

Hierarchical surprise signals in naturalistic violation of expectations

Vincent Pliakat^{1,2,3,*,+}, Pablo R. Grassi^{1,2,3,*,+}, Julius Frack⁴, Andreas Bartels^{1,2,3,+}

¹ Department of Psychology, University of Tübingen, Tübingen, Germany

² Centre for Integrative Neuroscience, Tübingen, Germany

³ Max-Planck Institute for Biological Cybernetics, Tübingen, Germany

⁴ Rocket Magic, Tübingen, Germany

*, equal contribution

+, corresponding authors

Abstract

Surprise responses signal both, high-level cognitive alerts that information is missing, and increasingly specific back-propagating error signals that allow updates in processing nodes. Studying surprise is hence central for cognitive neuroscience to understand internal world representations and learning. Yet, only few prior studies used naturalistic stimuli targeting our high-level understanding of the world. Here, we use magic tricks in an fMRI experiment to investigate neural responses to violations of core assumptions held by humans about the world. We showed participants naturalistic videos of three types of magic tricks, involving objects appearing, changing color, or disappearing, along with control videos without any violation of expectation. Importantly, the same videos were presented with and without prior knowledge about the tricks' explanation. Results revealed generic responses in frontoparietal areas, together with responses specific to each of the three trick types in posterior sensory areas. A subset of these regions, the midline areas of the default mode network (DMN), showed surprise activity that depended on prior knowledge. Equally, sensory regions showed sensitivity to prior knowledge, reflected in differing decoding accuracies. These results suggest a hierarchy of surprise signals involving three processing stages: first, generic processing of violation of expectations in frontoparietal areas. Second, surprise signals in sensory regions that are specific to the processed feature, demonstrating propagation of surprise signals to mid-level visual areas. Third, prior-knowledge dependent processing of surprise signals in parts of the DMN, showing its involvement in internal state monitoring and in making sense of events as they unfold in time.

Keywords: surprise, magic, predictive coding, violation of expectation, intuitive physics, fMRI

Word count: Abstract: 250, **Introduction:** 900, **Discussion:** 1818, **Total:** ca. 13000 (w/o refs)

1 Introduction

2 Prior experience and “intuitive” knowledge about the physical world guide our perception and allow
3 for a meaningful interaction with the environment. They set up constraints on our expectations
4 based on what we believe to be possible in the world. For example, prior world-knowledge informs
5 us that objects do not vanish of existence if occluded (object permanency), that objects tend to
6 keep their features (feature constancy), that objects cannot pass through other objects (solidity),
7 that objects to not appear out of the blue, and so forth. Informed expectations based on these
8 intuitive physical priors allow us to quickly make sense of incoming sensory information. Such
9 expectations have been shown to strongly modulate perception (de Lange et al., 2018). For
10 example, a “light-from-above” prior constrains depth perception from shading (Adams et al., 2004),
11 and knowledge that objects cannot occupy the same place at the same time explains our inability
12 to perceive two objects simultaneously in bistable perception (Hohwy et al., 2008).

13 Accordingly, perception can be understood as an inferential process in which top-down
14 informed expectations are matched with incoming sensory information (Friston, 2010, 2005; Lee
15 and Mumford, 2003; Rao and Ballard, 1999). In this “predictive processing” framework deviations
16 between expectations and incoming data (i.e., prediction errors) are used to update an internal
17 model of the world at different levels of complexity and abstraction. Higher-level priors (like object
18 permanency) represent more abstract aspects of the world and constrain lower-level inferences
19 that represent more immediate features of the world (like a particular object) (Clark, 2013; Hohwy,
20 2014). In this context, intuitive physical knowledge can be thought of as an internal model
21 representing different aspects of the causal structure of the physical world, not unsimilar to a
22 “physics engine” in a virtual environment (Battaglia et al., 2013).

23 This prior knowledge of physical principles can be studied using events that seemingly
24 violate them. For example, evidence from developmental psychology using so-called violation of
25 expectation (VOE) paradigms suggest that infants acquire important aspects about the workings
26 of the physical world during the first year of life, such as object permanency and solidity (Hespos
27 et al., 2009; Wang, 2004; Wynn, 1992). However, it is largely unclear how surprise-responses
28 relating to intuitive physical principles are represented in the human brain. Most previous
29 neuroimaging studies investigated responses to lower-level VOE, e.g. using paradigms involving
30 infrequent (and thus unexpected) stimuli (e.g., Egner et al., 2010; Kok et al., 2012a; Todorovic et
31 al., 2011; Wessel et al., 2012), or omission of expected stimuli (Kok et al., 2012b; SanMiguel et
32 al., 2013; e.g., Wacongne et al., 2011). These studies revealed lower-level stimulus-specific
33 prediction errors in modality- and feature-specific areas. In contrast, studies investigating VOE of
34 higher-level physical principles using more complex stimuli such as computer generated
35 animations (Bardi et al., 2017; Liu et al., 2024) or naturalistic videos showing magic tricks (Danek
36 et al., 2015; Parris et al., 2009) revealed higher-level surprise signals in frontoparietal areas. This
37 is in line with a frontoparietal role in the representation of physical concepts (Fischer et al., 2016;

1 Schwettmann et al., 2019). However, higher-level VOE studies have hitherto not observed VOE-
2 related activity in sensory areas.

3 Here, we set out to unify and resolve this longstanding divergence of results between high-
4 level and low-level VOE paradigms by using a novel VOE paradigm. Our paradigm was designed
5 to uncover hierarchical VOE signals of the brain's internal world model. It contained a battery of
6 standardized magic videos that we presented to human participants while measuring fMRI
7 responses to investigate: 1) which regions are generally involved when viewing natural videos that
8 violate physical principles, 2) whether specific types of violations, like appearance of objects, or
9 changes of color, modulate sensory areas known to process the feature concerned, 3) whether
10 knowledge of the explanation of a given magic trick modulates the observed VOE activity.

11 To this aim, we created and validated videos for a naturalistic VOE paradigm showing
12 either dedicated magic tricks (to create the illusion of seemingly impossible events to actually
13 occur, cf. Grassi and Bartels, 2021) or matched control actions that involved no violation of
14 physical principles. The VOE videos were designed to evoke surprise responses related to
15 unexpected object appearance, disappearance of objects, and feature change (color-changing
16 objects). They were performed by a professional magician (Julius Frack), and each trick-type was
17 performed using three common objects (balls, playing cards and pencils). Each trick was
18 presented before and after revealing the method of the tricks. This allowed us to compare
19 responses with and without VOE using identical videos.

20 Univariate analyses revealed a hierarchy of surprise signals: frontoparietal areas, including
21 areas of the default mode network (DMN), were involved when perceiving events violating physical
22 principles regardless of the type of trick used. In contrast, posterior sensory areas were modulated
23 specifically by the type of expectation-violation, such as color-processing medial fusiform cortex
24 by color change, and object selective LOC by the appearance of objects. Controls indicate that
25 their modulation is due to the feature-specific surprise and not due to the feature-change.
26 Multivariate analyses extended the results: information about the specific types of expectation-
27 violations was exclusively encoded in posterior regions, and significantly decodable down to the
28 earliest levels of cortical visual processing (V1-V3). Additionally, decoding accuracy significantly
29 decreased with prior knowledge, revealing a reduction in surprise signals. Together, our results
30 demonstrate a generic response in frontoparietal areas to violation of physical principles, along
31 with concurrent representations of specific expected information in early sensory areas.
32

Methods

Participants

We performed fMRI on 27 subjects. Three subjects were excluded from data analysis due to excessive movement and/or sleepiness during the scanning sessions. Data from a total of 24 subjects were analyzed (16 female; 8 male; mean age 24.4 ± 4.3 SD years). All subjects had normal or corrected to normal vision, no history of neurological impairments nor contraindication for fMRI. Participants had no expertise as magicians and were naive to the magic tricks used. Participants provided written informed consent prior to the experiment. The study was approved by the ethics committee of the University Clinic Tübingen and was conducted in accord with the Declaration of Helsinki.

VOE stimuli

Video recordings. A set of 63 videos were created for the VOE study. In accord with previous work (Parris et al., 2009), we created videos for three different conditions: the videos showed either magic tricks (magic condition, 18 videos), similar actions without a magic event (control condition, 18 videos) or unusual actions with the objects used in the magic tricks (unusual condition, 9 videos). We included the unusual condition to investigate neural correlates of surprise in a similar setting while not violating any physical concept (cf. Parris et al., 2009). We further created explanation videos showing how each of the magic tricks was achieved (18 videos).

All videos were performed by professional illusionist Julius Frack in a standardized setting consisting of a black background and a black table (see Figure 1A for an example). To investigate the effect of specific VOE, we presented magic tricks showing three different violations of physical principles: appearances (A) (i.e., a red object appears), color changes (C) (i.e., a red object changes color to blue) and vanishes (V) (i.e., a red object disappears) (see Figure 1C). To generalize across objects, we used three easily distinguishable objects: balls, playing cards and thick pencils. Each object was used equally often in each trick type. Finally, we created two versions of each type of VOE for each object (e.g., two different videos showing a red ball changing to a blue ball using different methods).

Full set of stimuli. Each of these magic tricks had a matched control video that showed the same sequences of actions as the trick, but without VOE (e.g., following the same actions a ball would not change its color). Thus, we had a total of 18 magic videos (3 type of VOEs \times 3 objects \times 2 methods), with 18 matching control videos and nine unusual videos (3 unusual actions per used object). As the solutions to the magic tricks were revealed in distinct sets throughout the experiment, we used a variety of different methods for each type of VOE. This ensured that participants were only able to infer trick solutions they were intended to understand.

Luminance and durations. Tricks were recorded in a standardized setting under the same lighting conditions. To balance out remaining inequalities, custom MATLAB (MathWorks, Natick, MA) scripts were used to standardize the videos (resolution of 1920 x 1080 with 25 frames per second), which were filmed on different days and had different lengths. We manually applied white-balance using 5 to 10 selected white-pixels for all videos of a same day, matched the luminance and contrast of the videos based their first frame, and shortened the videos to be no longer than 14 seconds. The final duration of the videos was $12.8 \text{ s} \pm 1.08 \text{ s}$ (mean \pm SD).

Stimulus presentation. Stimuli were presented using MATLAB 2019b using Psychtoolbox3 (version 3.0.16 <http://psychtoolbox.org/>) on a Linux computer and back-projected to a translucent screen mounted at the rear of the scanner bore using a VPIXX Pro projector (VPiXX Technologies, Saint-Bruno-de-Montarville, Canada) at a frame rate of 144 Hz. Participants viewed the screen (26.1×14.7 visual degrees) via a mirror mounted on the 64-channel head coil (Siemens, Erlangen, Germany) at a distance of 105 cm. To center the stimulus presentation, we cropped 160 pixels from the left and right side of the videos that only showed a black background. Accordingly, the shown part of the videos covered 21.8×14.7 visual degrees (1600 x 1080 pixels).

Behavioral evaluation of stimuli. Prior to the fMRI experiment we performed two psychophysics experiments with a total of 18 subjects (nine subjects in each experiment) to ensure the suitability of our stimuli and to select the magic tricks to be used in the fMRI experiment (see results in Supplementary Section *Behavioral evaluation of stimuli*). Participants were shown the videos using the same design as in the fMRI experiment (see Figure 1B). Participants were surprised when viewing the magic videos, even though they knew they were observing tricks (Grassi et al., 2024) and systematically reported to be more surprised when viewing the magic videos compared to the matched controls. Moreover, surprise responses to the VOE shown in the magic videos were reduced after providing an explanation to the tricks.

Experimental design and procedure

fMRI runs. Each run consisted of a total of 24 video trials: two repetitions of the six unique magic videos of one object (three types of VOE, each type recorded in two versions), resulting in 12 magic trick presentations; the six matching control videos (shown once); and two repetitions of three additional unusual videos, resulting in six unusual video presentations (see Figure 1D). The resulting 24 video trials were presented in a pseudo-random order that avoided repetitions of identical videos (different randomizations across runs).

Paradigm. The fMRI experiment consisted of three sets with four experimental fMRI runs each. Each set presented videos of only one object (i.e., balls, playing cards or pencils). The

1 presentation order of the object sets was counterbalanced across subjects. After the first two fMRI
2 runs in a set (pre-revelation runs) the method behind each magic trick in the set was revealed by
3 showing the participants each of the magic tricks again together with the matching revelation
4 video. Subjects could watch the videos as often as they wished and were asked to confirm per
5 button press that they understood how each of the tricks was achieved. Thereafter, two more runs
6 showing the same videos were performed (post-revelation runs). Together, each set consisted of
7 two fMRI runs before and two fMRI runs after the explanation of the tricks. A visualization of the
8 experimental design is shown in Figure 1E. We hypothesized that providing the explanation of the
9 tricks would decrease VOE responses because participants would adjust their expectations (e.g.,
10 knowing where the seemingly disappearing objects were concealed). Consequently, our
11 experimental design allowed us to compare responses with and without VOE using identical
12 stimuli.

13
14 **Individual trials.** Each video trial was presented for 14 seconds. In case a video was shorter (e.g.,
15 13 seconds), the last frame of the video was shown for the remaining time (e.g., for 1 second).
16 Participants were informed about this. Moreover, the first 500 ms of each video showed a central
17 cue (1.3 visual degrees) indicating whether the video was going to be a magic (M) or a non-magic
18 (X) video (i.e., control or unusual actions). After each video presentation participants were asked
19 to rate from 1 to 5 how surprising the content of the video was (1 = not surprising, 5 = very
20 surprising). They had two seconds to respond, after which the next trial started (total trial duration
21 = 16 seconds, see Figure 1B). Behavioral surprise ratings were given by the subjects using a
22 button-box with five keys.

