PREPRINT

# Distrust Before First Sight: Knowledge- and Appearance-Based Effects of Trustworthiness on the Visual Consciousness of Faces

Anna Eiserbeck[1], Alexander Enge[1], Milena Rabovsky[2], and Rasha Abdel Rahman[1]

[1] *Humboldt-Universität zu Berlin, Department of Psychology*
[2] *University of Potsdam*

Not all visual stimuli processed by the brain reach the level of conscious perception. Previous research has shown that the emotional value of a stimulus is one of the factors that can affect whether it is consciously perceived. Here, we investigated whether social-affective knowledge influences a face's chance to reach visual consciousness. Furthermore, we took into account the impact of facial appearance. Faces differing in facial trustworthiness (i.e., being perceived as more or less trustworthy based on appearance) were associated with neutral or negative socially relevant information. Subsequently, an attentional blink task was administered to examine whether the manipulated factors affect the faces' chance to reach visual consciousness under conditions of reduced attentional resources. Participants showed enhanced detection of faces associated with negative as compared to neutral social information. In event-related potentials (ERPs), this was accompanied by effects in the time range of the early posterior negativity (EPN) component. These findings indicate that social-affective person knowledge is processed already before or during attentional selection and can affect which faces are prioritized for access to visual consciousness. In contrast, no clear evidence for an impact of facial trustworthiness during the attentional blink was found. This study was pre-registered using the Open Science Framework (OSF).

*Keywords:* social-affective knowledge, facial trustworthiness, attentional blink, social judgments, visual consciousness

Of the wealth of visual information available at any given moment in time, we consciously perceive only a fraction. Besides a stimulus' inherent attributes in terms of visual salience, our experiences and our knowledge about the world are crucial in defining which input is selected for conscious processing (Alpers et al., 2005; Maier & Abdel Rahman, 2018). In this event-related potential (ERP) study, we examined how the access of faces to visual consciousness depends on social factors. Specifically, we investigated trustworthiness, which is an important dimension for social interaction, and can be evaluated based on very different types of information such as verbally transmitted knowledge about persons and visual impressions of their faces. We focused on the question whether social-affective knowledge about a person affects the access of a face to visual consciousness. In pursuing this question, we simultaneously controlled for and compared the influence of visually derived trustworthiness impressions based on facial appearance.

Our knowledge about the character and past behavior of individuals represents an important basis for evaluations and reactions in social contexts (e.g., Abdel Rahman, 2011;

Rezlescu et al., 2012). Even minimal information, such as a single sentence, is sufficient to influence person judgements on dimensions as valence, likeability, and trustworthiness (e.g., Baum et al., 2020; Baum & Abdel Rahman, 2020; Bliss-Moreau et al., 2008; Falvello et al., 2015). Furthermore, social-affective knowledge modulates face perception, as reflected in ratings of attractiveness, facial features or emotional expressions (Hassin & Trope, 2000; Nisbett & Wilson, 1977; Paunonen, 2006; Suess et al., 2015), and perception-related components of event-related potentials (ERPs), for instance, the early posterior negativity (EPN; Abdel Rahman, 2011; Luo et al., 2016; Suess et al., 2015; Wieser et al., 2014; Xu et al., 2016). The EPN, a relative negativity occurring at around 200 to 350 ms after stimulus onset at posterior sites, has been linked to enhanced attention to and facilitated processing of affective visual stimuli (Abdel Rahman, 2011; Schacht & Sommer, 2009; Schupp et al., 2003, 2004; Suess et al., 2015; for a review of EPN effects for emotional faces, see Schindler & Bublatzky, 2020).

A second factor determining face perception and person evaluation relates directly to facial appearance (for a review, see e.g. Todorov et al., 2015). Specifically, trustworthiness impressions based on facial features represent a central dimension underlying evaluations that closely corresponds to the general perceived valence of

faces with neutral expressions (Oosterhof & Todorov, 2008). Despite its questionable validity, facial trustworthiness has been found to influence impressions and character judgements, even when explicit knowledge about the person is available (Todorov & Olson, 2008; Verosky et al., 2018), and both factors may interact with each other (Rule et al., 2012). ERP studies indicate that effects of facial trustworthiness occur during visual processing in approximately the same time range as affective knowledge effects at around 200 ms after stimulus onset or possibly even earlier (Dzhelyova et al., 2012; Marzi et al., 2014; Rudoy, 2009; Shore et al., 2017; however, also see Lischke et al., 2018).

Can social-affective knowledge, possibly in interaction with appearance, influence whether we consciously perceive a face? While the discussed evidence demonstrates robust knowledge-based and appearance-induced influences on person perception and evaluation, the question whether conscious perception is necessary for the integration of this information, or whether the processing takes place or begins already beforehand, thereby influencing what we consciously perceive in the first place, is as yet unresolved. Concerning facial trustworthiness, some studies report effects on the access to visual consciousness (e.g., Getov et al., 2015; Stewart et al., 2012), while others indicate that observed effects may be due to conscious rather than pre-conscious processing and/or due to low-level visual differences (Abir et al., 2017; Stein et al., 2018).

With respect to face-related affective knowledge, initial evidence for this notion came from a study using binocular rivalry (E. Anderson et al., 2011) in which faces previously associated with negative socially relevant information were found to dominate longer in visual consciousness than faces associated with positive or neutral information. However, this finding might not necessarily indicate a prioritized access to consciousness since the measure of visual dominance could also reflect later (i.e., conscious) prioritization (see Stein et al., 2017), and indeed, no effect has been observed for the first percept to be reported. In subsequent studies using binocular rivalry or breaking continuous flash suppression, no evidence for an influence of affective knowledge on the access to visual consciousness was found (Rabovsky et al., 2016; Stein et al., 2017; also see Stein & Verosky, 2020, who find no effects of value learning on awareness of faces). However, the findings of a recent study (Eiserbeck & Abdel Rahman, 2020) using the attentional blink paradigm (Raymond et al., 1992) did provide additional support for the hypothesis of an impact of social-affective knowledge on visual consciousness. In the attentional blink paradigm, participants are instructed to detect two target stimuli—T1 and T2—among a series of distractor images in a Rapid Serial Visual Presentation (RSVP) stream. Successful detection of T1 thereby often impairs the detection of T2 when it follows in close temporal succession of approximately 200 to 500 ms (*short lag*), whereas detection is largely unimpaired for longer intervals (*long lag*). This *attentional blink* has been ascribed to the occupa-

tion of a capacity-limited processing stage by T1 leading to a disruption of attentional processing (and/) or of working memory encoding for the T2 stimulus (for a recent review of assumed underlying mechanisms, see Zivony & Lamy, 2020). This makes it possible to investigate which attributes of a stimulus determine access to conscious perception when attentional resources are limited. In line with an often observed detection advantage for emotional stimuli in the attentional blink (e.g., A. K. Anderson & Phelps, 2001; Schwabe et al., 2011), enhanced detection was observed for faces associated with negative as compared to neutral social behavioral information, whereas no effect of facial trustworthiness was found (Eiserbeck & Abdel Rahman, 2020). However, the results of this study have left questions open: The null effect of facial trustworthiness might be due to the fact that this factor comprised only two levels—average-trustworthy and low-trustworthy faces. Effects might depend on the inclusion of a broader range from low- to high-trustworthy faces, which was implemented in the current study (see below). Furthermore, no clear all-or-none pattern was found for influences of affective knowledge on visual consciousness, but rather a modulation of the strength or quality of the resulting percept—which raises the question whether the differences occurred at the time of attentional selection for visual consciousness, or at a later point in time. Although the result is in line with accounts that assume graded consciousness in the attentional blink (e.g., Fazekas & Overgaard, 2018; also see Eiserbeck et al., 2021), more direct evidence on the time course of the processing of affective knowledge in regard to the access to conscious perception is needed.

ERP studies on the neural correlates of the access to visual consciousness in the attentional blink have revealed a first (larger) divergence between detected and undetected stimuli in the time range of the N2 component at around 250 ms after stimulus onset, with enhanced negative amplitudes for detected stimuli over posterior regions (Koivisto & Revonsuo, 2008; Sergent et al., 2005). A similar early negative deflection has been observed using other paradigms as part of the broader *visual awareness negativity* (VAN) (for a review, see Koivisto & Revonsuo, 2010), which can occur as early as 100 ms after stimulus onset and last up to about 350 ms, including the time span of the N1 and N2 components. The VAN has been described as the correlate of visual awareness most consistently found across studies and is assumed to be indicative of the subjective experience of seeing (Koivisto & Revonsuo, 2010). Other perspectives (Sergent et al., 2005) suggest that modulations during the N2/VAN time range reflect preconscious differences and that later components (possibly in the P3 time range) mark the access to visual consciousness. These approaches can be roughly summarized by assuming that the selection for access to visual consciousness occurs or begins at around 250 ms after stimulus onset or possibly even later. Crucially, the early ERP markers related to consciousness coincide in time and space with the early ERP correlates (i.e., the described EPN effects) of social-

affective knowledge and visually derived trustworthiness. Due to this overlap it appears plausible to assume that affective knowledge and facial trustworthiness are integrated or processed before or while selection for conscious perception occurs, and thereby have the potential to influence which stimuli receive access to visual consciousness. This idea is further supported by an overlap in the functional significance of the components, taken to reflect enhanced attention towards and prioritized (conscious) processing of certain stimuli.

In the present study, we investigated effects of knowledge- and appearance-based trustworthiness on the access to conscious perception by utilizing the attentional blink paradigm combined with ERPs extracted from the EEG to examine the time course of the effects. Faces differing in facial trustworthiness (covering a range from low to high facial trustworthiness) were associated with negative or neutral social information, with manipulation checks for both factors included in the experiment. Subsequently, they were presented as T2-stimuli in an attentional blink task. As outlined above, facial trustworthiness represents perceptually salient affective information which may interact with social-affective knowledge. Taking into account and comparing the impact of facial trustworthiness as a second, more strongly visually based source of affective value may be informative in regard to the mechanisms underlying access of faces to visual consciousness: Is visual consciousness influenced by the overall affective/trustworthiness value ascribed to a face or does it depend on the type of information? Which type of information has a stronger impact and do the different factors interact?

In behavioral data we expected higher detection rates for faces associated with negative as compared to neutral knowledge (Eiserbeck & Abdel Rahman, 2020), and for less trustworthy as compared to more trustworthy looking faces (Abir et al., 2017). Based on reported congruency effects of affective knowledge and facial trustworthiness (e.g., in memory: Rule et al., 2012), we furthermore expected an interaction between both factors, with enhanced detection of faces with congruent negative information (negative knowledge combined with less trustworthy facial appearance). ERP analyses for face processing in the attentional blink focused on the N2/VAN as the earliest correlate of visual consciousness observed in attentional blink tasks in order to investigate overall differences between detected versus missed stimuli, as well as on the EPN component for knowledge- and appearance-specific effects of trustworthiness, also examining connections to the behavioral outcomes.

