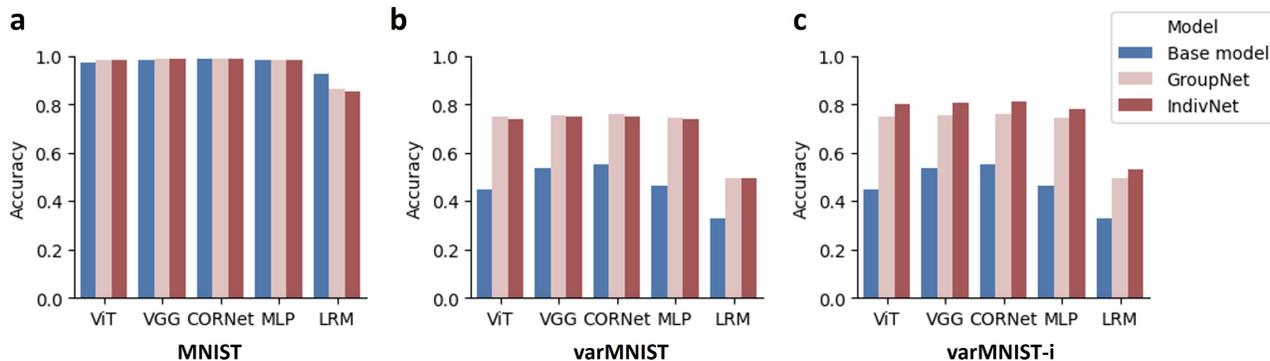


Figure A.8. **Correlation between model entropy and human behavior.** (a) Prediction accuracy on MNIST before and after fine-tuning for different classifiers. ViT and VGG show slight improvements in accuracy after group/individual fine-tuning, CORNet and MLP exhibit no significant changes, while LRM experiences a decrease in accuracy post-fine-tuning. (b) Prediction accuracy on variMNIST before and after fine-tuning. All five classifiers demonstrate substantial improvements in accuracy after group/individual fine-tuning. (c) Prediction accuracy on variMNIST-i before and after fine-tuning. All five classifiers show moderate improvements in accuracy after group/individual fine-tuning.

the ‘ToTensor’ transformation. For individual-level fine-tuning, the dataset is a mixture of variMNIST-i, variMNIST and MNIST at a ratio of 2:1:1 and the initial model is the group model. Training and testing sets were loaded with a batch size of 128, and the models were implemented with 5 different configurations (see Table. A.1) to map input images to 10 output classes. Training was performed on an NVIDIA GPU using the AdamW optimizer ($lr = 1 \times 10^{-3}$) for 16 epochs, and CrossEntropyLoss function was used to compute the classification loss.

C.3. Clustering analysis

There is a large variability in the subject’s digit recognition behaviors, since participants differ in high-level factors such as culture, ethnicity, educational background, regional customs, and psychological states. We hypothesize that participants could be grouped into several clusters, with participants within the same cluster likely to exhibit similar perceptual variability. To test this hypothesis, we used each participant’s subject-finetuned model to predict the behavior of all participants, and we calculated *inter-subject similarity matrix* based on the prediction results. The better the prediction performance, the higher the inter-subject similarity. As shown in Figure. A.19a, the similarity matrix between participants revealed the existence of eight distinct clusters. Furthermore, we observed that the subject-finetuned models performed better in predicting the behavior of participants within the same cluster (in-cluster) compared to those outside the cluster (out-cluster), as shown in Figure. A.19b. Our results indicate that the clustering is valid and that there are indeed high-level percept differences between participants.

D. Details of manipulating human perceptual variability

We present the result of further analysis (Prediction accuracy, success rate and targeted ratio) on each subject pair in Figure A.18. To illustrate the examples on different positions in the perceptual space and demonstrate the manipulation results, we present some examples of the manipulation stimuli in Figure A.27, A.28 and A.29.

Table A.2. Stimuli Counts before Experiment

Guidance Strategy	Classifier	Stimuli Count
Controversial	CORNet_LRM	1000
	CORNet_MLP	1000
	MLP_LRM	1000
	VGG_CORNet	1000
	VGG_LRM	1000
	VGG_MLP	1000
	ViT_CORNet	1000
	ViT_LRM	1000
	ViT_MLP	1000
	ViT_VGG	1000
Uncertainty	CORNet	2000
	LRM	2000
	MLP	2000
	VGG	2000
	ViT	2000
Sum		20000

Table A.3. Stimuli and Trial Counts after Experiment

Guidance Strategy	Classifier	Stimuli Count	Trial Count
Controversial	CORNet_LRM	997	5766
	CORNet_MLP	996	5688
	MLP_LRM	997	5684
	VGG_CORNet	995	5806
	VGG_LRM	994	5767
	VGG_MLP	999	5823
	ViT_CORNet	997	5865
	ViT_LRM	995	5949
	ViT_MLP	999	5811
	ViT_VGG	999	5881
Uncertainty	CORNet	1994	11631
	LRM	1992	11668
	MLP	1997	11849
	VGG	1996	11710
	ViT	1996	11817
Sum		19943	116715

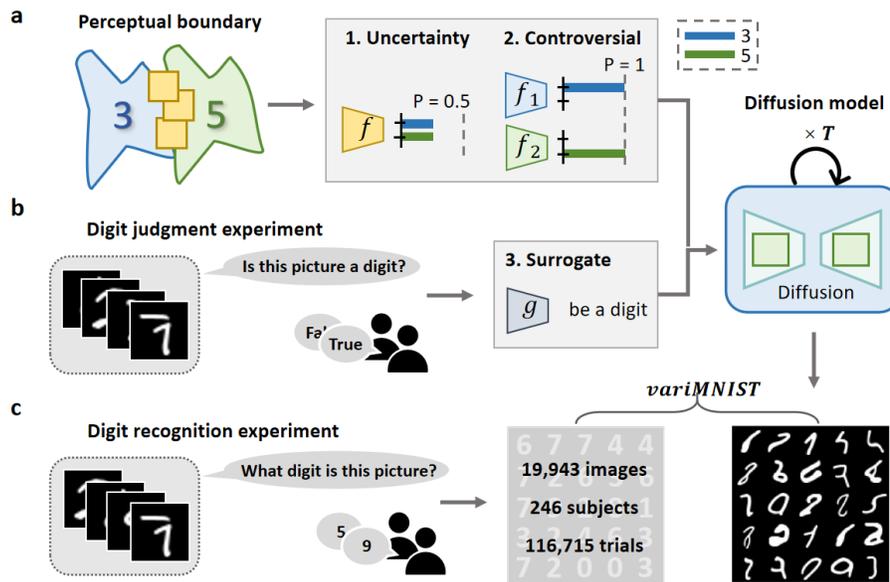


Figure A.9. **Generating images to elicit human perceptual variability.** (a) The example illustrates two guidance methods for sampling from the perceptual boundary between “3” and “5” in ANN: *uncertainty guidance* and *controversial guidance*. Specifically, *Uncertainty guidance* aims to make the ANN model f assign equal probabilities to “3” and “5,” while *controversial guidance* generates images classified as “3” by f_1 but as “5” by f_2 . One of these guidance methods is incorporated into the image generation process. (b) The synthetic images were used in a digit judgment experiment where participants answered, “Is this picture a digit?” We trained a *digit judgment surrogate* based on human responses and used it as a classifier to guide the image generation process. (c) We used the images synthesized using the two guidance methods, ANN perceptual boundary sampling and digit judgment surrogate, for the digit recognition human experiment. Participants were asked “What digit is this picture?” A total of 19,952 images were used, with 123,000 trials conducted across 246 participants, resulting in the high perceptual variability dataset variMNIST.