23
24 **Trial cues.** We included the cue (M or X) at the beginning of the trial to prevent participants from
25 being surprised by contextual or serial effects. For example, since we tried to create control and
26 magic videos that are visually as similar as possible, one could easily mistake a control video for
27 a magic trick at the beginning of the video. The missing magic trick could then be surprising for
28 participants expecting one. Yet, this was not the kind of surprise we wanted to investigate. We
29 hypothesized that prior knowledge about the content of a trial (magic vs. no-magic) should prevent
30 this from happening. Moreover, to reduce predictability of video content we flipped each video on
31 every second presentation of the same video horizontally. For example, if the magician performed
32 a color changing card trick on his left-hand side in the first presentation of the video, the next
33 presentation of the same video showed the color changing card trick on his right-hand side.

34
35 **Pre-scan instructions.** Before scanning, we instructed subjects about the task and design of the
36 experiment (i.e., the set design, meaning of trial cues, explanation videos, etc.). Participants were
37 informed that videos would show either magic tricks, control actions or unusual actions. Moreover,
38 participants were shown a magic and matching control example video (not used in the actual

experiment) to give them an impression about the kind and duration of the videos they were about to see. Further, to prevent participants from watching the videos in a “problem-solving” attitude, we instructed them to passively watch and enjoy the videos without trying to get behind the method of the tricks, as these would be explained during the experiment.

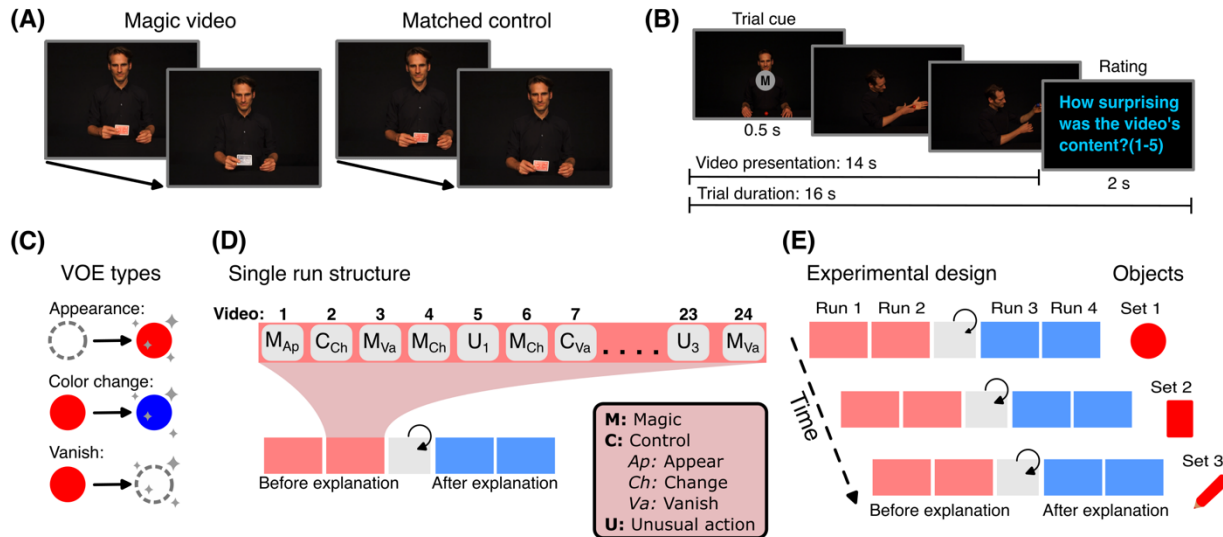


Figure 1. **A**, example of a color-changing card magic trick (left) and its corresponding matched control video (right). **B**, shown is the timeline of a single trial. Every video started with a central cue lasting 500 ms indicating whether the video will show a magic trick (M) or not (X). Each video presentation lasted 14 seconds, if a video happened to be shorter than 14 s, the last frame of the video was shown until the total duration was 14 s as a static image. After the video, subjects had 2 seconds to answer how surprising the video's content was. Note that the text was displayed in one line and did not cover the whole screen. **C**, shown are the three different types of VOE used (i.e., “magic events”). We created videos showing magic tricks using three objects (balls, playing cards and pencils) that showed either an unexpected object appear (Ap), change color (from red to blue) (Ch) or vanish (Va). For each object-VOE combination, we created two tricks (i.e., we had two color-changing card tricks). **D**, schematic example of a single experimental set. One set consisted of four runs – two before and two after explanation runs (i.e., pre and post revelation). Between pre- and post-revelation runs we showed participants videos explaining the magic tricks performed during the corresponding set. Each run showed 24 videos with a pseudo-randomized order ensuring that the same video was not presented in two consecutive trials (twelve magic videos [M], six matched controlled videos [C] and six unusual actions [U]). **E**, complete experimental design showing all three sets. The experiment was divided into three sets, one for each of the objects. Each set was divided into four experimental runs, where each run showed all the videos of an object but in a randomized order. After the second fMRI run the methods behind the tricks shown in the set were revealed, dividing the runs into *before* and *after* the explanation of the tricks (pre- and post-revelation, respectively).

fMRI data acquisition

fMRI data were acquired in a 3 Tesla Siemens Prisma scanner with a 64-channel head coil (Siemens, Erlangen, Germany). Functional images were acquired using an accelerated T2*-

1 weighted gradient-echo echoplanar imaging (EPI) sequence (multiband factor = 2, repetition time
2 (TR) = 2000 ms, echo time (TE) = 30 ms, flip angle (FA) = 75°, 62 slices with an isotropic voxel
3 size of 2 × 2 × 2 mm) using GRAPPA (GRAPPA = 2). Each run consisted of 198 images (total
4 duration = 6 minutes and 36 seconds). Moreover, a high-resolution T1-weighted structural scan
5 with whole-brain coverage was performed for each participant (TR = 2000 ms, TE = 3.06 ms,
6 inversion time (TI) = 1100 ms, FA = 9°, 192 slices and an isotropic voxel size of 1×1×1 mm). The
7 structural scan was measured during the explanation of the tricks in the first set.

8 **fMRI data preprocessing**

9 Functional MRI images were preprocessed first by removing thermal noise from the magnitude
10 EPI images using NORDIC, a PCA-based algorithm (Vizioli et al., 2021) in MATLAB. Then, we
11 discarded the first five volumes of each run to allow for T1 equilibration effects. Using SPM12 we
12 further performed motion correction (realigned to the first image), slice-time correction (using
13 middle slices as reference) and co-registration to the structural scan. Finally, functional MRI data
14 for whole-brain analyses were normalized to the Montreal Neurological Institute template brain
15 (MNI152) and spatially smoothed with a Gaussian kernel of full width at half-maximum of 6 mm
16 for univariate whole brain analyses. Region-of-interest (ROI) analyses were performed on
17 unsmoothed data in native space. Moreover, for the generation of subject-specific ROIs, we
18 generated inflated individual brain surfaces using Freesurfer 7.1.1 (Dale et al., 1999) using a
19 dedicated docker container (<https://hub.docker.com/r/freesurfer/freesurfer>).

20 **Univariate whole-brain data analysis**

21 To investigate differences in neural activity during high-level VOE in the human brain, we
22 performed whole-brain and ROI univariate analyses. We had three specific aims: 1) to investigate
23 neural responses to violations of physical principles in naturalistic stimuli, comparable to previous
24 reports using magic tricks (Danek et al., 2015; cf. Parris et al., 2009), 2) to investigate differential
25 responses between different VOE (i.e., magic trick types) and 3) to investigate the role of prior
26 knowledge in the perception of events violating physical principles by comparing responses to the
27 very same videos before and after revelation.

28 For these analyses, we created two event-related general linear models (GLM) using the
29 canonical hemodynamic response function and high-pass filtered data with a cut-off at 128s in
30 SPM12. In the first GLM (aim 1), we modelled responses using only three regressors of interest in
31 each run (magic, control, and unusual stimuli conditions). In a second, extended GLM, we
32 investigated possible differential responses to different forms of violation of expectations (aim 2),
33 using seven regressors of interest that modelled BOLD responses of the three magic types, three
34 matching controls and unusual actions for each run. As in previous reports, we modelled individual
35 trials within a regressor using discrete event times based on the moment of surprise (Danek et al.,
36 2015; cf. Parris et al., 2009). The definition of the event times was done by averaging the

1 independent selection of a suitable frame done by two authors (VP and PRG). For magic videos,
2 timing was decided based on when the “magic” (i.e., VOE) happened. For matching control videos,
3 the corresponding time point was selected, i.e., the moment in which one would expect to see a
4 magic event in the magic videos. For the videos showing unusual actions, timing was selected
5 based on the onset of the unusual actions. Moreover, we included participants’ response times,
6 six movement parameters and a column of ones as nuisance regressors in both GLMs. We
7 additionally computed the same analyses modelling the whole video durations (14 s) and report
8 these in the supplementary Tables S11-13 and supplementary Figure S7). We used these models
9 to address our three aims as follows.

10 Aim 1: to identify brain areas preferentially involved in signaling high-level VOE we
11 compared responses to magic and matching control videos in all six pre-revelation runs (i.e., two
12 runs from each object set) using the pooled regressors (from the first GLM). First-level contrasts
13 $Magic_{pre} > Control_{pre}$ were used for second-level random-effect analyses with non-parametric
14 permutations tests using the non-parametric mapping toolbox SnPM13 (Nichols and Holmes,
15 2002) for SPM12.

16 Aims 1 and 2: to investigate generic and violation-specific responses to high-level VOE
17 before the explanation of the tricks, we used a second-level, 2 (magic, control) x 3 (appear,
18 change, vanish) repeated measures ANOVA to conduct four conjunction analyses (Friston et al.,
19 2005). To find common responses to all VOE, we contrasted each magic trick type with its
20 corresponding control condition before the explanation of the tricks (i.e., $Magic_{A_pre} > Control_{A_pre}$)
21 and used those contrasts to perform a conjunction (across all three VOE types) with a threshold
22 of $p_{unc} \leq 0.001$ and a cluster threshold of $k = 10$. For each VOE type separately, we further
23 performed a conjunction analysis on the responses of one VOE type against the others. For
24 example, responses specific to an appearing object were investigated by means of the conjunction
25 of the contrasts $Magic_{A_pre} > Magic_{C_pre} \cap Magic_{A_pre} > Magic_{V_pre}$.

26 Aim 3: to identify brain areas generally involved in the perception of magic and affected by
27 prior knowledge we compared responses to magic after providing the explanation of the tricks
28 $Magic_{post} > Control_{post}$ and additionally compared the interaction between the revelation and magic
29 $(Magic_{pre} > Control_{pre}) > (Magic_{post} > Control_{post})$ using permutations tests. Finally, we compared the
30 same contrasts for each magic type separately using the second, extended GLM, e.g., $(Magic_{A_pre}$
31 $> Control_{A_pre}) > (Magic_{A_post} > Control_{A_post})$. Please note that these contrasts are controlling for
32 potential time confounds.

33 All non-parametric permutation tests were performed using 5000 permutations, a cluster-
34 forming threshold of $p_{unc} = 0.001$ and family-wise error correction (FWE, $\alpha = 0.05$). We report
35 clusters that do not survive FWE-correction but exceed $k = 10$ voxels. Conjunction analyses based
36 on parametric statistics were performed using an uncorrected threshold of $p \leq 0.001$ and a cluster
37 threshold of $k = 10$. Brain areas in whole brain analyses were first identified using the atlasreader

1 toolbox (Notter et al., 2019), applying the Automated Anatomical Labeling atlas 3 (AAL3) (Rolls et
2 al., 2020).

3 **Univariate whole-brain control analyses**

4 We performed two control analyses for the whole-brain univariate results. First, it is possible that
5 the conjunction analyses examining differences between the magic trick types are confounded by
6 differences in the visual content of the videos at specific time-points. This is because we are not
7 only comparing videos showing “appearances”, “color changes” and “vanishes”, but also videos
8 showing “red objects”, “blue objects” and “no objects” at a specific time point, respectively (see
9 Figure 1C). To control for this possible stimulus-driven confound we compared responses between
10 magic and control videos showing similar visual contents at specific time points. We contrasted
11 responses to magic appearances with the control videos for vanishing tricks (as both videos show
12 a “red object” at the specific moments), and magic vanishes with the control videos for tricks
13 showing something appear (as both videos show no objects at the specific moments). No similar
14 match was possible for the color-changing videos. Responses similar to that of the original
15 contrasts, which used control videos showing similar actions, would be indicative that the observed
16 differential activity reflects specific surprise responses and not differences in the visual content of
17 the videos.

18 Second, we looked for condition-independent time effects in our whole-brain results. We
19 performed a mixed-effects model using the average beta estimates of the magic tricks of each
20 run, using the significant clusters of the prior-knowledge dependent contrast ($Magic_{pre} > Control_{pre}$)
21 $> (Magic_{post} > Control_{post})$. We hypothesized that any prior-knowledge dependent modulation
22 should show high activity in the runs before the explanation of the tricks and decreased activity
23 afterwards. In contrast, any results driven by either a general or set-wise drop in alertness should
24 show a gradual drop in the magic estimates across all runs of the experiment or within a set,
25 respectively. Accordingly, we used four predictors in the mixed model to explain the 12 magic
26 estimates based on their presentation order: one predictor modelled the revelation condition (i.e.,
27 1, 1, -1, -1, etc.), one predictor modelled the run number (i.e., 5.5 to -5.5 in 12 steps) and two
28 predictors modelled a decay in activity either before (i.e., 1, -1, 0, 0, etc.) or after the revelation of
29 the tricks (i.e., 0, 0, 1, -1, etc.). We expect the last three regressors to account for any time-
30 dependent variance.