The hypotheses and methods of this study were preregistered using the Open Science Framework (OSF) and can be accessed under https://osf.io/us754 (pre-registration 1; for the aspect of affective knowledge) and https://osf.io/2yspe (pre-registration 2; for the aspect of facial trustworthiness and its interaction with affective knowledge).

## Methods

The methods used in this study were largely based on those used in a previous behavioral study (Eiserbeck & Abdel Rahman, 2020) and extended to include the recording and analyses of event-related potentials.

## Participants

Thirty-two native German speakers (21 female) with a mean age of 26.1 years ($SD = 6.65$) and normal or corrected-to-normal vision participated. Nineteen additional datasets were discarded based on pre-defined criteria described below. Participants provided written informed consent prior to participation. The study was conducted according to the principles expressed in the Declaration of Helsinki and was approved by the local Ethics Committee. Participants received either course credit or monetary compensation.

Planning of the sample size was based on a behavioral pilot test ($N = 5$). We used a generalized linear mixed model predicting T2 detection by affective knowledge (negative vs. neutral) and appearance[1] (continuous predictor), including by-participant and by-item random intercepts. The resulting effect size for the interaction of affective knowledge and appearance ($b = 0.16$) was entered in an a priori power analysis in R with the SIMR package (Green & Macleod, 2016). We aimed for a power of at least 80% as conventionally deemed adequate (see Green & MacLeod, 2016). After running 1,000 randomizations given different sample sizes, results indicated that we would need to test 20 participants to detect an effect with an expected power of 83.90%, 95% CI [81.47, 86.13].[2] For a balanced experimental design with a multiple of four participants and to have enough power to detect ERP effects, which may be smaller in size than the behavioral effects, we decided to test 32 participants.

Data sets were excluded and replaced if one or more of the following exclusion criteria applied, which were selected to ensure that the attentional blink manipulation was successful for all participants and that person knowledge was learned sufficiently well: (1) T1-performance below 80% (8 participants), (2) false alarm rate in T2-absent trials in the short lag condition above 50% (10 participants), (3) correct information recall (specific information or at least the valence of the information) for less than two thirds of the 24 T2-faces - as assessed by the retrieval at the end of the experiment (8 participants). In regard to ERPs during

---

[1] Please note that for reporting analyses and results throughout this article, we use the term *appearance* (rather than *facial trustworthiness*) to refer to the predictor in order to avoid confusion with the single-trial trustworthiness evaluations during the rating phases of the main experiment.

[2] A further power analysis was run to estimate the sample size needed to detect a main effect of affective knowledge ($b = 0.26$) which yielded a similar result of 23 participants.

the attentional blink task, these rigorous exclusion criteria served to ensure that enough (since only T1-correct trials were included in the ERP analysis) and informative (without too many guess trials that would dilute the analyses) ERP trials could be obtained.

## Materials

### Pictures

Stimuli were presented on a 19-inch LCD monitor with a 75-Hz refresh rate. During all phases of the experiment, the images were displayed on a grey background with a size subtending 5.8° vertical visual angle and 4.3° horizontal visual angle (viewing distance: 70 cm).

T2 target stimuli consisted of 24 portraits of faces (12 female) with Caucasian appearance, displaying neutral emotional expressions, taken from the Chicago Face Database (CFD; Ma et al., 2015). Based on the rating data of the database, faces were chosen to cover a range of trustworthiness evaluations from low to high perceived trustworthiness. The pictures were converted to greyscale images and cropped so that no hair and ears are visible. The outer shape of the face was retained (instead of, e.g., applying an oval mask), because shape may be a factor affecting trustworthiness impressions (Kleisner et al., 2013). To minimize low-level confounds, histograms (i.e., the distributions of brightness values) of the images were equated using the SHINE toolbox (Willenbockel et al., 2010) in MATLAB R2016a.

Six additional faces (three female) from the CFD, processed in the same way as the T2 faces, served as fillers, associated with positive knowledge during learning. They were not presented in the attentional blink task.

To serve as distractor images in the attentional blink task, 12 additional faces (6 female) with average trustworthiness ratings were chosen from the database and processed in the same way as described above. Additionally, the facial features were cut out, rotated and randomly placed in a different position within the face, thus creating abstract looking faces. For each distractor, features of two faces of the same sex were "mixed" to further contribute to an abstract impression. Our aim was to create distractors that are visually similar to the T2 targets (see Müsch et al., 2012, for the importance of target-distractor-similarity) but sufficiently distinguishable.

T1 target stimuli consisted of 36 images displaying either the face of a dog or a similarly looking blueberry muffin, all converted to greyscale and cropped to the same oval shape.

### Person-Related information

Twenty-four sentences, describing negative or neutral social behavior were recorded by a male speaker (mean duration = 2.63 s) and rated in a web-based questionnaire ($N = 20$) on valence (negative: $M = 1.67$, $SD = 0.36$; neutral: $M = 4.04$, $SD = 0.10$; difference: $t(11) = -31.3$, $p < .001$), and arousal (negative: $M = 4.94$, $SD = 0.54$; neutral: $M = 1.37$,

$SD = 0.13$); difference: $t(11) = 22.2$, $p < .001$). Six additional sentences describing a positive behavior (valence: $M = 6.39$, $SD = 0.15$; arousal: $M = 4.58$, $SD = 0.33$) served as fillers during learning. Sentences always started with "she"/"he" or "this woman"/"this man", followed by the description of a social behavior, e.g. "threatened a shop assistant with a knife" (negative knowledge condition) or "asked a waiter for the menu" (neutral knowledge condition). For a full list of sentences, see Supplement Table S1.

## Procedure

A graphical overview of the different experimental phases can be found in Figure 1B.

### Learning Phase

**Pre-Learning Ratings of Trustworthiness and Facial Expression.** Participants rated the trustworthiness and facial expression of all 30 faces (T2 target faces as well as filler faces associated with positive information in the learning phase) prior to knowledge acquisition. Ratings were completed block-wise with a counterbalanced order across participants. Faces were presented in random order within the blocks. At the beginning of each trial, a fixation cross was presented for 500 ms. Subsequently, a face was displayed for 1 s, followed by a short instruction and a 7-point scale, e.g. "Please rate the trustworthiness (/facial expression) of this woman". The ends of the scales were labeled, in case of the trustworthiness rating as "not at all trustworthy" and "very trustworthy," and for the facial expression rating as "negative" and "positive". The direction of the scales (left to right or right to left) was counterbalanced across participants. Participants used the left mouse button to indicate their choice. There were no time constraints for responses.

**Knowledge Acquisition.** After completion of the ratings, participants acquired knowledge about the persons. To this end, each of the 30 faces was presented together with the accompanying auditory information. During each trial, first a fixation cross was shown for 500 ms. Subsequently, the face was displayed for 6 s. Beginning at 1 s after face onset, the auditory information was presented via loudspeakers. Assignment of faces to affective knowledge conditions was counterbalanced across participants, such that each of the 24 T2-target faces was associated equally often with negative and neutral information. The filler-faces were accompanied by the same information for all participants. To foster learning, each face was presented together with the accompanying information for a total of five times in blocks of gradually increasing numbers of faces (4, 6, or 12 faces from each affective knowledge condition plus 2, 3, or 6 filler faces) and simple judgement tasks related to the presented behaviors were included (e.g., "Is this person's behavior common?"; Abdel Rahman, 2011; Baum et al., 2020; Suess et al., 2015).

After learning, the EEG was prepared and then recorded during the attentional blink task, the subsequent rating task,

and an eye movement calibration procedure at the end of the experiment.

### Test Phase

**Attentional Blink Task.** For each trial of the attentional blink task, first a fixation cross was presented for 500 ms. Then, 13 pictures were shown in rapid succession (for illustration, see Fig. 2A), with a presentation time of 107 ms each and without a time interval between pictures. Regular trials contained 11 distractor images, which were presented in randomized order, and two targets: a dog or muffin (T1) and a face (T2). T2 (if present) was always presented as the 10th stimulus whereas T1 position varied: It was either presented as the 3rd stimulus (entailing a lag of 7 items between T1 and T2; long lag) or as the 7th stimulus (entailing a lag of 3 items; short lag). The task comprised 696 trials in total. As T1, in 50% of cases a dog was shown and in 50% a muffin. All T2-faces were presented equally often—resulting in an equal number of trials for the two affective knowledge conditions (144 trials per affective knowledge condition for each short and long lag). In order to estimate the false alarm rate for each participant, within each lag, T2 was absent in 60 trials (17%) and instead another distractor was presented. All trial types (short or long lag, T2 present or absent, neutral or negative knowledge) were presented in randomized order.

Participants were instructed to look for the dog/muffin and the face. They were informed that both targets are equally important, but that not every sequence contains a face. After each trial, participants indicated via response keys (a) whether they saw the image of a dog or a muffin as T1 (options: *dog* / *muffin* / *I don't know*), (b) whether they saw a male or a female face as T2 (options: *male* / *female* / *I don't know*), and (c) how clear their subjective impression of T2 was on a four-point perception awareness scale (PAS; Ramsøy & Overgaard, 2004; options: *not seen* / *slight impression* / *strong impression* / *seen completely*).

**Post-Learning Ratings of Trustworthiness and Facial Expression.** The procedure of the second rating phase was identical to the first rating phase except that the tasks (trustworthiness and facial expression rating) were repeated three times. This was done in order to obtain enough trials for the ERP analyses.

After the experiment, successful acquisition of person-related information was checked via a computerized survey. Participants indicated which kind of behavior (negative or neutral) was associated with each target face and what they recalled that particular behavior to be.

### EEG Recording and Preprocessing

During the attentional blink task and post-learning ratings, the EEG was recorded with Ag/AgCl electrodes at 62 scalp sites according to the extended 10–20 system at a sampling rate of 500 Hz and with all electrodes referenced to the left mastoid. An external electrode below the left eye was used to measure electrooculograms. During recording, low- and high-cut-off filters (0.016 Hz and 1000 Hz) were applied and all electrode impedances were being kept below

10 kΩ. After the experiment, a calibration procedure was used to obtain prototypical eye movements for later artifact correction. Processing and analyses of the data were based on the EEG-processing pipeline by Frömer et al. (2018). An offline pre-processing was conducted using MATLAB (Version R2016a) and the EEGLAB toolbox (Version 13.5.4b; Delorme & Makeig, 2004). The continuous EEG was re-referenced to a common average reference and eye movement artifacts were removed using a spatio-temporal dipole modeling procedure with the BESA software (Ille et al., 2002). The corrected data were low-pass filtered with an upper pass-band edge at 40 Hz. Subsequently, they were segmented into epochs of -200 to 1,000 ms relative to T2 onset, and baseline-corrected using the 200 ms pre-stimulus interval. Segments containing artifacts (absolute amplitudes over ±150 µV or amplitudes changing by more than 50 µV between samples) were excluded from further analysis.

