

Value spillover: How contextually irrelevant values influence choice and vmPFC activity in humans

Nir Moneta^{1,2} (moneta@mpib-berlin.mpg.de), Hauke R. Heekeren² (hauke.heekeren@fu-berlin.de) & Nicolas W. Schuck^{1,3} (schuck@mpib-berlin.mpg.de)

¹Max Planck Research Group NeuroCode, Max Planck Institute for Human Development, 14195 Berlin, Germany

²Department of Education and Psychology, Freie Universität Berlin, 14195 Berlin, Germany

³Max Planck UCL Centre for Computational Psychiatry and Ageing Research, Berlin, 14195 Berlin, Germany

Abstract

Objects we choose between often have multiple features. Research has shown that decision-making is guided by attentional selection of contextually-relevant features, while ventromedial prefrontal-cortex (vmPFC) represents the expected outcomes associated with these features. Yet, if this selective value retrieval is not entirely perfect, irrelevant features, and values they carry in other contexts, will influence neural value representations as well as choice. We tested this idea by utilizing a context-dependent random-dot motion paradigm. Forty humans made decisions between two clouds of moving dots, each consisting of two features (motion direction and dot color). First, participants learned to associate each color and motion with specific rewards. During subsequent decision making, a context cue indicated the trials relevant feature-type (color/motion) and choices led to outcomes associated with the relevant feature. In line with our hypothesis, the more values of the irrelevant features agreed with the relevant, the faster participants reacted. fMRI analyses showed parametric modulation of the vmPFC/OFC signal by both (1) the relevant feature value and (2) the value difference between irrelevant features. These results indicate that contextually-irrelevant features influence value representation and suggest that the brain's decision system computes values in the presence of partial activation of irrelevant context or task states.

Keywords: decision-making; ventromedial-prefrontal cortex ; attention

Introduction

In the field of value-based decision-making, it is believed that people usually follow the principle of value maximization (e.g., Kahneman & Tversky, 1979), i.e. the subjective value is first evaluated for each alternative, and the option with the highest expected value is chosen. This representation of expected value is believed to reside in the ventromedial prefrontal cortex (vmPFC, e.g. Bartra, McGuire, & Kable, 2013).

In everyday life, we often have to make decisions between objects with multiple features (or attributes) that predict outcomes in different contexts. For example, when we choose a medicine, the active ingredients, but not the color of the pills, will predict its value. When picking a fruit during the same shopping trip, however, the color might be the best predictor of its value. There have been many investigations to

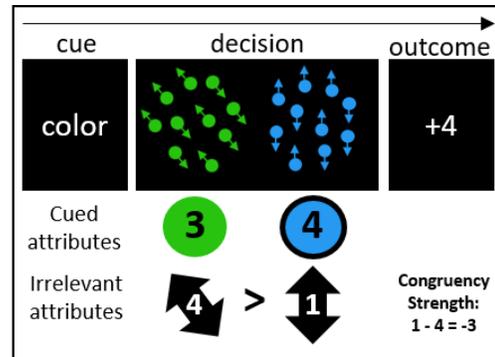


Figure 1: **Experimental paradigm.** Each trial started with a cue of the relevant dimension (left). Each cloud had two features (motion and color) and participants made a decision between the two clouds (middle). After a choice, participants received the outcome associated with the chosen cloud's cued feature (right). In the figure, an incongruent trial is presented, i.e., the irrelevant and cued dimensions disagree. Congruency strength is -3, see text in Figure.

how humans and animals make decisions between objects in the presence of distracting information (e.g., Li, Michael, Balaguer, Castan, & Summerfield, 2018; Mante, Sussillo, Shenoy, & Newsome, 2013) or integrate values associated with different features of the same object in order to elicit a decision (e.g., Pelletier & Fellows, 2018; Basten, Biele, Heekeren, & Fiebach, 2010). Previous research has shown that during decision-making the brain's attentional control network enhances the processing of features that are relevant in the current context (Niv et al., 2015). Yet, if such attentional filtering is not perfect, then outcome expectations associated with the features that are irrelevant in the current context might influence subjective value and choice.

