Viewpoint dependence and scene context effects generalize to depth rotated three-dimensional objects

Aylin Kallmayer

Melissa L.-H. Võ

Dejan Draschkow

Department of Psychology, Goethe University Frankfurt, Frankfurt am Main, Germany

Department of Psychology, Goethe University Frankfurt, Frankfurt am Main, Germany

> Department of Experimental Psychology, University of Oxford, Oxford, UK Oxford Centre for Human Brain Activity, Wellcome Centre for Integrative Neuroimaging, Department of Psychiatry, University of Oxford, Oxford, UK

Viewpoint effects on object recognition interact with object-scene consistency effects. While recognition of objects seen from "noncanonical" viewpoints (e.g., a cup from below) is typically impeded compared to processing of objects seen from canonical viewpoints (e.g., the string-side of a guitar), this effect is reduced by meaningful scene context information. In the present study we investigated if these findings established by using photographic images, generalize to strongly noncanonical orientations of three-dimensional (3D) models of objects. Using 3D models allowed us to probe a broad range of viewpoints and empirically establish viewpoints with very strong noncanonical and canonical orientations. In Experiment 1, we presented 3D models of objects from six different viewpoints (0°, 60°, 120°, 180° 240°, 300°) in color (1a) and grayscaled (1b) in a sequential matching task. Viewpoint had a significant effect on accuracy and response times. Based on the viewpoint effect in Experiments 1a and 1b, we could empirically determine the most canonical and noncanonical viewpoints from our set of viewpoints to use in Experiment 2. In Experiment 2, participants again performed a sequential matching task, however now the objects were paired with scene backgrounds which could be either consistent (e.g., a cup in the kitchen) or inconsistent (e.g., a guitar in the bathroom) to the object. Viewpoint interacted significantly with scene consistency in that object recognition was less affected by viewpoint when consistent scene information was provided, compared to inconsistent information. Our results show that scene context supports object recognition even when using extremely noncanonical orientations of depth rotated 3D objects. This supports the important role object-scene processing plays for

object constancy especially under conditions of high uncertainty.

Introduction

Object recognition happens fast, automatically, and in most cases seems effortless to us. Because our environment is highly dynamic, especially when interacting with it, one and the same object will produce a range of different images on the retina. In fact, it is very unlikely that an object would produce the same retinal image twice owing to changes in viewpoint, lighting, reflections, or viewing distance. Still, our visual system is able to flexibly transform this variable visual input in a way that object identity can successfully be read out from the resulting abstract representations in higher areas of visual cortex (see DiCarlo & Cox, 2007).

Whether object recognition is viewpoint dependent (recognition performance is sensitive to changes in viewpoints as indicated by accuracy and response time [RT] data) or viewpoint invariant (recognition performance is largely unaffected by changes in viewpoint) has been a debated topic (Biederman & Gerhardstein, 1993; Bülthoff & Edelman, 1992; Burgund & Marsolek, 2000; Leek, Atherton, & Thierry, 2007; Edelman, 1995; Graf, 2006; Hayward, 2003; Hayward & Tarr, 1997; Jolicoeur, 1990; Leek et al., 2007; Lowe, 1987; Marr, Nishihara, & Brenner, 1978; Ratan Murty & Arun, 2015; Stankiewicz, 2002; Tarr & Bülthoff, 1995; Tarr & Pinker, 1989; Wilson & Farah, 2003). Since the early debates, there has been

Citation: Kallmayer, A., Võ, M. L.-H., & Draschkow, D. (2023). Viewpoint dependence and scene context effects generalize to depth rotated three-dimensional objects. *Journal of Vision*, 23(10):9, 1–13, https://doi.org/10.1167/jov.23.10.9.

Received November 16, 2022; published September 14, 2023

ISSN 1534-7362 Copyright 2023 The Authors



 \searrow

overwhelming consensus that object recognition is neither solely viewpoint dependent nor solely viewpoint invariant and that evidence for both can be observed depending on experimental task and stimuli (Foster & Gilson, 2002; Hamm & McMullen, 1998; Jolicoeur, 1990; Leek et al., 2007; Ratan Murty & Arun, 2015; Sastyin, Niimi, & Yokosawa, 2015; Stankiewicz, 2002; Vanrie, Béatse, Wagemans, Sunaert, & Van Hecke, 2002).

Past research has made great advances toward understanding the mechanisms that underly invariant object recognition, when objects are presented in isolation (i.e., DiCarlo & Cox, 2007). More recently, however, researchers have started to investigate the viewpoint problem in the context of object-scene processing. Object recognition rarely occurs in isolation where the only available information are the objects' features. In our everyday lives, we encounter objects within certain contexts, which provides us with a pool of complex visual and multimodal information that is integrated during object recognition. Past research has shown that context facilitates object recognition (Biederman, Mezzanotte, & Rabinowitz, 1982; Oliva & Torralba, 2007; for a recent review see Lauer et al., 2021). Evidence from behavioral as well as neurophysiological studies (e.g., Brandman & Peelen, 2017) suggest an interactive processing of objects and scenes. For instance, objects placed in semantically consistent contexts are recognized faster and more accurately, often referred to as the scene-consistency effect (Davenport & Potter, 2004; Palmer, 1975). Accordingly, models of object recognition have been updated to incorporate the integration of contextual information (Bar, 2004). Further, frameworks incorporating object-scene and object-object relations (e.g., the so-called *scene* grammar) describe a set of internalized rules based on regularities found in real-world scenes that facilitate scene and object perception and guide our attention during different visual cognitive tasks (Draschkow & Võ, 2017; Josephs, Draschkow, & Võ, 2016; Võ, 2021; Võ, Boettcher, & Draschkow, 2019; Võ & Henderson, 2009; Võ & Wolfe, 2013a; Võ & Wolfe, 2013b).

