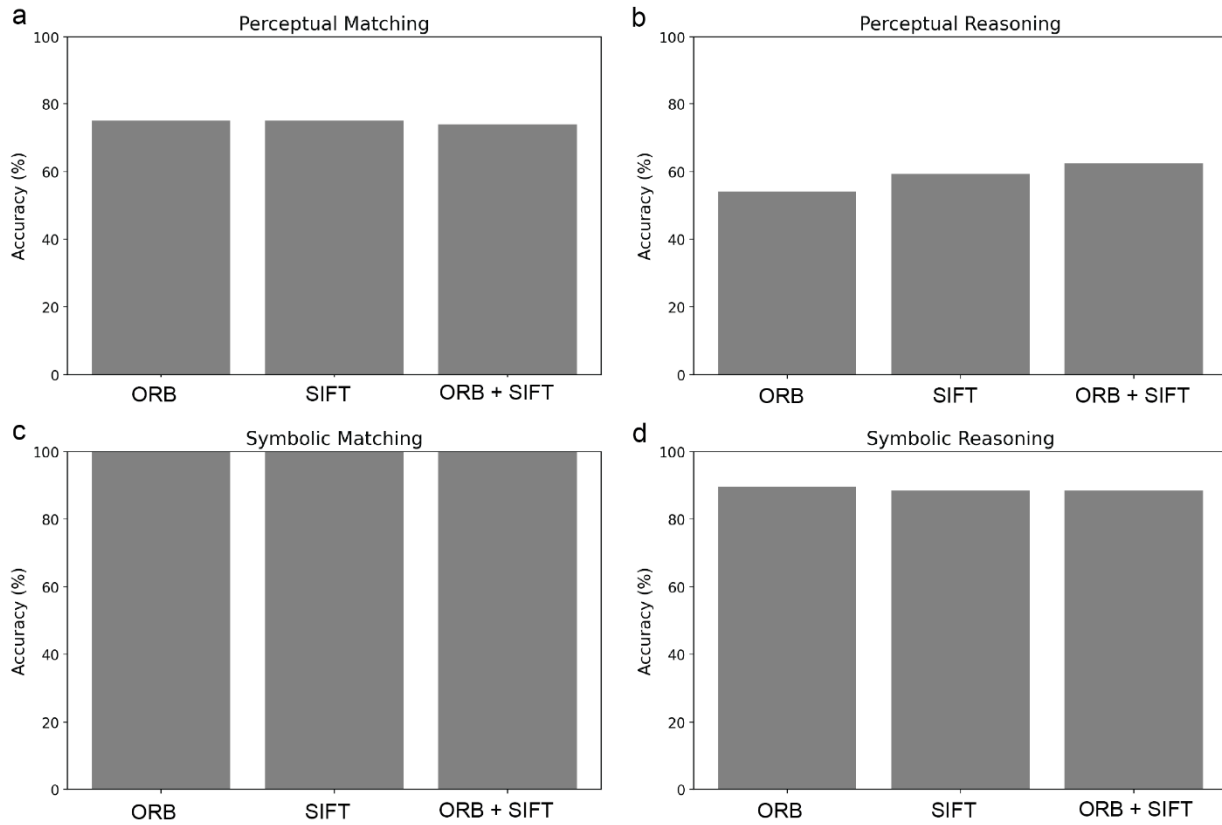


1 **Supplements**

2



3

4

5 **Supplemental Figure 1: Direct comparison of SIFT and ORB through the full model. a)**