## Data Analyses and Results

Analyses were conducted in R (Version 3.5.1, R Core Team, 2018) using the lme4 package (Version 1.1-17; Bates et al., 2015) and the lmerTest package (Version 3.1-1; Kuznetsova et al., 2017) to calculate *p*-values via the Satterthwaite approximation in the case of linear mixed models (while in case of generalized linear mixed models, *p*-values were based on the Wald *z*-test implemented in lme4). In all (G)LMM analyses, we aimed to include the maximal random effects structures justified by the design (Barr et al., 2013). If models failed to converge or yielded a singular fit, random effects were excluded based on least explained variance as indicated by the singular value decomposition of the non-converging model. Facial appearance was treated as a continuous predictor. To this end, the mean rating value of trustworthiness in the first rating phase across participants served as *appearance* score for each face.

### Manipulation Checks

As manipulation checks, we examined evaluations during both rating phases as well as ERPs during the second rating phase after the attentional blink task. This enabled us to assess the effects of affective knowledge and facial trustworthiness under conditions of conscious perception.

### Trustworthiness and Facial Expression Ratings

Rating data were analyzed with linear mixed models (LMMs) including crossed random effects for subjects and items, with trustworthiness or expression rating serving as dependent variable. The models include the fixed factors phase (before learning / after learning), affective knowledge (neutral / negative) and appearance. The predictors affective knowledge and appearance were nested within phase to specifically test effects before learning and after learning. Effect coding was applied for the factor affective knowledge (neutral: -0.5, negative: 0.5); the continuous predictor appearance was mean-centered.
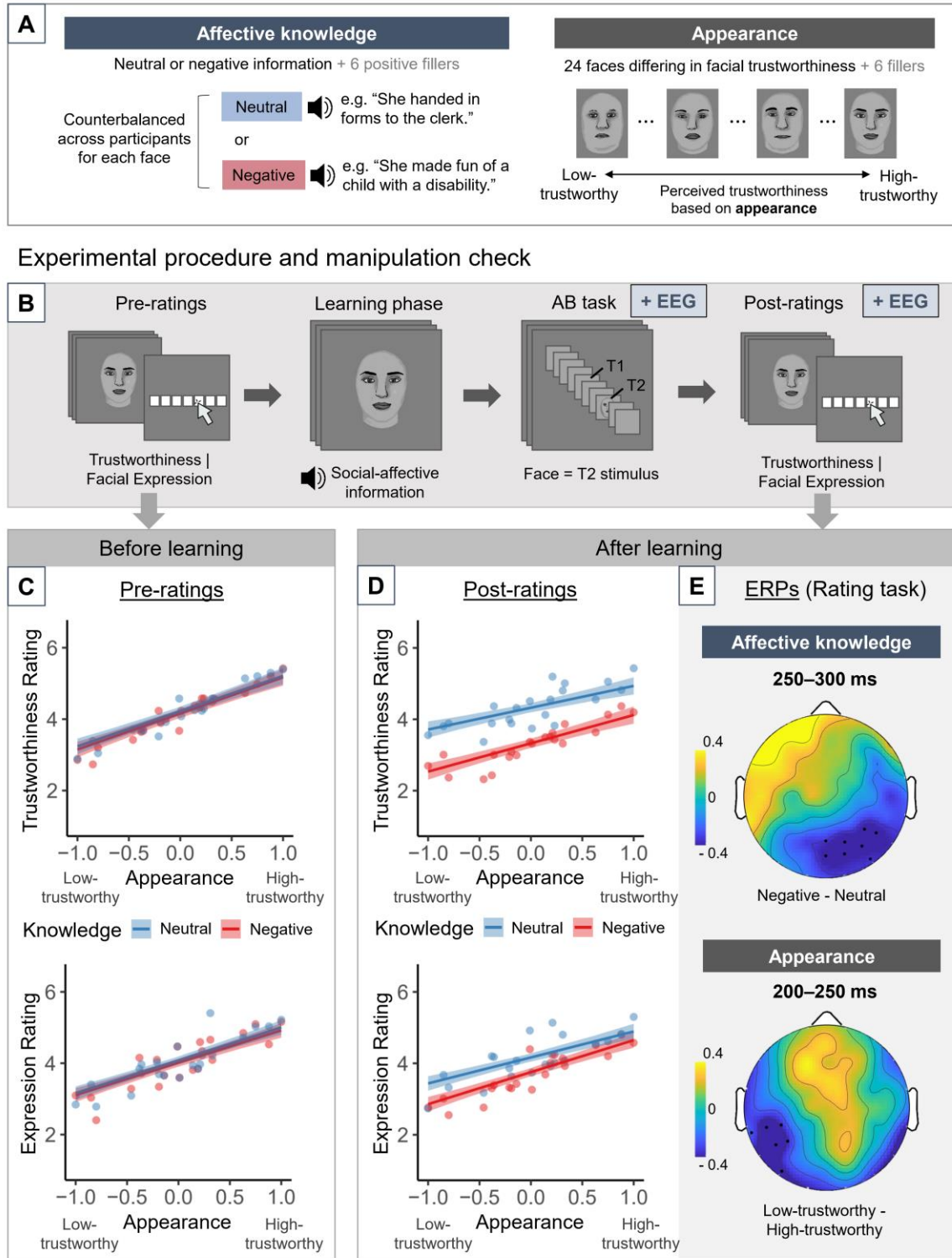
*Figure 1.* Experimental paradigm and manipulation check. Please note that in the illustration, the face stimuli used in the experiment (edited photographs of real faces; see description in the Methods section) have been replaced by roughly similar looking drawings. (A) Overview of the experimental manipulations. (B) The procedure consisted of pre-ratings of trustworthiness and facial expression (with task order counterbalanced across participants), learning of affective person knowledge, the attentional blink task, and post-ratings. (C), (D) Trustworthiness and facial expression ratings show a successful manipulation of affective knowledge and appearance: Before learning (C), only appearance influenced ratings; after learning, there were additive effects of both affective knowledge and appearance (D). The plots show partial effects from the LMM analyses; dots illustrate (descriptive) mean ratings for individual faces; error bands depict 95% confidence intervals. (E) Topographies of grand average ERPs time-locked to face onset illustrate independent effects during the EPN time range for affective knowledge (top) and appearance (bottom). For plotting, the continuous variable of appearance was converted into two levels (low- und high-trustworthy) by means of a median split.

Figures 1B and 1C provide an overview of the ratings before and after knowledge acquisition. A highly similar pattern of results was observed for trustworthiness and facial expression ratings: Before learning, there was a main effect of appearance (trustworthiness: $b = 1.00$, $t[41.08] = 8.47$, $p < .001$; expression: $b = 0.92$, $t[42.88] = 6.43$, $p < .001$), with more positive ratings for more trustworthy compared to less trustworthy faces (i.e. mean appearance scores—representing "consensus" tendencies across all participants—predicted individual participants' trustworthiness and expression ratings). Ratings for faces associated with neutral and negative knowledge did not differ significantly before knowledge acquisition (trustworthiness: $b = -0.06$, $t[51.33] = -0.54$, $p = .593$; expression: $b = -0.07$, $t[62.36] = -0.75$, $p = .457$), and there was no interaction between appearance and affective knowledge (trustworthiness: $b = 0.07$, $t[56.64] = 0.51$, $p = .610$; expression: $b = 0.02$, $t[173.09] = 0.13$, $p = .895$). After learning, appearance still predicted trustworthiness and expression ratings (trustworthiness: $b = 0.70$, $t[41.08] = 5.93$, $p < .001$; expression: $b = 0.80$, $t[42.88] = 5.64$, $p < .001$). Furthermore, there was a main effect of affective knowledge in both trustworthiness and expression task: Faces associated with negative information were rated as less trustworthy ($M = 3.29$) than faces associated with neutral information ($M = 4.32$; $b = -1.00$, $t[51.33] = -8.68$, $p < .001$), and the expressions of faces associated with negative information were rated as more negative ($M = 3.70$) than the expressions of faces associated with neutral information ($M = 4.13$; $b = -0.41$, $t[62.36] = -4.63$, $p < .001$). There was no significant interaction of appearance and affective knowledge (trustworthiness: $b = 0.13$, $t[56.64] = 0.96$, $p = .340$; expression: $b = 0.15$, $t[173.09] = 1.25$, $p = .215$). For full LMM output, see Supplement Table S2.

### ERPs During Second Rating Phase

As specified in the pre-registration, analysis of the EPN component during the second rating phase was based on previous evidence regarding effects of affective knowledge (Abdel Rahman, 2011; Suess et al., 2015) and facial trustworthiness (Marzi et al., 2014), and focused on a time range from 200 to 350 ms including electrodes PO7, PO8, PO9, PO10, TP9 and TP10. In a LMM analysis, the single-trial mean amplitudes across this time range and electrodes were predicted by the fixed effects affective knowledge, appearance and their interaction. Further model specifications were the same as described above. The analysis yielded no significant differences for either affective knowledge ($b = -0.11$, $t[76.97] = -0.77$, $p = .442$), appearance ($b = 0.12$, $t[25.88] = 0.77$, $p = .450$) or an interaction of knowledge and appearance ($b = -0.03$, $t[416.65] = -0.14$, $p = .891$). Since visual inspection of separate 50 ms time windows indicated more narrow effects with differential topographical and temporal distributions for affective knowledge and appearance within the EPN time range (see Figure 1D), we conducted additional exploratory analyses with restricted time ranges and electrode sites. A first analysis was focused on the time range of 250 to 300 ms and electrodes P4, P6, P8, POz, PO4, PO8, PO10, Oz, and O2. A

significant main effect of affective knowledge was observed ($b = -0.43$, $t[20.50] = -2.45$, $p = .024$), with enhanced negative amplitudes for faces associated with negative as compared to neutral person knowledge, whereas neither a main effect of appearance ($b = -0.19$, $t[22.38] = -1.19$, $p = .276$) nor an interaction of appearance and affective knowledge ($b = -0.21$, $t[21.58] = -0.76$, $p = .456$) was found (for full statistical output, see Supplement Table S3). A second analyses was focused on the time range of 200 to 250 ms and electrodes PO9, P5, TP7, CP5, P7, and TP9. A main effect of appearance was observed ($b = 0.32$, $t[20.62] = -2.46$, $p = .023$), with enhanced negative amplitudes for less trustworthy as compared to more trustworthy looking faces, whereas neither a main effect of affective knowledge ($b = -0.07$, $t[172.70] = -0.51$, $p = .608$) nor an interaction between appearance and affective knowledge ($b = -0.14$, $t[70.43] = -0.60$, $p = .549$) was found (for full statistical output, see Supplement Table S4). In the LPP time range, no significant differences were observed (all $ps \geq .232$; for full statistical output, see Supplement Table S5). Additional exploratory analyses indicated no differences in the P1 or N170 time range.