Recent work has also looked at influences of object and scene orientation on the scene consistency effect (Lauer, Schmidt, & Võ, 2020; Sastyin et al., 2015). Sastyin et al. (2015) conducted a series of experiments investigating the interaction between viewpoint and scene consistency on object and scene recognition. They used photographic images of objects from canonical and noncanonical viewpoints and paired them with consistent and inconsistent scenes. They evaluated viewpoints in a relative manner with canonical viewpoints containing more canonical characteristics than noncanonical viewpoints as determined by rating the stimuli. Others have defined canonical viewpoints as the viewpoint from which one would photograph an object or the viewpoint from which one sees the object when imagining it, mostly finding off-axis views to be preferred (Blanz, Tarr, & Bülthoff, 1999; Cutzu & Edelman, 1994; Palmer, Rosch, & Chase, 1981). It has been shown that using these criteria leads to relatively consistent results between participants. Sastyin et al. (2015) found a significant interaction between viewpoint and consistency, where the viewpoint effect was weaker when consistent scene information was provided. The authors concluded that object recognition relied more on context information if the object was presented from a noncanonical viewpoint.

These results are an impressive example of how contextual scene information can support object recognition. Here, we investigated if the contextual modulation of viewpoint effects generalizes to strongly noncanonical object orientations. That is, investigate object-scene processing under conditions that produce high uncertainty. This is an important test of the visual system's ability to flexibly rely more on recurrent top-down modulation from scene context when objects are difficult to recognize. In our study, we used three-dimensional (3D) models of objects to create our stimulus set. The use of 3D models to test conditions of object constancy has led to valuable insights such as uncovering the stages of shapeand size-invariant object recognition in the visual system (Isik, Meyers, Leibo, & Poggio, 2014), as well as investigating the features and computational transformations that support 3D object recognition (Biederman & Gerhardstein, 1993; Gauthier et al., 2002; Isik et al., 2014; Logothetis, Pauls, Bülthoff, & Poggio, 1994; Poggio & Edelman, 1990; Ratan Murty & Arun, 2018; Zisserman et al., 1995). In our case, the use of 3D models is motivated by the ability to create a set of highly noncanonical viewpoints in a controlled manner while retaining naturalistic properties, such as the 3D structure of the objects from each viewpoint. Recent work using 3D immersive environments has highlighted the importance of studying vision under more naturalistic constraints in order to investigate cognitive processes in the context of natural behavior (Draschkow, Kallmayer, & Nobre, 2021; Helbing, Draschkow, & Võ, 2020; Helbing, Draschkow, & Võ, 2022; Kristjánsson & Draschkow, 2021).

In the present study, we conducted three behavioral experiments. In our first two experiments (Experiments 1a and 1b), we presented 3D models of real-world objects from six different angles (0°, 60°, 180°, 120°, 240°, and 300°) rotated around the pitch axis in a word–picture verification task. Because rotating the objects around the pitch axis results in highly atypical viewpoints, we expected to find viewpoint-dependent recognition indicated by lower accuracy and slower RTs. In Experiment 1b, we wanted to replicate Experiment 1a with grayscale versions of the images, expecting similar effects of viewpoint as for Experiment 1a (Hayward & Williams, 2000). Experiments 1a and 1b also served to identify viewpoints that produced highest (canonical) and lowest (noncanonical) recognition performance, which we then used in Experiment 2.

In Experiment 2, we paired 3D objects presented in canonical (0° rotation) and noncanonical (120° rotation) viewpoints with semantically consistent and inconsistent scenes. Our aim was to test if viewpoint dependence and object–scene processing effects (Sastyin et al., 2015) generalize to depth rotated 3D models of objects.

General methods

Participants

Participants were recruited at Goethe-University Frankfurt am Main. The sample consisted of 12 participants who completed Experiment 1a (6 women, M age = 23.92 years, range = 19–29 years), 12 different participants who completed Experiment 1b (8 women, M age = 19 years, range = 18–22 years), and another set of 32 participants who completed Experiment 2 (25 women, M age = 24.28 years, range = 18–51 years). The sample size of Experiment 2 was a priori chosen to be higher compared to previous studies which found reliable effects across multiple experiments with 20 participants (e.g., Sastyin et al., 2015). In Experiment 1a, all except for six participants were psychology students who were compensated with course credits; the remaining participants volunteered for the experiment without any compensation. All had normal or corrected-to-normal vision, were native German speakers, and were unfamiliar with the stimulus materials. Written informed consent was obtained before participation, data collection and analyses were carried out according to guidelines approved by the Human Research Ethics Committee of the Goethe University Frankfurt.

Stimulus material

For Experiments 1a and 1b, we collected 100 3D models of objects from a broad range of categories such as furniture, foods, vehicles, plants, and electrical devices. Eighty-two of the 3D models were purchased from CG Axis Complete packages I, II, III, and V, and 18 additional models were obtained free of charge from sources like TurboSquid and free3D. Each model was rotated around its pitch axis by 0°, 60°, 120°, 180°, 240°, and 300° degrees and sized to fit a 60 cm \times 60 cm \times 60 cm box using the free 3D animation software

Blender. Importantly, we chose the most frontal view for the 0° label. Not necessarily because it was the most canonical out of all possible views (usually off-axis views are perceived as more canonical; e.g., Palmer et al., 1981), but because it did not include any additional in-plane rotations or rotations around other cardinal axes. Crucially, it still allowed us to determine the most canonical and noncanonical views out of the chosen set of views. A snapshot from each angle was systematically recorded in front of a gray background using the virtual reality software Vizward5 to create our final stimulus set of 600 images. Additionally, we created grey-scaled versions of these images for Experiment 1b using the GrayscaleEffect function in Vizard5 (https://docs.worldviz.com/vizard/latest/postprocess color.htm).

For Experiment 2, we used the same 3D models as in Experiment 1, adding an additional 56 models collected from the CGAxis packages, resulting in a total of 156 models. Instead of creating snapshots of all six angles, we chose the two viewpoints that had previously produced the highest (canonical viewpoint; 0°) and lowest (noncanonical viewpoint; 120°) recognition performance averaged over Experiments 1a and 1b. We gray-scaled the images using the previously described method.

Additionally, we collected 312 photographic images of scenes, one consistent and one inconsistent scene for each object. We defined a consistent scene as one in which we would expect the object to appear naturally. In both cases, the target object was not present in the scene. Most of the photographs were obtained from the SCEGRAM database (Öhlschläger & Võ, 2017), as well as from Google images. In Experiment 2, objects were presented as templates superimposed on scenes. This was done in line with previous work investigating the influence of object and scene orientation on scene-consistency effects (Lauer et al., 2020; Lauer & Võ, 2022).