31 **ROI definition**

32 For the ROI analyses, we defined 26 hypothesis-driven ROIs, separated into two groups.

33 First, we defined a set of 16 ROIs based on significant responses to magic videos from
34 previous experiments (Danek et al., 2015; Parris et al., 2009). The authors from Danek et al.,
35 (2015) kindly provided us with the corresponding parametric maps from their study which were

used to guide the ROI definition. We combined individual labels from a multi-modal parcellation of the human cortex (Glasser et al., 2016) to define the following 14 frontal and parietal ROIs that showed increased activity when viewing magic tricks: a posterior part of the dorsal anterior cingulate cortex (pdACC), an anterior part of the dorsal ACC (adACC), ventral ACC (vACC), inferior frontal junction (IFJ), left inferior frontal sulcus (IFS), Brodman area 6 (BA6), inferior premotor subdivision (6r), 8BM, anterior insula (AI), anterior ventral insula (AVI), inferior temporal gyrus temporo-occipital division (PH), left BA 46, left BA 8 and left inferior parietal cortex (IPC). The remaining two subcortical ROIs (caudate nucleus and left amygdala) were defined using the Freesurfer automatic parcellation (Fischl et al., 2002). A detailed list of the parcellations used to define these ROIs is in the supplementary section *Surprise-related region of interest definition*.

Second, we used a probabilistic map of visual fields (Wang et al., 2015) to define ten visual ROIs: primary visual cortex (V1), secondary visual cortex (V2), V3, V3A, V3B, human V4 (hV4), lateral occipital and ventral complex (LO and VO, respectively), intraparietal sulcus (IPS) and frontal eye-fields (FEF). All ROIs were defined in native space.

We included these visual ROIs to use in the decoding of the VOE type (i.e., appear, change and vanishing, see below) and to test for differences evoked by the different VOE types and the effect of prior knowledge. For example, areas of the ventral visual cortex are known to be responsive to color (Bartels and Zeki, 2000), while the lateral occipital complex is responsive to objects (Grill-Spector et al., 2001). As predictive coding approaches predict feature-specific prediction errors in functionally specialized regions, we expect unexpected color changes to affect color-responsive ROIs (e.g., hV4, VO and PH), and unexpected object appearances to affect object-responsive areas LO and VO, in line with recent imaging evidence (Jiang et al., 2016; Richter et al., 2018; Stefanics et al., 2019). Early visual areas (V1, V2, V3) were included to investigate possible top-down effects of prior knowledge in lower-level areas, when comparing the exact same videos before and after revelations, while the parietal (IPS) and prefrontal ROIs (FEF) were included due to their involvement in top-down voluntary attention (Corbetta and Shulman, 2002). See Figure 5A for a depiction of all 26 ROIs in an exemplary subject.

Multivariate pattern analysis (MVPA)

Apart from the univariate analyses testing for net signal differences, we further wanted to investigate which areas of the brain carry pattern information about the different types of VOE (unexpected appearance, feature change, and omission). To do so, we performed a series of multivariate pattern analyses (MVPA) on the 26 hypothesis-driven ROIs and a control ROI.

For the decoding analyses we computed a GLM in which every trial (i.e., video presentation) was modeled as a separate regressor to increase the number of data points for training and testing. All analyses were performed using a shrinkage linear discriminant analysis (LDA) on the de-meaned beta estimates of the individual trials (by the mean over all estimates, i.e., all trials, within each voxel) using the Python (version 3.8.13) package scikit-learn's class

LinearDiscriminantAnalysis (Pedregosa, 2011). To examine if any of the ROIs contained information about the different types of VOE (appear, color change, and vanish) we trained and tested a shrinkage LDA to predict the VOE types following a three-fold cross-validation scheme to ensure generalization across objects. We trained on the data of two objects (i.e., estimates from two sets, 48 trials) and tested on the third object (i.e., estimates from the third set, 24 trials). Significance testing of the decoding accuracies was done using a permutation analysis (1000 permutations) implementing the max statistic correction to correct for multiple comparisons (Nichols and Holmes, 2002). A control ROI (third ventricle) was included in the analysis, which should carry no information and thus reflect chance level.

Decoding analyses were performed separately for data before and after explanations of the magic tricks. Permutation-based corrected significance thresholds were 36.98% and 36.92% before revelation and after revelation, respectively. As both analyses were conducted using estimates based on the very same videos, we hypothesized that any significant difference in decoding accuracies between data before and after revelation would be indicative of decodable prior-knowledge dependent surprise signals. We tested for differences in decoding accuracies using paired t-tests between decoding using pre-revelation data and decoding using post-revelation data, only in those ROIs that showed significant decoding (corrected) using pre-revelation data.

Additionally, as an exploratory approach we performed a similar whole-brain searchlight analysis using a sphere to decode the magic types (4 mm radius), separately for data before and after explanation of the tricks using the SearchLight class implemented in nilearn (Abraham et al., 2014). The searchlight whole-brain accuracy maps were spatially smoothed with a 4 mm Gaussian kernel. A permutation-bootstrap hybrid method (in which each randomly generated accuracy map was also smoothed with a 4 mm Gaussian kernel) was used for significance testing and correction for multiple comparisons (Stelzer et al., 2013) using custom-made Python code.

Behavioral data

To test for differences in surprise ratings between videos before and after the explanation of the tricks we performed a 2 (before/after explanation) x 3 (magic, unusual and control videos) repeated measures ANOVA. We expected to see higher surprise ratings for magic videos compared to control videos and higher ratings for magic videos before compared to after the revelation of the methods. We also wanted to test if videos of magic and of unusual actions led to similarly high surprise ratings. Only this would allow us to use the unusual action videos as secondary comparison points for magic. Moreover, we tested for differences in surprise ratings of the magic videos for different objects and VOE types in 2 x 3 repeated-measures ANOVAs (rmANOVA) with the factors revelation (before/after) and object (ball/card/pencil) or VOE type (appear/change/vanish), respectively.

1 Eye tracking

2 **Data acquisition and preprocessing.** Gaze positions were measured using an MR-compatible
3 Eyelink 1000 (SR-Research, Ottawa, Canada) positioned at the rear end of the scanner bore at a
4 1000 Hz recording rate. Calibration of the eye tracker was performed at the beginning of the
5 experiment and drift correction was performed at the beginning of each run. If necessary, re-
6 calibration was performed at the beginning of a new run.

7 Eye tracking data (i.e., gaze position, blinks, and saccades) from 23 participants were
8 analyzed. Data from one participant was excluded because of technical problems during data
9 acquisition. Monocular gaze path data (x, y coordinates) were cleaned by removing values 150
10 ms before and after blinks, linearly interpolating the missing data (with an interpolation limit of 500
11 ms) and downsampling the data from 1000 Hz to 25 Hz (i.e., the framerate of our videos).
12 Identification of blinks and saccades was performed automatically using the Eyelink online parser
13 with default parameters (saccade detection threshold was 22 degrees/s).

14 Eye tracking data was used to examine if fMRI responses could be confounded by
15 systematic differences in gaze traces, saccades, or blinks during viewing of the videos. We
16 focused our eye tracking analyses to values within a -1 to +2 s window centered around the event
17 time from each video. We performed two types of tests. First, we used rmANOVAs to test for
18 differences in blink or saccade numbers across conditions. Second, we used correlation analyses
19 to test for differences in gaze traces.

20
21 **Gaze trace analysis.** To test for differences in gaze traces before and after the revelation of the
22 methods used in the magic tricks, we correlated the x and y positions of each video presentation
23 (eight presentations per video – two per run, two runs pre and post revelation) within a subject, for
24 each video separately (resulting in two 8×8 correlation matrices per video – one for x, one for y.).
25 Then, the correlation coefficients were transformed using the Fisher-z transformation and
26 averaged for the x and y traces. We pooled all values from comparisons *between* presentations
27 before and after the revelation of the tricks (i.e., the bottom left quadrant of the matrix) and those
28 *within* presentation before and after the revelation (see Figure 2E for an example) of the tricks
29 within a subject. The pooled correlation coefficients were then compared using paired t-tests.

30
31 **Blink and Saccade analysis.** Finally, to test for differences in blinks and saccades as a factor of
32 video condition, prior knowledge, and VOE, we compared the mean number of blinks and
33 saccades using two rmANOVAs, for blinks and saccades separately. The first rmANOVA had
34 video (magic and control) and revelation (before and after) condition as factors. The second
35 rmANOVA used magic data and had the VOE types (appear, color change, and vanish) and
36 revelation (before and after) as factors.

Inference statistics

Effect sizes for repeated measures ANOVAs and paired tests are presented as partial eta squared (η^2) and Cohen's d , respectively. Sphericity of rmANOVAs was tested using the Mauchly test. If sphericity was violated, degrees of freedom were adjusted using the Greenhouse-Geisser correction, the corresponding ϵ -correction factor is provided. Normality of data was assessed using the Shapiro-Wilk test (for paired tests and post-hoc tests). In case that data was normally distributed, we performed paired t-tests, otherwise we performed non-parametric Wilcoxon signed-rank tests instead. Correction for multiple comparison was performed using a step-down Holm-Bonferroni correction. Please note that in the post-hoc tests we corrected for each hypothesis separately (i.e., p-values for post-hoc tests of one factor are corrected independent of another factor or an interaction). In general, corrected p-values (p_{corr}) are reported in text, uncorrected p-values (p_{unc}) are reported in the corresponding tables. The threshold for statistical significance was set to 0.05 for all tests.

Results

Behavioral surprise ratings

Behavioral data was first tested for differences in surprise ratings for video condition (magic, control, and unusual videos) and revelation condition (before and after revelation) using a 2 x 3 rmANOVA (see Figure 2A). In sum, this analysis revealed that magic tricks were perceived as more surprising than the control videos, and that surprise ratings dropped after explanation of the tricks. In detail: the analysis revealed significant main effects for both factors (video condition: $F(2,46) = 57.494$, $p_{unc} < 0.001$, $\eta^2 = 0.478$, $\varepsilon = 0.932$, revelation: $F(1,23) = 103.242$, $p_{unc} < 0.001$, $\eta^2 = 0.211$, $\varepsilon = 1$) and interaction: $F(2,46) = 43.081$, $p_{unc} < 0.001$, $\eta^2 = 0.134$, $\varepsilon = 0.641$). As expected, post-hoc Wilcoxon signed-rank tests revealed that magic videos were more surprising than control and unusual videos, before ($W = 1$, $p_{corr} < 0.001$, Cohen's $d = 4.07$ and $W = 8$, $p_{corr} < 0.001$, Cohen's $d = 2.06$, respectively) and after explanation of the tricks ($W = 5$, $p_{corr} < 0.001$, Cohen's $d = 1.55$ and $W = 26$, $p_{corr} < 0.001$, Cohen's $d = 0.812$, respectively), and that all video conditions were more surprising before compared to after the revelation (all three video conditions: $p_{corr} \leq 0.001$). Unusual videos were more surprising than control videos pooled across runs ($W = 42$, $p_{corr} = 0.006$, Cohen's $d = 0.59$), but significantly so only in the first two runs ($W = 26$, $p_{corr} = 0.002$, Cohen's $d = 0.71$) (see Supplementary Table S1 for a detailed report of all post-hoc tests). Because of the large and significant differences in surprise ratings between the magic and unusual videos that preclude a meaningful comparison of corresponding neural responses, we report here no further analysis of the data from the unusual videos. For completeness, we show corresponding contrasts in the supplementary results (see supplementary Figure S3 and Tables S3-6).

Average surprise ratings of magic videos were consistently high before the explanation of the tricks (all group means > 3) and decreased afterwards (all group means < 2.5) (see Figure 2A, B and C). We further performed two rmANOVAs using only ratings from magic videos, to test for possible differences in VOE types and the objects used. Both rmANOVAs included the revelation condition as a factor and showed a significant decrease in surprise rating after the revelation of the tricks (as expected from the previous analysis). We additionally found a significant main effect for the magic type ($F(2,46) = 13.8$, $p_{unc} < 0.001$, $\eta^2 = 0.032$, $\varepsilon = 0.848$) and an interaction of magic type and revelation condition ($F(2,46) = 11.02$, $p_{unc} < 0.001$, $\eta^2 = 0.018$, $\varepsilon = 0.795$) (see Figure 2B). Post-hoc tests showed that appearances were rated less surprising than color changes and disappearances pre revelation ($W = 12$, $p_{corr} < 0.001$, Cohen's $d = -0.835$ and $W = 35.5$, $p_{corr} = 0.002$, Cohen's $d = -0.631$, respectively). No difference between color change and vanish magic tricks was observed (see Supplementary Table S2). No main effect for objects nor an interaction of object and revelation condition were found (all F -values < 2 , all $p_{unc} > 0.15$ and $\eta^2 < 0.012$) (see Figure 2C).

Eye-tracking results

Eye-tracking data were tested for systematic differences in gaze traces, saccades, and blinks during viewing of the videos (see also Supplementary section *Eye-tracking results* and Figure S2). First, for each magic video we computed the correlation of gaze path between video presentations in all subjects. We then averaged the transformed values (using Fisher-z transformation) of correlations between pre- and post-revelation presentations and within pre- and post-revelation presentations and tested for differences using paired t-tests. Gaze traces between pre- and post-revelation runs were similar (see an example visualization in Figure 2D) and no significant difference of correlations of gaze traces was observed (all $p_{unc} > 0.2$) (see Figure 2E for an example). Moreover, the number of saccades and blinks around the VOE times were similar between experimental conditions and only revealed small differences we deem unlikely to have affected the imaging results.