6 Performance of the model on the Perceptual Matching condition using ORB, SIFT and hybrid

7 ORB-SIFT features. b) Performance of the model on the Perceptual Reasoning condition using

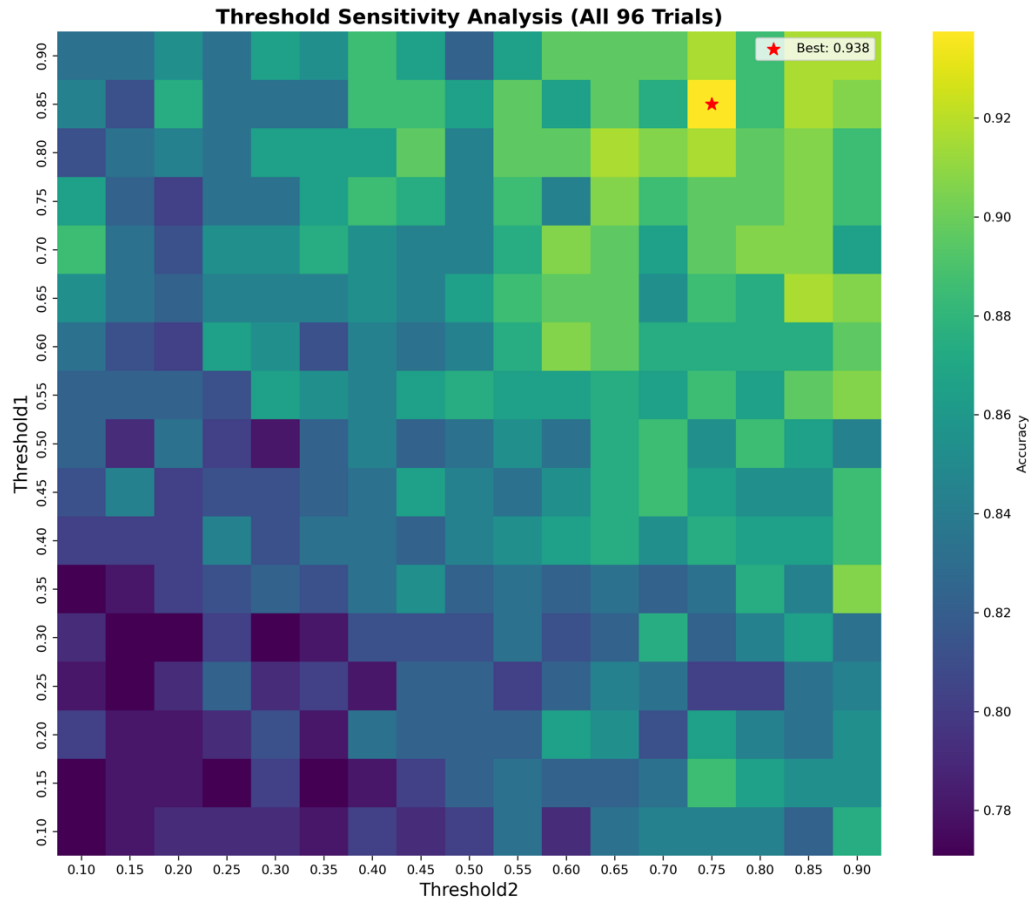
8 ORB, SIFT and hybrid ORB-SIFT features. c) Performance of the model on the Symbolic Matching

9 condition using ORB, SIFT and hybrid ORB-SIFT features. d) Performance of the model on the