All stimuli are available at https://github.com/ aylinsgl/2022-Viewpoint_and_Context.

Procedure

To investigate the speed and accuracy of object recognition, while keeping the procedure comparable with previous studies, a word-picture verification task was used for all experiments (Figure 1). Participants were instructed on screen as well as through standardized verbal instructions to decide as quickly and accurately as possible whether the object on screen matched the basic level category label presented to them at the beginning of the trial using a corresponding match or mismatch key. Participants were not made aware of the different viewpoint conditions beforehand.



Figure 1. Exemplary overview of a subselection of stimuli used in Experiment 1a and the viewpoints used when presenting them (A). Trial procedures for the matching task in Experiments 1a and 1b (B) and Experiment 2 (C). The object was presented in color in Experiment 1a and grayscaled in Experiment 1b. Note that the depicted labels are in English for visualization purposes (German in the original experiment). Participants had to press the "c" key on their keyboard to indicate a match between label and image, and the "m" key to indicate a mismatch.

Each experiment consisted of three practice trials during which the instructor stayed in the room with the participant. More detailed procedure and trial sequences are described in the individual Procedure sections of each experiment. Experiments 1a and 1b lasted approximately 30 minutes, and Experiment 2 lasted approximately 12 minutes.

Design

Experiments 1a and 1b consisted of six blocks with 100 trials each. In each block, the object was presented from a different angle (0°, 60°, 120°, 180°, 240°, or 300°) chosen randomly and counterbalanced between participants. The order of objects within each block was randomized. Each object appeared three times in the match condition (object image matched basic level category label) and three times in the mismatch condition (object image did not match basic level category label), randomized between blocks.

In the mismatch condition, the basic level category label stemmed from a different superordinate category than the object image (e.g., the label "chair" as part of the superordinate category "furniture" was paired with an image of a "car" as part of the superordinate category "vehicle").

Because there was no effect of viewpoint in the mismatch condition in Experiments 1a and 1b, most trials in Experiment 2 were match trials (n = 120), with 23% mismatch trials (n = 36) that were later excluded from analysis. In Experiment 2, each object was presented to each participant once, and we counterbalanced consistency (consistent vs. inconsistent) and viewpoint (canonical vs. noncanonical) between participants.

Data analysis

In Experiments 1a and 1b, we were interested in the effects of viewpoint (how far the object was rotated away from its canonical 0° angle) and match (whether the object matched the basic level category label as part of the experimental design) on reaction times (time between the onset of the object image and keypress response) and accuracy. In Experiment 2, we were interested in the interaction between viewpoint (canonical vs. noncanonical viewpoint), and scene consistency (consistent scene versus inconsistent scene) on reaction times and accuracy.

Raw data were preprocessed and analysed using R (R Core Team, 2021). Objects that produced accuracy ratings that deviated more than 2.5 standard deviations (SD) from the mean (computed for each condition separately) were excluded from analysis. Based on this criterion, we excluded four objects in Experiment 1a, one in Experiment 1b, and two in Experiment 2. We based our reaction time analysis on correctly matched trials only (percent trials removed: Experiment 1a = 4.45%, Experiment 1b = 10.16%, Experiment 2 = 8.55%).

In our data analysis, we used (generalized) linear mixed-effects models ((G)LMMs) using the lme4 package (Bates, Mächler, Bolker, & Walker, 2015). We chose this approach because of its potential advantages over analysis of variance, because it allows us to simultaneously estimate by-participant and by-stimulus variance (Baayen, Davidson, & Bates, 2008; Bates, Mächler, Bolker, & Walker, 2014; Kliegl, Wei, Dambacher, Yan, & Zhou, 2011). The random effects structure of each model was determined using a drop-one procedure starting with the full model including by-participant and by-stimulus varying intercepts and slopes for the main effects in our design. We then subsequently removed random slopes that did not contribute significantly to the goodness of fit as determined by likelihood ratio tests. This strategy allowed us to avoid overparameterization and produce converging models that are supported by the data. Details about the individual analysis and models are described in the Data analyses sections of each experiment. For each GLMM, we report β regression coefficients together with the z statistic and apply a two-tailed 5% error criterion for significance testing. The *p* values for the binary accuracy variable are based on asymptotic Wald tests. Additionally, reaction times were transformed following the Box–Cox procedure (Box & Cox, 1964) to correct for deviation from normality as to better meet LMM assumptions (see

individual Data analysis sections for further details). For the LMMs, regression coefficients are reported with the t-statistic and *p* values were calculated with the ImerTest package (Kuznetsova, Brockhoff, & Christensen, 2017). We defined sum contrasts for match (match vs. mismatch), and consistency (consistent vs. inconsistent) where slope coefficients represent differences between factor levels and the intercept is equal to the grand mean.

We used the ggplot2 package (Wickham, 2016) for graphics and emmeans (Lenth, 2023) for post hoc comparisons. Data and code are openly available at https://github.com/aylinsgl/2022-Viewpoint_and_Context.

Apparatus

All experimental sessions were carried out in the same six experimental cabins of the department of psychology at Goethe-University Frankfurt am Main, containing the same experimental set up (computers running OS Windows 10). Stimulus presentation, RTs and accuracy were systematically controlled and recorded by OpenSesame (Mathôt, Schreij, & Theeuwes, 2012), presented on a 19-in monitor (resolution = $1,680 \times 1,050$, refresh rate = 60 Hz, viewing distance = approximately 65 cm, subtending approximately $11.13^{\circ} \times 9.28^{\circ}$ of visual angle for the object images and approximately $19.0^{\circ} \times 15.84^{\circ}$ of visual angle for the background images).

Experiments 1a and 1b

In Experiments 1a and 1b, we investigated the effect of viewpoint on object recognition RT and accuracy using 3D models of objects rotated around the pitch axis (0°, 60°, 120°, 180°, 240°, and 300°). The only difference between the experiments was that 3D models were presented either in color (Experiment 1a) or a grayscale version of the model was used (Experiment 1b). Participants had to indicate whether the object matched the previously presented basic level category label.

