# RAPID REPORT | *Sensory Processing*

# Real-world structure facilitates the rapid emergence of scene category information in visual brain signals

● **Daniel Kaiser,**[1] **Greta Häberle,**[2,3,4] **and Radoslaw M. Cichy**[2,3,4,5]

[1]*Department of Psychology, University of York, York, United Kingdom;* [2]*Department of Education and Psychology, Freie Universität Berlin, Berlin, Germany;* [3]*Charité — Universitätsmedizin Berlin, Einstein Center for Neurosciences Berlin, Berlin, Germany;* [4]*Berlin School of Mind and Brain, Humboldt-Universität zu Berlin, Berlin, Germany; and* [5]*Bernstein Center for Computational Neuroscience Berlin, Berlin, Germany*

Kaiser D, Häberle G, Cichy RM. Real-world structure facilitates the rapid emergence of scene category information in visual brain signals. *J Neurophysiol* 124: 145–151, 2020. First published June 10, 2020; doi:10.1152/jn.00164.2020.—In everyday life, our visual surroundings are not arranged randomly but structured in predictable ways. Although previous studies have shown that the visual system is sensitive to such structural regularities, it remains unclear whether the presence of an intact structure in a scene also facilitates the cortical analysis of the scene's categorical content. To address this question, we conducted an EEG experiment during which participants viewed natural scene images that were either "intact" (with their quadrants arranged in typical positions) or "jumbled" (with their quadrants arranged into atypical positions). We then used multivariate pattern analysis to decode the scenes' category from the EEG signals (e.g., whether the participant had seen a church or a supermarket). The category of intact scenes could be decoded rapidly within the first 100 ms of visual processing. Critically, within 200 ms of processing, category decoding was more pronounced for the intact scenes compared with the jumbled scenes, suggesting that the presence of real-world structure facilitates the extraction of scene category information. No such effect was found when the scenes were presented upside down, indicating that the facilitation of neural category information is indeed linked to a scene's adherence to typical real-world structure rather than to differences in visual features between intact and jumbled scenes. Our results demonstrate that early stages of categorical analysis in the visual system exhibit tuning to the structure of the world that may facilitate the rapid extraction of behaviorally relevant information from rich natural environments.

**NEW & NOTEWORTHY** Natural scenes are structured, with different types of information appearing in predictable locations. Here, we use EEG decoding to show that the visual brain uses this structure to efficiently analyze scene content. During early visual processing, the category of a scene (e.g., a church vs. a supermarket) could be more accurately decoded from EEG signals when the scene adhered to its typical spatial structure compared with when it did not.

EEG; multivariate pattern analysis; real-world structure; scene representation; visual processing

## INTRODUCTION

IN EVERYDAY SITUATIONS, the input to our visual system is not random; rather, it rather arises from highly organized scenes,

Correspondence: D. Kaiser (danielkaiser.net@gmail.com).

which follow a predictable structure. In practically every real-word scene, visual information (such as the scene's layout properties or the objects contained in a scene) is distributed in meaningful ways across space (Bar 2004; Kaiser et al. 2019a; Oliva and Torralba 2007; Võ et al. 2019; Wolfe et al. 2011). Neuroimaging studies have shown that the visual system is sensitive to this structure, with cortical responses differing when scene elements do or do not adhere to typical real-world structure (Abassi and Papeo 2020; Baldassano et al. 2017; Bilalić et al. 2019; Kaiser et al. 2014; Kaiser and Peelen 2018; Kim and Biederman 2011; Roberts and Humphreys 2010). Although such studies suggest that the presence of real-world structure aids efficient scene representation, it is unclear how real-world structure impacts the representation of scene content. Specifically, does the presence of real-world structure facilitate the extraction of categorical information from a scene?

Evidence for an increase of visual category information in the presence of real-world regularities has already been reported for individual object processing. Several studies showed that typical real-world positioning enhances the neural representation of object category (Chan et al. 2010; de Haas et al. 2016; Kaiser and Cichy 2018; Kaiser et al. 2018); for example, neural responses to an airplane are better discriminable from responses to other objects when the airplane is shown in the upper visual field, where it is typically encountered in the real world. Does the presence of real-world structure similarly facilitate the representation of categorical scene content in scenes?