The aim of this study is to systematically test this influence. We hypothesized that contextually irrelevant features would have an impact on participants' choice behavior (i.e. reaction time, RT) as well as bias the expected value signal in the vmPFC.

Methods

Experimental procedure

Forty young participants performed a random dot motion paradigm (18 women, $\mu_{age} = 27.6$, $\sigma_{age} = 3.35$). Psychophys-

ical properties of four colors and four motion directions were first titrated using a staircasing procedure. Before fMRI scanning, participants learned the outcomes associated with each of these eight features. During fMRI scanning, participants were asked to make decisions between two random dot kinematograms, each of which had one color and one direction from the same set. Importantly, in each trial participants were first cued whether a decision should be made based on the color or motion features, and received only the outcome associated with the contextually relevant feature of their choice.

Staircasing procedure

In order to ensure that reaction times (RTs) depended mainly on the associated values (e.g., Hunt et al., 2012) and not on other stimulus properties, such as salience, we used a staircasing procedure that was conducted *prior to value learning*. In this procedure, motion coherence and color saturation were adjusted for each participant in order minimize between-feature RT differences in a perceptual detection task. Motion coherence and color saturation were adjusted for 168 trials, while luminance in YCbCr color space remained fixed throughout (cf. Abbott, Griffiths, & Regier, 2016). As can be seen in Figure 2A, RT differences between the eight used stimulus features (indexed by variability) was markedly reduced after compared to before the procedure ($\mu = 0.019$ vs. $\mu = 0.12$, $t_{(39)} = 5.83$, $p < .0001$). This ensured that any later changes in RT can be associated with the learned value information.

Outcome learning

Prior to the main task, participants learned to associate each of the four colors and four motion directions with deterministic outcomes. Outcomes associated with the four features on both dimensions were 10, 30, 50 and 70 credit-points (for simplicity, henceforth: 1,2,3,4). For this part, we used single-feature clouds only (henceforth: 1D), i.e. no coherent motion or no color (gray) for color and motion clouds, respectively. Therefore each cloud in this part only represented a single feature. Eighty 'forced-choice' trials in which only one cloud was presented were followed by a blocks of 72 free-choice trials in which two 1D clouds were shown. Free-choice blocks repeated until participants reached a minimum of 85% accuracy of choosing the higher valued cloud in a block (minimum 2 blocks, maximum 4). Crucially, in free-choice trials the two 1D clouds could be of the same dimension (e.g. color and color) but also from different dimensions (e.g. color and motion). This was done to encourage mapping of the values for each dimension on similar scales. This type of trials did not repeat in the decision making task.

Decision Making Task

In order to investigate how contextually irrelevant features influence the expected value of choices, we utilized a context-dependent random dot motion paradigm (e.g. Pilly & Seitz, 2009), in which a choice between two clouds of moving dots was presented. As can be seen in Figure 1, each cloud had

two features (a motion direction and a color) that were associated with specific rewards in the previous phase (henceforth: 2D trials). At the start of each trial, participants were cued to focus only on one dimension, and had to choose the stimulus with the higher valued feature on the cued dimension. Following a choice, the outcome associated with the contextually relevant feature of the chosen cloud was presented. Thus, in each trial only the two features on the cued dimension were relevant for determining the choice. This part included four blocks, each of 96 trials (36 1D trials and 60 2D trials).

Importantly, four features were present during 2D trials: the two relevant features, one on the left and one on right, and the two irrelevant features. Each of these four features had one associated outcome, and the main purpose of this paradigm was to study how these values influence choices and brain activity. To constrain complexity, the two features on the cued dimension always had a value difference of 1, i.e. the choices on the cued dimension were only between outcomes of 1 vs. 2, 2 vs. 3 or 3 vs. 4. The features on the irrelevant dimension could have bigger value differences (see congruency strength below). In some trials, no irrelevant stimulus features were present, i.e. both clouds were 1D clouds.