Uncertainty without MSE



Uncertainty with MSE



Controversial without MSE



Controversial with MSE



Figure A.10. Comparison of generated images with/without MSE loss. The losses here are all untargeted.

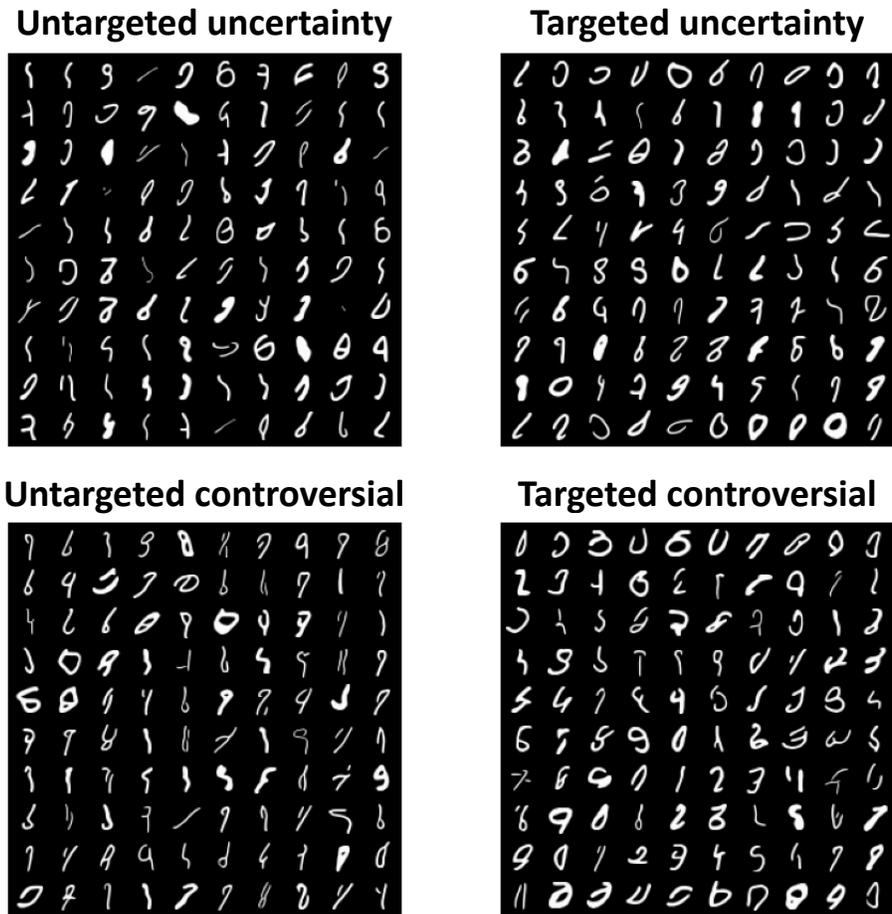


Figure A.11. Examples of generated stimuli.

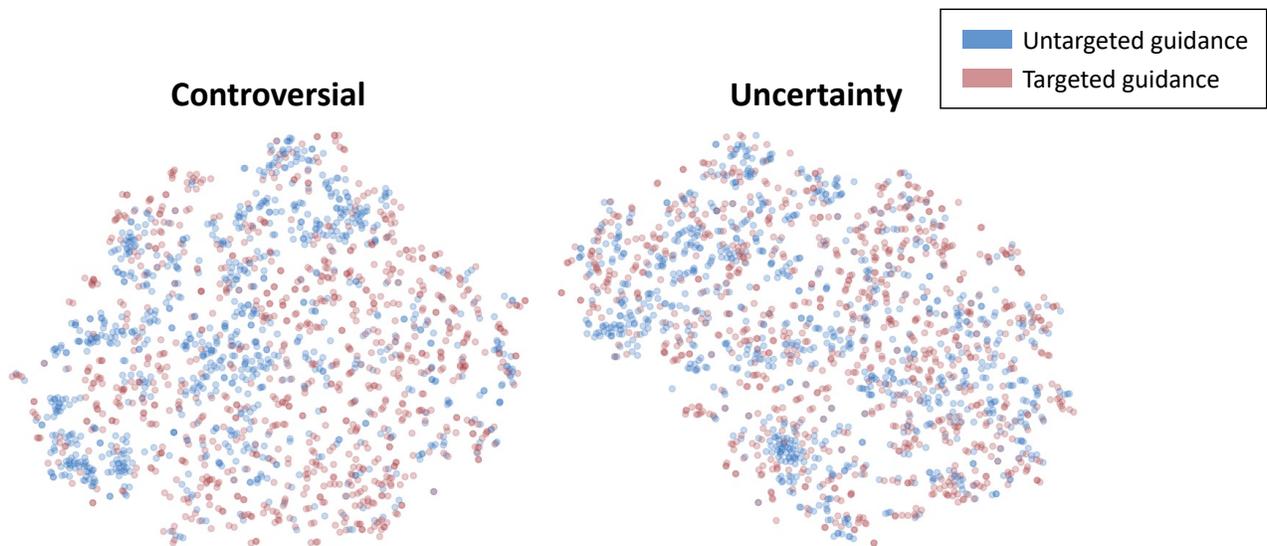


Figure A.12. t-SNE analysis on the pixel space for the generated images of different guidance methods. From the figure it is obvious that stimuli from targeted sampling are distributed more uniformly.