10 Symbolic Reasoning condition using ORB, SIFT and hybrid ORB-SIFT features.

11

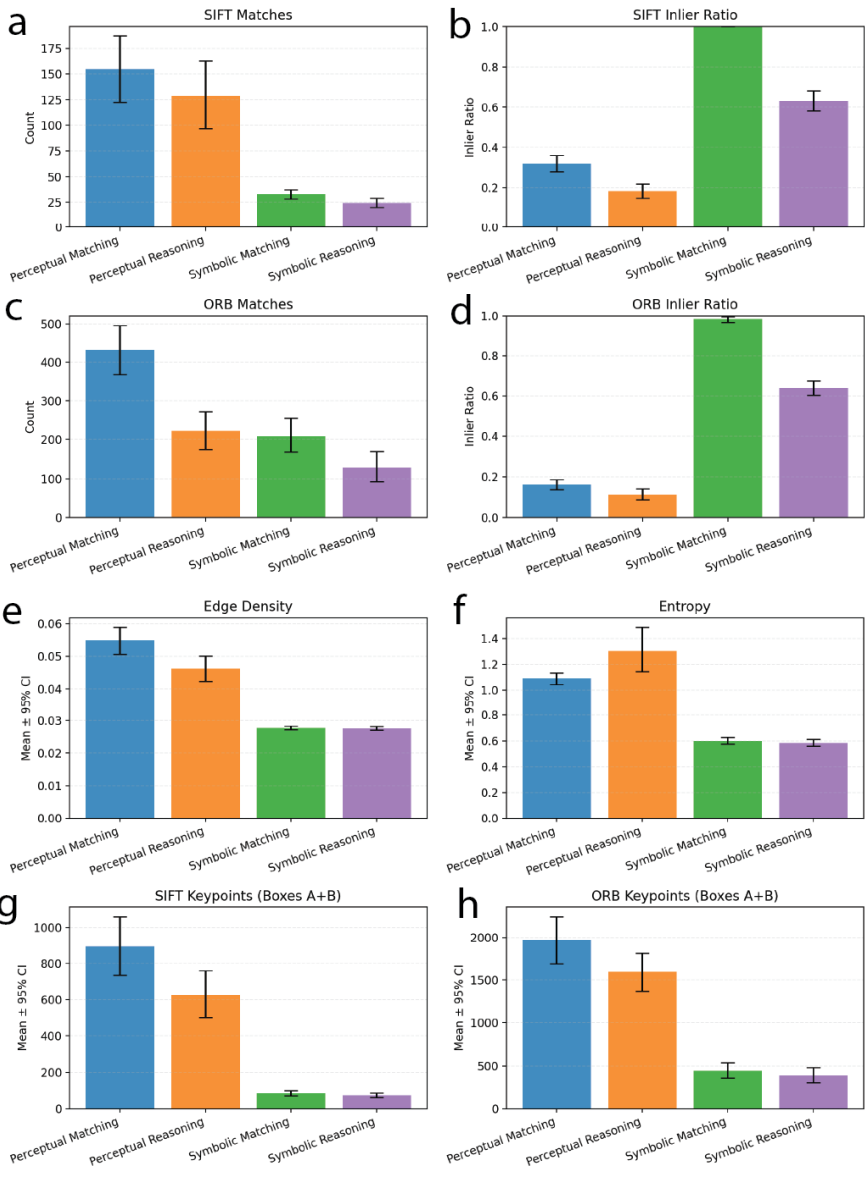
12



13

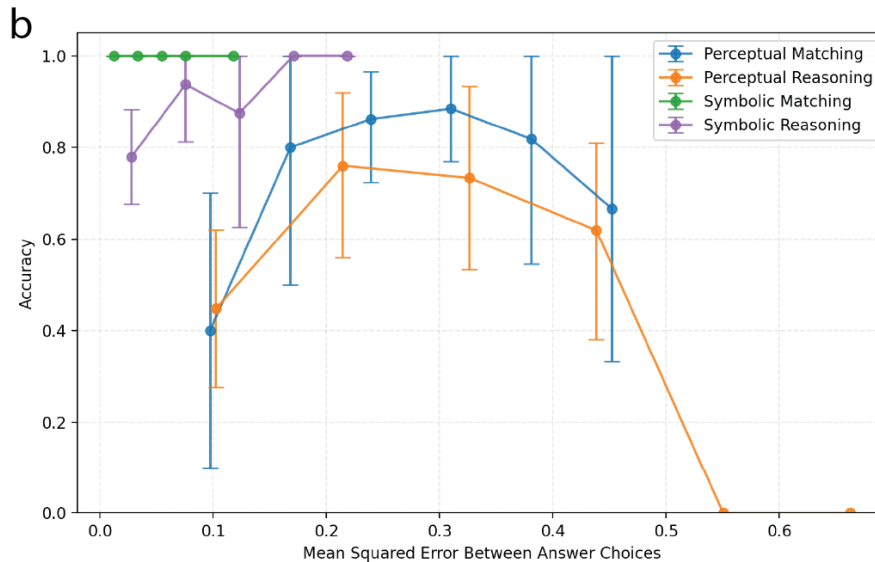
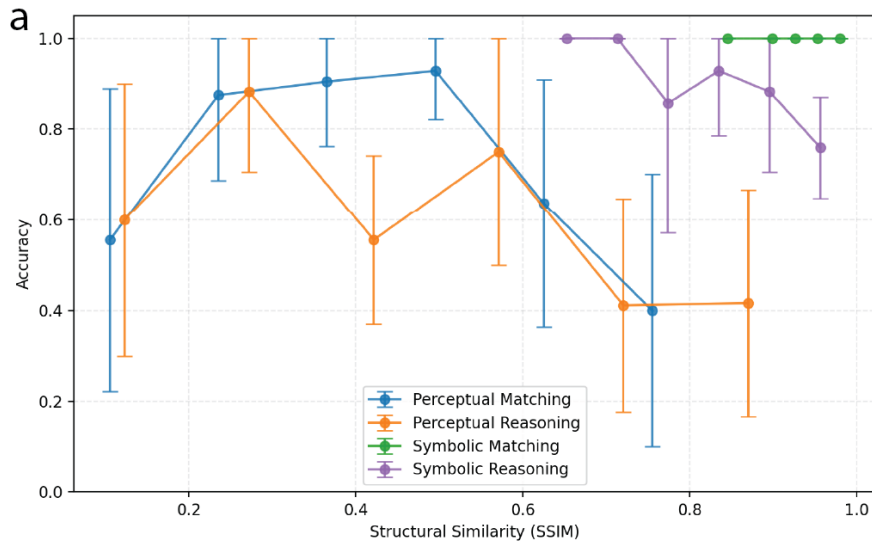
14 **Supplemental Figure 2: Sensitivity analysis of threshold parameters and their effect on model**  
 15 **performance.** The heatmap shows model accuracy (color-coded from dark blue = low to yellow  
 16 = high) as a function of two active threshold parameters, Threshold 1 and Threshold 2. Each  
 17 parameter was systematically varied from 0.10 to 0.90 in 17 increments, resulting in 289 total  
 18 combinations (17 × 17 grid). Each configuration was tested across 96 trials in the symbolic  
 19 reasoning condition. Overall accuracy ranged from 77.1% to 93.8% (mean = 84.5% ± 3.7%),  
 20 revealing a 16.7-point performance spread across the parameter space. The optimal  
 21 configuration occurred at Threshold 1 = 0.85 and Threshold 2 = 0.75, achieving 93.8% accuracy.