To address this question, we used a jumbling paradigm (Biederman 1972; Biederman et al. 1974) that manipulates natural scenes' spatial structure. Individual parts of the scene could either appear in their typical, "intact" positions or in atypical, "jumbled" positions (Fig. 1). In a recent neuroimaging study (Kaiser et al. 2020a), we employed this paradigm to show that in scene-selective visual cortex (fMRI) and after 250 ms of vision (EEG), spatially intact scenes were represented differently from jumbled scenes. Here, we analyzed the EEG data from this jumbling paradigm to investigate whether the typical real-world structure, in contrast to an atypical structure, facilitates the visual representation of scene category.
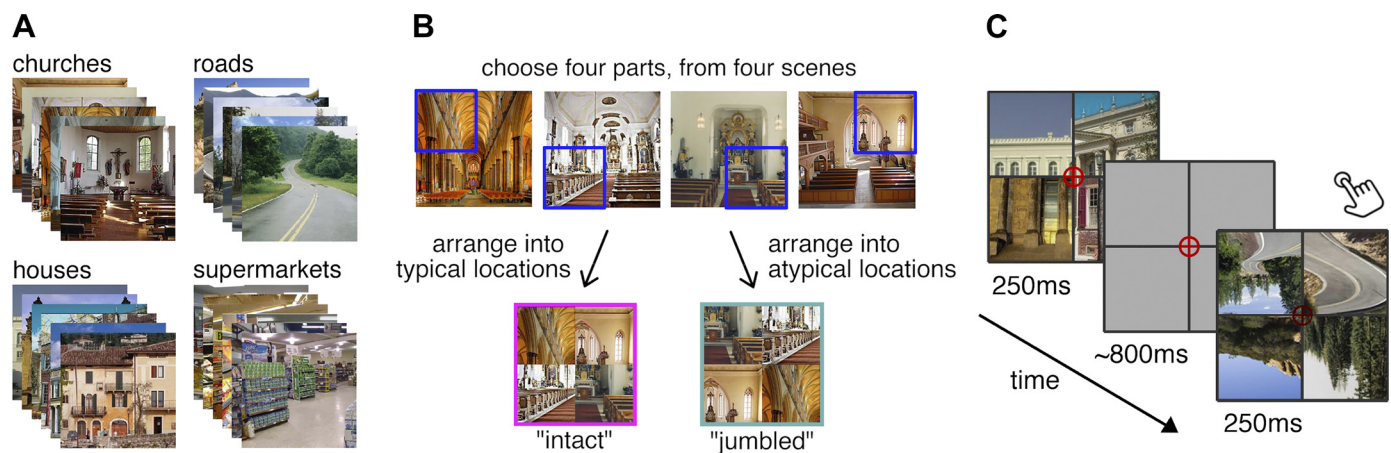
Fig. 1. Experimental design. *A*: stimulus set was constructed from natural scene photographs of 4 categories. *B*: intact and jumbled scenes were created by combining parts of 4 different scenes of the same category in either typical locations or in atypical locations (with positions swapped in a crisscrossed way). *C*: during the EEG experiment, participants viewed the scenes in upright and inverted orientation for 250 ms each, in random order. Participants performed an orthogonal task, where they responded whenever the fixation cross darkened.

To extract differences in category information between intact and jumbled scenes with high sensitivity, we used a cumulative multivariate decoding approach (Ramkumar et al. 2013), which maximizes the amount of data available at every time point along the processing cascade. In line with previous reports (Dima et al. 2018; Kaiser et al. 2019b, 2020b; Lowe et al. 2018), this analysis showed that scene category information emerges rapidly (within the first 100 ms of vision). Critically, the early emergence of scene category information was facilitated for intact compared with jumbled scenes. This benefit was only present for upright but not inverted scenes, indicating that the early facilitation of scene analysis is related to the presence of real-world structure rather than differences in basic visual features.

### MATERIALS AND METHODS

*Participants.* Twenty healthy adults (mean age 26.6 yr, SD = 5.8; 9 female) participated. All participants had normal or corrected-to-normal vision. Participants provided written informed consent and received either monetary reimbursement or course credits. All procedures were approved by the ethical committee of the Department of Psychology at Freie Universität Berlin and were in accordance with the Declaration of Helsinki.

*Stimuli.* Stimuli were scenes from four different categories: churches, houses, roads, and supermarkets (Fig. 1*A*). The stimuli were taken from an online resource (Konkle et al. 2010). For each category, six different exemplars were used. To manipulate scenes' adherence to real-world structure, we first split each original image into quadrants. We then systematically recombined parts (quadrants) from different scenes such that the scenes' spatial structure was either intact or jumbled (Fig. 1*B*). For the intact scenes, four parts from four different scenes of the same scene category were combined in their correct spatial locations. For the jumbled scenes, four parts from four different scenes of the same scene category were combined, but their spatial locations were arranged in a crisscrossed way. This jumbling manipulation simultaneously disrupted multiple structural regularities in the scene, such as visual feature distributions, scene geometry, absolute and relative object positions, and cues to three-dimensional structure. Additionally, the stimulus set entailed scenes that were jumbled in their categorical content (with the individual scene parts stemming from different categories); these scenes were created to answer a different research question (see Kaiser et al. 2020a) and not used in the analyses reported in this paper. In both

conditions relevant for this paper, we used parts from four different scenes to equate the presence of visual discontinuities between fragments. Separately for each participant, 24 unique intact and 24 unique jumbled stimuli were generated by randomly drawing suitable fragments from different scenes. Each scene was presented upright and upside down. Although the key manipulation was the positioning of the individual scene parts relative to each other, it is worth noting that stimuli from the four resulting conditions adhered to, or violated, real-world structure on different levels: *1*) upright intact scenes featured typical orientation of the individual parts, typical absolute locations of the parts, and typical relative positions of the parts; *2*) upright jumbled scenes featured typical orientation of the individual parts, atypical absolute locations of the parts, and atypical relative positions of the parts; *3*) inverted intact scenes featured atypical orientation of the individual parts, atypical absolute locations of their individual parts, and typical relative positions of the parts; and *4*) inverted jumbled scenes featured atypical orientation of the individual parts, typical absolute locations of the parts, and atypical relative positions of the parts.

*Paradigm.* During the EEG experiment, the different stimuli were randomly intermixed within a single session. Within each trial, a scene appeared for 250 ms. Stimuli appeared in a black grid (4.5° visual angle), which served to mask visual discontinuities between quadrants (Fig. 1*C*). Each trial was followed by an intertrial interval that varied randomly between 700 ms and 900 ms. For this paper, only parts of the collected data (spatially intact and spatially jumbled scenes in upright and upside-down orientation) were analyzed. Each of these four conditions covered 384 trials (96 trials per scene category). Additionally, 1,152 target trials were measured. During the target trials, the crosshair changed into a slightly darker red at the same time the scene was presented. When detecting a target, participants had to press a button; additionally, they were asked to blink during the target trials, making it easier for them to refrain from blinking during nontarget trials. Target detection was purposefully made challenging to ensure sufficient attentional engagement (mean accuracy 78.1%, SE = 3.6%). Target trials were not included in subsequent analyses. Furthermore, 1,536 trials where the scenes' categorical structure was altered were measured. This data has been analyzed elsewhere (see Kaiser et al. 2020a). Furthermore, participants were instructed to maintain central fixation throughout the experiment. Stimulus presentation was controlled using the Psychtoolbox (Brainard 1997).