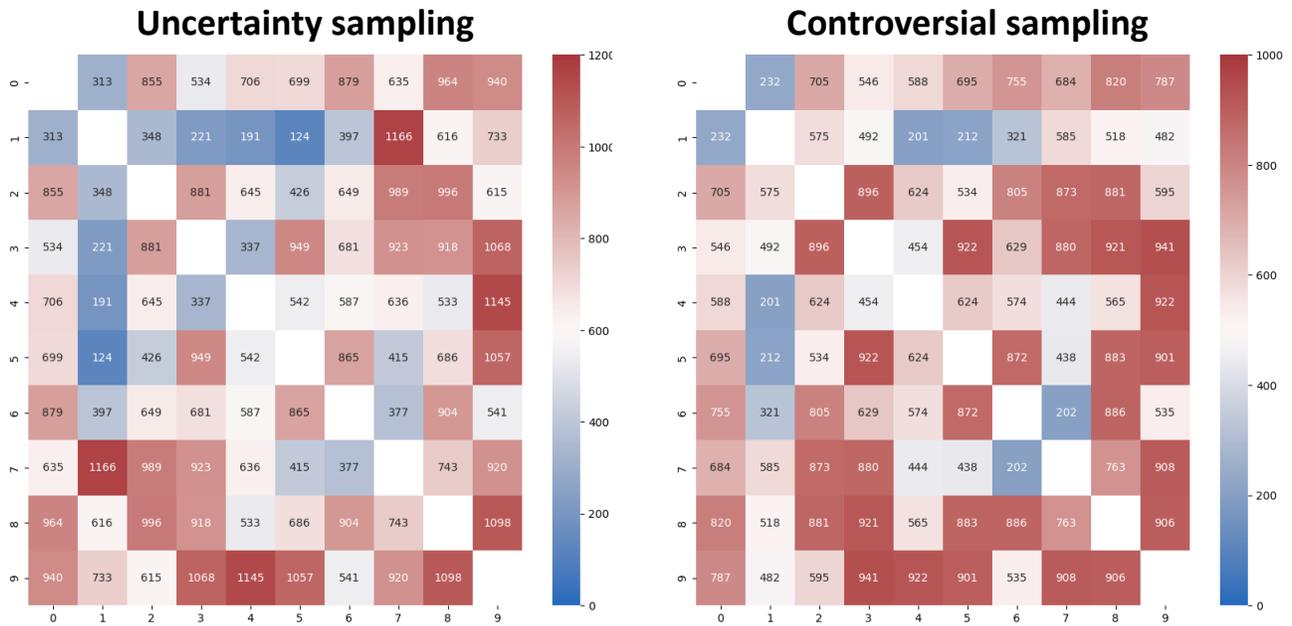


Figure A.13. **Category distribution in uncertainty sampling and controversial sampling.** Certain target pairs yield a higher number of stimuli that successfully pass the filtering criteria. For instance, target pairs such as (4, 9) and (8, 9) consistently produce more valid stimuli under both uncertainty and controversial sampling. In contrast, pairs like (1, 5) and (1, 3) result in significantly fewer stimuli meeting the filtering requirements.

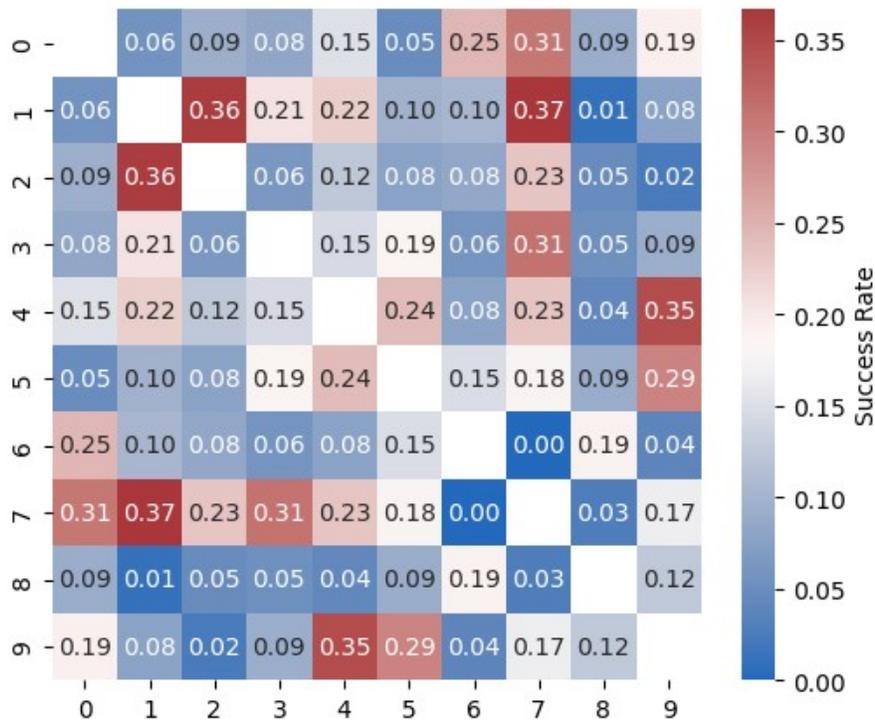


Figure A.14. **The success rate (proportion of successful stimuli) varies significantly across different guidance targets.** Pairs such as (1, 7), (1, 2), and (4, 9) achieve the highest success rates, exceeding 0.35. In contrast, pairs such as (1, 8), (2, 9), and (7, 8) have the lowest success rates, falling below 0.03.