## Main Task: Attentional Blink

### Behavioral Data

Behavioral data of the attentional blink task were analyzed with binomial generalized linear mixed models (GLMMs), only including trials in which T1 was correctly identified in order to ensure that attention was paid to the first target as a pre-requisite for the attentional blink to occur. Hit (encoded as 1 for hit and 0 for miss) served as the dependent variable. As specified in the pre-registration, analyses were conducted separately for two criteria defining trials as T2 hit or miss (see Figures 2B and 2C). To count as a hit trial, for both criteria, the gender of T2 needed to be classified correctly. Furthermore, participants had to indicate either at least a *slight impression* (liberal hit criterion), or at least a *strong impression* (strict hit criterion) as the subjectively rated visibility of T2. These two different criteria were implemented since previous findings indicate that it may be important to take into account the threshold for considering a trial as hit or miss (see Eiserbeck & Abdel Rahman, 2020). To verify the presence of an (overall) attentional blink effect, GLMMs with the fixed factor lag (short lag / long lag) were computed. To test the hypotheses concerning affective knowledge and appearance-based trustworthiness, analyses were confined to short lag trials (as specified in the pre-registration), including the fixed factors affective knowledge (neutral / negative) and appearance. Effect coding was applied for the factors lag (short: -0.5, long: 0.5) and affective knowledge (neutral: -0.5, negative: 0.5); the continuous predictor appearance was mean-centered.

Mean T1 recognition rate was 91.06% (CI $\pm$ 0.38). Mean correct rejection rate in T2 absent trials was 85.37% (CI $\pm$ 1.12). The GLMM analyses confirmed the presence of an attentional blink effect, i.e., an effect of lag, for both liberal ($b = 0.95$, $z = 5.78$, $p < .001$) and strict criterion ($b = 0.76$, $z = 4.52$, $p < .001$) and with higher hit rates in the

long as compared to the short lag condition (liberal criterion: 83.32%, CI ± 1.09 vs. 69.52%, CI ± 1.25; strict criterion: 48.64%, CI ± 1.25 vs. 36.74%, CI ± 1.22; for full statistical output, see Supplement Table S6).

Table 1 contains estimates (regression coefficients $b$) of the fixed effects, standard errors, $z$- and $p$-values for the analyses of short lag trials for both hit criteria. Graphical illustrations of descriptive by-participant and by-item hit rates (i.e., the proportion of T2 hits in T1-correct trials) and distributions can be found in Figures 2D and 2E. For the liberal criterion, neither main effects of affective knowledge or appearance nor their interaction reached statistical significance. In contrast, for the strict criterion, a main effect of affective knowledge was found. Mean hit rates were higher in the negative (38.23%, CI ± 1.42) relative to the neutral

knowledge condition (35.27%, CI ± 1.39). The main effect of appearance did not reach statistical significance ($p$ = .082) but a trend for an enhanced detection of low- as compared to high-trustworthy faces was observed. The interaction effect of affective knowledge and appearance did not reach statistical significance ($p$ = .102). Nonetheless, in a further exploratory analysis, we were interested in testing influences of appearance separately for each knowledge condition and observed a significant effect of appearance in the neutral knowledge condition ($b$ = -0.40, $z$ = -2.44, $p$ = .015) but not in the negative knowledge condition ($b$ = -0.17, $z$ = -0.88, $p$ = .378); for full model output, see Table 2. Graphical illustrations of predicted hit rates by appearance and knowledge can be found in Figures 2E (liberal hit criterion) and 2F (strict hit criterion).

Table 1

*GLMM Statistics for Analysis of Short Lag Trials*

| Variable | Liberal hit criterion | | | | Strict hit criterion | | | |
|---|---|---|---|---|---|---|---|---|
| | $b$ | SE | $z$ | $p$ | $b$ | SE | $z$ | $p$ |
| Intercept | 0.99 | 0.20 | 4.87 | **<.001** | -1.12 | 0.38 | -2.96 | **.003** |
| Knowledge (Neg-Neu) | 0.01 | 0.07 | 0.13 | .897 | 0.21 | 0.09 | 2.25 | **.025** |
| Appearance | -0.19 | 0.12 | -1.60 | .111 | -0.29 | 0.17 | -1.74 | .082 |
| Knowledge:Appearance | 0.09 | 0.11 | 0.78 | .437 | 0.23 | 0.14 | 1.64 | .102 |
| Random effects | Var. | SD | | | Var. | SD | | |
| Participants (Intercept) | 1.20 | 1.10 | | | 4.16 | 2.04 | | |
| Knowledge | 0.02 | 0.15 | | | 0.04 | 0.21 | | |
| Appearance | 0.12 | 0.34 | | | 0.11 | 0.34 | | |
| Knowledge:Appearance | 0.02 | 0.13 | | | 0.06 | 0.25 | | |
| Items (Intercept) | 0.61 | 0.78 | | | 0.20 | 0.45 | | |
| Knowledge | 0.03 | 0.18 | | | 0.05 | 0.21 | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

Table 2

*GLMM Statistics for Analysis of Short Lag Trials With Appearance Nested Within Knowledge Condition (for Strict Hit Criterion)*

| Variable | Strict hit criterion | | | |
|---|---|---|---|---|
| | $b$ | SE | $z$ | $p$ |
| Intercept | -1.12 | 0.38 | -2.96 | **.003** |
| Knowledge (Neg-Neu) | 0.21 | 0.09 | 2.25 | **.025** |
| Neg/Appearance | -0.17 | 0.19 | -0.88 | .378 |
| Neu/Appearance | -0.41 | 0.19 | -2.44 | **.015** |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; "/" indicates nesting of fixed factors. Variances and standard deviations of the random effects are equivalent to those provided for the strict hit criterion in Table 1.
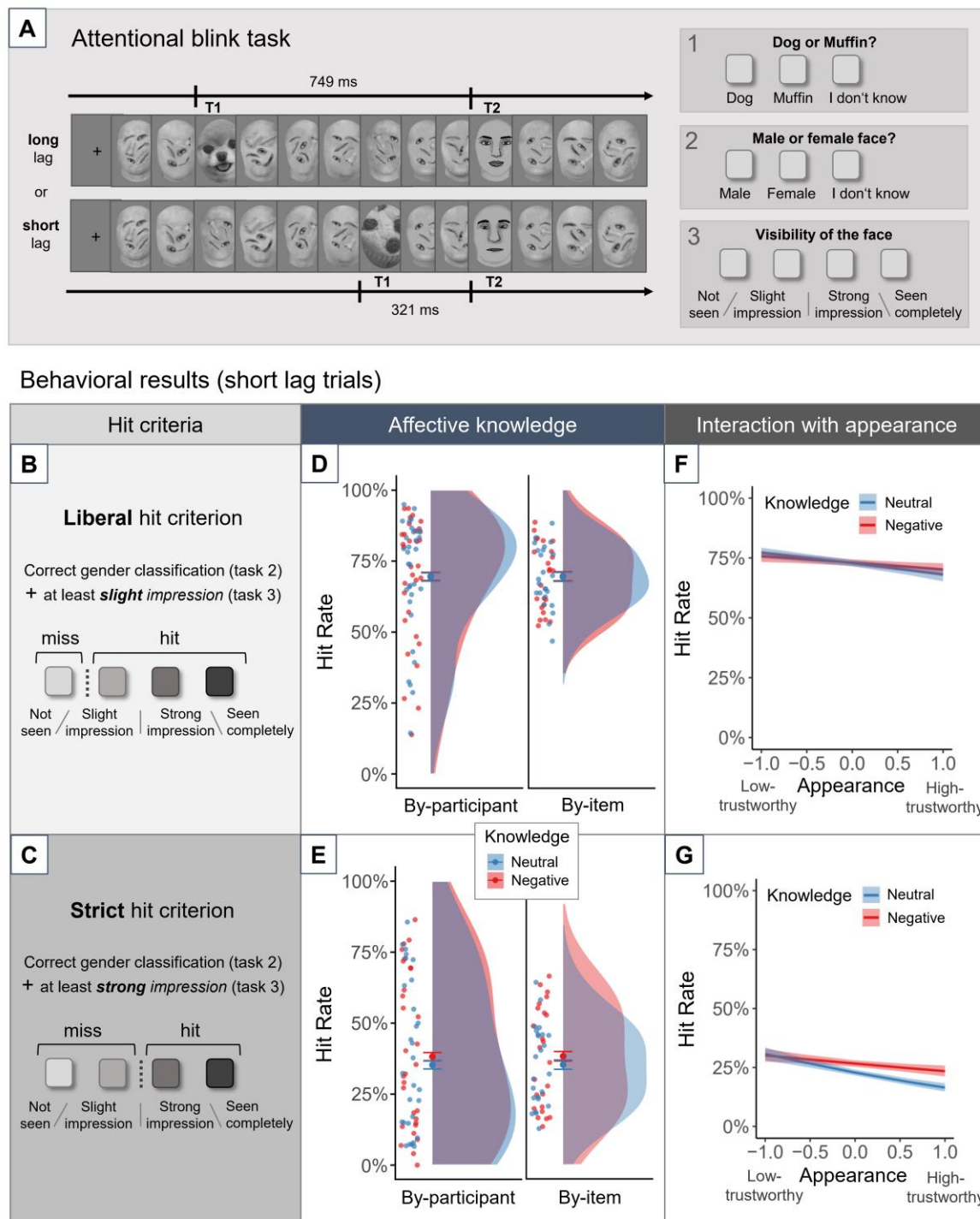
*Figure 2.* Design and behavioral results of the attentional blink task. Please note that in the illustration, the face stimuli used in the experiment (edited photographs of real faces; see description in the Methods section) have been replaced by roughly similar looking drawings. (A) Illustration of T2-present trials in the attentional blink task. 13 images were shown in rapid succession, with a presentation time of 107 ms each. In long lag trials, there was a 749 ms interval between T1-onset and T2-onset. In short lag trials, there was a 321 ms interval. T2 faces differed in associated knowledge (neutral or negative, counterbalanced across participants) and facial trustworthiness. After each trial, participants answered three questions via button press regarding the identity of T1, the gender of the T2 face, and the visibility of T2. (B), (C) Overview of the two hit criteria implemented in the analyses, representing two differently stringent thresholds for considering a trial as a T2 hit or miss. (D), (E) By-participant and by-item hit rates (ratio of T2 hits in T1-correct trials) and distributions in short lag T2-present trials depending on knowledge condition, for liberal (D) and strict hit criterion (E). Error bars depict 95% confidence intervals. Analyses on a single trial basis using GLMMs revealed a significant main effect of affective knowledge for the prediction of hit or miss for the strict but not for the liberal hit criterion. (F), (G) Predicted hit rates in short lag trials depending on knowledge condition and appearance for liberal (F) and strict (G) hit criterion, based on partial effects from GLMM analyses. Error bands depict 95% confidence intervals.