*EEG recording and preprocessing.* The EEG data were the same as in Kaiser et al. (2020a). EEG signals were recorded using an EASYCAP 64-electrode system and a Brainvision actiCHamp amplifier. For two participants, only data from 32 electrodes were recorded because of

technical problems. Electrodes were arranged in accordance with the 10–10 system. EEG data was recorded at 1,000 Hz sampling rate and filtered online between 0.03 Hz and 100 Hz. All electrodes were referenced online to the Fz electrode. Offline preprocessing was performed using FieldTrip (Oostenveld et al. 2011). EEG data were epoched from −200 ms to 800 ms relative to stimulus onset and were baseline corrected by subtracting the mean prestimulus signal. Channels and trials containing excessive noise were removed based on visual inspection. Blink and eye movement artifacts were removed using independent components analysis and visual inspection of the resulting components (Jung et al. 2000). The epoched data were downsampled to 200 Hz.

*EEG decoding.* Decoding analyses were performed using CoSMo-MVPA (Oosterhof et al. 2016). To track cortical representations across time, we used a cumulative classification approach that takes into account all time points before the current time point for each time point across the epoch (Ramkumar et al. 2013). This classification technique uses larger amounts of data at each subsequent time point while maintaining temporal precision in the forward direction (i.e., it only collapses across information backward in time but not forward). Cumulative decoding may thus provide increased sensitivity for detecting decoding onsets compared with standard timeseries decoding (Grootswagers et al. 2017).

We used such cumulative classifiers to discriminate between the four scene categories. This analysis was done separately for the intact and jumbled scenes. Classification analyses were performed repeatedly, with the amount of information available to the classifier accumulating across time (Fig. 2); that is, for the first time point in the epoch, the classifier was trained and tested on response patterns across the electrodes at this time point. At the second time point in the epoch, the classifier was trained and tested on response patterns across the

electrodes at the first and second time point in this epoch. Finally, at the last time point in the epoch, the classifier was trained on response patterns across all electrodes and at all time points in this epoch.

The richer information contained in these cumulative response patterns comes at the expense of a higher dimensionality of the data, which potentially harms classification. To reduce the dimensionality of the data at each time point, we performed principal component analyses (PCAs). These PCAs were always done on the classifier training set, and the PCA solution was projected onto the testing set (Grootswagers et al. 2017). For each PCA, we retained as many components as needed to explain 99% of the variance in the training set data (average number of components retained at example time points; at 0 ms: 225, SE = 11; at 200 ms: 250, SE = 10; at 800 ms: 269, SE = 10).

For classification, we used linear discriminant analysis classifiers. For each classifier, the covariance matrix was regularized by adding the identity matrix scaled by 1% of the mean of the diagonal elements (as implemented in the *cosmo_classify_lda* function in CoSMo-MVPA; Oosterhof et al. 2016). Classification was performed in a cross-validation scheme with 12 distinct folds. Classifiers were trained on data from 11 of these folds and tested on data from the left-out fold. The amount of data in the training set was always balanced across the four categories. Classification was done repeatedly until every fold was left out once. Classification accuracies were averaged across these repetitions. These analyses resulted in separate decoding timeseries for intact and jumbled scenes, which reflect the temporal accrual of category information (i.e., how well the four categories are discriminable from the neural data).

*Statistical testing.* To compare decoding timeseries against chance level and the different conditions' decoding timeseries against each other, we used a threshold-free cluster enhancement (TFCE) proce-
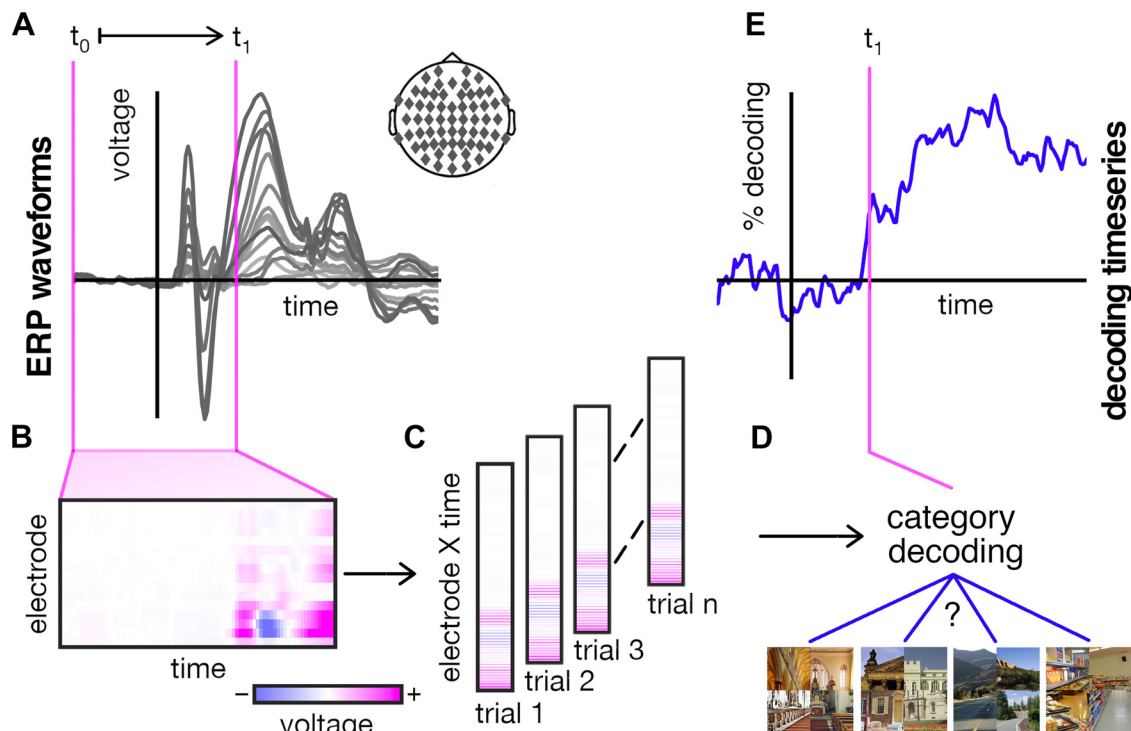