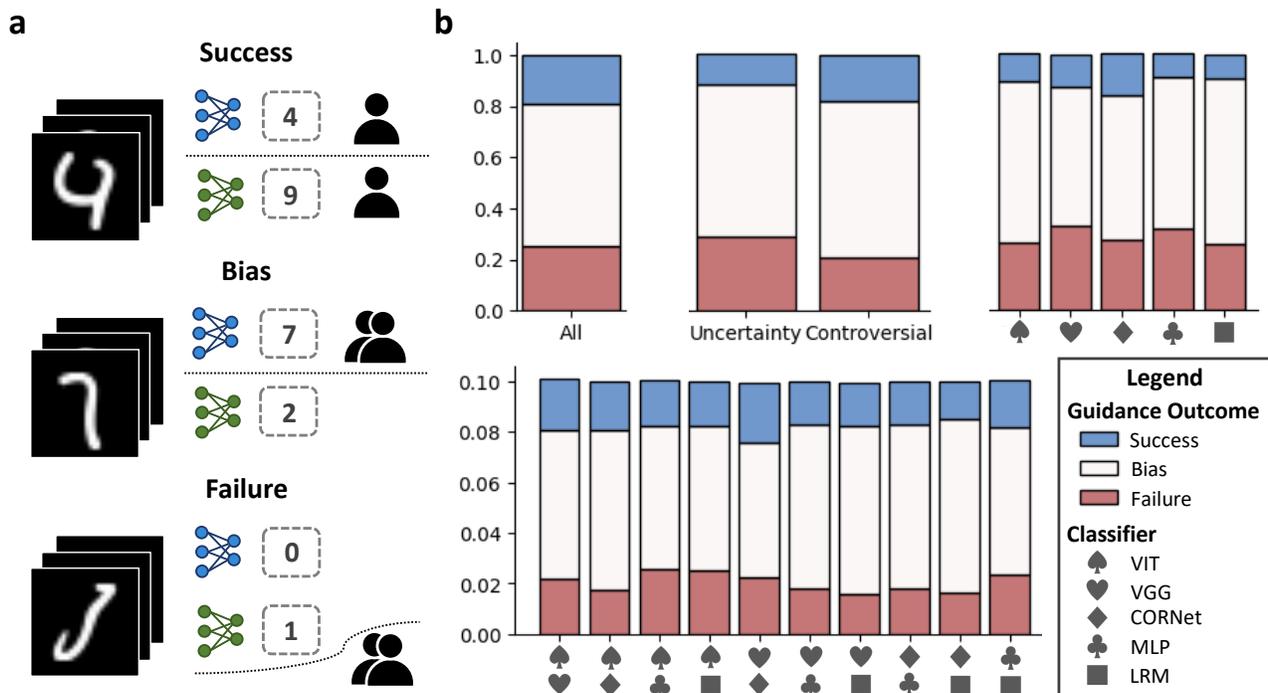


Figure A.15. **Quantitative analysis of variMNIST.** (a) Examples of three types of *guidance outcome*: *success*, *bias*, and *failure*. (b) Guidance outcomes across strategies and classifiers. The average sum of overall *success* and *bias* rates approaches 80%. Controversial guidance achieves a higher *success* rate than uncertainty guidance, with similar *bias* rates. CORNet performs best in uncertainty guidance, while LRM performs worst. In controversial guidance, combinations of VGG and CORNet achieve the highest *success* rates and lowest *bias* rates, but exhibit relatively high *bias* rates when paired with other classifiers.

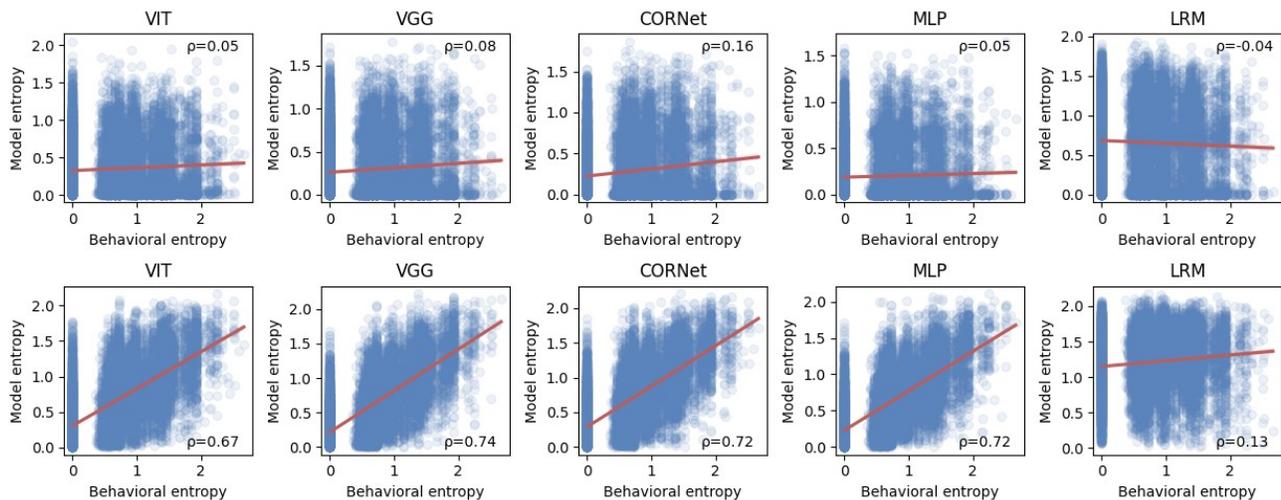


Figure A.16. **Correlation between model entropy and human behavior.** (a) Positive correlation between the entropy calculated from participants' behavior and the entropy predicted by the model for visual stimuli across five models. Each blue dot represents an image stimulus, and the red line shows the fitted result. (b) Significant improvement in the correlation between behavioral entropy and model-predicted entropy after fine-tuning on variMNIST, across five models.

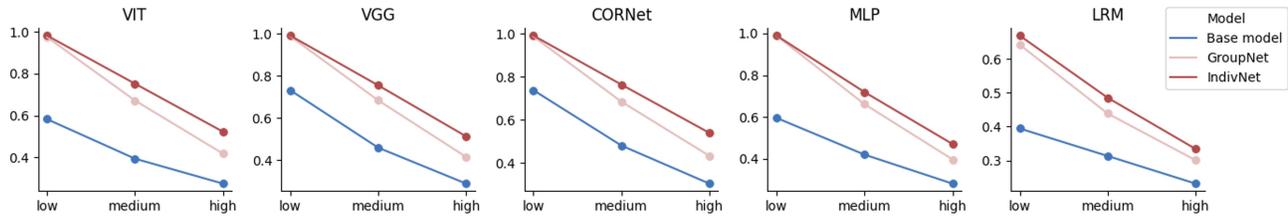


Figure A.17. Correlation between model entropy of different classifiers and human behavior. Prediction accuracy of different classifiers on images with varying entropy levels before and after fine-tuning. For all five classifiers, fine-tuned models show significant improvements in accuracy across all entropy levels compared to the baseline models. For the four classifiers other than LRM, the improvements of individual fine-tuned models over group-fine-tuned models are primarily observed on high-entropy images.

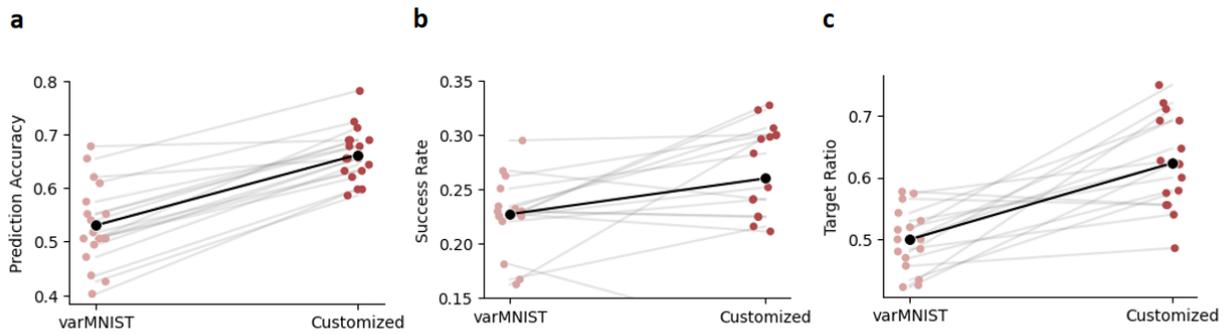


Figure A.18. Detailed analysis of each subject pair. (a) Prediction accuracy of individual models trained on in-lab participants, showing consistent improvement after fine-tuning. The original models are group models fine-tuned on variMNIST. (b) Comparison of guidance success rates between variMNIST and customized stimuli, indicating notable improvement for the majority of subject pairs. (c) Target ratio comparisons on variMNIST and customized stimuli, demonstrating an increase across nearly all subject pairs.