Lastly, to prevent context confusion caused by frequent switching, the cued dimension stayed the same for 5-7 trials and only then switched (in a non-predictive manner). Participants were instructed that they collect all the reward points throughout the experiment and that all collected points would be translated to monetary reward at the end of the experiment (1 Euro payment per 600 credit points).

Results

Behavioral analysis

We included only accurate trials in our analysis, i.e. trials in which participants choose the highest value based on the cued dimension. Overall accuracy was well above chance ($\mu = 0.896$, $\sigma = 0.054$, $t_{(39)} = 46.14$, $p < .0001$). Replicating previous findings (e.g. Hunt et al., 2012), participants reacted faster when choosing higher valued options (linear mixed effects model: $\chi^2(1) = 370.282$, $p < .0001$). In line with our hypothesis, we additionally found that RTs were faster when the irrelevant dimension was congruent (i.e. indicated the same choice), compared to incongruent (i.e. indicated the other option, $t_{(39)} = 4.5737$, $p < .001$, Figure 2B.). We further quantified congruency strength of a choice by taking the value difference of the irrelevant features (chosen minus unchosen options, 1 minus 4 in Figure 1). The congruency strength ranged from -3 to 3, with negative numbers representing more incongruent trials. Crucially, the congruency strength represents only the outcome-irrelevant values of the trial. As can be seen in Figure 2C, the more positive the congruency strength, the faster participants reacted (linear mixed effects model: $\chi^2(1) = 11.623$, $p < .0001$). I.e., the larger the value associated with the irrelevant feature on the chosen side was, relative to the unchosen side, the faster were RTs.

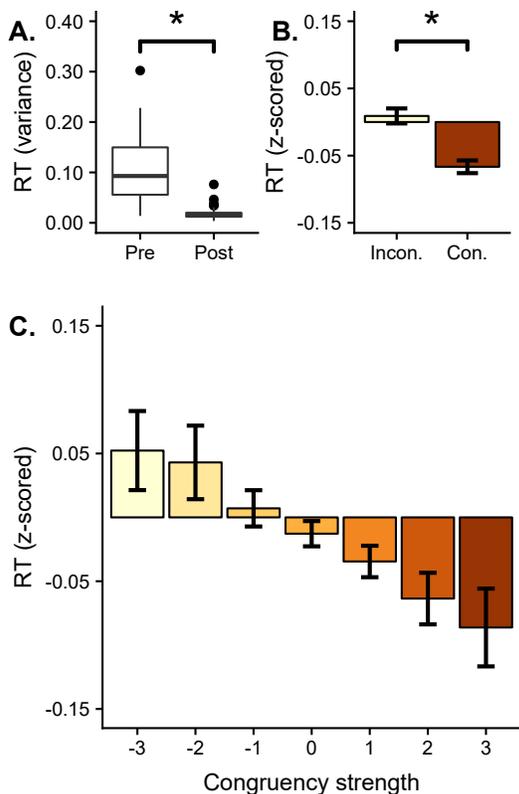


Figure 2: **Behavioral results.** (A) Staircasing procedure. RT variance (y-axis) for the different features was significantly reduced. Variance was computed for each participant as the difference between the median RT per feature and their overall median RT. (B) Main congruency effect. Participants reacted faster for congruent trials compared to incongruent trials (C) Effect of congruency strength. Participants responded faster with increasing congruency and slower with increasing incongruency. All effects presented are significant at $p < 0.001$, $N=40$. In B and C, Y-axis represents the RT z-scored for each participant separately and error bars in B and C represent standard error of the mean.

fMRI analysis

Preprocessing We acquired fMRI data using a multi-band sequence (acceleration factor 4, TR 1250ms, TE 26ms, FA 71, 64 2mm slices, 2x2mm in plane resolution). A tilt angle of 30 degrees from AC-PC was used in order to maximize signal from the orbitofrontal cortex (OFC, see Weiskopf, Hutton, Josephs, & Deichmann, 2006). Preprocessing was done using fMRIPrep (Esteban et al., 2018). In short, based on the estimated susceptibility distortion generated from the deformation field estimated based on two echo-planar imaging (EPI) references with opposing phase-encoding directions, an unwarped BOLD reference was calculated and used for a more accurate co-registration with the anatomical image. BOLD runs were slice-time corrected and the BOLD time-series were resampled to MNI152NLin2009cAsym standard space. Six

head-motion estimates were calculated and used for correction in later analysis.