22  
 23  
 24  
 25  
 26  
 27  
 28

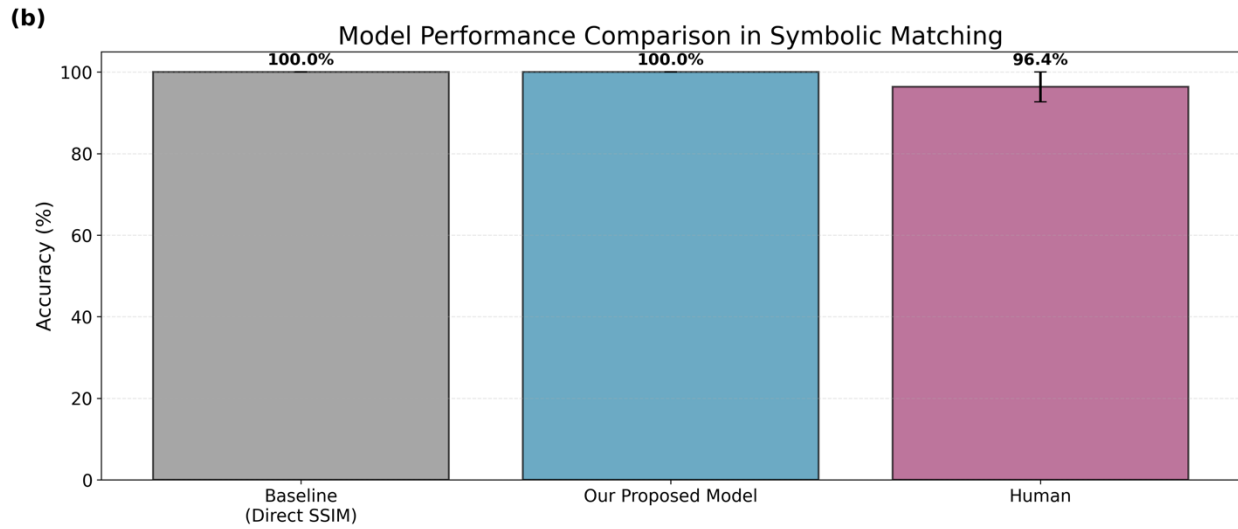
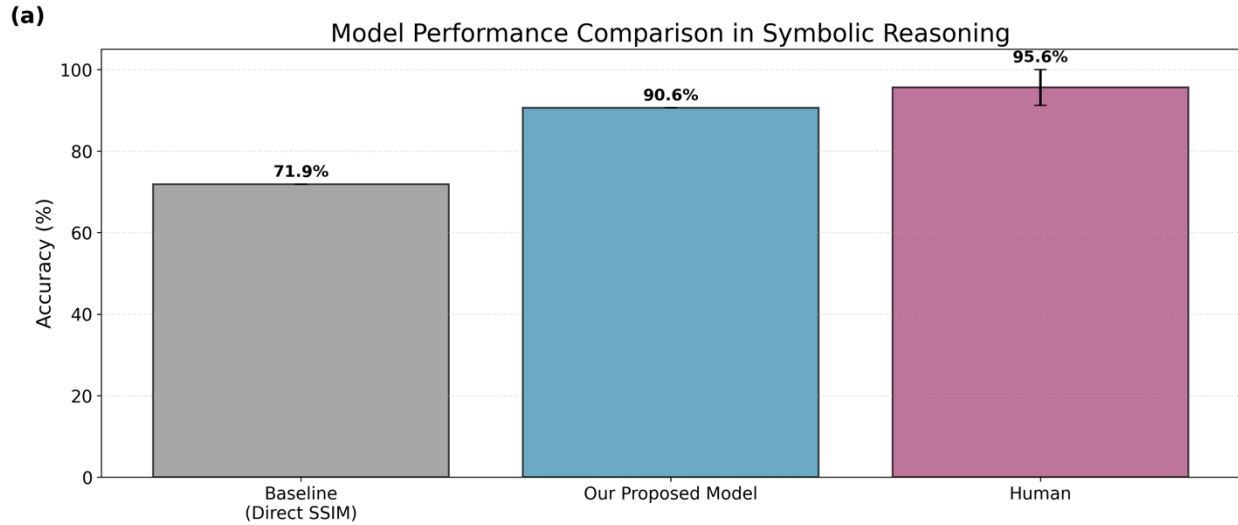