Fig. 2. Schematic depiction of the cumulative decoding approach. *A*: for each time point $t_1$ across the epoch, a separate decoding analysis was performed. *B*: for each of these analyses, we aggregated event-related potential waveforms across all EEG electrodes and all time points between $t_1$ and the beginning of the epoch ($t_0$). *C*: for each trial, we then unfolded these two-dimensional response patterns across electrodes and time into a one-dimensional response pattern. *D*: these one-dimensional response patterns were first subjected to principal component analysis to reduce dimensionality (see MATERIALS AND METHODS) and then fed to linear discriminant analysis classifiers, which were trained to discriminate the 4 scene categories. Decoding accuracy was computed by repeatedly assessing classifier performance on single trials left out during classifier training. *E*: repeating this analysis across time yielded a decoding timeseries with 200 Hz resolution. Importantly, the cumulative nature of this analysis allowed us to increase power by increasing the amount of data available to the classifier without losing temporal precision regarding the onset of category information.

dure (Smith and Nichols 2009). Multiple-comparison correction was based on a sign-permutation test (with null distributions created from 10,000 bootstrapping iterations) as implemented in CoSMoMVPA (Oosterhof et al. 2016). The resulting statistical maps were thresholded at $z > 1.96$ (i.e., $P_{corr} < 0.05$). However, the onset of statistical significance for TFCE methods may be biased by the presence of strong clusters following the onset (as expected from the cumulative decoding performed here) and can therefore not be directly interpreted (Sassenhagen and Draschkow 2019). We thus additionally provide statistics for conventional one-sample $t$ tests, which we corrected for multiple comparisons using false discovery rate (FDR) corrections. For all tests, only clusters of at least 4 consecutive significant time points (i.e., more than 20 ms) were considered.

*Data availability.* Data are publicly available on OSF (https://doi.org/10.17605/OSF.IO/ECMA4).

## RESULTS

We first analyzed data from the upright scenes, where we expected a facilitation of category information for spatially

intact, compared with jumbled, scenes. We found that EEG signals conveyed robust scene category information. Categories were discriminable for both intact scenes (significant decoding obtained from TFCE statistics: between 75 ms and 800 ms; significant decoding obtained from FDR-corrected statistics: between 75 ms and 800 ms) and jumbled scenes (TFCE: between 120 ms and 800 ms; FDR: between 135 ms and 800 ms) (Fig. 3A). Crucially, we found significantly enhanced decoding for the spatially intact scenes compared with the jumbled scenes (TFCE: between 105 ms and 800 ms; FDR: between 105 ms and 800 ms) (Fig. 3C).

The inclusion of inverted scenes allowed us to investigate whether the effects of scene structure were genuinely related to the scenes adhering to real-world structure rather than differences in their low-level visual attributes. If the enhanced category information for spatially intact scenes is indeed related to their adherence with real-world structure, then no effects should be seen when the same scenes are viewed upside
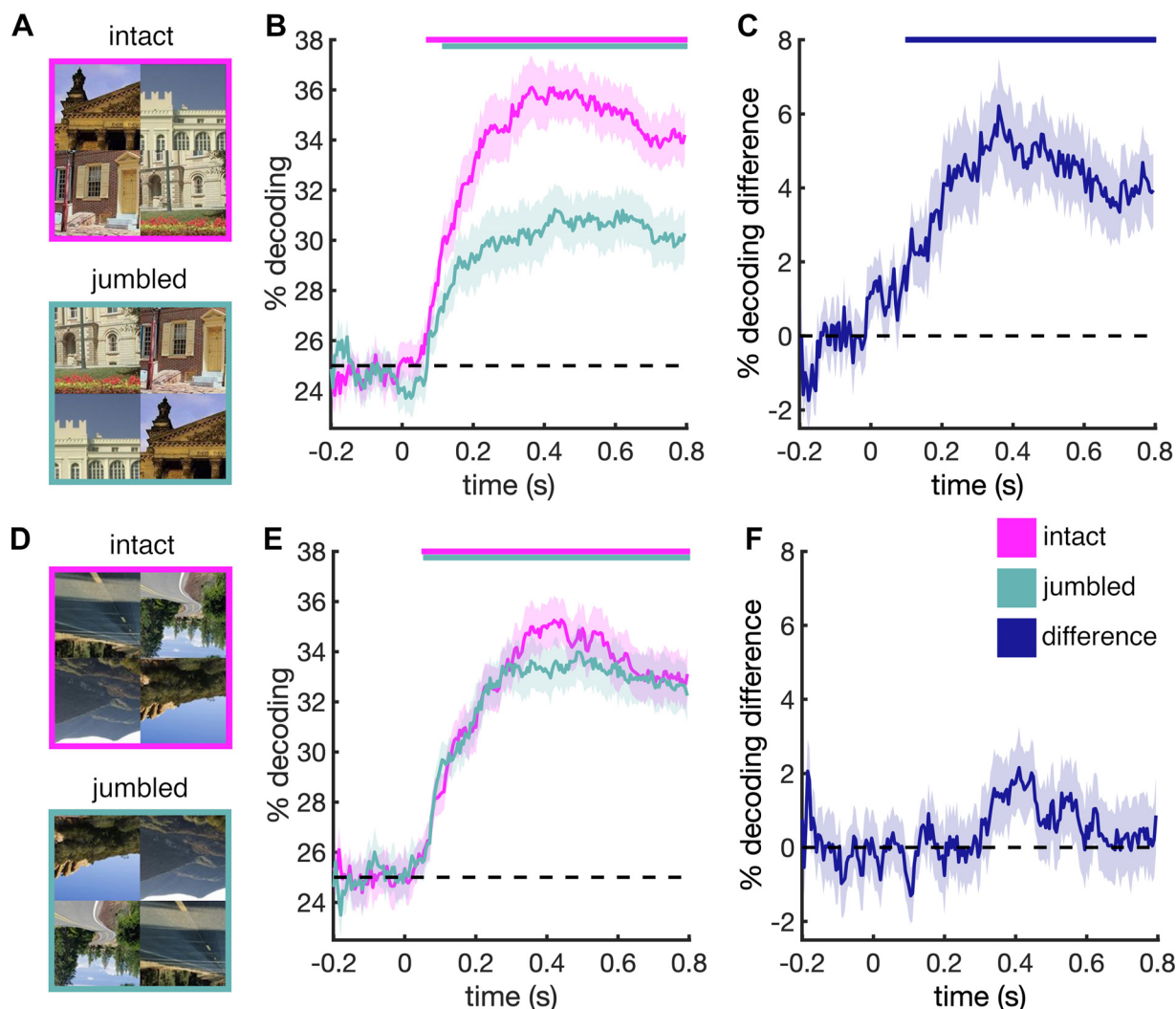