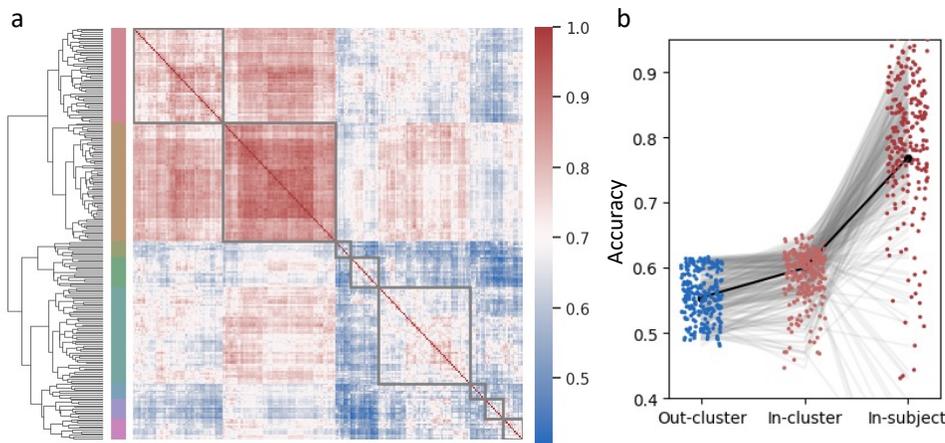


Figure A.19. Subject clustering analysis. (a) Subject similarity matrix and clustering results. The subject-finetuned model was used to predict the entire variMNIST dataset, and similarity between subjects was computed based on their prediction results. The left axis and gray boxes indicate subjects belonging to the same cluster, with a total of eight clusters. (b) Performance of the subject-finetuned model in predicting data from different groups: out-cluster, in-cluster, and in-subject correspond to different clusters, the same cluster, and the subject itself, respectively. Each point represents the average prediction performance of a subject on data from the corresponding group, and The black line represents the average of all subjects.

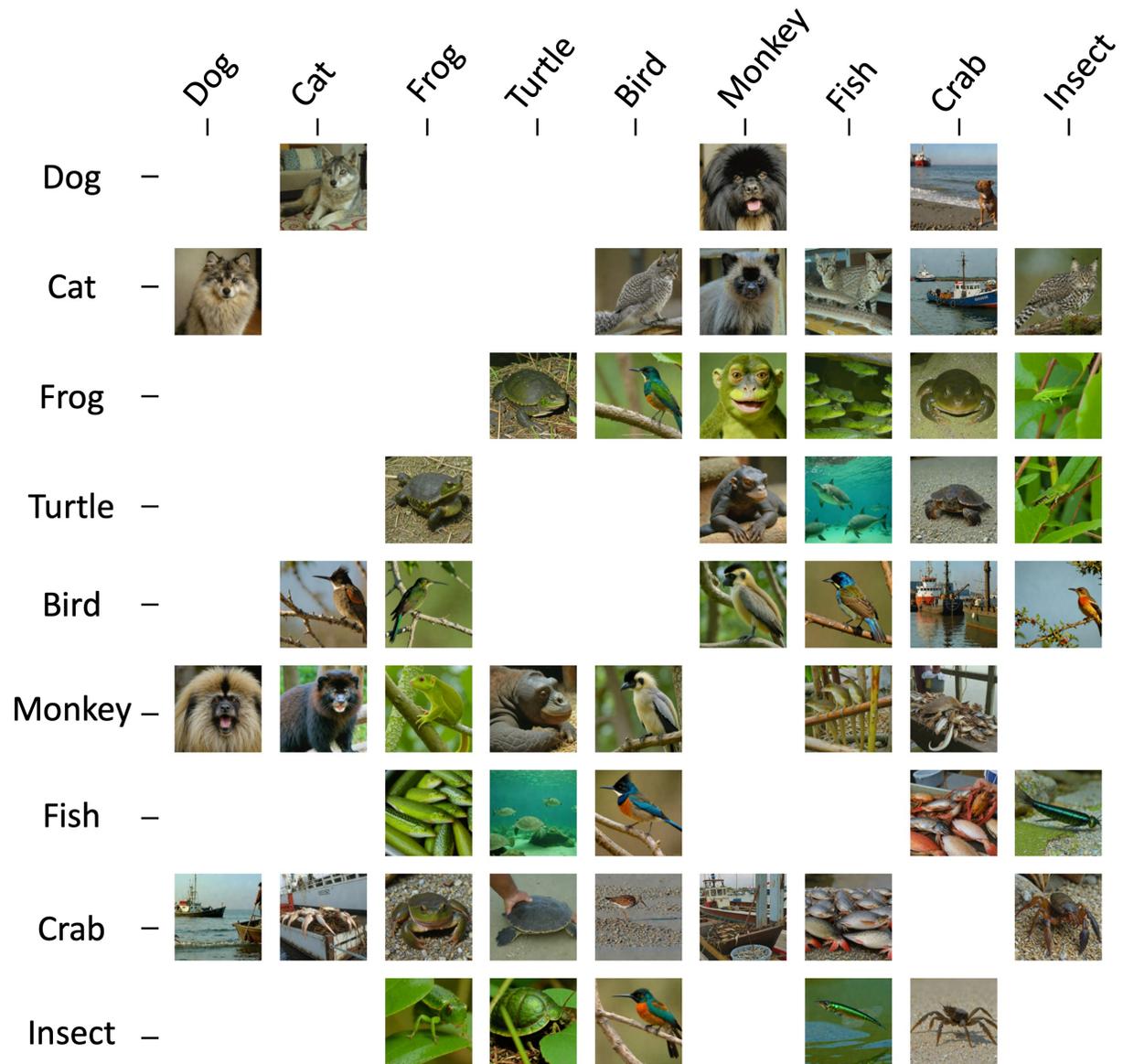


Figure A.20. Examples of stimuli for the natural image based experiment. We only use the categories from the Restricted ImageNet (Engstrom et al., 2019).

Examples of label pairs with the highest Success rate



Figure A.21. Examples of stimuli for the natural image based experiment that have the highest value of uncertainty. We only use the categories from the Restricted ImageNet (Engstrom et al., 2019).

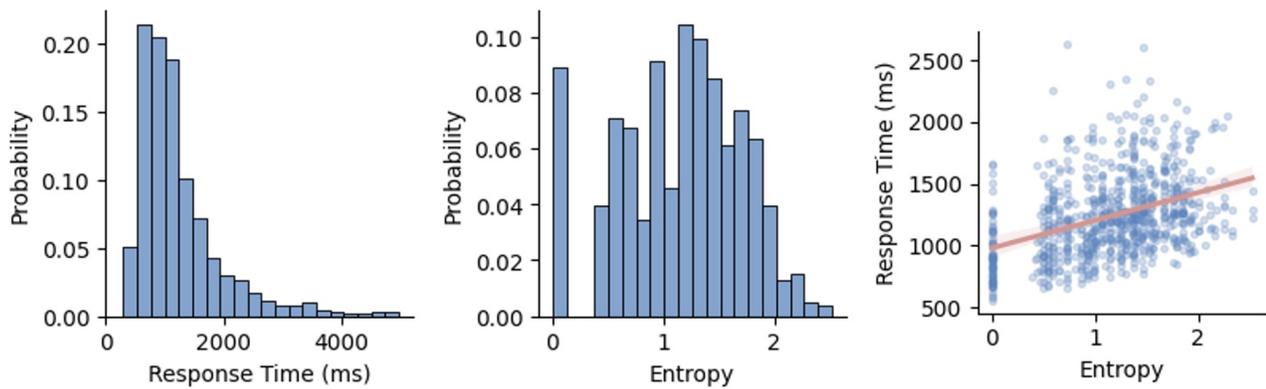


Figure A.22. Further analysis on the samples. On the left we show the distribution of the response time of the samples, in the middle we show the entropy distribution of the samples. On the right, we show the correlation between the entropy and the response time. It can be observed that the entropy and the response time are correlated.