- 9 -

### ERPs

**Neural Correlates of Consciousness.** As specified in the pre-registration (2), a first analysis focused on the N2 component. Based on previous evidence (e.g. Del Cul et al., 2007; Sergent et al., 2005) as well as visual inspection of the differences between T2-present and T2-absent short lag trials (representing a contrast which is independent of the analyses relevant to our hypotheses), analyses included the time range of 220 to 300 ms after T2 onset and electrodes TP9, TP10, P7, P8, PO9, PO10, O1, and O2. Single-trial mean amplitudes across this time range and region of interest were centered and entered as a single fixed effect in a GLMM predicting detection (hit / miss), also including by-participant and by-item random intercepts.[3] A significant main effect of mean N2 amplitude in the prediction of hits was observed for both hit criteria (liberal hit criterion: $b = -0.09$, $z = -13.66$, $p < .001$; strict hit criterion: $b = -0.09$; $z = 12.81$, $p < .001$), with enhanced negative amplitudes for hit as compared to miss trials (for full model outputs, see Supplement Table S7).[4] In line with evidence regarding the visual awareness negativity (Förster et al., 2020; Koivisto & Revonsuo, 2010), visual inspection indicated that the observed N2 differences are part of a broader pattern of posterior negativity ranging from approximately 150 to 400 ms (for topographical differences between hit and miss trials across the whole time range, see Supplement Figures S1 and S2). The prediction of detection by mean amplitudes obtained by averaging across this broader time range was confirmed in additional analyses (liberal hit criterion: $b = -0.1$, $z = -13.21$, $p < .001$; strict hit criterion: $b = -0.09$; $z = 11.91$, $p < .001$; for full model outputs, see Supplement Table S8).

**Effects of Affective Knowledge and Appearance.** Following the idea of using the rating phase as a localizer task (as planned in the pre-registration), we analyzed the same regions of interest and time ranges as previously identified and reported above in order to examine the presence of EPN effects. A first analysis was focused on the time range of 200 to 250 ms and electrodes PO9, P5, TP7, CP5, P7, and TP9 for which a main effect of appearance had been found in the rating task. No significant effects were observed for either affective knowledge ($b = -0.06$, $t[20.58] = -0.51$, $p = .613$), appearance ($b = 0.20$, $t[23.14] = 1.78$, $p = .089$) or an interaction between affective knowledge and appearance ($b = 0.19$, $t[21.35] = 0.99$, $p = .334$, for full statistical output, see Supplement Table S9), albeit in case of appearance a trend for enhanced negative amplitudes for less trustworthy faces was apparent (see statistical values above). A second analysis was focused on the time range of 250 to 300 ms and electrodes P4, P6, P8, PO3, POz, PO4, PO8, PO10, Oz, and O2 for which a main effect of knowledge had been found in the rating task. Again, no significant effects were observed for affective knowledge ($b = 0.15$, $t[38.21] = 1.06$, $p = .298$), appearance ($b = -0.14$, $t[22.28] = -1.22$, $p = .234$) or an interaction between appearance and affective knowledge ($b = -0.09$, $t[410.87] = -0.52$, $p = .601$; for full statistical output, see Supplement Table S10).

However, the observed null effects could be due to two influencing factors: Firstly, EPN activity may be temporally or topographically shifted compared to activity during ratings due to task differences. Secondly, effects of affective knowledge and/or appearance might not have been observed due to an interaction with detection, i.e., the described N2/VAN effect. Based on these considerations, we conducted additional exploratory analyses: We utilized cluster-based permutation tests (CBPTs) as a systematic method to examine potential connections between the detection of faces and the factors affective knowledge and appearance within the EPN time range and region of interest. Furthermore, we again also examined potential main effects of affective knowledge and appearance as well as an interaction of knowledge and appearance. A detailed description of the CBPT procedure can be found in Supplementary Information S1. The CBPT approach enabled us to limit the time range and electrode sites to where and when EPN effects can be expected while at the same time allowing for variance in regard to the location of effects and controlling for multiple testing. No differences were observed for the comparison of affective knowledge conditions, appearance, an interaction of affective knowledge and appearance, or an interaction of appearance and detection.

However, the test indicated a connection between affective knowledge and detection. Based on the topographical and temporal distribution of the corresponding cluster—approximately 180 to 240 ms at electrodes P7, P5, P3, Pz, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, O1, Oz, and O2—single-trial mean amplitudes were obtained and mean-centered to be entered into GLMM analyses. Thereby, as planned in the pre-registration, we extended the previously specified GLMM described in the behavioral analyses by the predictor *mean ROI amplitude*: Hits (0/1) were predicted by knowledge (neutral / negative), appearance (continuous predictor) and mean ROI amplitude (continuous predictor), including all interactions between the predictors, as well as the previously specified random effects structure. The models revealed a significant interaction between affective knowledge and mean amplitude for the liberal hit criterion ($b = -0.05$, $z = -3.65$, $p < .001$) as well as for the strict hit criterion ($b = -0.03$, $z = -2.01$, $p = .044$; for full model outputs, see supplement Table S11). As graphically illustrated in Figure 3C and 3D, nested models showed that the mean amplitude has a higher predictive value for

---

[3] Since neural activity precedes the behavioral outcomes, we deemed it more plausible to predict detection by neural activity (mean ROI amplitude) rather than predicting neural activity by detection.

[4] For completeness, we would like to note that N2 activity based on the same data was investigated in the context of a different research topic and with different model specifications in another manuscript (Eiserbeck et al., 2021), where we examined N2 activity for each level of visibility (*not seen / slight impression / strong impression / seen completely*) and observed graded amplitude differences between successive levels.
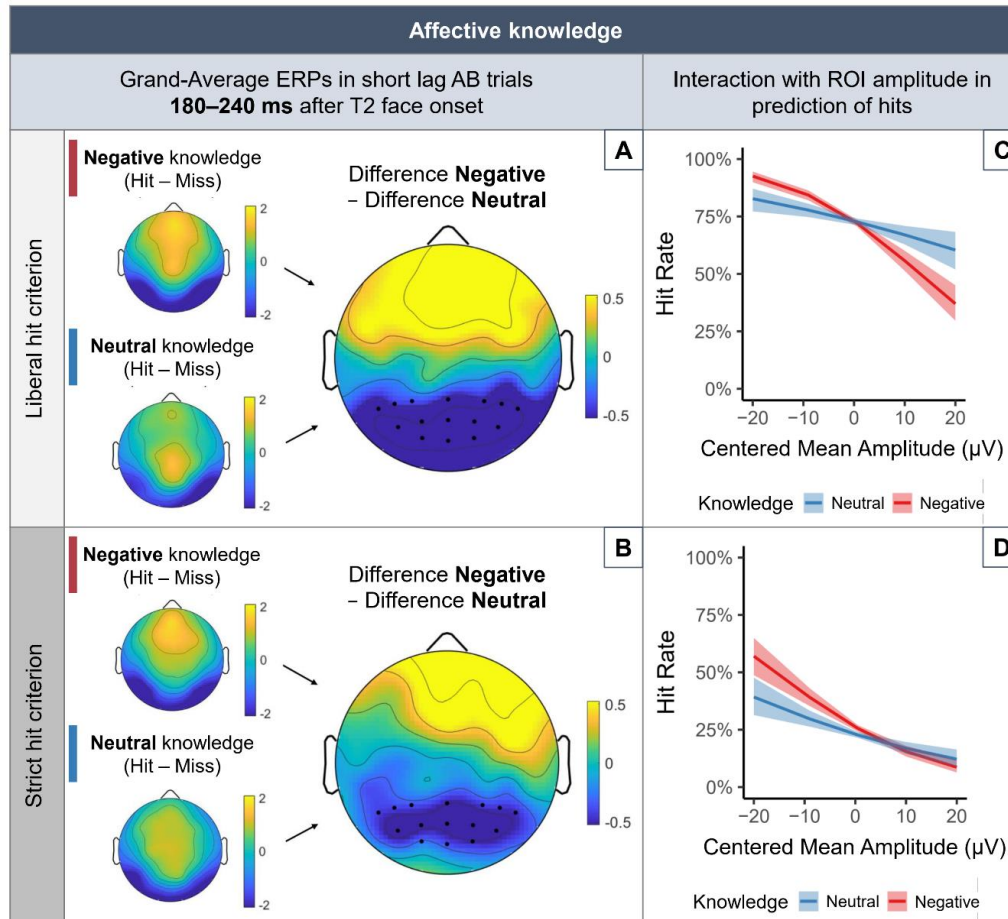
*Figure 3.* Connection of ERP activity between 180 to 240 ms to behavioral outcomes in short lag attentional blink trials. (A), (B) Topographies of grand average ERPs time-locked to T2 onset illustrate a more pronounced difference between T2 hit and miss trials in the negative as compared to the neutral knowledge condition, for both liberal (A) and strict hit criterion (B). Electrode markers indicate sites taken into account in the GLMM analyses. (C), (D) Interaction of affective knowledge and mean amplitude between 180 and 240 ms at electrodes P3, Pz, P4, P6, PO7, PO3, POz, PO4, O1, and O2 for the prediction of hits. GLMM analyses revealed a stronger predictive value of ROI activity in the negative as compared to the neutral knowledge condition.

determining a trial as a hit or miss in the negative knowledge condition (liberal criterion: $b = -0.08$, $z = -8.09$, $p < .001$; strict criterion: $b = -0.07$, $z = -6.75$, $p < .001$) as compared to the neutral knowledge condition (liberal criterion: $b = -0.03$, $z = -3.03$, $p = .002$; strict criterion: $b = -0.04$, $z = -3.83$, $p < .001$; for full model output, see supplement Table S12).

Additional LMM analyses were conducted to examine whether the observed amplitude differences between the negative and neutral knowledge condition could be due to an effect in hit trials specifically: No significant effect of affective knowledge on mean ROI amplitude during hit trials was observed (liberal criterion: $b = -0.19$, $t[31.36] = -1.29$, $p = .207$; strict criterion: $b = -0.19$, $t[23.25] = -1.13$, $p = .271$; for full model output, see supplement Table S13).