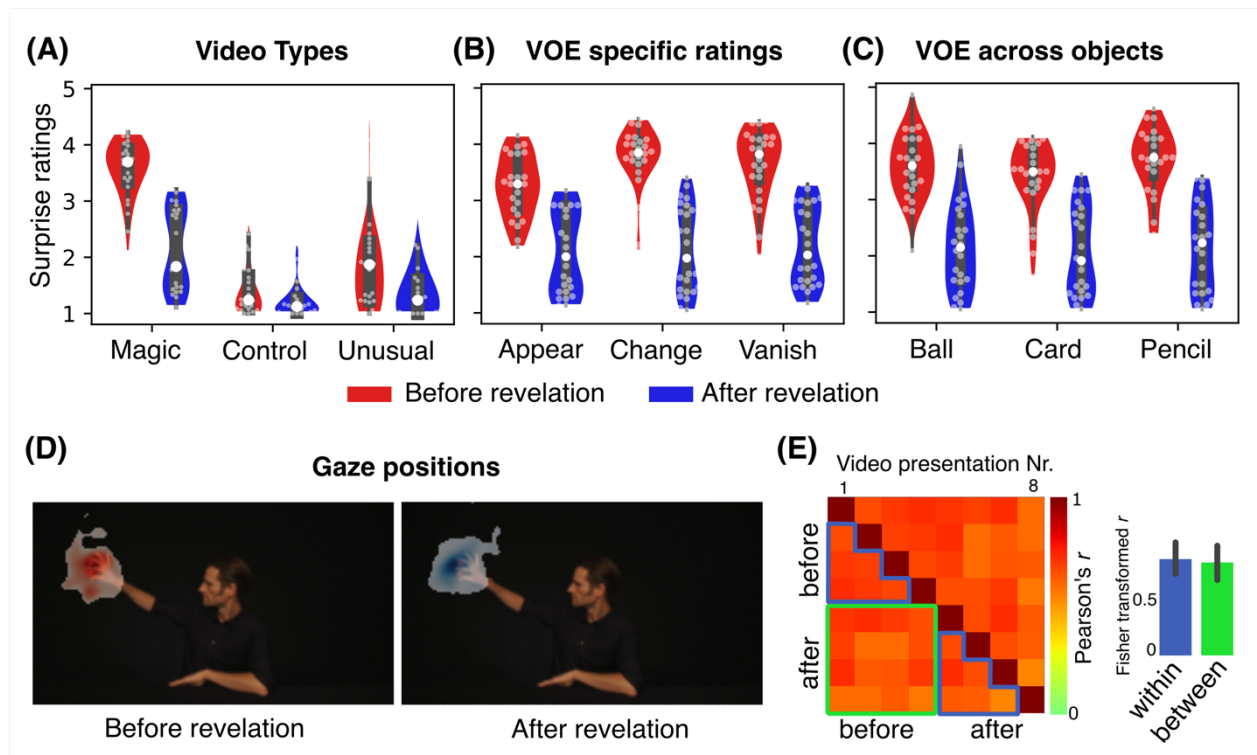


Figure 2. Shown are behavioral surprise ratings separated for the different video types (A), for magic videos across VOE types (B) and across objects (C). Surprise ratings before the revelations of the magic tricks (red) were consistently higher than after the revelations (blue). D, Exemplary group gaze positions for the same magic video in the moment a ball appears, before (left) and after (right) the revelation of the trick. E, Correlation matrix for gaze positions around the moment of magic (-1 s and +2 s) from an exemplary video. To test for differences between gaze paths before and after the revelation of the tricks, we compared the average Fisher z-transformed correlations between gaze paths of all

combinations of presentations (from presentation Nr. 1 to Nr. 8) *between* pre- and post-presentations (marked green) and those *within* pre- and *within* post-presentations (marked blue), as shown in the bar plot.

Univariate whole-brain analysis

Whole-brain surprise responses

To investigate neural correlates of high-level VOE when viewing seemingly impossible events, we first compared whole-brain responses to the magic and matched control videos before the revelation of the magic tricks ($Magic_{pre} > Control_{pre}$). This contrast revealed several clusters of activity in frontal and parietal cortices, as well as subcortical areas, such as the caudate nucleus, largely in line with previous studies (Danek et al., 2015; Parris et al., 2009) (see Figure 3A and Table 1). In particular, large clusters of activity were observed in the medial part of Brodmann area 8 (preSMA), the dorsal and ventral anterior cingulate cortex (dACC and vACC), caudate nucleus (CN), and the posterior parietal cortex (PPC, especially the superior parietal lobe and the precuneus). No lower-level sensory area was differentially modulated in view of the unexpected events.

Generic surprise responses

To specifically test for generic surprise responses in the brain showing a significant involvement of all three types of VOE, we performed a conjunction analysis combining all different VOE types before the revelation of the tricks ($Magic_{A_pre} > Control_{A_pre} \cap Magic_{C_pre} > Control_{C_pre} \cap Magic_{V_pre} > Control_{V_pre}$). While no cluster survived FWE correction, using a voxel-wise uncorrected threshold of $p \leq 0.001$ and $k = 10$ we found small clusters in the dorsal anterior cingulate cortex (dACC) bilaterally and posterior parietal cortex (precuneus) revealing generic responses (see Figure 3B, left and Table 1). Also, using a more liberal threshold of $p_{unc} \leq 0.005$ (and $k = 10$), we further observed generic activity in the vACC/mPFC.

Table 1. Significant clusters of activity from the whole-brain contrast comparing responses between magic videos and matched controls ($Magic_{pre} > Control_{pre}$, thresholded at $p \leq 0.001$ and $k = 10$, uncorrected) and the conjunction analysis testing for common areas involved in the processing of violation of expectations (thresholded at $p_{unc} \leq 0.001$ and $k = 10$, uncorrected). P-values show cluster statistics. k = cluster size, T = t statistic at peak voxel, x , y , z = peak voxel MNI coordinates [mm]. dAAC: dorsal anterior cingulate cortex; PFC: prefrontal cortex; vACC: ventral ACC; mPFC: medial PFC; PCC: posterior cingulate cortex; DLPFC: dorso-lateral prefrontal cortex; preSMA: supplementary motor area.

Brain region	AAL atlas labels	p(FWE)	p(unc)	k	T	x	y	z
$Magic_{pre} > Control_{pre}$								
Posterior parietal cortex	Occipital_Mid_L	0.0004	<0.0001	1871.0	7.16	-34	-82	34
	Precuneus_R	-	-	-	6.87	4	-68	52
	Precuneus_L	-	-	-	6.48	-10	-68	44
dAAC	Cingulate_Ant_L	0.0048	0.0005	886.0	6.46	-6	30	20

	Cingulate_Ant_R	-	-	-	6.05	4	32	24
	Cingulate_Mid_R	-	-	-	5.51	2	24	32
L. Parieto-occipital sulcus	Cuneus_L	0.0594	0.0082	162.0	6.18	-14	-62	22
	Calcarine_L	-	-	-	4.21	-12	-60	14
R. Intraparietal sulcus	Occipital_Mid_R	0.1752	0.0284	74.0	5.53	36	-80	32
L. anterior PFC	Frontal_Sup_2_L	0.1982	0.0332	67.0	5.41	-26	58	4
L. anterior intraparietal	Postcentral_L	0.0548	0.0074	172.0	5.2	-46	-36	54
	Parietal_Inf_L	-	-	-	4.47	-42	-38	44
vACC/mPFC	Cingulate_Ant_L	0.0548	0.0074	172.0	5.13	-4	48	2
L. Superior frontal gyrus	Frontal_Sup_2_L	0.0522	0.007	179.0	5.02	-22	12	58
	Frontal_Sup_2_L	-	-	-	3.7	-22	0	52
L. anterior middle frontal gyrus	Frontal_Mid_2_L	0.1722	0.028	75.0	5.01	-34	54	14
L. PCC	no_label	0.5014	0.1159	25.0	5.01	10	-36	2
L. Thalamus	Thalamus_L	0.5456	0.1334	22.0	4.86	-6	-6	10
L. Hippocampus	no_label	0.5158	0.1214	24.0	4.76	-18	-22	-10
L. Caudate nucleus	Caudate_L	0.5946	0.1558	19.0	4.73	-12	16	2
L. Postcentral gyrus	Postcentral_R	0.1134	0.0173	103.0	4.71	52	-16	54
	Postcentral_R	-	-	-	4.07	52	-18	46
	Postcentral_R	-	-	-	3.72	46	-22	42
R. Superior frontal gyrus	Frontal_Sup_2_R	0.2394	0.0417	57.0	4.61	28	10	62
L. PCC	no_label	0.2832	0.0511	49.0	4.48	-2	-30	26
R. Thalamus	Thalamus_R	0.7188	0.2273	13.0	4.17	12	-8	12
L. Fusiform gyrus	Fusiform_L	0.6338	0.1754	17.0	4.16	-24	-54	-16
	Cerebelum_6_L	-	-	-	3.75	-24	-62	-18
L. Ventral occipital cortex	Lingual_L	0.7188	0.2273	13.0	4.1	-2	-78	-8
L. Anterior intraparietal	Parietal_Inf_L	0.5306	0.1277	23.0	3.84	-40	-50	50
L. DLPFC	Frontal_Mid_2_L	0.762	0.2623	11.0	3.81	-36	26	32
Conjunction of $Magic_{pre}$ – $Control_{pre}$ per VOE type								
dACC/preSMA	Supp_Motor_Area_L	0.722	0.096	38	4.02	-4	16	50
Precuneus	Precuneus_L	0.885	0.162	26	3.99	-8	-68	48
	Precuneus_R	0.722	0.096	38	3.94	8	-68	48
dACC	Cingulum_Mid_L	0.966	0.254	17	3.7	-6	18	40
	Cingulum_Mid_L	-	-	-	3.31	-4	26	38
R. Postcentral gyrus	Precentral_R	0.986	0.318	13	3.5	46	-14	54

1 Specific surprise responses

2 After establishing what areas are generally involved in the processing of violation of expectations
3 (i.e., commonly active in seemingly impossible appearances, disappearances, and color
4 changes), we looked for VOE type specific differential activity in the brain (e.g., areas responsive
5 to something unexpected appearing but not disappearing or changing color). Beyond the
6 systematic generic activation of frontoparietal areas described above, all three types of VOE

evoked differential responses in posterior visual areas, but in a segregated and even partially opposite manner (see Figure 3C and Supplementary Table S8). For example, while appearances and color changes evoked an increase of activity in ventral visual areas, the disappearance of objects reduced activity in overlapping regions. To better visualize this diverse modulation of sensory areas by the different VOE, we tested which areas were significantly more activated by one type of VOE than by the other two using a conjunction test (e.g., test for appear-specific responses: $Magic_{A_pre} > Magic_{C_pre} \cap Magic_{A_pre} > Magic_{V_pre}$) (see Figure 3D and E and Supplementary Table S7). As hypothesized, activity related to objects appearing were observed in early visual areas and in higher-level visual areas of the lateral occipito-temporal cortex (peak coordinates: $x = 30, y = -88, z = 22$ and $x = -22, y = -78, z = 44$), whilst activity related to color changes were observed specifically in color-responsive ventral areas of the fusiform gyrus (FFG, peak coordinates: $x = -30, y = -50, z = -16$ and $x = 32, y = -54, z = -14$). Finally, vanishing objects evoked significant responses – among others – in the anterior parts of the calcarine sulcus, close to the parieto-occipital sulcus (peak coordinates: $x = -20, y = -64, z = 12$ and $x = 18, y = -76, z = 8$).

Control analysis for visual content

As these VOE type-specific patterns of activity are located predominantly in visual processing areas, they are potentially related to general visual differences between the conditions tested. The compared videos are not only showing “appearances”, “color changes” and “vanishes”, but also “red objects”, “blue objects” or “no objects”, respectively. To test if these VOE type-specific responses are confounded by different visual input, we performed a control analysis comparing neural responses between 1) videos showing red objects, either as the product of a magic trick (appearances) or as a control to other tricks (vanishes) ($Magic_{A_pre} > Control_{V_pre}$) and 2) between videos showing no object, either as the product of a magic trick (vanishes) or as a control to the other tricks (appearances) ($Magic_{V_pre} > Control_{A_pre}$). Differential responses in posterior visual areas to these control contrasts were similar to the results of the corresponding conjunction analyses (for appearances and vanishes), suggesting that violation-specific responses observed in posterior visual areas are unlikely to be driven merely by visual content (see overlays in Supplementary Figure S5). However, if these signals in visual areas reflect specific prediction errors based on different VOE, we would additionally expect them to be modulated by prior knowledge and to show decreased responses after the explanation of the tricks.