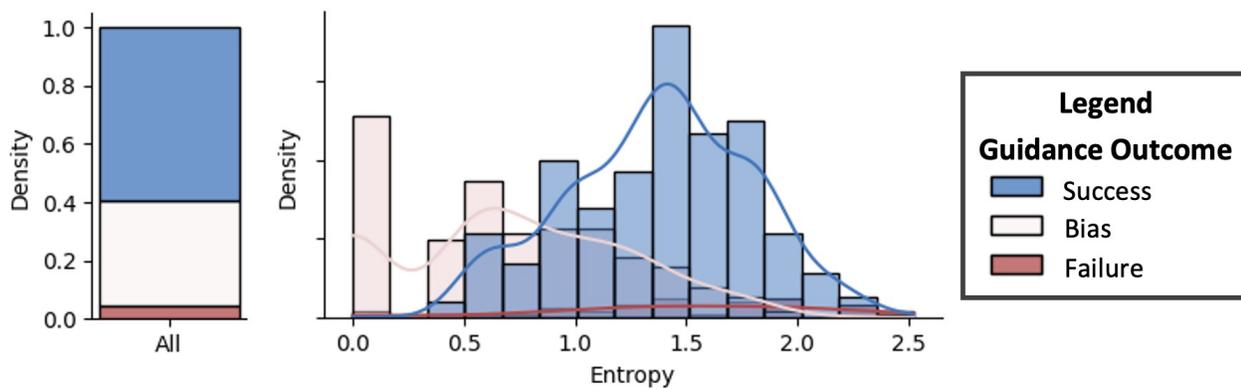


Figure A.23. Guidance outcome of the Imagenet based experiment. The left part shows the proportion of the Success, Bias and Failure cases. The right part shows the entropy distribution of the samples.

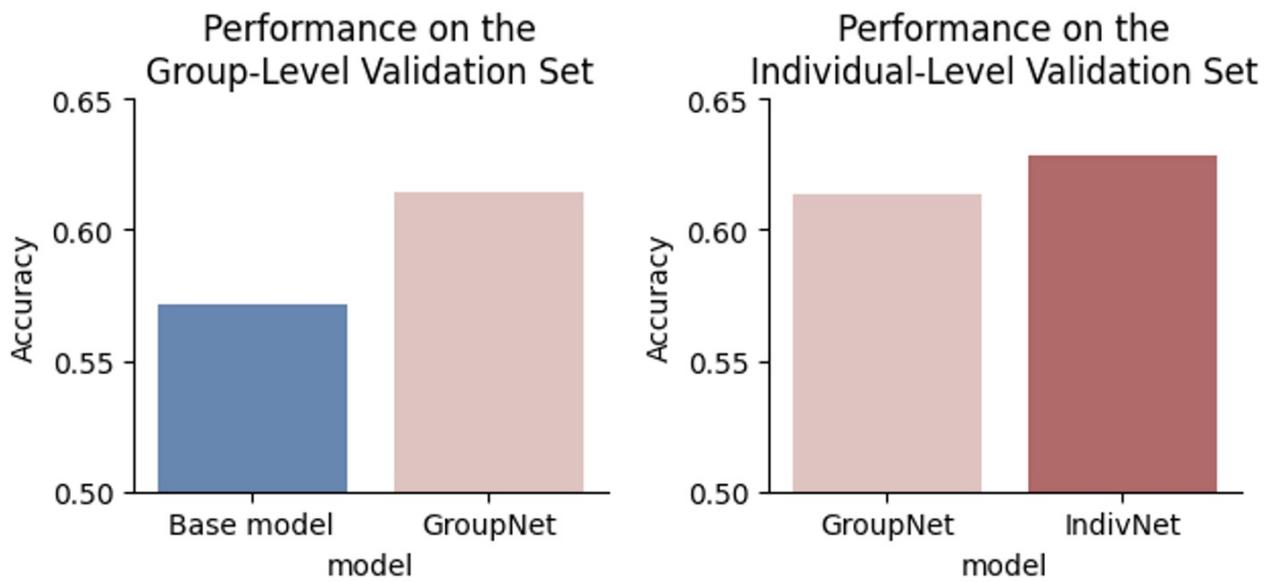


Figure A.24. **Performance of models on different validation sets.** The accuracy of the group model on the group-level validation set increased by about 4 percent and the individual-level finetune increased the accuracy by about 2 percent, which aligns with the results on MNIST.