29

30 **Supplemental Figure 3: Task difficulty and visual complexity metrics across four task**  
 31 **conditions. a–d)** Bar charts show mean ± 95% bootstrap confidence intervals (2,000 resamples;  
 32 fixed RNG seed) for **SIFT Matches**, **SIFT Inlier Ratio**, **ORB Matches**, and **ORB Inlier Ratio**  
 33 across the four conditions (*Perceptual Matching*, *Perceptual Reasoning*, *Symbolic Matching*, *Symbolic*  
 34 *Reasoning*). “Matches” denotes the number of descriptor correspondences; **Inlier**  
 35 **Ratio** represents the fraction of matches deemed geometrically consistent by RANSAC during  
 36 similarity-transform estimation (bounded 0–1; y-axis limited accordingly). **e–h)** Bar charts show  
 37 mean ± 95% bootstrap confidence intervals (2,000 resamples; fixed RNG seed) for **visual-**  
 38 **complexity metrics** computed on the cue images: **Edge Density**, **Shannon Entropy**, and feature  
 39 counts (**SIFT Keypoints** and **ORB Keypoints**) aggregated over Boxes A + B. Error bars reflect  
 40 percentile confidence intervals; x-axis labels denote the four task conditions (*Perceptual*  
 41 *Matching*, *Perceptual Reasoning*, *Symbolic Matching*, *Symbolic Reasoning*).



44 **Supplemental Figure 4. Model accuracy as a function of answer-choice similarity. a)** Model  
 45 accuracy (y-axis) plotted against the **structural similarity index (SSIM)** between the two answer  
 46 choices (x-axis). SSIM ranges from 0 (dissimilar) to 1 (identical) and reflects perceptual  
 47 similarity. **b)** Model accuracy (y-axis) plotted against the **mean squared error (MSE)** between  
 48 the two answer choices (x-axis). MSE represents the normalized per-pixel difference between  
 49 candidate answers; lower values indicate higher similarity. For both panels, points show mean  
 50 accuracies within uniformly spaced bins (e.g., 20 bins), connected to emphasize trends. Shaded  
 51 regions denote **95% bootstrap confidence intervals** (2,000 resamples) around the mean  
 52 accuracy in each bin.