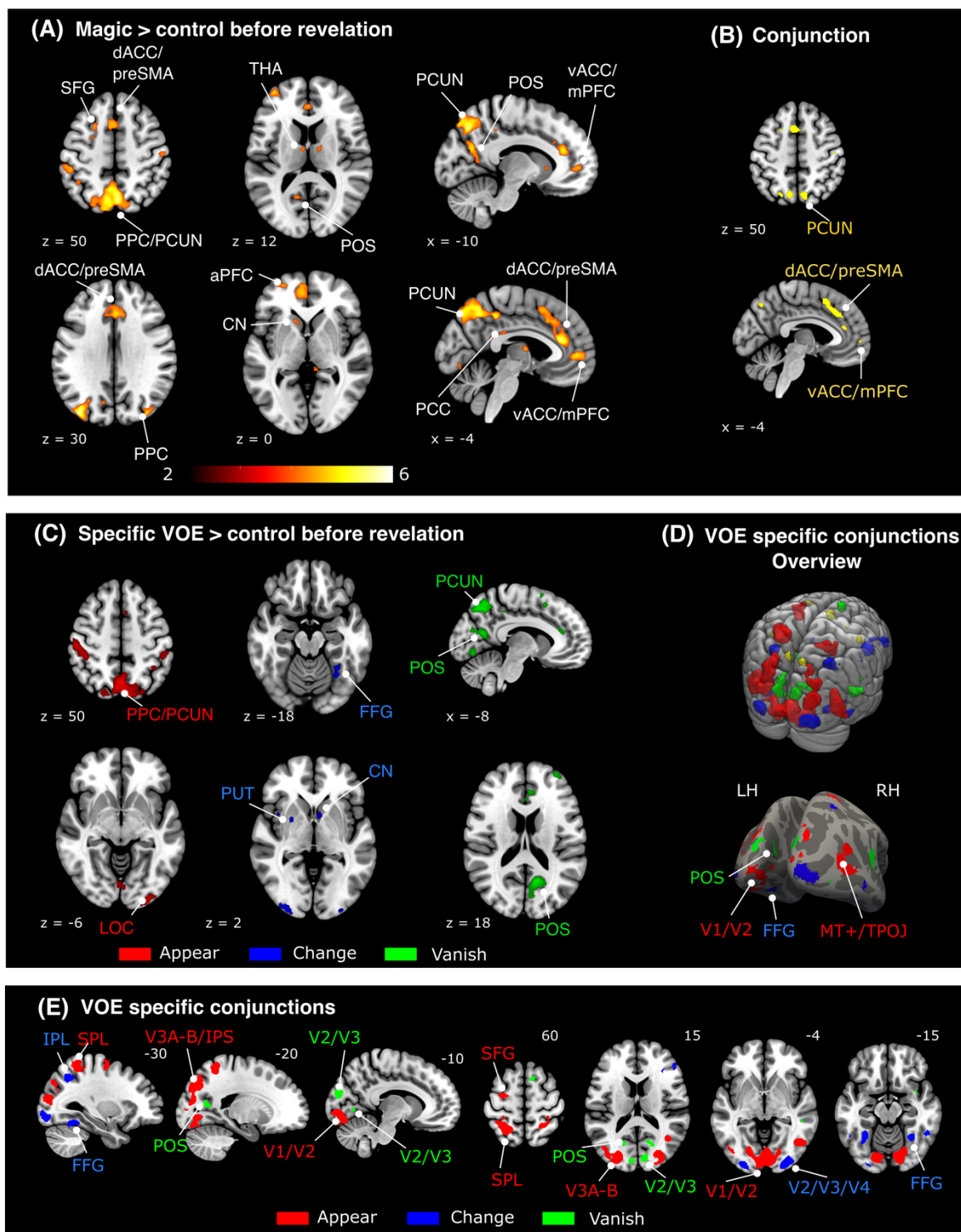


Figure 3. **A**, shown are active regions during viewing magic tricks compared to matched control videos before the participants knew how the tricks were performed (thresholded at $p_{unc} < 0.001$ and $k = 10$, uncorrected). **B**, shown are results from the conjunction analyses testing for generic activity (left, yellow) (at $p_{unc} < 0.005$). Only the dorsal ACC and preSMA and the precuneus (PCUN) revealed generic responses to high-level VOE at $p_{unc} < 0.001$. **C**, VOE specific activity revealing several posterior visual areas responsive to appearances (red), color changes (blue) and objects disappearing (green). **D**, overview of all conjunctions results in the MNI152 template volume (upper) and of the VOE-specific responses projected on an average surface (fsaverage) (lower). **E**, shown are VOE-specific conjunction results. **Abbreviations:** ACC: anterior cingulate cortex; dACC: dorsal ACC; vACC: ventral ACC; CN: caudate nucleus; aPFC: anterior prefrontal cortex; FFG: fusiform gyrus; IPL: inferior parietal lobe; IPS: intraparietal sulcus; MT+: motion area; MT: mPFC: ventral prefrontal cortex; PCC: posterior cingulate cortex; PCUN: precuneus; POS: parieto-occipital sulcus; PPC: posterior parietal cortex; SFG: superior frontal gyrus; SMA: supplementary motor area; SPL: superior parietal lobe; THA: thalamus; TPOJ: temporo-parietal-occipital junction. LH: left hemisphere; RH: right hemisphere.

Prior-knowledge dependent whole-brain responses

To investigate the effect of prior knowledge on brain responses, we provided participants with the methods behind the magic tricks. Surprisingly, neural activity after explanation of the tricks ($Magic_{post} > Control_{post}$) were similar to that before the explanation of the tricks (see Figure 4A and Supplementary Table S10). Thus, areas related to the processing of surprising events remained significantly active in view of VOE also after the explanation of the tricks (i.e., dACC, caudate nucleus, anterior insula). The interaction contrast comparing surprise responses before and after providing the explanations ($Magic_{pre} > Control_{pre}$) $>$ ($Magic_{post} > Control_{post}$) revealed only a small number of areas showing prior-knowledge dependent modulations (see Figure 4B and C and Table 2). We observed higher activation in a large cluster of the medial prefrontal cortex (mPFC) and ventral ACC (peak coordinates: $x = 2$, $y = 46$, $z = -10$) and right posterior cingulate cortex (PCC) with FWE-correction (peak coordinates: $x = 6$, $y = -46$, $z = 8$). These regions hence decreased magic-related activity after the revelation. These areas overlap with and constitute of subset of the activity observed with $Magic_{pre} > Control_{pre}$ before the revelation of the tricks.

Interestingly, the observed decrease of activity in the vACC/mPFC and PCC after revelation of the tricks coincides with the midline core areas of the default mode network (DMN), whilst parietal areas of the dorsal attention network (DAN) showed increased activity. This indicates that our findings are unlikely a result of a decrease in attention, as the DAN is known to direct top-down attention (Corbetta and Shulman, 2002). A visualization of these patterns of activity, together with the DMN and DAN is shown in Figure 4C.

No prior-knowledge dependency of trick-specific modulations

While prior-knowledge driven signals overlapped with neural activations to unspecific VOE ($Magic_{pre} > Control_{pre}$), none of the previously observed trick-specific visual areas showed a modulation as factor of prior knowledge (even when using a more liberal threshold of $p < 0.001$ uncorrected, and a cluster size of $k = 10$, see Table 2 and Figure 4B and C). Further, the contrasts investigating the effect of prior knowledge for each VOE type separately revealed similar response

patterns to those reported in the main contrasts (see supplementary Figure S4). This indicates that visual areas responsive to a specific VOE (e.g., an object changing color) were not modulated by prior knowledge.

Controlling for time effects in net responses

To rule out possible time confounds in results comparing responses before and after revelation of the tricks, we inspected how the individual values changed across time as a control analysis. We extracted betas estimates corresponding to all magic and matched control presentations from the suprathreshold clusters from the prior-knowledge interaction contrast, averaged them over subjects and calculated the same contrast within each run separately. We then ran a mixed-effects model on the contrast values with four predictors: a pre-post predictor (i.e. 1, 1, -1, -1 for each set), a constant decay predictor over all 12 fMRI runs (i.e. 5.5 to -5.5 in 12 steps) and two time-decay predictors for values before and after revelation separately (i.e. 1, -1, 0, 0 and 0, 0, 1, -1, for each set). We hypothesized that areas showing prior-knowledge dependent activity should have a significant pre-post predictor. The pre-post predictor was significant in all clusters (all $p_{corr} < 0.05$). However, in some clusters a significant amount of variance was also explained by the time-decay regressors, showing that the activity in some clusters had a contribution of time. Yet, the fact that in all clusters the pre-post regressor remained significant despite inclusion of the time-regressor shows that knowledge-dependent effects were true (see supplementary Figure S6 and Table S9).

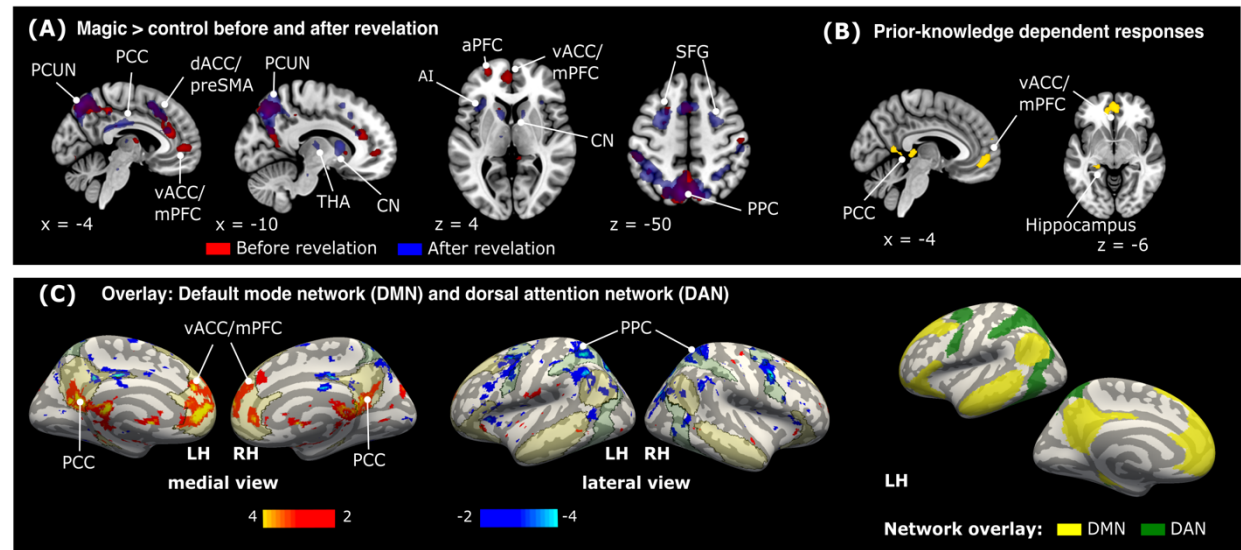


Figure 4. Prior knowledge dependent modulation of brain responses. **A**, overlay of *Magic* > *Control* before (red) and after (blue) explanation of the magic tricks. Both contrasts reveal a consistent activation of surprise related areas, such as the dorsal ACC (dACC), caudate nucleus (CN), and posterior parietal cortex (PPC). **B**, the difference of both contrasts (i.e., the interaction between video condition and revelation) revealed only few areas that were significantly more active before the explanation of the tricks: the posterior cingulate cortex (PCC), the ventral ACC/medial prefrontal cortex (vACC/mPFC) and left Hippocampus. Threshold at $p_{unc} < 0.001$ and $k = 10$. **C**, shown are the interaction testing for prior-knowledge dependent responses (thresholded from $t = 2$ to $t = 4$, red to yellow and $t = -2$ to $t = -4$, dark blue to light blue) and the default mode network (DMN) and the dorsal attention network (DAN) as transparent overlay projected on to an average brain surface (fsaverage). The network overlays are provided by (Yeo et al., 2011). **Abbreviations:** AI: anterior insula; aPFC: anterior prefrontal cortex; LH: left hemisphere; PCUN: precuneus; preSMA: supplementary motor area; RH: right hemisphere; SFG: superior frontal gyrus; THA: thalamus.

Table 2: Results of contrasts comparing neural responses to VOs before and after explanation of the tricks (thresholded at $p_{unc} < 0.001$ and $k = 10$, uncorrected). P-values show permutation-based cluster statistics. K = cluster size, T = t statistic at peak voxel, x, y, z = peak voxel MNI coordinates [mm]. ACC: anterior cingulate cortex; vACC: ventral ACC; mPFC: medial prefrontal cortex; POS: parieto-occipital sulcus; PCC: posterior cingulate cortex; dACC: dorsal ACC.

Brain region	AAL atlas labels	p(FWE)	p(unc)	k	T	x	y	z
vACC/mPFC	Frontal_Med_Orb_R	0.0068	0.0007	538.0	5.72	2	46	-10
	Frontal_Med_Orb_L	-	-	-	4.96	-10	42	-6
R. Cuneus (V3)	Cuneus_R	0.25	0.0385	52.0	5.23	8	-86	34
L. POS	Precuneus_L	0.1082	0.014	101.0	4.99	-8	-56	8
L. PCC	Calcarine_L	-	-	-	4.31	-6	-48	4
R. PCC	Precuneus_R	0.0424	0.005	185.0	4.91	6	-46	8
	no_label	-	-	-	4.59	-4	-30	6
	no_label	-	-	-	4.56	8	-36	4
L. Temporal pole	Temporal_Pole_Sup_L	0.7558	0.2264	12.0	4.67	-32	14	-28
L. dACC	Cingulate_Ant_L	0.6462	0.1596	17.0	4.53	-10	38	14
L. Hippocampus	Hippocampus_L	0.665	0.1695	16.0	4.21	-22	-32	-6
L. Hippocampus	Hippocampus_L	0.5714	0.1262	21.0	4.14	-24	-18	-14

Univariate ROI analysis

To complement the whole-brain analysis, we defined a set of 26 hypothesis-driven regions-of-interest (ROIs) from which we extracted and analyzed parameter estimates (10 visual and 16 surprise-related ROIs based on prior literature) from unsmoothed data in native space and tested them with paired tests (correcting for the number of ROIs). Overall, our ROI-based analyses confirmed our above whole-brain findings. None of the lower-level visual cortices showed a significant increase of neural responses during the magic compared to the matched control condition before the explanation of the tricks (all $p_{unc} > 0.24$). In contrast, significant responses to magic were observed in the visual parietal ROI IPS ($t(23) = 2.746$, $p_{unc} = 0.012$, Cohen's $d = 0.53$),

as well as in several higher-level surprise-related ROIs, especially in adACC ($t(23) = 5.32$, $p_{corr} = 0.001$, *Cohen's d* = 0.917), as well as in vACC, pvACC, BA6, BA46 and caudate nucleus (all $p_{unc} < 0.05$, see detailed results in supplementary Table S14-17). Only two ACC ROIs (adACC and vACC) and 8BM (directly superior to the adACC) showed a decrease of activity after the revelation of the tricks ($t(23) = 3.43$, $p_{unc} = 0.002$, *Cohen's d* = 0.41, $W = 79$, $p_{unc} = 0.042$, *Cohen's d* = 0.42 and $t(23) = 2.54$, $p_{unc} = 0.01$, *Cohen's d* = 0.41, respectively). None of the visual ROIs showed a similar response pattern. In contrast, VOE-specific responses were largely constrained to visual ROIs (see supplementary section *univariate ROI results*).

Multivariate pattern analysis

Our univariate analyses revealed generic responses (modulated by prior knowledge) in fronto-parietal areas, while showing trick-specific responses in posterior sensory areas (unaffected by prior knowledge). However, differences in specific VOE and prior-knowledge modulations might also be reflected in activation patterns and not only in net signal differences. Accordingly, we complemented our univariate analysis by a multivariate pattern analysis to investigate if specific information about VOE types were present in the activity patterns of our surprise-related fronto-parietal ROIs (which showed only generic responses to magic). Further, we performed this analysis separately for data before and after the explanation of the tricks, to investigate if posterior sensory areas (which showed trick-specific effects) might show a difference in decoding accuracies before and after revelation of the tricks.