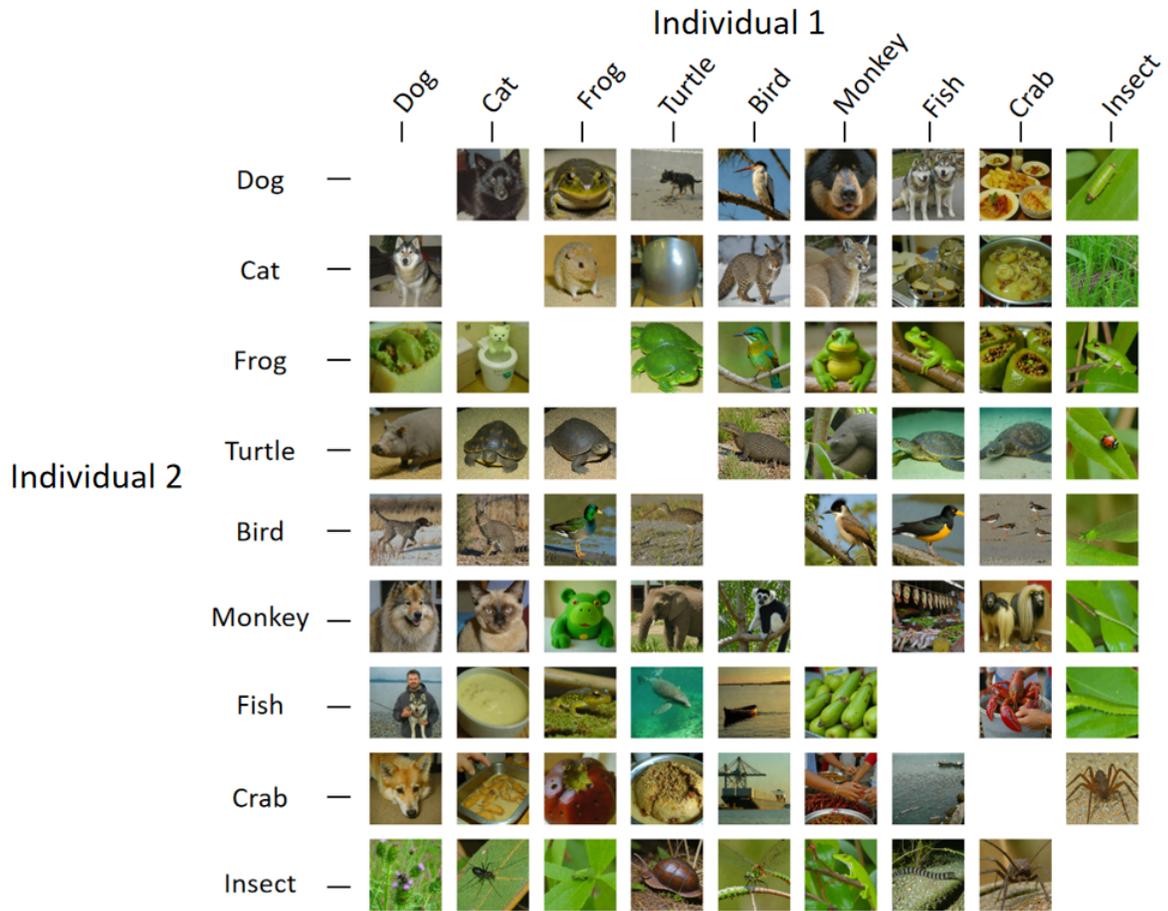


Figure A.25. Controversial generation results on individually finetuned models.

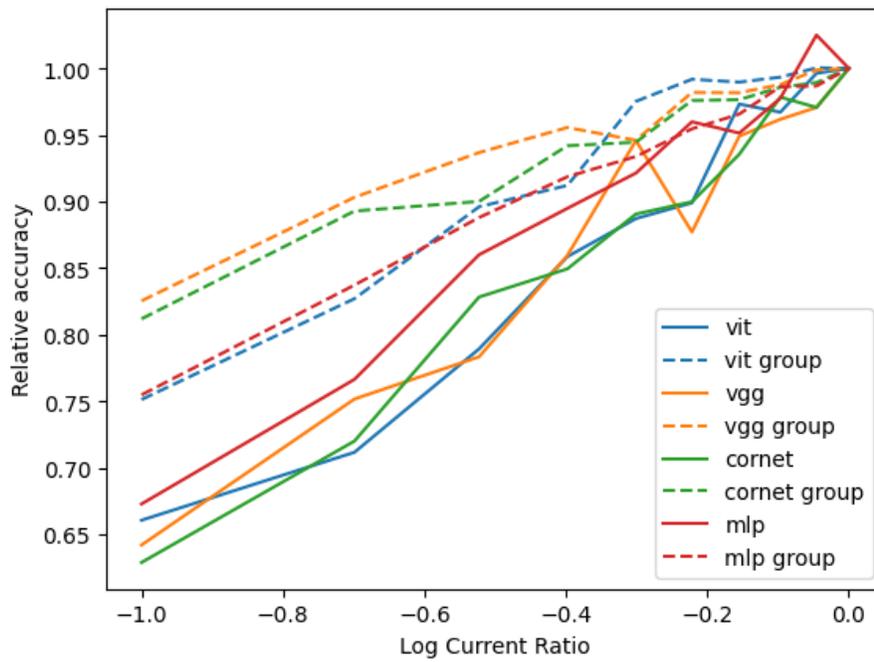


Figure A.26. Relative accuracies (improvements at a reduced data volume/improvement at full data volume) v.s. data volume. It can be observed that increasing data volume has diminishing returns on GroupNet accuracy, whereas IndivNet accuracy continues to improve significantly.

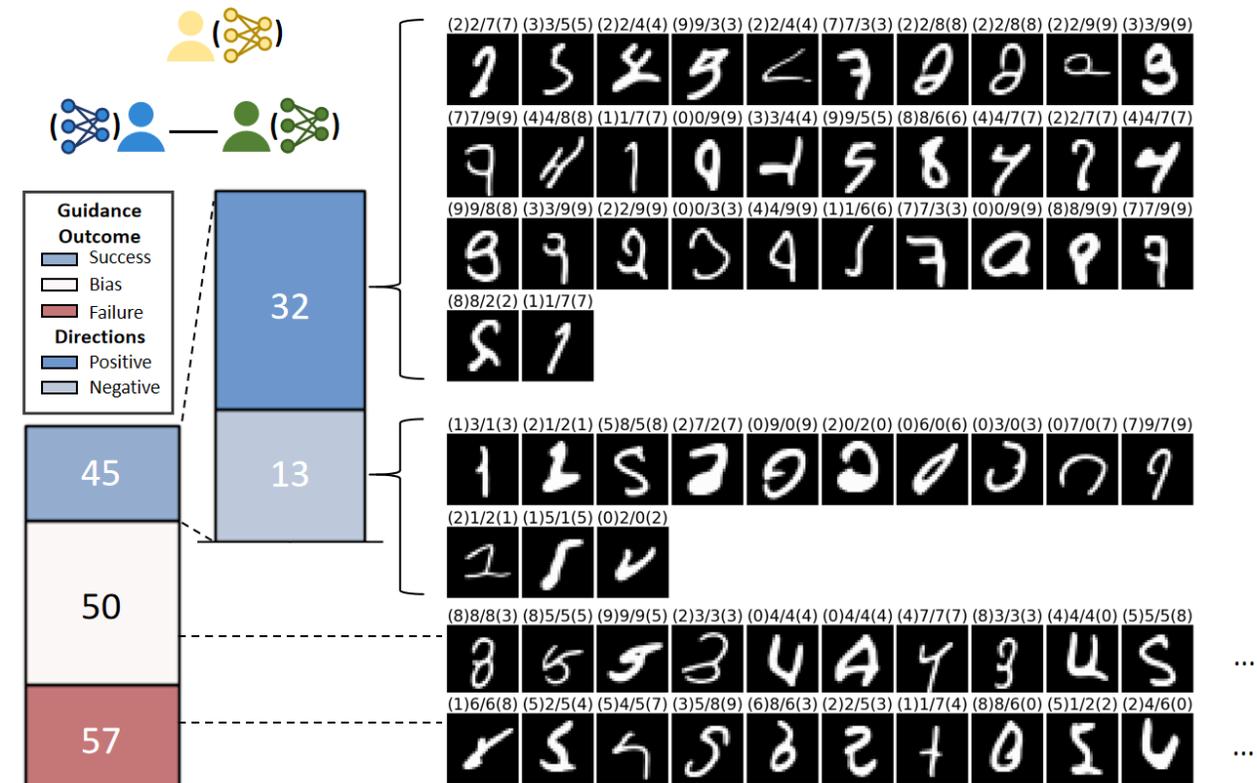


Figure A.27. Examples of manipulation stimuli for subject 1 and subject 2. The left part of the figure shows the actual numbers of each category of stimuli. The real stimuli used to manipulate the subjects are shown on the right. The choices of the subjects are in the middle, with the guidance label marked in parentheses. All positive and negative examples are presented, along with 10 typical bias and failure cases.