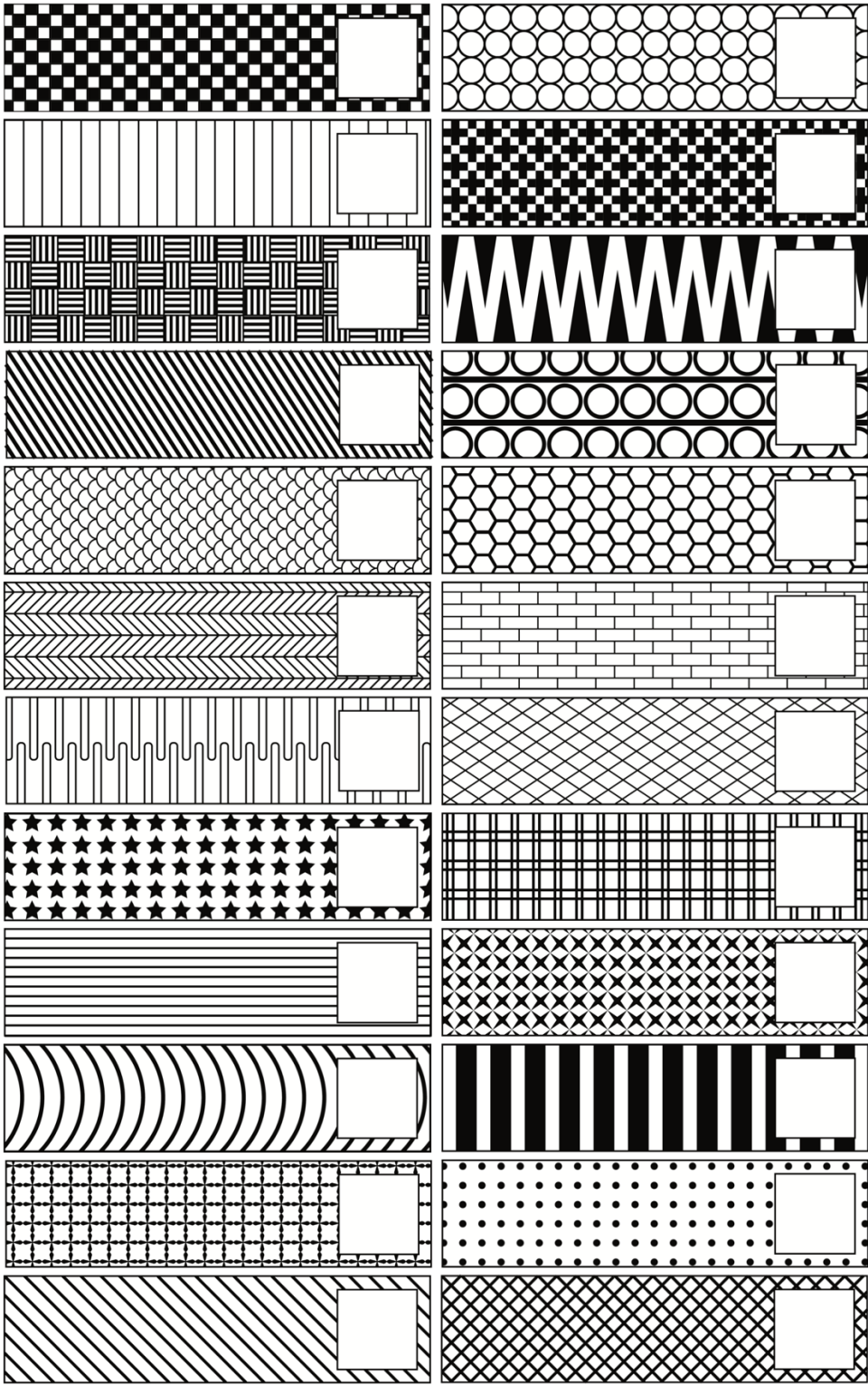
54

55

56 **Supplemental Figure 5: Comparison with a template matching strategy.** Our Proposed Model  
 57 substantially outperforms a direct SSIM baseline on symbolic reasoning and closely approaches  
 58 human performance. Top row (Symbolic Reasoning): Baseline achieves 71.9%, Our Proposed  
 59 Model reaches 90.6%, and Humans achieve 95.6% (error bars denote 95% confidence interval,  
 60 capped within 0–100%). Bottom row (Symbolic Matching): Baseline and Our Proposed Model  
 61 both achieve 100%, while Humans achieve 96.36% ( $\pm 18.73\%$ , capped). Bars show mean  
 62 accuracies; human error bars reflect inter-subject variability.