We followed a three-fold cross-decoding approach to generalize each magic type across objects: we trained a linear classifier to classify specific VOE types, appear (A), color change (C), and vanishes (V), based on the estimates of the magic videos from two objects (48 trials, from two sets) and tested on the estimates of the remaining object (24 trials, from one set).

Visual ROIs

Decoding of the specific VOE types across objects was possible in all posterior visual ROIs (V1, V2, V3, hV4, V3A, V3B, LO, VO, IPS) before the explanation of the magic tricks (all corrected p -values < 0.001 ; corrected using a permutation maximal statistic, Nichols and Holmes, 2002). After revelation of the method behind the magic tricks, decoding accuracies significantly dropped in most ROIs (V1, V2, V3, LO and IPS) (all $p_{corr} \leq 0.05$; Holm-Bonferroni corrected for the number of visual ROIs), being below threshold in LO and IPS (see Figure 5B and supplementary Table S20). Additionally, we found uncorrected significant differences in hV4 and V3B ($p_{unc} = 0.05$ and $p_{unc} = 0.016$, respectively).

1 **Surprise related ROIs**

2 In contrast, no surprise-related ROI (nor the control ROI) showed significant above-chance
3 decoding accuracies between VOE types using the permutation-max statistic for correction, except
4 for the PH ROI ($p_{corr} < 0.001$), an area located in the inferior temporal sulcus (temporo-occipital
5 division), using data before revelation of the tricks. Uncorrected significant above chance decoding
6 (chance level = 33%) was observed in IFJ, AI, pdACC, BA6, BA8, and BA46 using data before
7 revelation. Differences in decoding accuracies before and after revelation were observed only in
8 PH ($p_{corr} = 0.005$, Holm-Bonferroni corrected for the number of ROIs that could significantly decode
9 before revelation) (see Figure 5B and supplementary Table S20) and in BA8 (corrected for the
10 number of ROIs that could *not* significantly decode the VOE type using pre-revelation data, i.e.,
11 16 ROIs) (see Figure 5B and supplementary Table S21).

12 Therefore, while frontoparietal areas showed a generic and surprise-dependent
13 involvement in processing VOE in the univariate analyses, they carried no or only weak information
14 (i.e., in BA8) as to what exactly happened. In contrast, information about specific VOE were
15 observed in all posterior sensory areas across the visual hierarchy. An in-depth view of the
16 informative ROIs is shown in the confusion matrices in Figure 5C. They reveal that appearances
17 and color-changes were consistently decodable above chance and that most confusions
18 happened between them and less so with vanishing events.

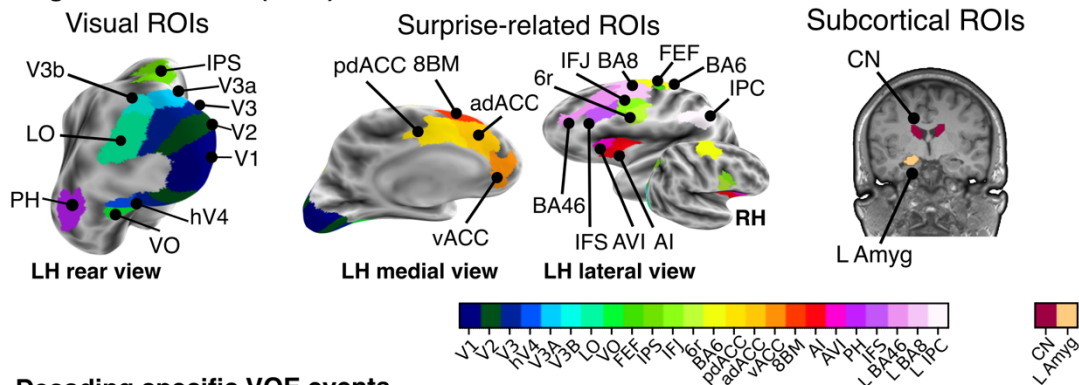
19 **Searchlight analysis**

20 To complement the ROI decoding analysis that used hypothesis-driven visual and surprise-related
21 ROIs, we performed a whole-brain searchlight analysis to better visualize the results and explore
22 if additional brain areas differentially respond to specific VOE. As in the ROI decoding analysis,
23 we trained linear classifiers to classify specific VOE (i.e., appear, change, vanish) in a three-fold
24 approach separately for both, magic estimates before and after revelation of the tricks. Results of
25 the whole-brain searchlight analysis revealed that only posterior visual areas of the brain could
26 significantly decode the magic type in both conditions (using a permutation-bootstrap hybrid
27 correction method, Stelzer et al., 2013). As shown in Figure 5D, significant decoding was possible
28 in large areas of the visual cortex, including most of the occipital cortex and small parts of the
29 temporal and parietal cortex. Crucially, decoding accuracies before the explanation of the magic
30 tricks were significant in more voxels and larger clusters compared to after explanation. Significant
31 decoding before the explanation of the tricks extends to parts of temporal and parietal cortex,
32 whereas significant decoding after revelation is largely restricted to posterior visual areas. In sum,
33 and in contrast to the univariate results that showed no net modulation as a function of prior
34 knowledge (and surprise) in visual areas, the differences in decoding accuracies using data before
35 and after explanation of the tricks suggests that visual areas are indeed sensitive to changes in
36 knowledge and encode specific expectations.

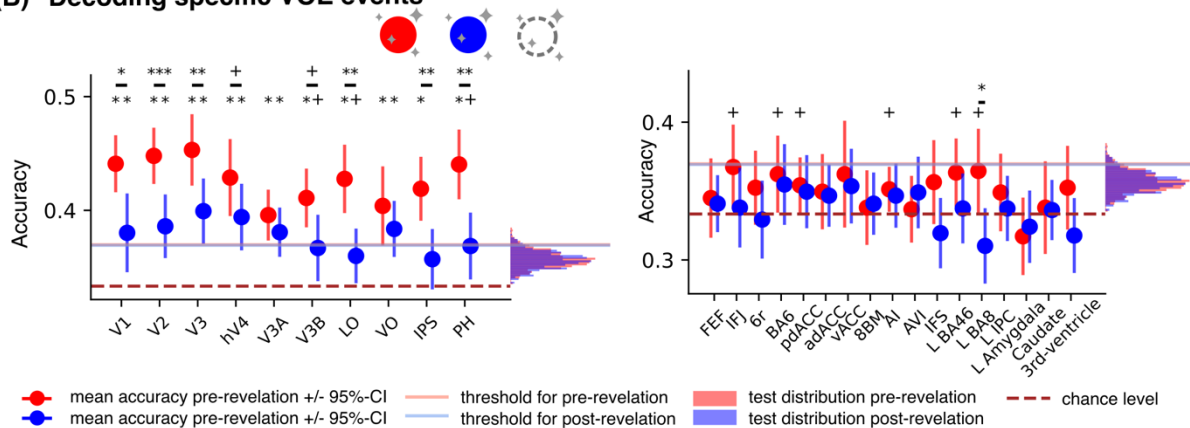
Controlling for time effects in pattern activity

Arguably, suprathreshold decoding of magic effects in posterior sensory areas may reflect general stimulus differences in the moment of magic between the different magic trick types independent of the object used, as appear videos systematically showed red objects, color changing videos blue objects and vanishing videos no object. However, the observed significant differences in decoding accuracy before and after revelation suggest surprise-dependent modulations of activity patterns, as these differences are present in view of the very same videos. Yet, these differences could reflect time- and/or design-related confounds, such as a general decrease of attention and alertness over time. However, since we did not find any significant changes in univariate comparisons in posterior visual areas and we observed a general increase in parts of the dorsal attention network, we believe that our results are not confounded by time and/or design related factors. Nonetheless we performed control analyses comparing decoding accuracy of objects present or absent in control videos in the pre vs post revelation phase. These analyses showed no modulation of time in control videos (detailed results can be found in the Supplementary Materials Section *MVPA control analyses*).

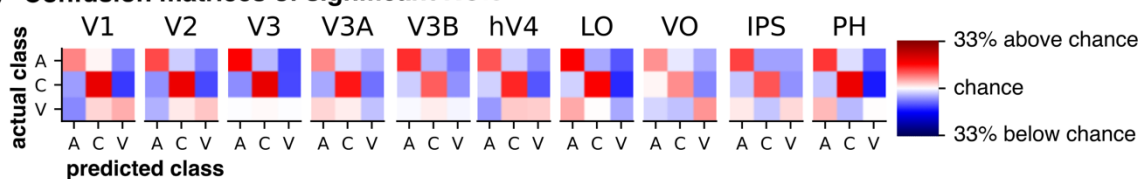
(A) Regions of Interest (ROIs)



(B) Decoding specific VOE events



(C) Confusion matrices of significant ROIs



(D) Searchlight analysis

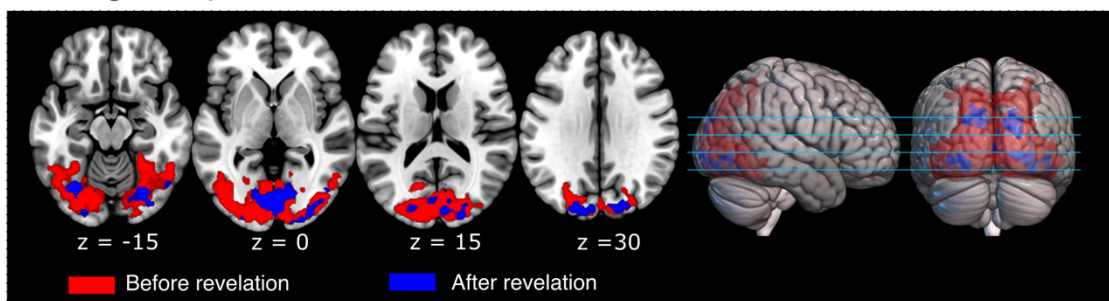


Figure 5. Results of decoding analyses **A**, regions of interest (ROIs) used in the experiment. Visual ROIs are shown left, all of which (except for the inferior temporal area PH) were defined using a probabilistic map of visual areas (Wang et al., 2015). Surprise-related ROIs presented on the right panel were defined based on previous results (Danek et al., 2015; Parris et al., 2009) using a multimodal parcellation atlas (Glasser et al., 2016), except for frontal eye field (FEF), which was also defined using the probabilistic atlas from Wang et al., (2015) (see supplementary section *Surprise-related region of Interest definition* for more information). The two subcortical ROIs, caudate nucleus, and left amygdala were defined using the Freesurfer automatic parcellation (Fischl et al., 2002). **B**, shown are the decoding accuracies for decoding the VOE types over objects in a three-fold cross-decoding approach in our theory-driven ROIs (left: ROIs that significantly decoded the VOE type using pre-revelation data, right: ROIs that did *not* significantly decode the VOE type using pre-revelation data). Decoding was performed with data before (red) and after (blue) revelation. All statistics were corrected for multiple tests by using the max-statistic correction across all ROIs (Nichols and Holmes, 2002) **C**, confusion matrices from decoding VOE types before revelation for ROIs that could significantly decode VOE types before revelation. It appears that VOEs due to unexpected color changes are predicted best, while the VOEs due to objects vanishing is predicted less than VOEs due to objects appearing and changing. **D**, whole-brain searchlight decoding results. We can significantly decode VOE types in the majority of visual cortex before revelation and less so after revelation (correcting using a permutation-bootstrap hybrid method, Stelzer et al., 2013).

1 Discussion

2 In this fMRI study, we used naturalistic video stimuli showing magic tricks and matched control
3 actions to investigate responses to violation of expectations (VOE) of deeply held beliefs about
4 the physical world. We used three distinct magic types (object appearance, object disappearance,
5 and feature-change) that were presented with and without prior knowledge about the underlying
6 deceptive methods (i.e., sleights-of-hand). Each magic type was presented using three distinct
7 objects to allow for object-invariant classification of magic types. We looked for 1) generic
8 prediction error responses to perceived violation of physical principles, 2) specific responses to
9 the different magic types and, 3) effects of the viewers' prior knowledge on prediction error
10 processing, for both generic and specific responses.

11 Our results revealed a hierarchy of surprise signals. First, we observed generic effects of
12 world-model VOE (i.e. common to all magic types) in several clusters of the prefrontal and parietal
13 cortex (such as the dorsal and ventral ACC and the posterior parietal cortex). Then, differential
14 activity specific to the different types of magic was evoked predominately in posterior visual areas
15 of the occipital and parietal cortex. These specific prediction error signals were evident in the
16 univariate analyses and in decoding of the magic types, both of which were confined to posterior
17 areas across the visual hierarchy. Finally, following explanation of the tricks, responses were
18 largely unaffected by participants knowledge and only decreased in select parts of the network
19 showing generic effects of VOE (midline areas of the default mode network). While net activity in
20 visual areas was not significantly modulated by the prior knowledge, decoding of VOE type-
21 specific signals was sensitive to changes in the participants knowledge, showing decreased
22 decoding when participants knew the tricks. These results suggest that higher-level predictive
23 information affects even the earliest levels of cortical visual processing (V1-V3).

24 Generic responses to violation of expectations

25 Witnessing magic events that violate intuitive physical principles evoked activity in a large network
26 of frontoparietal (dACC, vACC/mPFC and posterior parietal cortex) and subcortical (caudate
27 nucleus) areas, with no involvement of lower-level sensory areas. This pattern of activity is
28 consistent with that observed in a recent large meta-analysis of surprising events (Fouragnan et
29 al., 2018) and with previous experiments investigating surprise responses using naturalistic
30 videos, such as magic tricks (Danek et al., 2015), computer generated animations (Bardi et al.,
31 2017) or learned sequences of movements (Schiffer and Schubotz, 2011). Accordingly, our results
32 add to prior evidence showing the key role of the dACC in processing incongruent information
33 (Alexander and Brown, 2019, 2011) and of the caudate nucleus in signaling unexpected and
34 rewarding events (Schultz et al., 1997; Wittmann et al., 2008; Zink et al., 2003). Most importantly,
35 our results suggest that higher-level VOE in view of seemingly impossible events are processed
36 similarly to breaches of lower-level expectations, such as the presence of infrequent stimuli or

unlikely events (cf. Grassi and Bartels, 2021). This suggests the existence of a dedicated frontoparietal network signaling the detection of incongruent information in the human brain.