Procedure

Participants were presented with a fixation point in the middle of the screen followed by a basic level object category label (in German, font. Droid Sans Mono; font size. 26; color. black). This presentation was followed by the target object presented in the middle of the screen, which could either match or mismatch the label, until the participant gave



Figure 2. Partial effect plots of the interactions of viewpoint (0°, 60°, 120°, 180° 240°, and 300°) and match (match vs. mismatch) on accuracy for Experiment 1a (colored; **A**), and Experiment 1b (grayscaled; **C**), and the effect of viewpoint on RT for Experiment.

a response (Figure 1A). Participants were given feedback on screen if their answer was incorrect. The next trial automatically started with a new fixation point.

Data analysis

After data preprocessing, we used a binomial GLMM to examine the effects of viewpoint and match on accuracy. As fixed effects we included viewpoint (0°, 60°, 120°, 180°, 240°, or 300°) as a first- and second-degree polynomial, the match versus mismatch comparison, and the interactions between these terms. The second-degree polynomial viewpoint term was added, because we expected viewpoint to affect recognition in a nonlinear manner (symmetry around 180°). Our final model included random intercepts for participants and stimuli, as well as a by-stimuli random slope for the match versus mismatch effect for Experiment 1a, and random intercepts for participants and stimuli, as well as a by-stimuli and by-participant random slope for the match effect for Experiment 1b.

Based on the power coefficient output of the Box– Cox procedure ($\lambda = 0.22$), RTs were log transformed. We used the same fixed effects structure for the RT–LMMs as for the accuracy–GLMMs. As random effects, we entered random intercepts for participants and stimuli, as well as by-participant and by-stimuli random slopes for the effect of match for Experiments 1a and 1b.

Results

Accuracy

The average accuracy in Experiment 1a was quite high (M = 0.95, SD = 0.21) and slightly lower in Experiment 1b (M = 0.9, SD = 0.3). In line with our hypothesis, the GLMM yielded a significant main effect for the second-degree polynomial viewpoint term in both experiments (Experiment 1a: $\beta = 16.67$, SD =5.61, z = 2.97, p = 0.003; Experiment 1b: $\beta = 18.82$, SE = 3.79, z = 4.97, p < 0.001), meaning that the effect of viewpoint on accuracy can be well-described by a quadratic function (Figure 2A and 2C). There was also a significant interaction between the second-degree polynomial of viewpoint and the match condition in both experiments, Experiment 1a: $\beta = 23.62$, SE =5.69, z = 4.15, p < 0.001; Experiment 1b: $\beta = 15.23$, SE = 3.82, z = 3.98, p < 0.001. Comparing the viewpoint trend for the match and mismatch conditions, we found

that the second-degree viewpoint trend was significant in the match condition (Experiment 1a: $\beta = 0.19$, SE = 0.03, 95% CI = [0.13–0.25]; Experiment 1b: $\beta =$ 0.16, SE = 0.02, 95% CI = [0.12–0.21), but not in the mismatch condition, Experiment 1a: $\beta = -0.03$, SE = 0.04, 95% CI = [-0.12 to 0.05]; Experiment 1b: $\beta =$ -0.02, SE = 0.03, 95% CI = [-0.04 to 0.07]. A detailed overview of performance for each object and viewpoint is provided in the Appendix (Figure A1).

RTs

Participants were slightly faster on average in Experiment 1b (M = 685 ms, SD = 358 ms) than Experiment 1a (M = 738 ms, SD = 299 ms). In line with our hypothesis, the LMM revealed a significant main effect for the second-degree polynomial viewpoint term in both experiments: Experiment 1a: $\beta = -2.2$, SE =0.29, t = -7.48, p < 0.001; Experiment 1b: $\beta = -1.42$, SE = 0.29, t = -4.99, p < 0.001 (Figure 2B and 2D). In both experiments, there was no interaction between viewpoint and match, Experiment 1a: $\beta = -0.12$, SE =0.29, t = -0.4, p = 0.69; Experiment 1b: $\beta = -0.38$, SE =0.29, t = -1.34, p = 0.18.

Discussion

In Experiment 1a, we found viewpoint-dependent object recognition for objects rotated around the pitch axis. This effect can best be described by a quadratic curve that approximates symmetry around 120° rotation. We also found that in our sequential matching task, only the match condition produced viewpoint-dependent behavior, whereas mismatch trials seemed unaffected by viewpoint. Finding a mismatch might rely more on the analysis of global, viewpoint-invariant features, whereas matching might be more dependent on the analysis of local, viewpoint-dependent features (e.g., Jolicoeur, 1990a) (e.g., deciding a shape is not a car might require less viewpoint-dependent information than identifying the shape as a chair). In Experiment 1b, we were able to replicate our results from Experiment 1a. Grayscaling the images seemed to have made the overall task slightly more difficult, while still producing similarly viewpoint-dependent behavior. Although some studies report mirror confusion effects for rotations around 180° (e.g., Gregory & McCloskey, 2010), we did not encounter this phenomenon in our study. In our case, rotating around the pitch axis produced views such as "upside-down, from behind" which is untypical for images that usually produce mirror confusions. The canonical (0°) and noncanonical (120°) viewpoints we used in Experiment 2 represented viewpoints that produced the best and worst recognition performance derived from average accuracy ratings obtained from Experiments 1a and 1b.

Experiment 2

In Experiment 2, we paired canonical (0°) and noncanonical (120°) viewpoints with consistent and inconsistent scene contexts. We were specifically interested in the interaction between viewpoint and consistency, with the expectation that meaningful scene context information would decrease the effect of viewpoint on object recognition.

Procedure

In Experiment 2, we used the same word-picture verification task as in Experiments 1a and 1b (Figure 1B). Scene context was provided by first previewing the consistent or inconsistent scene for 300 ms and then overlaying the target object on top of the scene background until a response was given.

Data analysis

For both the accuracy–GLMM and RTs LMM, we entered interaction terms between viewpoint and consistency as fixed effects. The GLMM included random intercepts for participants and stimuli, as well as a by-stimuli random slope for the effect of viewpoint. RT data were log transformed.