General Linear Model analysis We analyzed the fMRI data using a general linear model (GLM). We included 5 regressors of interest including two onset regressors of the two trial-types (1D trials/2D trials). Each trial-type had an additional parametric modulator representing the value of the relevant feature of the chosen side. For the 2D trials, we included an additional parametric modulators of the congruency strength of the trial, i.e. positive numbers for congruent and negative for incongruent trials. Additional regressors reflected cue and outcome onset, in addition to six motion nuisance regressors. The parametric regressors representing the chosen values of both trial types captured a main effect of chosen value in orbitofrontal cortex, the ventral striatum, and hippocampus (Figure 3A.). The parametric regressors representing the congruency strength of the 2D trials captured a main effect in the anterior insular/posterior OFC (Figure 3B, $p < 0.001$, uncorrected for both effects). These results indicate that the signal in the vmPFC is influenced by the contextually outcome-irrelevant values, both of the chosen as well as the unchosen objects (since the congruency strength is dependent on both)

Discussion

In daily life, selecting which features should lead our decisions is a very challenging task. In this study, we showed that even when the relevant features are cued explicitly and the reward that appears immediately after choice is not influenced by the irrelevant feature, this selective value retrieval process is less than perfect. Participants RTs were influenced by all the valued-features that were associated with each object. Furthermore, this influence was not merely based on

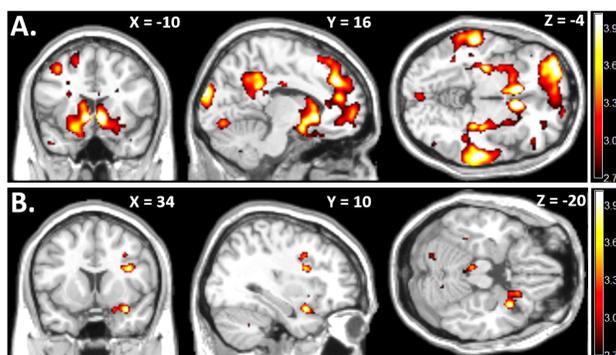


Figure 3: **Second-level group analysis** (A) Effect of the chosen value. Visualization of BOLD activation for the chosen value parametric modulation in the vmPFC. (B) Effect of congruency strength. Visualization of BOLD activation for the ConS parametric modulation in the anterior insular/posterior OFC. All t-value maps are thresholded at $p < 0.005$ uncorrected for illustration purposes. The upper right corner insets denote the MNI coordinate of the respective slice.

the agreement of the ignored and cued features, but rather on the specific strength of the agreement (or disagreement). This shows that participants meaningfully processed the cued and ignored features of both objects, and both influenced their choice.

Moreover, the strength in which the outcome-irrelevant features agree with the decision parametrically modulated the signal in the vmPFC. The vmPFC plays a crucial role in guiding behavior based on the representation of rewards (Rushworth, Kolling, Sallet, & Mars, 2012). The signal measured in this region is believed to represent the expected-value of a choice, irrespective of the stimuli that are associated with the reward (Lebreton, Jorge, Michel, Thirion, & Pessiglione, 2009). It is likely that this requires placing all options' values on a common scale (Chib, Rangel, Shimojo, & O'Doherty, 2009; McNamee, Rangel, & O'doherty, 2013).

These results indicate that contextually-irrelevant features can influence participants expected value of their choices – and highlight the economy of adding irrelevant but otherwise valuable features to products. The influence on the expected value representation in the vmPFC elucidate the potential neural mechanisms underlying choices made based on features that are known to be irrelevant for the outcome. Future analysis could shed light on the exact nature of the relations between these effects and attentional mechanisms using multivariate analysis, and whether relevant and irrelevant values are processed in parallel or are integrated in vmPFC representations during choices.