Fig. 3. Decoding of scene category for intact and jumbled scenes. *A*: first, we decoded the category of intact and jumbled scenes when they were presented upright. *B*: this analysis revealed widespread clusters of category decoding for both intact and jumbled scenes. *C*: critically, we found more accurate decoding of scene category when the scene was intact, suggesting that adherence to real-world structure boosts early visual category information. *D*: second, we decoded the category of upside-down scenes. *E*: for upside-down scenes, category could be similarly decoded from the EEG signals. *F*: however, there was no benefit of intact scene structure when the scenes were inverted, suggesting that adherence to real-world structure, rather than low-level differences, explains the enhanced category decoding for structured scenes when they are upright. Error margins indicate standard errors of the difference. Significance markers (colored horizontal lines) indicate $P < 0.05$, corrected for multiple comparisons using threshold-free cluster enhancement.

down, as both types of inverted scenes do not adhere to real-world structure in the same way as upright scenes: (*1*) although their individual parts appear in typical relative positions, the inverted intact scenes have parts that are themselves inverted and each appear in atypical absolute locations, and *2*) although their individual parts appear in typical absolute positions, the inverted jumbled scenes have parts that are themselves inverted and each appear in atypical relative positions.

Performing the category decoding analysis on the inverted scenes (Fig. 3*D*) revealed a qualitative difference to the upright scenes. The effect of scene structure was significantly stronger for the upright scenes (TFCE: between 170 ms and 800 ms; FDR: between 95 ms and 115 ms and between 185 ms and 800 ms). Indeed, no significant differences between intact and jumbled scenes were observed for the inverted scenes, although the category of both intact scenes (TFCE: between 55 ms and 800 ms; FDR: between 60 ms and 800 ms) and jumbled scenes (TFCE: between 60 ms and 800 ms; FDR: between 75 ms and 800 ms) could be decoded from the EEG signals (Fig. 3, *E* and *F*). This indicates that the early facilitation of scene category information for spatially structured scenes can be attributed to the scenes adhering to typical real-world structure, rather than to low-level features differing between the intact and jumbled scenes.

Our results establish that for processing of upright scenes, scene structure matters more than for processing inverted scenes. Additionally, one can also ask how robustly category information emerges as a function of whether the scene is presented upright or upside down. To answer this question, we directly compared category information for the intact upright scenes, the jumbled upright scenes, and the inverted scenes (Fig. 4*A*). For the inverted scenes we averaged across the intact and jumbled conditions, because there were no statistical differences between them. We found that category decoding accuracy for the inverted scenes was numerically in between the intact and jumbled upright scenes (Fig. 4*B*). When directly comparing the decoding time courses (Fig. 4*C*), we found that

category decoding was not significantly stronger in the intact upright scenes compared with the inverted scenes. By contrast, category decoding for the upright jumbled scenes was significantly weaker than for the inverted scenes (TFCE: between 170 ms and 800 ms; FDR: between 200 ms and 800 ms). This result suggests that for the inverted scenes, category can be decoded similarly as for the intact upright scenes. However, once the structure of an upright scene is destroyed, only weaker categorical representations emerge in the visual system.

## DISCUSSION

Our results provide evidence that real-world regularities facilitate the extraction of scene category information during visual analysis. We show that this facilitation of category information emerges within the first 200 ms of vision. Our findings highlight the pervasive role of real-world structure in perceptual processing, suggesting that already at relatively early processing stages cortical scene representations are tightly linked to the typical composition of our daily surroundings.