63

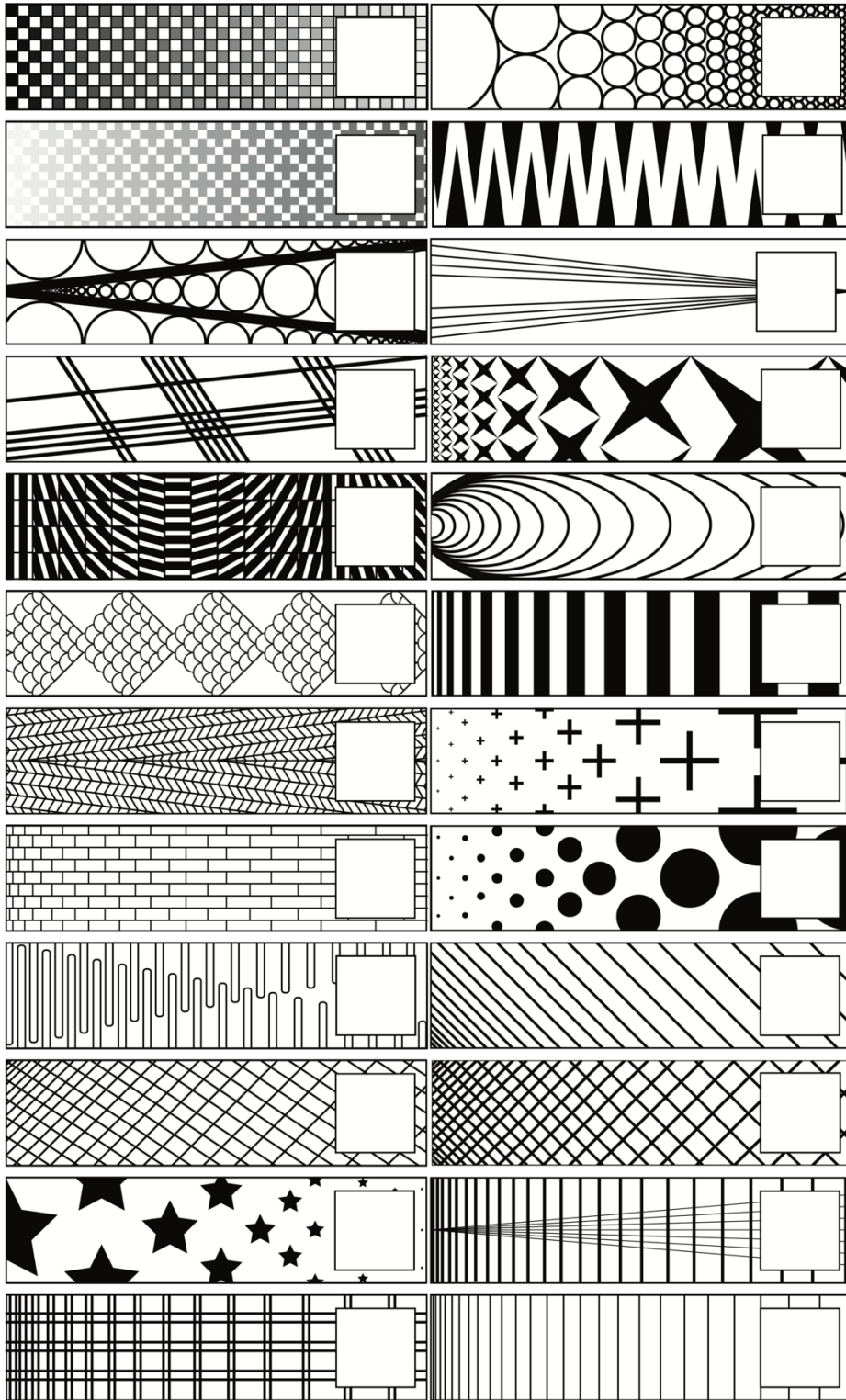
Perceptual Matching Stimuli



64  
65

Supplemental Figure 6: Additional Examples for Perceptual Matching Stimuli.

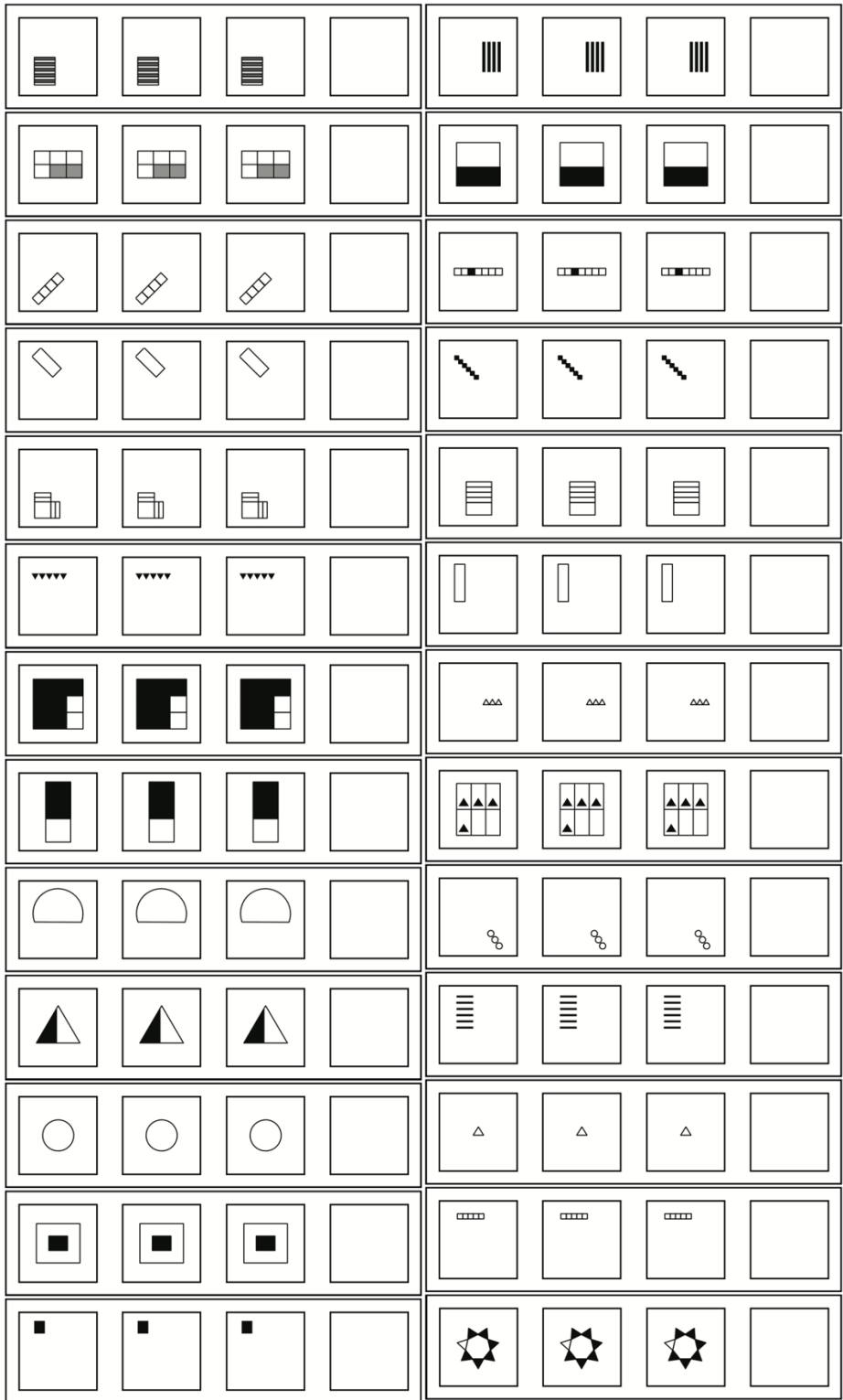
Perceptual Reasoning Stimuli



66  
67  
68  
69

Supplemental Figure 7: Additional Examples for Perceptual Reasoning Stimuli.