Specific responses to violation of expectations

Complementing the generic frontoparietal involvement in processing naturalistic violations of physical principles, we further looked into the specific effects of the different types of VOE (appear, change, vanish). Specific responses evoked by the different VOE in net activity and multivariate activation patterns (i.e., allowing for a distinction between trick-types) were observed exclusively in posterior sensory areas. In contrast, frontal areas revealed no or only weak specific VOE responses. Since this divergent pattern of net activity was only observable when looking into the individual VOE, it is possible that previous studies failed to report the involvement of sensory areas because of pooling responses to different types of VOE (Danek et al., 2015; Liu et al., 2024; Parris et al., 2009).

Previous results using dynamically occluded stimuli report the neural representation of occluded objects in posterior visual areas (Erlikhman and Caplovitz, 2017; Hulme and Zeki, 2007; Olson et al., 2004) and in neurons of the inferotemporal cortex of macaque monkeys (Puneeth and Arun, 2016) at different levels of complexity (e.g., occluded faces selectively engaged the fusiform face areas, Hulme and Zeki, 2007). The observed differential activity in visual areas using naturalistic stimuli in the present study are likely VOE type-specific surprise signals when violating said representations.

The following reasons support this interpretation: First, net responses were not driven by differences in visual content (because they were evident across distinct objects, and additional control contrasts ruled content-driven responses out). Second, decoding did not work in the absence of VOE when we used similar sensory occurrences using data from the matched control videos. Third, prior knowledge significantly reduced decoding accuracies. And finally, the observed responses occurred in functionally specialized regions of the visual cortex. For example, the unexpected appearance of objects evoked activity in the object-responsive LOC and the perception of unexpected colors evoked activity in the ventral color areas of the fusiform gyrus. Both of which are compatible with prior evidence showing increased activity in the LOC upon perceiving unexpected objects (Richter et al., 2018) and in color areas when viewing unexpected colors (Jiang et al., 2016; Stefanics et al., 2019). Hence, our results suggest that memory-based expectations related to higher-level principles affect visual processing already at the earliest cortical visual processing areas: they encode information about the presence, absence, and features of objects.

Please note that the specific responses observed may additionally be related to attentional mechanisms engaged in processing the feature-specific VOE. Indeed, attentional mechanisms are likely entangled with processing of VOE, to enhance model updating (Hohwy, 2012). However, differences in attention are unlikely the sole account of our VOE-specific results, as we observed

an enhanced involvement of the dorsal attention network (DAN) post-explanation compared to pre-explanation, while the VOE-type specific activity remained constant.

Hierarchical prediction errors in naturalistic perception

The observation of surprise-related information in posterior visual areas is in line with a variety of higher-level memory-based signals that have been reported in visual areas, such as memory color (Bannert and Bartels, 2013), scene context (Muckli et al., 2015), scene segmentation (Grassi et al., 2018, 2017; Scholte et al., 2008), expected visual stimuli (Ekman et al., 2017) and working memory (Harrison and Tong, 2009). Importantly, these signals have been interpreted as evidence of recurrent predictive signals from higher-level areas, as they encoded information that is not thought to originate from V1 (such as memory color, 3D or Gestalt and scene information). Consistent with this, further studies located corresponding signals in superficial and/or deeper layers of the cortex using laminar fMRI (Aitken et al., 2020; Lawrence et al., 2018; Muckli et al., 2015) or electrophysiological measurements in monkeys (e.g., Papale et al., 2022; Self et al., 2013).

Together, the current results fall in line with predictive coding theories and extend them to VOE regarding higher-level world models (Friston, 2005; Lee and Mumford, 2003; Rao and Ballard, 1999). We show a clear dissociation: generic responses to VOE in higher-level frontoparietal areas and segregated surprise responses in functionally specialized lower-level sensory areas (involved in the processing of the expected information). This reflects the hierarchical structure of our internal world model (Clark, 2013; Hohwy, 2014), with frontoparietal areas involved in representing more abstract aspects of the world (such as object permanency), while sensory areas represent lower-level inferences about the immediate and detailed features (such as color and shapes). Accordingly, the observed surprise-related responses in lower-level sensory areas can be thought of the product of a mismatch between top-down predictions (“a red ball”) fed back to lower-level areas to be compared with incoming sensory evidence (“a blue ball”) (Grassi and Bartels, 2021).

Knowledge-dependent modulations

To further investigate these hierarchical VOE responses, we probed how prior knowledge affected them. We hypothesized that providing the participants with knowledge about the mechanics of the trick for each video would avert VOE: why should we be surprised when viewing a disappearing ball when we know how and that the magician is actually hiding it behind his hands?

Surprisingly, while participants subjective surprise ratings were significantly reduced after providing them with the explanation of the tricks, net brain responses were almost indistinguishable: areas involved in the processing of unexpected events, such as the dACC, anterior insula and caudate nucleus (Fouragnan et al., 2018) were systematically active when

1 observing the magic videos even after participants had rational explanations for the tricks. As
2 neural responses were reminiscent to those signaling surprise, it suggests that repeated viewing
3 of explainable events did not prevent VOE. This intriguing observation likely reflects that people
4 can be moved by things they know to be unreal, such as fictions (i.e., the "paradox of fiction", see
5 Radford and Weston, 1975) or magic illusions (i.e., the "paradox of theatrical magic", see Grassi
6 et al., 2024). This is akin to how we still perceive visual illusions, even if we know how they work.
7 For example, when a magician convincingly saws someone in half on stage, the audience is
8 genuinely moved by the illusion (i.e., surprised), but do not attempt to prevent it nor call the police
9 (because they know it is unreal). Our results suggest that the compelling perceptual illusions that
10 magic provides are initially appraised as surprising, even with existing prior-information (cf. Grassi
11 et al., 2024).

12 In turn, the only areas whose activity decreased following the explanation of the tricks were
13 two midline core areas of the default mode network (DMN), the ventral ACC/mPFC and the
14 posterior cingulate cortex. The modulation of midline core areas of the DMN by prior knowledge
15 is consistent with recent reports showing their involvement in processing surprising events in
16 movies (Brandman et al., 2021), jokes (Jääskeläinen et al., 2016) and structured events unfolding
17 in time (Baldassano et al., 2018; Regev et al., 2013; Simony et al., 2016). Based on these findings,
18 it has been suggested that the DMN is not to be understood exclusively as an "intrinsic" network
19 (as originally proposed, cf. Raichle, 2015), but as a dynamic "sense-making" network involved in
20 the creation of rich models of events by integrating incoming information with prior knowledge as
21 they unfold over time instead (Stawarczyk et al., 2021; Yeshurun et al., 2021).

22 Here, we show that areas of the "sense-making" network may be sensitive to the rational
23 explanation of magic tricks, whereas all other identified surprise-related regions continued to be
24 sensitive to the VOE (even once the tricks were understood). The decreased involvement of areas
25 of the DMN, together with an increase of activity in frontoparietal areas of the dorsal attention
26 network (DAN), is consistent with a reduction in prediction error (surprise) signals related to
27 narrative understanding (engaging the DMN) and an increase of top-down attention after
28 explanation of the tricks (engaging the DAN).

29 **Conclusion**

30 We used a naturalistic paradigm to violate deeply held beliefs of our physical world, involving three
31 types of expectation violations (object appearance, color change, and object disappearance). Our
32 results show a hierarchy of surprise signals: generic responses to unexpected events in
33 frontoparietal areas, and responses specific to the type of VOE in distinct functionally specialized
34 sensory areas. Our results suggest that world-model VOE are processed similarly to other
35 surprising events in dedicated areas of the prefrontal cortex and striatum, and that core midline
36 areas of the default-mode network decrease their involvement once rational understanding is

1 established. Most importantly, we show that early and functionally specialized areas of the visual
2 cortex encode memory-based predictions about the presence, absence, and features of objects.
3

Data accessibility

Code (for preprocessing, analysis and visualisations) and preprocessed data for group analyses can be found at https://osf.io/kn2af/?view_only=067b698a4567441e93b01518a88860a0. Raw MRI data can be provided upon reasonable request.

CRedit author statement

Vincent Plikat: conceptualization, formal analysis (lead), investigation (equal), visualisation, writing – original draft, writing – review and editing. **Pablo R. Grassi:** conceptualization (lead), formal analysis (supporting), investigation (equal), methodology, supervision (equal), visualization, writing – original draft, writing – review and editing. **Julius Frack:** resources. **Andreas Bartels:** conceptualization, methodology, supervision (equal), writing – review and editing.

Acknowledgements

This work was supported by Barbara-Wengeler-Foundation, the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) project number 465409366, and by the Max Planck Society.

References

- Abraham A, Pedregosa F, Eickenberg M, Gervais P, Mueller A, Kossaifi J, Gramfort A, Thirion B, Varoquaux G. 2014. Machine learning for neuroimaging with scikit-learn. *Front Neuroinform* **8**. doi:10.3389/fninf.2014.00014
- Adams WJ, Graf EW, Ernst MO. 2004. Experience can change the “light-from-above” prior. *Nature Neuroscience* **7**:1057–1058. doi:10.1038/nn1312
- Aitken F, Menelaou G, Warrington O, Koolschijn RS, Corbin N, Callaghan MF, Kok P. 2020. Prior expectations evoke stimulus-specific activity in the deep layers of the primary visual cortex. *PLoS Biol* **18**:e3001023. doi:10.1371/journal.pbio.3001023
- Alexander WH, Brown JW. 2019. The Role of the Anterior Cingulate Cortex in Prediction Error and Signaling Surprise. *Topics in Cognitive Science* **11**:119–135. doi:10.1111/tops.12307
- Alexander WH, Brown JW. 2011. Medial prefrontal cortex as an action-outcome predictor. *Nature Neuroscience* **14**:1338–1344. doi:10.1038/nn.2921
- Baldassano C, Hasson U, Norman KA. 2018. Representation of Real-World Event Schemas during Narrative Perception. *J Neurosci* **38**:9689–9699. doi:10.1523/JNEUROSCI.0251-18.2018
- Bannert MM, Bartels A. 2013. Decoding the yellow of a gray banana. *Current Biology* **23**:2268–2272. doi:10.1016/j.cub.2013.09.016
- Bardi L, Desmet C, Nijhof A, Wiersema JR, Brass M. 2017. Brain activation for spontaneous and explicit false belief tasks overlaps: New fMRI evidence on belief processing and violation of expectation. *Social Cognitive and Affective Neuroscience* **12**:391–400. doi:10.1093/scan/nsw143
- Bartels A, Zeki S. 2000. The architecture of the colour centre in the human visual brain: New results and a review. *European Journal of Neuroscience* **12**:172–193. doi:10.1046/j.1460-9568.2000.00905.x
- Battaglia PW, Hamrick JB, Tenenbaum JB. 2013. Simulation as an engine of physical scene understanding. *Proc Natl Acad Sci USA* **110**:18327–18332. doi:10.1073/pnas.1306572110

- 1 Brandman T, Malach R, Simony E. 2021. The surprising role of the default mode network in naturalistic
2 perception. *Commun Biol* **4**:79. doi:10.1038/s42003-020-01602-z
- 3 Clark A. 2013. Whatever next? Predictive brains, situated agents, and the future of cognitive science.
4 *Behavioral and Brain Sciences* **36**:181–204. doi:10.1017/S0140525X12000477
- 5 Corbetta M, Shulman GL. 2002. Control of goal-directed and stimulus-driven attention in the brain. *Nature*
6 *reviews Neuroscience* **3**:201–15. doi:10.1038/nrn755
- 7 Dale AM, Fischl B, Sereno MI. 1999. Cortical Surface-Based Analysis. *NeuroImage* **9**:179–194.
8 doi:10.1006/nimg.1998.0395
- 9 Danek AH, Öllinger M, Fraps T, Grothe B, Flanagan VL. 2015. An fMRI investigation of expectation violation
10 in magic tricks. *Frontiers in Psychology* **6**:1–11. doi:10.3389/fpsyg.2015.00084
- 11 de Lange FP, Heilbron M, Kok P. 2018. How Do Expectations Shape Perception? *Trends in Cognitive*
12 *Sciences* **22**:764–779. doi:10.1016/j.tics.2018.06.002
- 13 Egnér T, Monti JM, Summerfield C. 2010. Expectation and surprise determine neural population responses
14 in the ventral visual stream. *Journal of Neuroscience* **30**:16601–16608.
15 doi:10.1523/JNEUROSCI.2770-10.2010
- 16 Ekman M, Kok P, De Lange FP. 2017. Time-compressed preplay of anticipated events in human primary
17 visual cortex. *Nature Communications* **8**:1–9. doi:10.1038/ncomms15276
- 18 Erlichman G, Caplovitz GP. 2017. Decoding information about dynamically occluded objects in visual cortex.
19 *NeuroImage* **146**:778–788. doi:10.1016/j.neuroimage.2016.09.024
- 20 Fischer J, Mikhael JG, Tenenbaum JB, Kanwisher N. 2016. Functional neuroanatomy of intuitive physical
21 inference. *Proc Natl Acad Sci USA* **113**. doi:10.1073/pnas.1610344113
- 22 Fischl B, Salat DH, Busa E, Albert M, Dieterich M, Haselgrove C, van der Kouwe A, Killiany R, Kennedy D,
23 Klaveness S, Montillo A, Makris N, Rosen B, Dale AM. 2002. Whole Brain Segmentation. *Neuron*
24 **33**:341–355. doi:10.1016/S0896-6273(02)00569-X
- 25 Fouragnan E, Retzler C, Philiastides MG. 2018. Separate neural representations of prediction error valence
26 and surprise: Evidence from an fMRI meta-analysis. *Hum Brain Mapp* **39**:2887–2906.
27 doi:10.1002/hbm.24047
- 28 Friston K. 2010. The free-energy principle: A unified brain theory? *Nature Reviews Neuroscience* **11**:127–
29 138. doi:10.1038/nrn2787
- 30 Friston K. 2005. A theory of cortical responses. *Philosophical transactions of the Royal Society of London*
31 *Series B, Biological sciences* **360**:815–36. doi:10.1098/rstb.2005.1622
- 32 Friston KJ, Penny WD, Glaser DE. 2005. Conjunction revisited. *NeuroImage* **25**:661–7.
33 doi:10.1016/j.neuroimage.2005.01.013
- 34 Glasser MF, Coalson TS, Robinson EC, Hacker CD, Harwell J, Yacoub E, Ugurbil K, Andersson J,
35 Beckmann CF, Jenkinson M, Smith SM, Van Essen DC. 2016. A multi-modal parcellation of human
36 cerebral cortex. *Nature* **536**:171–178. doi:10.1038/nature18933
- 37 Grassi PR, Bartels A. 2021. Magic, Bayes and wows: A Bayesian account of magic tricks. *Neuroscience &*
38 *Biobehavioral Reviews* **126**:515–527. doi:10.1016/j.neubiorev.2021.04.001
- 39 Grassi PR, Plikat V, Wong HY. 2024. How can we be moved by magic? *British Journal of Aesthetics* **64**:187–
40 204. doi:10.1093/aesthj/ayad026
- 41 Grassi PR, Zaretskaya N, Bartels A. 2018. A Generic Mechanism for Perceptual Organization in the Parietal
42 Cortex. *The Journal of Neuroscience* **38**:7158–7169. doi:10.1523/JNEUROSCI.0436-18.2018
- 43 Grassi PR, Zaretskaya N, Bartels A. 2017. Scene segmentation in early visual cortex during suppression of
44 ventral stream regions. *NeuroImage* **146**:71–80. doi:10.1016/j.neuroimage.2016.11.024
- 45 Grill-Spector K, Kourtzi Z, Kanwisher N. 2001. The lateral occipital complex and its role in object recognition.
46 *Vision Research* **41**:1409–1422. doi:10.1016/S0042-6989(01)00073-6
- 47 Harrison SA, Tong F. 2009. Decoding reveals the contents of visual working memory in early visual areas.
48 *Nature* **458**:632–635. doi:10.1038/nature07832