## Discussion

Can social-affective knowledge about persons affect the access of faces to visual consciousness? In the present study, we investigated this question while additionally taking into account the impact of facial trustworthiness as a source of affective information derived from visual appearance. As a manipulation check, we tested and observed

effects of affective knowledge and facial appearance in explicit evaluations of trustworthiness and emotional expressions and associated ERP modulations. Knowledge- and appearance-based effects did not interact (see Fig. 1). In the attentional blink, replicating previous behavioral findings (Eiserbeck & Abdel Rahman, 2020), stimulus detection under conditions of reduced attention was affected primarily by affective knowledge about the person and not by his or her appearance (see Fig. 2). This influence of knowledge on the access to visual consciousness was also reflected in the ERPs: The mean amplitude at posterior regions between 180 to 240 ms after face onset—a location and time range corresponding to the EPN component—had a higher predictive value in determining whether a face was detected or not in the negative as compared to the neutral condition (see Fig. 3). No clear evidence for an impact of facial trustworthiness in the attentional blink was found.

## Manipulation Check: Evaluations of Trustworthiness and Facial Expressions

The trustworthiness and facial expression evaluations after learning demonstrate successful manipulations of affective knowledge and appearance. In line with previously reported effects of affective person knowledge, faces

associated with negative information were rated as less trustworthy and their emotional expression were rated as more negative than faces associated with relatively neutral knowledge. No interaction between affective knowledge and appearance was observed. ERPs confirmed this pattern of independent effects: While modulations in the EPN time range were found for both factors, their time course and topographical distributions differed. For affective knowledge, an effect with a pattern typical for the EPN component could be observed, with enhanced negative mean amplitudes for faces associated with negative as compared to neutral knowledge between approximately 250 to 300 ms after stimulus onset at posterior sites (see Abdel Rahman, 2011; Luo et al., 2016; Suess et al., 2015; Wieser et al., 2014; Xu et al., 2016). Activity in this region of interest and time range was not influenced by appearance and no interaction between affective knowledge and appearance was found. An effect of appearance was found slightly earlier, emerging around 200 ms, over left temporo-parietal sites. It corresponded to the expected direction of enhanced negative amplitudes for less trustworthy as compared to more trustworthy faces (Marzi et al., 2014). Contrary to findings of other studies (Lischke et al., 2018; Yang et al., 2011), no modulations during the LPP time range were observed. Since not directly relevant for interpretation of the attentional blink results as the main focus of this study, and because the evidence on ERP effects of facial trustworthiness is scarce, this is not further discussed. Overall, the behavioral and ERP data associated with trustworthiness and emotional expression evaluations indicate successful manipulations of both factors, with independent contributions of affective knowledge and facial trustworthiness at relatively early, perception-related stages.

## Effects of Affective Knowledge on the Access to Visual Consciousness

In the attentional blink task, detection under conditions of reduced attention was enhanced for faces associated with negative as compared to neutral knowledge. Since the assignment of faces to affective knowledge condition was counterbalanced across participants, this finding cannot be explained by low-level visual differences. Replicating a previous report (Eiserbeck & Abdel Rahman, 2020), the effect depended on the subjective visibility rating and was only observed for the strict hit criterion of at least a *strong impression*, not for the liberal criterion of at least a *slight impression*. This result pattern indicates that the intensity or quality of the percept—rather than the precision with which the objective (gender classification) task is solved—is influenced by affective knowledge (see Eiserbeck & Abdel Rahman, 2020; Fazekas & Overgaard, 2018).

The behavioral effect was accompanied by an early ERP modulation in the EPN time range from 180 to 240 ms after T2 onset: The mean amplitude in the posterior ROI had a higher predictive value for face detection in the negative compared to the neutral knowledge condition. This

difference was present for both liberal and strict hit criterion. This stronger ERP modulation compared to the behavioral data might be due to a higher sensitivity of the neurophysiological measure, making it possible to detect finer differences. Indeed, discrepancies between behavioral and neurophysiological measures (with differences only found for the later) have previously been reported in regard to an attentional bias for faces paired with negative social information (Xu et al., 2016) and ascribed to limitations of the behavioral measures which "reflect the combined effects of a sequence of many distinct neural processes" (Xu et al., 2016, p. 8) whereas ERPs enable tracking the continuous unfolding of processing over time.

Importantly, comparing the activity following T2-onset only for *detected* faces associated with neutral and negative knowledge did not yield significant differences—indicating that the observed effect cannot be explained merely by knowledge effects for faces that have already been detected. In combination with the observed enhanced detection of faces associated with negative knowledge, this suggests an intertwining of social-affective knowledge with the access to consciousness—rather than effects of knowledge depending on conscious perception.

The results point to slightly earlier differences between affective knowledge conditions (in interaction with detection) in the attentional blink task as compared to the evaluations. This pattern of results could be explained by the different task demands: During the evaluations, the images were shown comparatively long (1 s) and one by one. In the attentional blink task, the images were presented only very briefly (117 ms), preceded and followed by distractor stimuli. These increased task demands in the attentional blink task may have resulted in a faster integration of affective knowledge.

In order to examine the time course of the integration of affective knowledge in regard to the overall temporal course of processing during the attentional blink, general differences between detected and missed stimuli (as defined by the hit criteria) were also investigated in the ERPs. As reported in previous attentional blink studies (e.g., Sergent et al., 2005), an increased posterior negativity was observed for detected compared to missed stimuli in the N2 time range. However, the observed differences emerged already earlier and extended beyond this time span, in the form of one broad negative posterior cluster ranging from about 150 ms to 400 ms after T2 stimulus onset. In its topography and latency, this cluster corresponds to the VAN component (Förster et al., 2020; Koivisto & Revonsuo, 2010), which has been described as the most consistent ERP correlate of visual consciousness across different paradigms and is assumed to reflect phenomenal consciousness (i.e., the subjective experience of seeing). Our findings suggest that the processes underlying this ERP correlate are influenced by affective knowledge during the initial access to visual consciousness.

To conclude, the findings indicate that social-affective knowledge associated with a person affects perception-

related processing before or while attentional selection takes place. As a result, faces associated with negative compared to neutral information gain prioritized access to visual consciousness. These findings replicate and extend results from a previous behavioral study (Eiserbeck & Abdel Rahman, 2020). To the best of our knowledge, these studies are the first to report influences of social-affective knowledge on the visual consciousness of faces, whereas (apart from the discussed initial evidence from E. Anderson et al., 2011) this could not be shown in previous studies with other paradigms (binocular rivalry and breaking continuous flash suppression; Rabovsky et al., 2016; Stein et al., 2017). These differing results may be due to differences in the suppression techniques utilized in the paradigms (for a comparison of underlying mechanisms, see, e.g., Kanai et al., 2010): Binocular rivalry and continuous flash suppression rely on interocular suppression—a suppression of low-level sensory signals— and it is not yet clear whether or to what extent higher-level (e.g., emotional) processing of stimuli is possible under these conditions (Moors et al., 2017, 2019; Sklar et al., 2018). The attentional blink, on the other hand, is character- ized by attentional blindness (i.e., low-level signals cannot be accessed despite being available) and may enable pro- cessing up to a conceptual level (Martens & Wyble, 2010).

## Facial Trustworthiness and Visual Consciousness

In contrast to the affective knowledge effects we found little evidence that appearance-based facial trustworthiness has an impact on visual consciousness in the attentional blink. In the behavioral data, only a non-significant trend for a main effect of appearance was observed (whereas in one model additional including the ERP amplitude in the EPN time range, the $p$-value just exceeded the threshold for statistical significance, see Supplement Table S11). This trend was in the direction postulated in the hypotheses: Less trustworthy faces showed a tendency for enhanced detec- tion under conditions of reduced attention. Furthermore, while the interaction with knowledge did not reach signifi- cance, out of interest we nonetheless performed contrasts (for the strict hit criterion), which revealed an influence of appearance on detection in the neutral, but not in the negative knowledge condition. If replicated in future work, this pattern of results may indicate systematic differences based on the emotional value: Due to the counterbalanced design, all faces were shown equally often in both knowledge conditions and an effect of appearance in the neutral but not in the negative knowledge condition can therefore not be explained by low-level visual differences. Instead, the pattern could be explained by the diagnostic value of the information: Persons associated with negative knowledge can be assumed to pose a potential threat, and increased awareness can therefore be seen as appropriate in any case (regardless of facial trustworthiness). For persons associated with neutral knowledge, it remains unclear whether they might represent a danger—since neutral

knowledge does not exclude this possibility. It is therefore possible that trustworthiness assessments based on facial appearance may have a higher weight as an additional source of information in this case. However, no connection to electrophysiological activity could be found and the timing of the integration of this information therefore remains unclear. Furthermore, as noted above, even though we tentatively discuss possible interpretations, future work is needed to confirm the reliability of these results.

We conclude that, even though facial trustworthiness affects explicit and conscious evaluations of persons and expressions and may have been processed to a certain degree in the attentional blink task, its influence is very limited, in line with a previous report comparing knowledge- and appearance-based trustworthiness effects on visual consciousness (Eiserbeck & Abdel Rahman, 2020), and it is not clear whether the processing took place during the time range of attentional selection or afterwards.

## Conclusions

The present study demonstrates that social-affective knowledge has an influence on the access of faces to visual consciousness. Faces associated with negative information have a higher chance to be selected for enhanced conscious processing as compared to faces associated with relatively neutral information. ERPs revealed a connection between perception- and attention-related processing at a latency of about 180 ms. Specifically, starting at this time posterior ERP amplitudes had a higher predictive value for detection of faces associated with negative compared to neutral information. Our findings suggest that social affective knowledge about individuals can influence to what degree visual facial information becomes available for conscious processing, providing an important basis for social perception. Beyond descriptive trends in the behavioral data, no evidence for an effect of facial trustworthiness on the access to visual consciousness was observed.