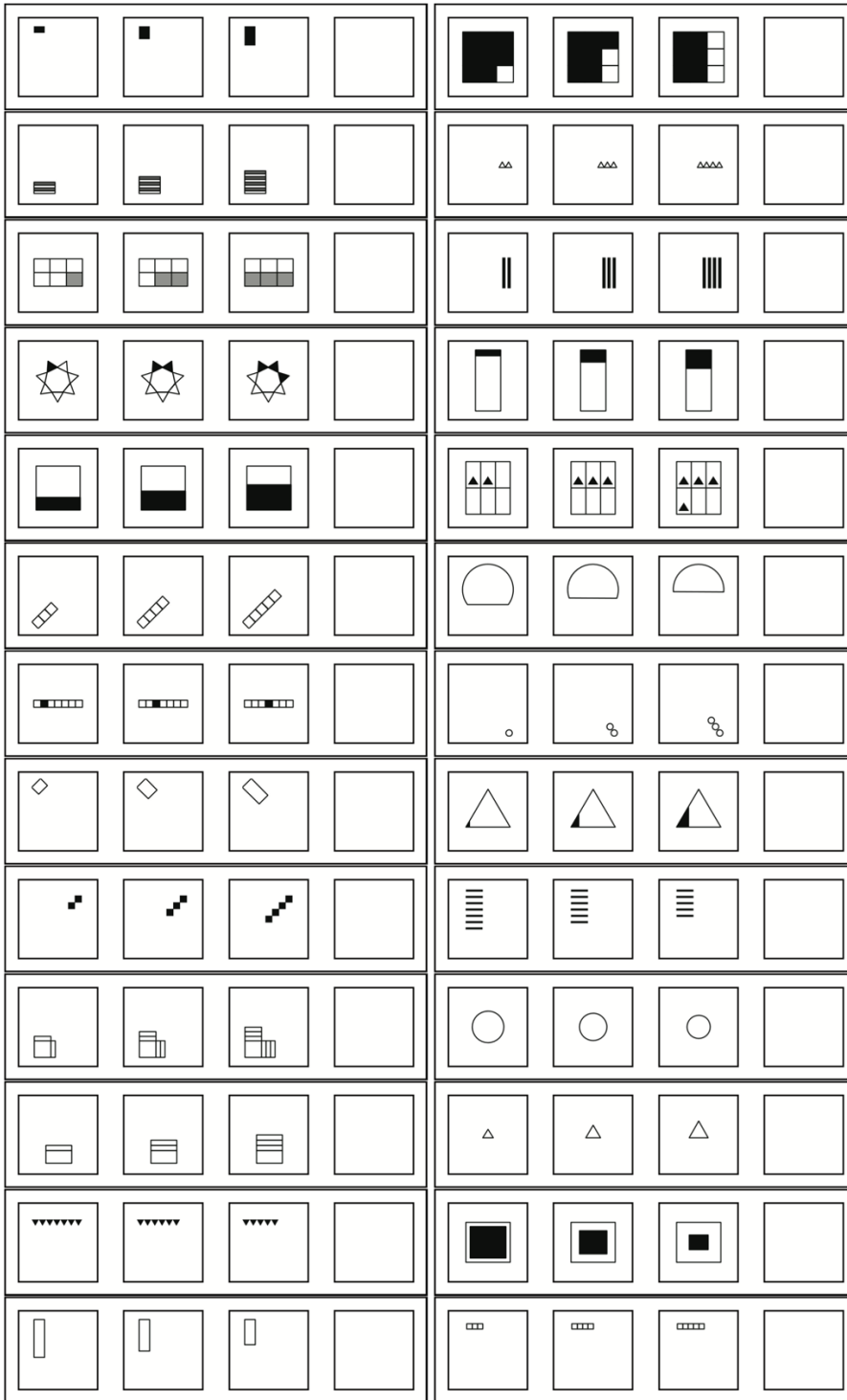
# Symbolic Matching Stimuli



70  
71

Supplemental Figure 8: Additional Examples for Symbolic Matching Stimuli.

### Symbolic Reasoning Stimuli



72  
73  
74  
75

**Supplemental Figure 9: Additional Examples for Symbolic Reasoning Stimuli.**