- 1 Hespos SJ, Ferry AL, Rips LJ. 2009. Five-Month-Old Infants Have Different Expectations for Solids and
2 Liquids. *Psychol Sci* **20**:603–611. doi:10.1111/j.1467-9280.2009.02331.x
- 3 Hohwy J. 2014. The Predictive Mind. Oxford: OUP Oxford.
4 doi:10.1093/acprof:oso/9780199682737.001.0001
- 5 Hohwy J. 2012. Attention and conscious perception in the hypothesis testing brain. *Frontiers in Psychology*
6 **3**:1–14. doi:10.3389/fpsyg.2012.00096
- 7 Hohwy J, Roepstorff A, Friston K. 2008. Predictive coding explains binocular rivalry: an epistemological
8 review. *Cognition* **108**:687–701. doi:10.1016/j.cognition.2008.05.010
- 9 Hulme OJ, Zeki S. 2007. The Sightless View: Neural Correlates of Occluded Objects. *Cerebral Cortex*
10 **17**:1197–1205. doi:10.1093/cercor/bhl031
- 11 Jääskeläinen IP, Pajula J, Tohka J, Lee H-J, Kuo W-J, Lin F-H. 2016. Brain hemodynamic activity during
12 viewing and re-viewing of comedy movies explained by experienced humor. *Sci Rep* **6**:27741.
13 doi:10.1038/srep27741
- 14 Jiang J, Summerfield C, Egner T. 2016. Visual Prediction Error Spreads Across Object Features in Human
15 Visual Cortex. *J Neurosci* **36**:12746–12763. doi:10.1523/JNEUROSCI.1546-16.2016
- 16 Kok P, Jehee JFM, de Lange FP. 2012a. Less Is More: Expectation Sharpens Representations in the
17 Primary Visual Cortex. *Neuron* **75**:265–270. doi:10.1016/j.neuron.2012.04.034
- 18 Kok P, Rahnev D, Jehee JFM, Lau HC, de Lange FP. 2012b. Attention Reverses the Effect of Prediction in
19 Silencing Sensory Signals. *Cerebral Cortex* **22**:2197–2206. doi:10.1093/cercor/bhr310
- 20 Lawrence SJD, van Mourik T, Kok P, Koopmans PJ, Norris DG, de Lange FP. 2018. Laminar Organization
21 of Working Memory Signals in Human Visual Cortex. *Current Biology*.
22 doi:10.1016/j.cub.2018.08.043
- 23 Lee TS, Mumford D. 2003. Hierarchical Bayesian inference in the visual cortex. *Journal of the Optical*
24 *Society of America A* **20**:1434. doi:10.1364/josaa.20.001434
- 25 Liu S, Lydic K, Mei L, Saxe R. 2024. Violations of physical and psychological expectations in the human
26 adult brain. *Imaging Neuroscience* **2**:1–25. doi:10.1162/imag_a_00068
- 27 Muckli L, De Martino F, Vizioli L, Petro LS, Smith FW, Ugurbil K, Goebel R, Yacoub E. 2015. Contextual
28 Feedback to Superficial Layers of V1. *Current Biology* **25**:2690–2695.
29 doi:10.1016/j.cub.2015.08.057
- 30 Nichols TE, Holmes AP. 2002. Nonparametric permutation tests for functional neuroimaging: A primer with
31 examples. *Human Brain Mapping* **15**:1–25. doi:10.1002/hbm.1058
- 32 Notter M, Gale D, Herholz P, Markello R, Notter-Bielser M-L, Whitaker K. 2019. AtlasReader: A Python
33 package to generate coordinate tables, region labels, and informative figures from statistical MRI
34 images. *JOSS* **4**:1257. doi:10.21105/joss.01257
- 35 Olson IR, Gatenby JC, Leung H-C, Skudlarski P, Gore JC. 2004. Neuronal representation of occluded
36 objects in the human brain. *Neuropsychologia* **42**:95–104. doi:10.1016/S0028-3932(03)00151-9
- 37 Papale P, Wang F, Morgan AT, Chen X, Gilhuis A, Petro LS, Muckli L, Roelfsema PR, Self MW. 2022.
38 Feedback brings scene information to the representation of occluded image regions in area V1 of
39 monkeys and humans. doi:10.1101/2022.11.21.517305
- 40 Parris BA, Kuhn G, Mizon GA, Benattayallah A, Hodgson TL. 2009. Imaging the impossible: An fMRI study
41 of impossible causal relationships in magic tricks. *NeuroImage* **45**:1033–1039.
42 doi:10.1016/j.neuroimage.2008.12.036
- 43 Pedregosa F. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*
44 **12**:2825–2830.
- 45 Puneeth NC, Arun SP. 2016. A neural substrate for object permanence in monkey inferotemporal cortex.
46 *Sci Rep* **6**:30808. doi:10.1038/srep30808
- 47 Radford C, Weston M. 1975. How Can We Be Moved by the Fate of Anna Karenina? *Aristot Soc Suppl Vol*
48 **49**:67–94. doi:10.1093/aristoteliansupp/49.1.67

- 1 Raichle ME. 2015. The Brain's Default Mode Network. *Annu Rev Neurosci* **38**:433–447.
2 doi:10.1146/annurev-neuro-071013-014030
- 3 Rao RPN, Ballard DH. 1999. Predictive coding in the visual cortex: A functional interpretation of some extra-
4 classical receptive-field effects. *Nature Neuroscience* **2**:79–87. doi:10.1038/4580
- 5 Regev M, Honey CJ, Simony E, Hasson U. 2013. Selective and Invariant Neural Responses to Spoken and
6 Written Narratives. *Journal of Neuroscience* **33**:15978–15988. doi:10.1523/JNEUROSCI.1580-
7 13.2013
- 8 Richter D, Ekman M, de Lange FP. 2018. Suppressed Sensory Response to Predictable Object Stimuli
9 throughout the Ventral Visual Stream. *J Neurosci* **38**:7452–7461. doi:10.1523/JNEUROSCI.3421-
10 17.2018
- 11 Rolls ET, Huang C-C, Lin C-P, Feng J, Joliot M. 2020. Automated anatomical labelling atlas 3. *NeuroImage*
12 **206**:116189. doi:10.1016/j.neuroimage.2019.116189
- 13 SanMiguel I, Widmann A, Bendixen A, Trujillo-Barreto N, Schroger E. 2013. Hearing Silences: Human
14 Auditory Processing Relies on Preactivation of Sound-Specific Brain Activity Patterns. *Journal of*
15 *Neuroscience* **33**:8633–8639. doi:10.1523/JNEUROSCI.5821-12.2013
- 16 Schiffer A-M, Schubotz RI. 2011. Caudate Nucleus Signals for Breaches of Expectation in a Movement
17 Observation Paradigm. *Front Hum Neurosci* **5**. doi:10.3389/fnhum.2011.00038
- 18 Scholte HS, Jolij J, Fahrenfort JJ, Lamme VAF. 2008. Feedforward and recurrent processing in scene
19 segmentation: Electroencephalography and functional magnetic resonance imaging. *Journal of*
20 *Cognitive Neuroscience* **20**:2097–2109. doi:10.1162/jocn.2008.20142
- 21 Schultz W, Dayan P, Montague PR. 1997. A neural substrate of prediction and reward. *Science* **275**:1593–
22 1599. doi:10.1126/science.275.5306.1593
- 23 Schwettmann S, Tenenbaum JB, Kanwisher N. 2019. Invariant representations of mass in the human brain.
24 *eLife* **8**:e46619. doi:10.7554/eLife.46619
- 25 Self MW, van Kerkoerle T, Supér H, Roelfsema PR. 2013. Distinct Roles of the Cortical Layers of Area V1
26 in Figure-Ground Segregation. *Current Biology* **23**:2121–2129. doi:10.1016/j.cub.2013.09.013
- 27 Simony E, Honey CJ, Chen J, Lositsky O, Yeshurun Y, Wiesel A, Hasson U. 2016. Dynamic reconfiguration
28 of the default mode network during narrative comprehension. *Nat Commun* **7**:12141.
29 doi:10.1038/ncomms12141
- 30 Stawarczyk D, Bezdek MA, Zacks JM. 2021. Event Representations and Predictive Processing: The Role
31 of the Midline Default Network Core. *Top Cogn Sci* **13**:164–186. doi:10.1111/tops.12450
- 32 Stefanics G, Stephan KE, Heinzle J. 2019. Feature-specific prediction errors for visual mismatch.
33 *NeuroImage* **196**:142–151. doi:10.1016/j.neuroimage.2019.04.020
- 34 Stelzer J, Chen Y, Turner R. 2013. Statistical inference and multiple testing correction in classification-based
35 multi-voxel pattern analysis (MVPA): Random permutations and cluster size control. *NeuroImage*
36 **65**:69–82. doi:10.1016/j.neuroimage.2012.09.063
- 37 Todorovic A, van Ede F, Maris E, de Lange FP. 2011. Prior expectation mediates neural adaptation to
38 repeated sounds in the auditory cortex: An MEG study. *Journal of Neuroscience* **31**:9118–9123.
39 doi:10.1523/JNEUROSCI.1425-11.2011
- 40 Vizioli L, Moeller S, Dowdle L, Akçakaya M, De Martino F, Yacoub E, Uğurbil K. 2021. Lowering the thermal
41 noise barrier in functional brain mapping with magnetic resonance imaging. *Nat Commun* **12**:5181.
42 doi:10.1038/s41467-021-25431-8
- 43 Wacongne C, Labyt E, Van Wassenhove V, Bekinschtein T, Naccache L, Dehaene S. 2011. Evidence for a
44 hierarchy of predictions and prediction errors in human cortex. *Proceedings of the National*
45 *Academy of Sciences of the United States of America* **108**:20754–20759.
46 doi:10.1073/pnas.1117807108
- 47 Wang L, Mruczek REB, Arcaro MJ, Kastner S. 2015. Probabilistic maps of visual topography in human
48 cortex. *Cerebral Cortex* **25**:3911–3931. doi:10.1093/cercor/bhu277

- 1 Wang S. 2004. Young infants' reasoning about hidden objects: evidence from violation-of-expectation tasks
2 with test trials only. *Cognition* **93**:167–198. doi:10.1016/j.cognition.2003.09.012
- 3 Wessel JR, Danielmeier C, Morton JB, Ullsperger M. 2012. Surprise and Error: Common Neuronal
4 Architecture for the Processing of Errors and Novelty. *Journal of Neuroscience* **32**:7528–7537.
5 doi:10.1523/JNEUROSCI.6352-11.2012
- 6 Wittmann BC, Daw ND, Seymour B, Dolan RJ. 2008. Striatal Activity Underlies Novelty-Based Choice in
7 Humans. *Neuron* **58**:967–973. doi:10.1016/j.neuron.2008.04.027
- 8 Wynn K. 1992. Addition and subtraction by human infants. *Nature* **358**:749–750. doi:10.1038/358749a0
- 9 Yeo BTT, Krienen FM, Sepulcre J, Sabuncu MR, Lashkari D, Hollinshead M, Roffman JL, Smoller JW, Zollei
10 L, Polimeni JR, Fischl B, Liu H, Buckner RL. 2011. The organization of the human cerebral cortex
11 estimated by intrinsic functional connectivity. *Journal of neurophysiology* **106**:1125–1165.
12 doi:10.1152/jn.00338.2011.
- 13 Yeshurun Y, Nguyen M, Hasson U. 2021. The default mode network: where the idiosyncratic self meets the
14 shared social world. *Nat Rev Neurosci* **22**:181–192. doi:10.1038/s41583-020-00420-w
- 15 Zink CF, Pagnoni G, Martin ME, Dhamala M, Berns GS. 2003. Human Striatal Response to Salient
16 Nonrewarding Stimuli. *J Neurosci* **23**:8092–8097. doi:10.1523/JNEUROSCI.23-22-08092.2003
- 17