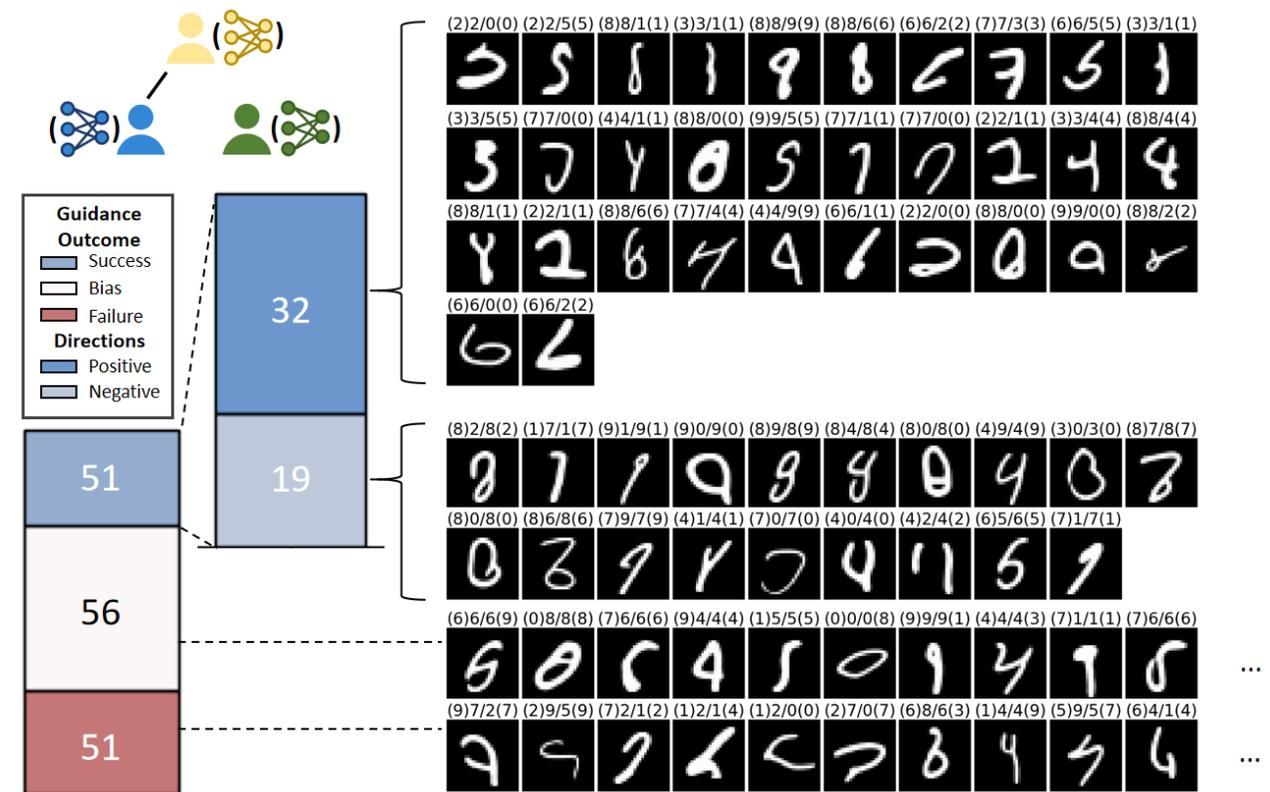


Figure A.28. Examples of manipulation stimuli for subject 1 and subject 3. The structure is the same as Figure A.27.

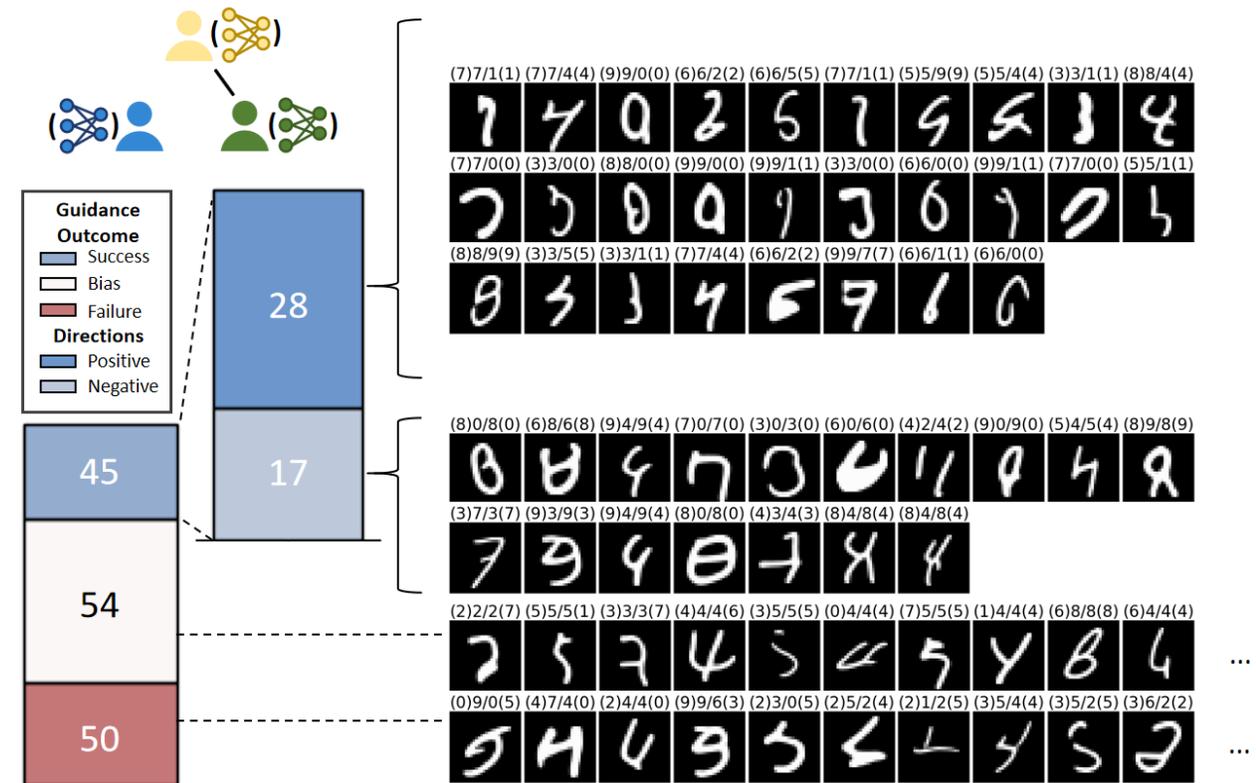


Figure A.29. Examples of manipulation stimuli for subject 2 and subject 3. The structure is the same as Figure A.27.