Acknowledgments

This work was funded by a research group grant awarded to NWS by the Max Planck Society (M.TN.A.BILD0004). We thank the Ernst Ludwig Ehrlich Studienwerk (ELES) for financial support through this study. We thank Gregor Caregnato for help with participant recruitment, Anika Löwe and Lena Maria Krippner, Sonali Beckmann and Nadine Taube for help with data acquisition and all participants for their participation.

References

Abbott, J. T., Griffiths, T. L., & Regier, T. (2016). Focal colors across languages are representative members of color categories. *Proceedings of the National Academy of Sciences*, *113*(40), 11178-11183.

Bartra, O., McGuire, J. T., & Kable, J. W. (2013). The valuation system: a coordinate-based meta-analysis of bold fmri experiments examining neural correlates of subjective value. *Neuroimage*, *76*, 412-427.

Basten, U., Biele, G., Heekeren, H. R., & Fiebach, C. J. (2010). How the brain integrates costs and benefits during decision making. *Proceedings of the National Academy of Sciences*, *107*(50), 21767-21772. doi: 10.1073/pnas.0908104107

Chib, V. S., Rangel, A., Shimojo, S., & O'Doherty, J. P. (2009). Evidence for a common representation of decision values for dissimilar goods in human ventromedial prefrontal cor-

tex. *Journal of Neuroscience*, *29*(39), 12315–12320. doi: 10.1523/JNEUROSCI.2575-09.2009

Hunt, L. T., Kolling, N., Soltani, A., Woolrich, M. W., Rushworth, M. F., & Behrens, T. E. (2012). Mechanisms underlying cortical activity during value-guided choice. *Nature neuroscience*, *15*(3), 470-476.

Kahneman, D., & Tversky, A. (1979). Prospect theory : An analysis of decisions under risk. *Econometrica*, *47*, 278.

Lebreton, M., Jorge, S., Michel, V., Thirion, B., & Pessiglione, M. (2009). An automatic valuation system in the human brain: Evidence from functional neuroimaging. *Neuron*, *64*(3), 431 - 439. doi: <https://doi.org/10.1016/j.neuron.2009.09.040>

Li, V., Michael, E., Balaguer, J., Castan, S. H., & Summerfield, C. (2018). Gain control explains the effect of distraction in human perceptual, cognitive, and economic decision making. *Proceedings of the National Academy of Sciences*, *115*(38), E8825-E8834.

Mante, V., Sussillo, D., Shenoy, K. V., & Newsome, W. T. (2013). Context-dependent computation by recurrent dynamics in prefrontal cortex. *nature*, *503*(7474), 78.

McNamee, D., Rangel, A., & O'doherty, J. P. (2013). Category-dependent and category-independent goal-value codes in human ventromedial prefrontal cortex. *Nature neuroscience*, *16*(4), 479-485.

Niv, Y., Daniel, R., Geana, A., Gershman, S. J., Leong, Y. C., Radulescu, A., & Wilson, R. C. (2015). Reinforcement learning in multidimensional environments relies on attention mechanisms. *Journal of Neuroscience*, *35*(21), 8145-8157.

Pelletier, G., & Fellows, L. K. (2018). A critical role for human ventromedial frontal lobe in value comparison based on multi-attribute configuration. *bioRxiv*, 483719.

Pilly, P. K., & Seitz, A. R. (2009). What a difference a parameter makes: A psychophysical comparison of random dot motion algorithms. *Vision research*, *49*(13), 1599-1612.

Rushworth, M. F., Kolling, N., Sallet, J., & Mars, R. B. (2012). Valuation and decision-making in frontal cortex: one or many serial or parallel systems? *Current Opinion in Neurobiology*, *22*(6), 946 - 955. (Decision making) doi: <https://doi.org/10.1016/j.conb.2012.04.011>

Weiskopf, N., Hutton, C., Josephs, O., & Deichmann, R. (2006). Optimal epi parameters for reduction of susceptibility-induced bold sensitivity losses: A whole-brain analysis at 3t and 1.5t. *NeuroImage*, *33*(2), 493 - 504. doi: <https://doi.org/10.1016/j.neuroimage.2006.07.029>