Here, we used a cumulative decoding technique to establish differences in the initial emergence of information in EEG signals. This technique uses all the available historical data (i.e., data before the current time point) for classification. Together with using PCA for dimensionality reduction, the availability of this larger amount of data promises high detection sensitivity. The availability of historical data at later time points may also hold true for the brain, where downstream regions have access to information coded earlier in upstream regions. However, as a note of caution, classifiers may also use temporally distinct information that is not necessarily available in the same way in the brain, particularly when looking at late processing stages. Cumulative decoding nonetheless provides a useful approach to reveal early differences in cortical information processing.

The early facilitation of category information is consistent with results from single-object processing, where representa-
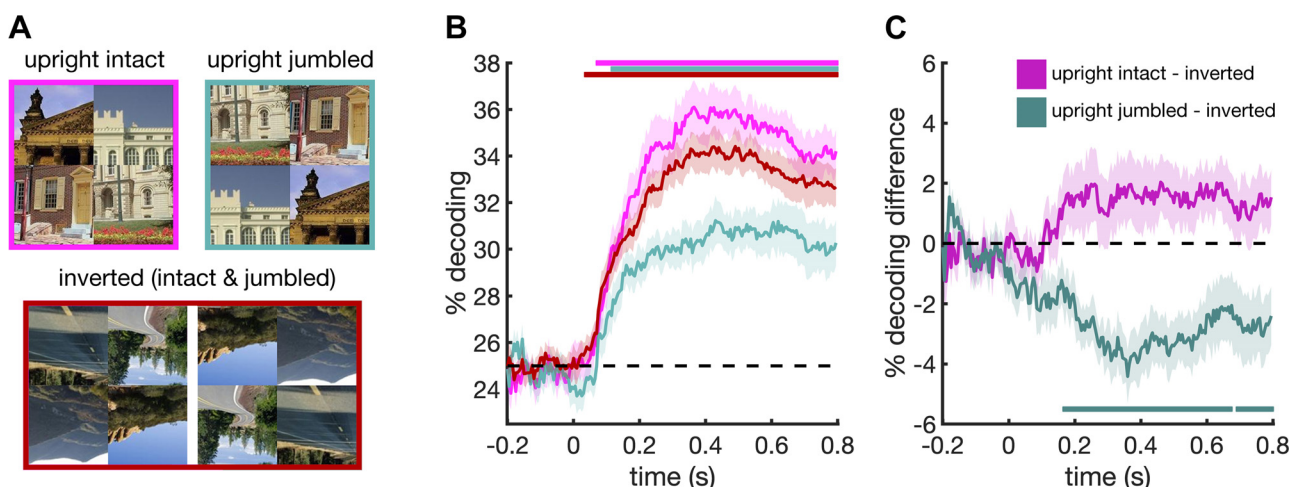


Fig. 4. Comparing category decoding between upright and inverted scenes. *A*: we compared the emergence of category information for the intact upright scenes, the jumbled upright scenes, and the inverted scenes; for the inverted scenes, we averaged across the intact and jumbled scenes, as no significant differences between the two were found. *B*: numerically, category decoding accuracy for the inverted scenes was in between the accuracies observed for the intact and jumbled upright scenes. *C*: when subtracting decoding in the inverted condition from decoding in the upright conditions, we found that statistically, category information was comparable for intact upright scenes and inverted scenes. By contrast, weaker category information was found for the jumbled upright scenes, compared with the inverted scenes, suggesting that jumbling specifically harms the emergence of category information in upright scenes. Error margins indicate standard errors of the difference. Significance markers (colored horizontal lines) indicate *P* < 0.05, corrected for multiple comparisons using threshold-free cluster enhancement.

tions of individual objects are rapidly enhanced (within the first 150 ms of vision) when the objects appear in their typical real-world locations, such as an eye in the upper visual field (Issa and DiCarlo 2012) or a shoe in the lower visual field (Kaiser et al. 2018). Together, these findings therefore support the idea that real-world structure can boost basic visual analysis across diverse stimuli and processing levels (Kaiser et al. 2019a).

When directly comparing neural category information in upright and inverted scenes, we found that it was equally pronounced when the scenes were intact and upright and when the scenes were inverted, regardless of their structural arrangement—only when the upright scenes were jumbled did we find significantly reduced category information. One interpretation of this result is that jumbling causes a specific disruption for upright scenes because for these scenes, the jumbling manipulation may be perceptually more salient. Alternatively, the pattern of results may be explained by an interaction of two different effects. The inverted intact scenes still retain the intact relative positioning of their parts, which may explain why they are better decodable than the upright jumbled scenes. The inverted jumbled scenes do not have this intact relative positioning, but by means of inversion they gain an intact absolute positioning of their parts (e.g., a piece of sky would be in the upper part of an inverted jumbled scene, which is where it belongs); this may explain why these scenes yield better category decoding than upright jumbled scenes. At this point, further studies are needed to fully understand this pattern of results. Challenges with interpreting inversion effects in the current paradigm may necessitate the inclusion of other low-level stimulus controls in these future studies.

Although our effects demonstrate an enhanced early representation of scenes that adhere to real-world structure compared with scenes that do not, studies on object-scene consistency suggest that EEG waveforms only become affected by typical object positioning after around 250 ms of vision (Coco et al. 2020; Draschkow et al. 2018; Ganis and Kutas 2003; Mudrik et al. 2010, 2014; Võ and Wolfe 2013). How do these early and late effects of scene structure relate to each other?

As one possibility, later effects may partly reflect increased responses to inconsistencies rather than an enhanced processing of consistent scene-object combinations (Faivre et al. 2019). Together with our results, these findings may suggest that early responses are biased toward scenes that predictably follow real-world structure, whereas later responses may be more biased toward violations of this structure. This idea is consistent with a recent proposal in predictive processing, which suggests a temporal succession of more general processing biases, first toward the expected and then toward the surprising (Press et al. 2020).