---

### Declaration of Conflicting Interests
None.

# References

Abdel Rahman, R. (2011). Facing good and evil: Early brain signatures of affective biographical knowledge in face recognition. *Emotion*, *11*(6), 1397–1405. https://doi.org/10.1037/a0024717

Abir, Y., Sklar, A. Y., Dotsch, R., Todorov, A., & Hassin, R. R. (2017). The determinants of consciousness of human faces. *Nature Human Behaviour*, *2*(3), 194. https://doi.org/10.1038/s41562-017-0266-3

Alpers, G. W., Ruhleder, M., Walz, N., Mühlberger, A., & Pauli, P. (2005). Binocular rivalry between emotional and neutral stimuli: A validation using fear conditioning and EEG. *International Journal of Psychophysiology*, *57*(1), 25–32. https://doi.org/10.1016/j.ijpsycho.2005.01.008

Anderson, A. K., & Phelps, E. A. (2001). Lesions of the human amygdala impair enhancedperception of emotionally salient events. *Nature*, *411*(17), 305–309. https://doi.org/10.1038/35077083

Anderson, E., Siegel, E. H., Bliss-Moreau, E., & Barrett, L. F. (2011). The visual impact of gossip. *Science*, *332*(6036), 1446–1448. https://doi.org/10.1126/science.1201574

Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, *68*(3), 255–278. https://doi.org/10.1016/j.jml.2012.11.001

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

Baum, J., & Abdel Rahman, R. (2020). Emotional news affects social judgments independent of perceived media credibility. *Social Cognitive and Affective Neuroscience*, *December*, 1–12. https://doi.org/10.1093/scan/nsaa164

Baum, J., Rabovsky, M., Rose, S. B., & Abdel Rahman, R. (2020). Clear judgments based on unclear evidence: Person evaluation is strongly influenced by untrustworthy gossip. *Emotion*, *20*(2), 248–260. https://doi.org/10.1037/emo0000545

Bliss-Moreau, E., Barrett, L. F., & Wright, C. I. (2008). Individual differences in learning the affective value of others under minimal conditions. *Emotion*, *8*(4), 479–493. https://doi.org/10.1037/1528-3542.8.4.479

Del Cul, A., Baillet, S., & Dehaene, S. (2007). Brain dynamics underlying the nonlinear threshold for access to consciousness. *PLoS Biology*, *5*(10), 2408–2423. https://doi.org/10.1371/journal.pbio.0050260

Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics. *Journal of Neuroscience Methods*, *13*, 9–21. https://doi.org/http://doi.org/10.1016/j.jneumeth.2003.10.009

Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of trustworthiness perception. *Brain Research*, *1435*, 81–90. https://doi.org/10.1016/j.brainres.2011.11.043

Eiserbeck, A., & Abdel Rahman, R. (2020). Visual consciousness of faces in the attentional blink : Knowledge- based effects of trustworthiness dominate over appearance-based impressions. *Consciousness and Cognition*, *83*, 102977. https://doi.org/10.1016/j.concog.2020.102977

Eiserbeck, A., Enge, A., Rabovsky, M., & Abdel Rahman, R. (2021). Graded visual consciousness during the attentional blink. *BioRxiv*. https://doi.org/10.1101/2021.01.15.426792

Falvello, V., Vinson, M., Ferrari, C., & Todorov, A. (2015). The robustness of learning about the trustworthiness of other people. *Social Cognition*, *33*(5), 368–386. https://doi.org/10.1521/soco.2015.33.5.368

Fazekas, P., & Overgaard, M. (2018). A multi-factor account of degrees of awareness. *Cognitive Science*, *42*(6), 1833–1859. https://doi.org/10.1111/cogs.12478

Förster, J., Koivisto, M., & Revonsuo, A. (2020). ERP and MEG correlates of visual consciousness: The second decade. *Consciousness and Cognition*, *80*, 102917. https://doi.org/10.1016/j.concog.2020.102917

Frömer, R., Maier, M., & Abdel Rahman, R. (2018). Group-level EEG-processing pipeline for flexible single trial-based analyses including linear mixed models. *Frontiers in Neuroscience*, *12*, 1–15. https://doi.org/10.3389/fnins.2018.00048

Getov, S., Kanai, R., Bahrami, B., & Rees, G. (2015). Human brain structure predicts individual differences in preconscious evaluation of facial dominance and trustworthiness. *Social Cognitive and Affective Neuroscience*, *10*(5), 690–699. https://doi.org/10.1093/scan/nsu103

Green, P., & Macleod, C. J. (2016). SIMR: An R package for power analysis of generalized linear mixed models by simulation. *Methods in Ecology and Evolution*, *7*(4), 493–498. https://doi.org/10.1111/2041-210X.12504

Hassin, R., & Trope, Y. (2000). Facing faces: Studies on the cognitive aspects of physiognomy. *Journal of Personality and Social Psychology*, *78*(5), 837–852. https://doi.org/10.1037/0022-3514.78.5.837

Ille, N., Berg, P., & Scherg, M. (2002). Artifact correction of the ongoing EEG using spatial filters based on artifact and brain signal topographies. *Journal of Clinical Neurophysiology*, *19*(2), 113–124. https://doi.org/10.1097/00004691-200203000-00002

Kanai, R., Walsh, V., & Tseng, C. H. (2010). Subjective discriminability of invisibility: A framework for distinguishing perceptual and attentional failures of awareness. *Consciousness and Cognition*, *19*(4), 1045–1057. https://doi.org/10.1016/j.concog.2010.06.003

Kleisner, K., Priplatova, L., Frost, P., & Flegr, J. (2013). Trustworthy-looking face meets brown eyes. *PLoS ONE*, *8*(1), 1–7. https://doi.org/10.1371/journal.pone.0053285

Koivisto, M., & Revonsuo, A. (2008). Comparison of event-related potentials in attentional blink and repetition blindness. *Brain Research*, *1189*(1), 115–126. https://doi.org/10.1016/j.brainres.2007.10.082

Koivisto, M., & Revonsuo, A. (2010). Event-related brain potential correlates of visual awareness. *Neuroscience and Biobehavioral Reviews*, *34*(6), 922–934. https://doi.org/10.1016/j.neubiorev.2009.12.002

Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). lmerTest package: Tests in linear mixed effects models. *Journal of Statistical Software*, *82*(13). https://doi.org/10.18637/jss.v082.i13

Lischke, A., Junge, M., Hamm, A. O., & Weymar, M. (2018). Enhanced processing of untrustworthiness in natural faces with neutral expressions. *Emotion*, *18*(2), 181–189. https://doi.org/10.1037/emo0000318

Luo, Q. L., Wang, H. L., Dzhelyova, M., Huang, P., & Mo, L. (2016). Effect of affective personality information on face processing: Evidence from ERPs. *Frontiers in Psychology*, *7*, 810. https://doi.org/10.3389/fpsyg.2016.00810

Ma, D. S., Correll, J., & Wittenbrink, B. (2015). The Chicago face database: A free stimulus set of faces and norming data. *Behavior Research Methods*, *47*(4), 1122–1135. https://doi.org/10.3758/s13428-014-0532-5

Maier, M., & Abdel Rahman, R. (2018). Native language promotes access to visual consciousness. *Psychological Science*, *29*(11), 1757–1772. https://doi.org/10.1177/0956797618782181

Martens, S., & Wyble, B. (2010). The attentional blink: Past, present, and future of a blind spot in perceptual awareness. *Neuroscience and Biobehavioral Reviews*, *34*(6), 947–957. https://doi.org/10.1016/j.neubiorev.2009.12.005

Marzi, T., Righi, S., Ottonello, S., Cincotta, M., & Viggiano, M. P. (2014). Trust at first sight: Evidence from ERPs. *Social Cognitive and Affective Neuroscience*, *9*(1), 63–72. https://doi.org/10.1093/scan/nss102

Moors, P., Gayet, S., Hedger, N., Stein, T., Sterzer, P., van Ee, R., Wagemans, J., & Hesselmann, G. (2019). Three criteria for evaluating high-level processing in continuous flash suppression. *Trends in Cognitive Sciences*, *23*(4), 267–269. https://doi.org/10.1016/j.tics.2019.01.008

Moors, P., Hesselmann, G., Wagemans, J., & van Ee, R. (2017). Continuous flash suppression: Stimulus fractionation rather than integration. *Trends in Cognitive Sciences*, *21*(10), 719–721. https://doi.org/10.1016/j.tics.2017.06.005

Müsch, K., Engel, A. K., & Schneider, T. R. (2012). On the blink: The importance of target-distractor similarity in eliciting an attentional blink with faces. *PLoS ONE*, *7*(7). https://doi.org/10.1371/journal.pone.0041257

Nisbett, R. E., & Wilson, T. D. (1977). The halo effect: Evidence for unconscious alteration of judgments. *Journal of Personality and Social Psychology*, *35*(4), 250–256. https://doi.org/10.1037/0022-3514.35.4.250

Oosterhof, N. N., & Todorov, A. (2008). The functional basis of face evaluation. *Proceedings of the National Academy of Sciences*, *105*(32), 11087–11092. https://doi.org/10.1073/pnas.0805664105

Paunonen, S. V. (2006). You are honest, therefore I like you and find you attractive. *Journal of Research in Personality*, *40*(3), 237–249. https://doi.org/10.1016/j.jrp.2004.12.003

R Core Team. (2018). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. https://www.r-project.org/

Rabovsky, M., Stein, T., & Abdel Rahman, R. (2016). Access to awareness for faces during continuous flash suppression is not modulated by affective knowledge. *PLoS ONE*, *11*(4), e0150931. https://doi.org/10.1371/journal.pone.0150931

Ramsøy, T. Z., & Overgaard, M. (2004). Introspection and perception. *Phenomenology and the Cognitive Sciences*, *3*(1), 1–23. https://doi.org/10.1007/BF00776206

Raymond, J. E., Shapiro, K. L., & Arnell, K. M. (1992). Temporary suppression of visual processing in an RSVP task: An attentional blink? In *Journal of Experimental Psychology: Human Perception and Performance* (Vol. 18, Issue 3, pp. 849–860). https://doi.org/10.1037/0096-1523.18.3.849

Rezlescu, C., Duchaine, B., Olivola, C. Y., & Chater, N. (2012). Unfakeable facial configurations affect strategic choices in trust games with or without information about past behavior. *PLoS ONE*, *7*(3). https://doi.org/10.1371/journal.pone.0034293

Rudoy, J. D., & Paller, K. A. (2009). Who can you trust? Behavioral and neural differences between perceptual and memory-based influences. *Frontiers in Human Neuroscience*, *3*(16), 1–6. https://doi.org/10.3389/neuro.09.016.2009

Rule, N. O., Slepian, M. L., & Ambady, N. (2012). A memory advantage for untrustworthy faces. *Cognition*, *125*(2), 207–218. https://doi.org/10.1016/j.cognition.2012.06.017

Schacht, A., & Sommer, W. (2009). Emotions in word and face processing: Early and late cortical responses. *Brain and Cognition*, *69*(3), 538–550. https://doi.org/10.1016/j.bandc.2008.11.005

Schindler, S., & Bublatzky, F. (2020). Attention and emotion: An integrative review of emotional face processing as a function of attention. *Cortex*, *130*, 362–386. https://doi.org/10.1016/j.cortex.2020.06.010

Schupp, H. T., Junghöfer, M., Weike, A. I., & Hamm, A. O. (2003). Attention and emotion: An ERP analysis of facilitated emotional stimulus processing. *NeuroReport*, *14*(8), 1107–1110. https://doi.org/10.1097/00001756-200306110-00002