For the RT-LMM, we had random intercepts for participants and stimuli, and a by-participant random slope for the effect of viewpoint and by-stimuli random slopes for the effects of viewpoint and consistency.

Results

Accuracy

Accuracy was significantly higher for canonical viewpoints than for noncanonical viewpoints as revealed by the GLMM ($\beta = 0.68$, SE = 0.14, z = 4.82, p < 0.001), but there was no significant main effect for consistency ($\beta = 0.06$, SE = 0.07, z = 0.75, p = 0.45). Critically, there was a significant interaction between viewpoint and consistency ($\beta = -0.21$, SE = 0.07, z = -2.84, p = 0.004) (Figure 3A). Post hoc interaction contrasts revealed that the viewpoint-dependence effect was significantly stronger in the inconsistent scene condition compared to the consistent scene condition $(\beta = -0.84, SE = 0.3, z = -2.84, p = 0.005)$. This finding is in line with our hypothesis that providing meaningful scene context can decrease the effects of viewpoint on object recognition. Additionally, the scene-consistency effect was only significant in the noncanonical condition ($\beta = 0.53$, SE = 0.15, z = 3.45,



Figure 3. Experiment 2 accuracy difference scores per participant (canonical vs. noncanonical) for consistent and inconsistent scene backgrounds (**A**). Adjusted response times (**B**) were obtained with the remef package (Hohenstein & Kliegl, 2023). *p < 0.05. ***p < 0.001.

p < 0.001), but not in the canonical condition ($\beta = -0.31$, SE = 0.25, z = -1.22, p = 0.22).

RTs

The LMM yielded a significant main effect for viewpoint ($\beta = -0.07$, SE = 0.01, t = -7.26, p < 0.001), where RTs were faster for canonical (M = 558 ms, SD = 255 ms) than for noncanonical viewpoints (M = 645 ms, SD = 333 ms) (Figure 3B). There was no significant interaction between viewpoint and consistency ($\beta = 0.004$, SE = 0.005, t = 0.83, p = 0.41).

Discussion

In general, object recognition accuracy was viewpoint dependent; however, there was a significant interaction between viewpoint and consistency. In line with our hypothesis, the viewpoint effect was significantly weaker for consistent scenes and the scene consistency effect was only observed for noncanonical viewpoints (Figure 3A). Noncanonical viewpoints were recognized significantly slower than canonical viewpoints. However, this result was unaffected by scene consistency.

General discussion

In the present study, we investigated how scene context information modulates viewpoint-dependent object recognition under conditions of high uncertainty using 3D models of everyday objects. Although providing meaningful context did not eradicate the viewpoint effect fully, it significantly decreased recognition accuracy costs. By extending previous findings (Sastyin et al., 2015), we provide further support for a model of object recognition that incorporates context (e.g., Bar, 2004), while dynamically adapting to the amount of available information based not only on visual features of the object (Burgund & Marsolek, 2000; Hayward & Tarr, 1997; Jolicoeur, 1990), but also context.

It is assumed that, when objects are presented in context, rapidly accessed low spatial frequency information is fed back to the occipito-temporal cortex facilitating high spatial frequency based analysis during object recognition (Bar, 2004; Kauffmann, Ramanoël, & Peyrin, 2014; Peyrin, Chauvin, Chokron, & Marendaz, 2003; Peyrin, Baciu, Segebarth, & Marendaz, 2004). The highly noncanonical viewpoints we used in our experiments produce high uncertainty in the initial set of possible target objects. We show that, under conditions where low spatial frequency analysis of the object leads to very ambiguous target candidates, the visual system relies more on top–down regulation modulated by recurrent processing of low spatial frequency information from the scene (Bar, 2004).

It further motivates models of object constancy—the visual system's ability to produce representations that are robust to changes in, for example, viewpoint or lighting (e.g., DiCarlo & Cox, 2007)—that efficiently integrate contextual information and can lead to both viewpoint-dependent and invariant behavior based on available information and the task at hand.

A key component of the present study was to generalize previous findings on object-scene processing effects and viewpoint dependence to depth-rotated 3D objects. We want to highlight the importance of generalizing findings from traditional two-dimensional settings to more naturalistic settings and stimuli. Kristjánsson and Draschkow (2021) have shown very illustratively for a variety of phenomena that, given more naturalistic constraints, a system is able to circumvent, for example capacity limits by drawing on the rich visual experience of natural environments. Although we did not use fully immersive environments, using 3D models offers a more realistic encounter of everyday objects and, therefore, a more precise measure of viewpoint dependence in real-world object recognition. It should be noted, however, that there is a trade-off between naturalistic looking stimuli (i.e., photographs) and stimuli that more precisely capture naturalistic properties (i.e., 3D structure of objects from different viewpoints) in a highly controlled manner, while not looking as naturalistic. Here, we opted for providing more naturalistic 3D properties of the displayed objects.

From the present study, it is unclear what kind of information contained in the scenes was responsible for decreasing the viewpoint costs. Rapidly accessed global information such as the gist of the scene (Oliva & Torralba, 2007) could be the main factor. At the same time, more local information such as the detection and recognition of certain objects in the scene preview could provide information about related possible target objects based on internalized scene–object and object–object regularities (Võ et al., 2019). Revealing the time course of when what kind of contextual information is integrated to buffer viewpoint effects would provide new insights into how the visual system so effortlessly achieves invariant object recognition.

Varying what information is presented during the task (i.e., providing meaningful context vs. showing objects in isolation) is one way to probe the visual system's ability to overcome processing limitations in viewpoint-dependent object recognition. Alternatively, one could keep the visual input constant, but vary the level at which participants have to perform the matching task (Hamm & McMullen, 1998). If there are object representations that contain more or less viewpoint-dependent or invariant information, how does this factor interact with the integration of contextual information in the form of scene context?

Finally, we would like to address that, on average, performance was high in the matching task throughout all our experiments. These ceiling effects are probably due to the type of task we chose, which are different from the tasks usually used to study scene consistency effects (Davenport & Potter, 2004; Sastyin et al., 2015). Despite these differences in difficulty, we were able to demonstrate a significant decrease in viewpoint costs by providing meaningful scene context.