Alternatively, the beneficial effects of real-world regularities may not immediately result in consistency signals. Whether visual inputs generally are consistent with our real-world experience may only be analyzed following more basic visual analysis. Supporting this idea, generic consistency signals in our data only emerge later than the enhanced category processing. As previously reported, intact and jumbled scenes (independent of their category) evoked reliably different responses only after 255 ms of processing (Kaiser et al. 2020a).

More broadly, the findings can add to our understanding of efficient everyday vision. Even under challenging real-world conditions, human vision is remarkably efficient; in fact, it is much more efficient than findings from simplified laboratory experiments would predict (Wolfe et al. 2011; Peelen and Kastner 2014). Behavioral studies using jumbling paradigms have suggested that typical scene structure contributes to this efficiency. When scenes are structurally intact, observers can better categorize them (Biederman et al. 1974), recognize objects within them (Biederman 1972), or detect visual changes in the scene (Varakin and Levin 2008). These perceptual benefits may be linked to the rapid facilitation of neural category information for typical scenes observed in the current study. However, our participants performed an orthogonal fixation task, which precludes directly linking brain and behavior here. Future studies combining neural recordings with naturalistic behavioral tasks may reveal that the early cortical tuning to real-world structure may be a crucial asset for solving complex real-world tasks.

## AUTHOR CONTRIBUTIONS

D.K. and R.M.C. conceived and designed research; G.H. performed experiments; D.K. and G.H. analyzed data; D.K., G.H., and R.M.C. interpreted results of experiments; D.K. prepared figures; D.K. drafted manuscript; D.K., G.H., and R.M.C. edited and revised manuscript; D.K., G.H., and R.M.C. approved final version of manuscript.

## REFERENCES

**Abassi E, Papeo L.** The representation of two-body shapes in the human visual cortex. *J Neurosci* 40: 852–863, 2020. doi:10.1523/JNEUROSCI.1378-19.2019.

**Baldassano C, Beck DM, Fei-Fei L.** Human-object interactions are more than the sum of their parts. *Cereb Cortex* 27: 2276–2288, 2017. doi:10.1093/cercor/bhw077.

**Bar M.** Visual objects in context. *Nat Rev Neurosci* 5: 617–629, 2004. doi:10.1038/nrn1476.

**Biederman I.** Perceiving real-world scenes. *Science* 177: 77–80, 1972. doi:10.1126/science.177.4043.77.

**Biederman I, Rabinowitz JC, Glass AL, Stacy EW Jr.** On the information extracted from a glance at a scene. *J Exp Psychol* 103: 597–600, 1974. doi:10.1037/h0037158.

**Bilalić M, Lindig T, Turella L.** Parsing rooms: the role of the PPA and RSC in perceiving object relations and spatial layout. *Brain Struct Funct* 224: 2505–2524, 2019. doi:10.1007/s00429-019-01901-0.

**Brainard DH.** The psychophysics toolbox. *Spat Vis* 10: 433–436, 1997. doi:10.1163/156856897X00357.

**Chan AW, Kravitz DJ, Truong S, Arizpe J, Baker CI.** Cortical representations of bodies and faces are strongest in commonly experienced configurations. *Nat Neurosci* 13: 417–418, 2010. doi:10.1038/nn.2502.

**Coco MI, Nuthmann A, Dimigen O.** Fixation-related brain potentials during semantic integration of object-scene information. *J Cogn Neurosci* 32: 571–589, 2020. doi:10.1162/jocn_a_01504.

**de Haas B, Schwarzkopf DS, Alvarez I, Lawson RP, Henriksson L, Kriegeskorte N, Rees G.** Perception and processing of faces in the human brain is tuned to typical feature locations. *J Neurosci* 36: 9289–9302, 2016. doi:10.1523/JNEUROSCI.4131-14.2016.

**Dima DC, Perry G, Singh KD.** Spatial frequency supports the emergence of categorical representations in visual cortex during natural scene perception. *Neuroimage* 179: 102–116, 2018. doi:10.1016/j.neuroimage.2018.06.033.

**Draschkow D, Heikel E, Võ ML, Fiebach CJ, Sassenhagen J.** No evidence from MVPA for different processes underlying the N300 and N400 incongruity effects in object-scene processing. *Neuropsychologia* 120: 9–17, 2018. doi:10.1016/j.neuropsychologia.2018.09.016.

**Faivre N, Dubois J, Schwartz N, Mudrik L.** Imaging object-scene relations processing in visible and invisible natural scenes. *Sci Rep* 9: 4567, 2019. doi:10.1038/s41598-019-38654-z.

**Ganis G, Kutas M.** An electrophysiological study of scene effects on object identification. *Brain Res Cogn Brain Res* 16: 123–144, 2003. doi:10.1016/S0926-6410(02)00244-6.

**Grootswagers T, Wardle SG, Carlson TA.** Decoding dynamic brain patterns from evoked responses: a tutorial on multivariate pattern analysis applied to time series neuroimaging data. *J Cogn Neurosci* 29: 677–697, 2017. doi:10.1162/jocn_a_01068.

**Issa EB, DiCarlo JJ.** Precedence of the eye region in neural processing of faces. *J Neurosci* 32: 16666–16682, 2012. doi:10.1523/JNEUROSCI.2391-12.2012.

**Jung TP, Makeig S, Humphries C, Lee TW, McKeown MJ, Iragui V, Sejnowski TJ.** Removing electroencephalographic artifacts by blind source separation. *Psychophysiology* 37: 163–178, 2000. doi:10.1111/1469-8986.3720163.