Past research has made strong advances toward understanding the computations that underly invariant object recognition (DiCarlo & Cox, 2007). Understanding these mechanisms in isolation is key to understanding object recognition in general. We argue that understanding how the visual system is able to make use of richly structured naturalistic environments to circumvent computational bottlenecks will ultimately lead to better, more robust models of object recognition and inspire approaches in fields such as computer vision (e.g., Bomatter et al., 2021).

To conclude, in the present study we show that scene context supports object recognition, even when using extremely noncanonical orientations of depth rotated 3D objects. We highlight the importance of testing capacity limits of object recognition in more naturalistic frameworks to build more robust and flexible models and move toward a better understanding of vision under naturalistic constraints.

Keywords: object recognition, viewpoint dependence, scene context effects

Acknowledgments

Supported by SFB/TRR 26 135 project C7 to Melissa L.-H. Võ and the Hessisches Ministerium für Wissenschaft und Kunst (HMWK; project 'The Adaptive Mind') and the Main-Campus-Doctus stipend awarded by the Stiftung Polytechnische Gesellschaft to Aylin Kallmayer.

The Wellcome Centre for Integrative Neuroimaging is supported by core funding from the Wellcome Trust (203139/Z/16/Z). The work is supported by the NIHR Oxford Health Biomedical Research Centre. The funders had no role in the decision to publish or in the preparation of the manuscript.

Commercial relationships: none. Corresponding author: Aylin Kallmayer. Email: kallmayer@psych.uni-frankfurt.de. Address: Department of Psychology, Goethe University Frankfurt, Theodor-W.-Adorno-Platz 6, Frankfurt am Main 60323, Germany.

References

- Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal* of Memory and Language, 59(4), 390–412, https://doi.org/10.1016/j.jml.2007.12.005.
- Bar, M. (2004). Visual objects in context. *Nature Reviews Neuroscience*, 5(8), Article 8, https: //doi.org/10.1038/nrn1476.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2014). *Fitting Linear Mixed-Effects*

Models using lme4 (arXiv:1406.5823). arXiv, https://doi.org/10.48550/arXiv.1406.5823.

Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67, 1–48, https://doi.org/10.18637/jss.v067.i01.

Biederman, I., & Gerhardstein, P. C. (1993). Recognizing depth-rotated objects: Evidence and conditions for three-dimensional viewpoint invariance. *Journal of Experimental Psychology: Human Perception and Performance*, 19(6), 1162– 1182, https://doi.org/10.1037/0096-1523.19.6.1162.

Biederman, I., Mezzanotte, R. J., & Rabinowitz, J. C. (1982). Scene perception: Detecting and judging objects undergoing relational violations. *Cognitive Psychology*, 14(2), 143–177, https://doi.org/10.1016/0010-0285(82)90007-X.

Blanz, V., Tarr, M. J., & Bülthoff, H. H. (1999). What object attributes determine canonical views? *Perception*, 28(5), 575–599, https: //doi.org/10.1068/p2897.

Bomatter, P., Zhang, M., Karev, D., Madan, S., Tseng, C., & Kreiman, G. (2021). When pigs fly: Contextual reasoning in synthetic and natural scenes. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2021*, 255–264, https://openaccess.thecvf.com/content/ ICCV2021/html/Bomatter_When_Pigs_Fly_ Contextual_Reasoning_in_Synthetic_and_ Natural_Scenes_ICCV_2021_paper.html.

Box, G. E. P., & Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society: Series B (Methodological), 26*(2), 211–243, https://doi.org/10.1111/j.2517-6161.1964.tb00553.x.

Brandman, T., & Peelen, M. V. (2017). Interaction between scene and object processing revealed by human fMRI and MEG decoding. *Journal* of Neuroscience, 37(32), 7700–7710, https: //doi.org/10.1523/JNEUROSCI.0582-17.2017.

Bülthoff, H. H., & Edelman, S. (1992). Psychophysical support for a two-dimensional view interpolation theory of object recognition. *Proceedings* of the National Academy of Sciences of the United States of America, 89(1), 60–64, https://doi.org/10.1073/pnas.89.1.60.

Burgund, E. D., & Marsolek, C. J. (2000). Viewpointinvariant and viewpoint-dependent object recognition in dissociable neural subsystems. *Psychonomic Bulletin & Review*, 7(3), 480–489, https://doi.org/10.3758/BF03214360.

Charles Leek, E., & Johnston, S. J. (2006). A polarity effect in misoriented object recognition: The role of polar features in the computation of orientation-invariant shape representations. *Visual Cognition*, *13*(5), 573–600, https://doi.org/10.1080/13506280544000048.

Cutzu, F., & Edelman, S. (1994). Canonical views in object representation and recognition. *Vision Research*, 34(22), 3037–3056, https: //doi.org/10.1016/0042-6989(94)90277-1.

Davenport, J. L., & Potter, M. C. (2004). Scene consistency in object and background perception. *Psychological Science*, 15(8), 559–564, https: //doi.org/10.1111/j.0956-7976.2004.00719.x.

DiCarlo, J. J., & Cox, D. D. (2007). Untangling invariant object recognition. *Trends in Cognitive Sciences*, 11(8), 333–341, https: //doi.org/10.1016/j.tics.2007.06.010.

Draschkow, D., Kallmayer, M., & Nobre, A. C. (2021). When natural behavior engages working memory. *Current Biology*, *31*(4), 869–874.e5, https://doi.org/10.1016/j.cub.2020.11.013.

Draschkow, D., & Võ, M. L.-H. (2017). Scene grammar shapes the way we interact with objects, strengthens memories, and speeds search. *Scientific Reports*, 7(1), Article 1, https://doi.org/10.1038/s41598-017-16739-x.

Edelman, S. (1995). Class similarity and viewpoint invariance in the recognition of 3D objects. *Biological Cybernetics*, 72(3), 207–220, https://doi.org/10.1007/BF00201485.

Foster, D. H., & Gilson, S. J. (2002). Recognizing novel three–dimensional objects by summing signals from parts and views. *Proceedings of the Royal Society* of London. Series B: Biological Sciences, 269(1503), 1939–1947, https://doi.org/10.1098/rspb.2002.2119.

Gauthier, I., Hayward, W. G., Tarr, M. J., Anderson, A. W., Skudlarski, P., & Gore, J. C. (2002). BOLD activity during mental rotation and viewpoint-dependent object recognition. *Neuron*, *34*(1), 161–171, https: //doi.org/10.1016/S0896-6273(02)00622-0.

Graf, M. (2006). Coordinate transformations in object recognition. *Psychological Bulletin*, *132*(6), 920–945, https://doi.org/10.1037/0033-2909.132.6.920.

Gregory, E., & McCloskey, M. (2010). Mirror-image confusions: Implications for representation and processing of object orientation. *Cognition*, 116(1), 110–129, https://doi.org/10.1016/j.cognition.2010. 04.005.

Hamm, J. P., & McMullen, P. A. (1998). Effects of orientation on the identification of rotated objects depend on the level of identity. *Journal of Experimental Psychology: Human Perception and Performance, 24*(2), 413–426, https://doi.org/10.1037/0096-1523.24.2.413.

Hayward, W. G. (2003). After the viewpoint debate: Where next in object recognition?

Trends in Cognitive Sciences, 7(10), 425–427, https://doi.org/10.1016/j.tics.2003.08.004.

- Hayward, W. G., & Tarr, M. J. (1997). Testing conditions for viewpoint invariance in object recognition. *Journal of Experimental Psychology: Human Perception and Performance*, 23(5), 1511–1521, https://doi.org/10.1037/0096-1523.23.5. 1511.
- Hayward, W. G., & Williams, P. (2000). Viewpoint dependence and object discriminability. *Psychological Science*, 11(1), 7–12, https://doi.org/ 10.1111/1467-9280.00207.
- Helbing, J., Draschkow, D., & Võ, M. L.-H. (2020). Search superiority: Goal-directed attentional allocation creates more reliable incidental identity and location memory than explicit encoding in naturalistic virtual environments. *Cognition*, 196, 104147, https://doi.org/10.1016/j.cognition.2019.104147.
- Helbing, J., Draschkow, D., & Võ, M. L.-H. (2022). Auxiliary Scene-Context Information Provided by Anchor Objects Guides Attention and Locomotion in Natural Search Behavior. *Psychological Science*, 33(9), 1463–1476, https://doi.org/10.1177/09567976221091838.
- Hohenstein, S, & Kliegl, R. (2023). remef: Remove partial effects. R package version 1.0.7, https://github.com/hohenstein/remef/.
- Isik, L., Meyers, E. M., Leibo, J. Z., & Poggio, T. (2014). The dynamics of invariant object recognition in the human visual system. *Journal of Neurophysiology*, 111(1), 91–102, https://doi.org/10.1152/jn.00394.2013.
- Jolicoeur, P. (1990). Identification of disoriented objects: A dual-systems theory. *Mind & Language*, 5(4), 387–410, https://doi.org/10.1111/j.1468-0017. 1990.tb00170.x.
- Josephs, E. L., Draschkow, D., Wolfe, J. M., & Võ, M. L.-H. (2016). Gist in time: Scene semantics and structure enhance recall of searched objects. *Acta Psychologica*, 169, 100–108, https://doi.org/10.1016/j.actpsy.2016.05.013.
- Kauffmann, L., Ramanoël, S., & Peyrin, C. (2014). The neural bases of spatial frequency processing during scene perception. *Frontiers in Integrative Neuroscience*, 8, 27, https://www.frontiersin.org/ articles/10.3389/fnint.2014.00037.
- Kliegl, R., Wei, P., Dambacher, M., Yan, M., & Zhou, X. (2011). Experimental effects and individual differences in linear mixed models: Estimating the relationship between spatial, object, and attraction effects in visual attention. *Frontiers in Psychology*, *1*, 238, https://www.frontiersin.org/article/10.3389/ fpsyg.2010.00238.

- Kristjánsson, Á., & Draschkow, D. (2021). Keeping it real: Looking beyond capacity limits in visual cognition. *Attention, Perception, & Psychophysics, 83*(4), 1375–1390, https://doi.org/ 10.3758/s13414-021-02256-7.
- Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). ImerTest package: Tests in linear mixed effects models. *Journal of Statistical Software*, 82, 1–26, https://doi.org/10.18637/jss.v082.i13.
- Lauer, T., Schmidt, F., & Võ, M. L.-H. (2021). The role of contextual materials in object recognition. *Scientific Reports, 11*(1), Article 1, https://doi.org/10.1038/s41598-021-01406-z.
- Lauer, T., & Võ, M. L.-H. (2022). The ingredients of scenes that affect object search and perception.
 In B. Ionescu, W. A. Bainbridge, & N. Murray (Eds.), *Human perception of visual information: psychological and computational perspectives* (pp. 1–32). New York: Springer International Publishing, https://doi.org/10.1007/978-3-030-81465-6_1.
- Lauer, T., Willenbockel, V., Maffongelli, L., & Võ, M. L.-H. (2020). The influence of scene and object orientation on the scene consistency effect. *Behavioural Brain Research*, 394, 112812, https://doi.org/10.1016/j.bbr.2020.112812.
- Leek, E. C., Atherton, C. J., & Thierry, G. (2007). Computational mechanisms of object constancy for visual recognition revealed by event-related potentials. *Vision Research*, 47(5), 706–713, https://doi.org/10.1016/j.visres.2006.10.021.
- Lenth, R. (2023). emmeans: Estimated marginal means, aka least-squares means. R package version 1.8.7, https://CRAN.R-project.org/package=emmeans.
- Logothetis, N. K., Pauls, J., Bülthoff, H. H., & Poggio, T. (1994). View-dependent object recognition by monkeys. *Current Biology*, 4(5), 401–414, https://doi.org/10.1016/S0960-9822(00)00089-0.
- Lowe, D. G. (1987). Three-dimensional object recognition from single two-dimensional images. *Artificial Intelligence*, *31*(3), 355–395, https://doi.org/10.1016/0004-3702(87)90070-1.
- Marr, D., Nishihara, H. K., & Brenner, S. (1978). Representation and recognition of the spatial organization of three-dimensional shapes. *Proceedings of the Royal Society of London. Series B. Biological Sciences, 200*(1140), 269–294, https://doi.org/10.1098/rspb.1978.0020.
- Mathôt, S., Schreij, D., & Theeuwes, J. (2012). OpenSesame: An open-source, graphical experiment builder for the social sciences. *Behavior Research Methods*, 44(2), 314–324, https://doi.org/10.3758/s13428-011-0168-7.
- Öhlschläger, S., & Võ, M. L.-H. (2017). SCEGRAM: An image database for semantic and syntactic

inconsistencies in scenes. *Behavior Research Methods*, 49(5), 1780–1791, https://doi.org/10.3758/ s13428-016-0820-3.