**Kaiser D, Cichy RM.** Typical visual-field locations enhance processing in object-selective channels of human occipital cortex. *J Neurophysiol* 120: 848–853, 2018. doi:10.1152/jn.00229.2018.

**Kaiser D, Häberle G, Cichy RM.** Cortical sensitivity to natural scene structure. *Hum Brain Mapp* 41: 1286–1295, 2020a. doi:10.1002/hbm.24875.

**Kaiser D, Inciuraite G, Cichy RM.** Rapid contextualization of fragmented scene information in the human visual system. *Neuroimage* 117045, 2020. doi:10.1016/j.neuroimage.2020.117045.

**Kaiser D, Moeskops MM, Cichy RM.** Typical retinotopic locations impact the time course of object coding. *Neuroimage* 176: 372–379, 2018. doi:10.1016/j.neuroimage.2018.05.006.

**Kaiser D, Peelen MV.** Transformation from independent to integrative coding of multi-object arrangements in human visual cortex. *Neuroimage* 169: 334–341, 2018. doi:10.1016/j.neuroimage.2017.12.065.

**Kaiser D, Quek GL, Cichy RM, Peelen MV.** Object vision in a structured world. *Trends Cogn Sci* 23: 672–685, 2019a. doi:10.1016/j.tics.2019.04.013.

**Kaiser D, Stein T, Peelen MV.** Object grouping based on real-world regularities facilitates perception by reducing competitive interactions in visual cortex. *Proc Natl Acad Sci USA* 111: 11217–11222, 2014. doi:10.1073/pnas.1400559111.

**Kaiser D, Turini J, Cichy RM.** A neural mechanism for contextualizing fragmented inputs during naturalistic vision. *eLife* 8: e48182, 2019b. doi:10.7554/eLife.48182.

**Kim JG, Biederman I.** Where do objects become scenes? *Cereb Cortex* 21: 1738–1746, 2011. doi:10.1093/cercor/bhq240.

**Konkle T, Brady TF, Alvarez GA, Oliva A.** Scene memory is more detailed than you think: the role of categories in visual long-term memory. *Psychol Sci* 21: 1551–1556, 2010. doi:10.1177/0956797610385359.

**Lowe MX, Rajsic J, Ferber S, Walther DB.** Discriminating scene categories from brain activity within 100 milliseconds. *Cortex* 106: 275–287, 2018. doi:10.1016/j.cortex.2018.06.006.

**Mudrik L, Lamy D, Deouell LY.** ERP evidence for context congruity effects during simultaneous object-scene processing. *Neuropsychologia* 48: 507–517, 2010. doi:10.1016/j.neuropsychologia.2009.10.011.

**Mudrik L, Shalgi S, Lamy D, Deouell LY.** Synchronous contextual irregularities affect early scene processing: replication and extension. *Neuropsychologia* 56: 447–458, 2014. doi:10.1016/j.neuropsychologia.2014.02.020.

**Oliva A, Torralba A.** The role of context in object recognition. *Trends Cogn Sci* 11: 520–527, 2007. doi:10.1016/j.tics.2007.09.009.

**Oostenveld R, Fries P, Maris E, Schoffelen JM.** FieldTrip: Open source software for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput Intell Neurosci* 2011: 156869, 2011. doi:10.1155/2011/156869.

**Oosterhof NN, Connolly AC, Haxby JV.** CoSMoMVPA: Multi-modal multivariate pattern analysis of neuroimaging data in Matlab/GNU Octave. *Front Neuroinform* 10: 27, 2016. doi:10.3389/fninf.2016.00027.

**Peelen MV, Kastner S.** Attention in the real world: toward understanding its neural basis. *Trends Cogn Sci* 18: 242–250, 2014. doi:10.1016/j.tics.2014.02.004.

**Press C, Kok P, Yon D.** The perceptual prediction paradox. *Trends Cogn Sci* 24: 13–24, 2020. doi:10.1016/j.tics.2019.11.003.

**Ramkumar P, Jas M, Pannasch S, Hari R, Parkkonen L.** Feature-specific information processing precedes concerted activation in human visual cortex. *J Neurosci* 33: 7691–7699, 2013. doi:10.1523/JNEUROSCI.3905-12.2013.

**Roberts KL, Humphreys GW.** Action relationships concatenate representations of separate objects in the ventral visual system. *Neuroimage* 52: 1541–1548, 2010. doi:10.1016/j.neuroimage.2010.05.044.

**Sassenhagen J, Draschkow D.** Cluster-based permutation tests of MEG/EEG data do not establish significance of effect latency or location. *Psychophysiology* 56: e13335, 2019. doi:10.1111/psyp.13335.

**Smith SM, Nichols TE.** Threshold-free cluster enhancement: addressing problems of smoothing, threshold dependence and localisation in cluster inference. *Neuroimage* 44: 83–98, 2009. doi:10.1016/j.neuroimage.2008.03.061.

**Varakin DA, Levin DT.** Scene structure enhances change detection. *Q J Exp Psychol (Hove)* 61: 543–551, 2008. doi:10.1080/17470210701774176.

**Võ ML, Boettcher SE, Draschkow D.** Reading scenes: how scene grammar guides attention and aids perception in real-world environments. *Curr Opin Psychol* 29: 205–210, 2019. doi:10.1016/j.copsyc.2019.03.009.

**Võ ML, Wolfe JM.** Differential electrophysiological signatures of semantic and syntactic scene processing. *Psychol Sci* 24: 1816–1823, 2013. doi:10.1177/0956797613476955.

**Wolfe JM, Võ ML, Evans KK, Greene MR.** Visual search in scenes involves selective and nonselective pathways. *Trends Cogn Sci* 15: 77–84, 2011. doi:10.1016/j.tics.2010.12.001.