- Oliva, A., & Torralba, A. (2007). The role of context in object recognition. *Trends in Cognitive Sciences*, 11(12), 520–527, https://doi.org/10.1016/j.tics.2007.09.009.
- Palmer, T. E. (1975). The effects of contextual scenes on the identification of objects. *Memory & Cognition*, 3(5), 519–526, https://doi.org/10.3758/BF0319 7524.
- Palmer, S. E., Rosch, E., & Chase, P. (1981). Canonical Perspective and the perception of objects. *Attention* and *Performance*, 9, 135–151.
- Peyrin, C., Baciu, M., Segebarth, C., & Marendaz, C. (2004). Cerebral regions and hemispheric specialization for processing spatial frequencies during natural scene recognition. An event-related fMRI study. *NeuroImage*, 23(2), 698–707, https://doi.org/10.1016/j.neuroimage.2004.06.020.
- Peyrin, C., Chauvin, A., Chokron, S., & Marendaz, C. (2003). Hemispheric specialization for spatial frequency processing in the analysis of natural scenes. *Brain and Cognition*, 53(2), 278–282, https://doi.org/10.1016/S0278-2626(03)00126-X.
- Poggio, T., & Edelman, S. (1990). A network that learns to recognize three-dimensional objects. *Nature, 343*(6255), Article 6255, https://doi.org/10.1038/343263a0.
- Ratan Murty, N. A., & Arun, S. P. (2015). Dynamics of 3D view invariance in monkey inferotemporal cortex. *Journal of Neurophysiology*, 113(7), 2180–2194, https://doi.org/10.1152/jn.00810.2014.
- Ratan Murty, N. A., & Arun, S. P. (2018). Multiplicative mixing of object identity and image attributes in single inferior temporal neurons. *Proceedings* of the National Academy of Sciences of the United States of America, 115(14), E3276–E3285, https://doi.org/10.1073/pnas.1714287115.
- R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria, https://www.R-project.org/.
- Sastyin, G., Niimi, R., & Yokosawa, K. (2015). Does object view influence the scene consistency effect? *Attention, Perception, & Psychophysics*, 77(3), 856– 866, https://doi.org/10.3758/s13414-014-0817-x.
- Stankiewicz, B. J. (2002). Empirical evidence for independent dimensions in the visual representation of three-dimensional shape. *Journal of Experimental Psychology: Human Perception and Performance, 28*(4), 913–932, https://doi.org/10.1037/0096-1523.28.4.913.

- Tarr, M. J., & Bülthoff, H. H. (1995). Is human object recognition better described by geon structural descriptions or by multiple views? Comment on Biederman and Gerhardstein (1993). *Journal of Experimental Psychology: Human Perception and Performance, 21*(6), 1494–1505, https://doi.org/10.1037/0096-1523.21.6.1494.
- Tarr, M. J., & Pinker, S. (1989). Mental rotation and orientation-dependence in shape recognition. *Cognitive Psychology*, 21(2), 233–282, https: //doi.org/10.1016/0010-0285(89)90009-1.
- Vanrie, J., Béatse, E., Wagemans, J., Sunaert, S., & Van Hecke, P. (2002). Mental rotation versus invariant features in object perception from different viewpoints: An fMRI study. *Neuropsychologia*, 40(7), 917–930, https: //doi.org/10.1016/S0028-3932(01)00161-0.
- Võ, M. L.-H. (2021). The meaning and structure of scenes. *Vision Research*, *181*, 10–20, https://doi.org/10.1016/j.visres.2020.11.003.
- Võ, M. L.-H., Boettcher, S. E., & Draschkow, D. (2019). Reading scenes: How scene grammar guides attention and aids perception in real-world environments. *Current Opinion in Psychology, 29*, 205–210, https://doi.org/10.1016/j.copsyc.2019.03. 009.
- Võ, M. L.-H., & Henderson, J. M. (2009). Does gravity matter? Effects of semantic and syntactic inconsistencies on the allocation of attention during scene perception. *Journal of Vision*, 9(3), 24, https://doi.org/10.1167/9.3.24.
- Võ, M. L.-H., & Wolfe, J. M. (2013a). Differential electrophysiological signatures of semantic and syntactic scene processing. *Psychological Science*, 24(9), 1816–1823, https://doi.org/10.1177/ 0956797613476955.
- Võ, M. L.-H., & Wolfe, J. M. (2013b). The interplay of episodic and semantic memory in guiding repeated search in scenes. *Cognition*, 126(2), 198–212, https://doi.org/10.1016/j.cognition.2012.09.017.
- Wickham, H. (2016). ggplot2: Elegant graphics for data analysis. New York: Springer-Verlag, https://ggplot2.tidyverse.org.
- Wilson, K. D., & Farah, M. J. (2003). When does the visual system use viewpoint-invariant representations during recognition? *Cognitive Brain Research*, 16(3), 399–415, https://doi.org/10.1016/ S0926-6410(03)00054-5.
- Zisserman, A., Forsyth, D., Mundy, J., Rothwell, C., Liu, J., & Pillow, N. (1995). 3D object recognition using invariance. *Artificial Intelligence*, 78(1), 239–288, https://doi.org/10.1016/0004-3702(95) 00023-2.

Appendix A



Figure A1. Accuracy for Experiment 1b (A). Each line represents average performance for each viewpoint of each object. Some objects were highlighted to demonstrate the range of performance differences between viewpoints for different objects (B).