Schupp, H. T., Junghöfer, M., Weike, A. I., & Hamm, A. O. (2004). The selective processing of briefly presented affective pictures: An ERP analysis. *Psychophysiology*, *41*(3), 441–449. https://doi.org/10.1111/j.1469-8986.2004.00174.x

Schwabe, L., Merz, C. J., Walter, B., Vaitl, D., Wolf, O. T., & Stark, R. (2011). Emotional modulation of the attentional blink: The neural structures involved in capturing and holding attention. *Neuropsychologia*, *49*(3), 416–425. https://doi.org/10.1016/j.neuropsychologia.2010.12.037

Sergent, C., Baillet, S., & Dehaene, S. (2005). Timing of the brain events underlying access to consciousness during the attentional blink. *Nature Neuroscience*, *8*(10), 1391–1400. https://doi.org/10.1038/nn1549

Shore, D. M., Ng, R., Bellugi, U., & Mills, D. L. (2017). Abnormalities in early visual processes are linked to hypersociability and atypical evaluation of facial trustworthiness: An ERP study with Williams syndrome. *Cognitive, Affective and Behavioral Neuroscience*, *17*(5), 1002–1017. https://doi.org/10.3758/s13415-017-0528-6

Sklar, A. Y., Deouell, L. Y., & Hassin, R. R. (2018). Integration despite fractionation: Continuous flash suppression. *Trends in Cognitive Sciences*, *22*(11), 956–957. https://doi.org/10.1016/j.tics.2018.07.003

Stein, T., Awad, D., Gayet, S., & Peelen, M. V. (2018). Unconscious processing of facial dominance: The role of low-level factors in access to awareness. *Journal of Experimental Psychology: General*, *147*(11), e1–e13. https://doi.org/10.1037/xge0000521

Stein, T., Grubb, C., Bertrand, M., Suh, S. M., & Verosky, S. C. (2017). No impact of affective person knowledge on visual awareness: Evidence from binocular rivalry and continuous flash suppression. *Emotion*, *17*(8), 1199–1207. https://doi.org/10.1037/emo0000305

Stein, T., & Verosky, S. (2020). No effect of value learning on awareness and attention for faces : Evidence from continuous flash suppression and the attentional blink. *PsyArXiv*. https://doi.org/10.31234/osf.io/zc2wr

Stewart, L. H., Ajina, S., Getov, S., Bahrami, B., Todorov, A., & Rees, G. (2012). Unconscious evaluation of faces on social dimensions. *Journal of Experimental Psychology: General*, *141*(4), 715–727. https://doi.org/10.1037/a0027950

Suess, F., Rabovsky, M., & Abdel Rahman, R. (2015). Perceiving emotions in neutral faces: Expression processing is biased by affective person knowledge. *Social Cognitive and Affective Neuroscience*, *10*(4), 531–536. https://doi.org/10.1093/scan/nsu088

Todorov, A., Olivola, C. Y., Dotsch, R., & Mende-Siedlecki, P. (2015). Social attributions from faces: Determinants, consequences, accuracy, and functional significance. *Annual Review of Psychology*, *66*(1), 519–545. https://doi.org/10.1146/annurev-psych-113011-143831

Todorov, A., & Olson, I. R. (2008). Robust learning of affective trait associations with faces when the hippocampus is damaged, but not when the amygdala and temporal pole are damaged. *Social Cognitive and Affective Neuroscience*, *3*(3), 195–203. https://doi.org/10.1093/scan/nsn013

Verosky, S. C., Porter, J., Martinez, J. E., & Todorov, A. (2018). Robust effects of affective person learning on evaluation of faces. *Journal of Personality and Social Psychology*, *114*(4), 516–528. https://doi.org/10.1037/pspa0000109

Wieser, M. J., Gerdes, A. B. M., Büngel, I., Schwarz, K. A., Mühlberger, A., & Pauli, P. (2014). Not so harmless anymore: How context impacts the perception and electrocortical processing of neutral faces. *NeuroImage*, *92*, 74–82. https://doi.org/10.1016/j.neuroimage.2014.01.022

Willenbockel, V., Sadr, J., Fiset, D., Horne, G. O., Gosselin, F., & Tanaka, J. W. (2010). Controlling low-level image properties: The SHINE toolbox. *Behavior Research Methods*, *42*(3), 671–684. https://doi.org/10.3758/BRM.42.3.671

Xu, M., Li, Z., Diao, L., Fan, L., & Yang, D. (2016). Contextual valence and sociality jointly influence the early and later stages of neutral face processing. *Frontiers in Psychology*, *7*, 1258. https://doi.org/10.3389/fpsyg.2016.01258

Yang, D., Qi, S., Ding, C., & Song, Y. (2011). An ERP study on the time course of facial trustworthiness appraisal. *Neuroscience Letters*, *496*(3), 147–151. https://doi.org/10.1016/j.neulet.2011.03.066

Zivony, A., & Lamy, D. (2020). *What processes are disrupted during the attentional blink? An integrative review of event-related potentials research*. https://doi.org/10.31234/osf.io/epfbt

# Supplement

## 1. Additional Information

### Information S1: Cluster-Based Permutation Tests

Cluster-based permutation tests (CBPTs) as implemented in FieldTrip (Maris & Oostenveld, 2007) with 10,000 randomizations were used to test for differences between conditions in the attentional blink task within a restricted time range and topographical location typical for the EPN component. Only short lag trials in which T2 was present and T1 was correctly identified were examined. A dependent samples *t*-statistic was used to evaluate the effect at the sample level to determine cluster inclusion, using an alpha level of .05 for a single test. To be included in a cluster, a minimum number of two below-threshold neighborhood channels was required. Thereby, an electrode's spatial neighborhood was (manually) defined as directly adjacent electrodes in the cap. As the test statistic on the cluster level, the maximum of the cluster-level statistics was used (i.e., the largest sum of sample-specific *t*-statistics for each of the different clusters produces the test statistic) and a two-tailed test was applied. To specifically investigate EPN differences, the test was restricted to a time range and (broadly defined) posterior topographical region typical for the component (Abdel Rahman, 2011; Schacht & Sommer, 2009; Schupp et al., 2004; Suess et al., 2015), namely 200 to 350 ms after T2 stimulus onset including electrodes TP9, TP7, TP8, TP10, P7, P5, P3, Pz, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, PO10, O1, Oz, and O2. Since CBPTs are based on the comparison of two conditions, the continuous variable of appearance was converted into a factor with two levels by using a median split to separate between low and high trustworthy looking faces. To investigate interactions between affective knowledge and appearance, we used a double subtraction procedure, comparing the difference between untrustworthy and trustworthy appearance in the negative and neutral knowledge condition ([Negative: untrustworthy – trustworthy] – [Neutral: untrustworthy – trustworthy]).

No main effects of affective knowledge or appearance and no interaction of affective knowledge and appearance were observed (for the analyses of all short lag trials regardless of detection). However, when additionally taking into account the detection of the T2 stimuli as a factor (i.e., whether a face was detected or not according to the hit criteria), a significant interaction between affective knowledge and detection—in form of a double-difference between *negative: hit – miss* and *neutral: hit – miss*—was found for the liberal hit criterion (*p* = .004). The corresponding negative cluster was observed between 200 and 236 ms at electrodes P7, P5, P3, Pz, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, O1, Oz, and O2. Since the cluster appeared to be cut-off at 200 ms due to the restricted time range, we repeated the same analysis for a slightly larger time range from 150 to 350 ms in order to be able to obtain its full extent. Again, a significant difference was found (*p* =.006) and the corresponding cluster ranged from 184 to 236 ms after stimulus onset including the same electrodes (P7, P5, P3, Pz, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, O1, Oz, and O2). No interaction between appearance and detection was observed in the CBPTs.

## 2. Tables

Table S1

*Sentences Containing Social-Affective Information Used in the Study (English Translations of Used German Sentences)*

| Knowledge category | Sentence | Valence rating (1-7) | SD | Arousal rating (1-7) | SD |
|---|---|---|---|---|---|
| Negative | He made a lot of money with child pornography. | 1.10 | 0.31 | 5.50 | 1.73 |
| Negative | This woman set fire to a refugee shelter. | 1.20 | 0.41 | 5.75 | 1.45 |
| Negative | He fired into a crowd near a church. | 1.35 | 0.67 | 5.75 | 1.59 |
| Negative | This man has killed baby monkeys in a zoo. | 1.37 | 0.60 | 5.25 | 1.62 |
| Negative | This man beat his wife regularly. | 1.70 | 1.08 | 5.20 | 1.40 |
| Negative | This woman cheated her best friend out of a lot of money. | 1.70 | 1.08 | 4.60 | 1.82 |
| Negative | He threatened a shop assistant with a knife. | 1.75 | 0.79 | 4.70 | 1.53 |
| Negative | She sabotaged the brake lines on her neighbors' bikes. | 1.80 | 0.70 | 4.60 | 1.64 |
| Negative | She ran over a man and committed a hit-and-run. | 1.85 | 0.99 | 4.85 | 1.63 |
| Negative | She made fun of a child with a disability. | 1.85 | 0.81 | 4.55 | 1.73 |
| Negative | This man slapped his colleague for no reason. | 2.20 | 0.83 | 4.35 | 1.84 |
| Negative | This woman pretends to be a cleaning lady to steal from her customers. | 2.21 | 0.71 | 4.15 | 1.50 |
| Neutral | She took the change from the shop assistant. | 3.85 | 0.88 | 1.45 | 1.19 |
| Neutral | She described the haircut to the hairdresser. | 3.95 | 0.60 | 1.35 | 0.49 |
| Neutral | He showed a technician the connection. | 3.95 | 0.76 | 1.30 | 0.57 |
| Neutral | She handed in forms to the clerk. | 4.00 | 0.46 | 1.50 | 1.19 |
| Neutral | This man made an appointment with an eye doctor. | 4.00 | 0.79 | 1.20 | 0.41 |
| Neutral | This man asked a waiter for the menu. | 4.00 | 0.56 | 1.15 | 0.37 |
| Neutral | This man spoke to his employer. | 4.05 | 0.22 | 1.60 | 1.39 |
| Neutral | This woman took the elevator with a neighbor. | 4.05 | 0.76 | 1.40 | 0.75 |
| Neutral | He was watching an old western on night television. | 4.10 | 0.45 | 1.45 | 1.00 |
| Neutral | He received a registered letter from the mailman. | 4.10 | 0.72 | 1.20 | 0.62 |
| Neutral | This woman left her dress at the tailor's. | 4.15 | 0.49 | 1.40 | 0.99 |
| Neutral | This woman opened the door for a salesman. | 4.25 | 0.85 | 1.40 | 0.60 |
| Positive (Filler) | This man has taken stranded seals back into the sea. | 6.20 | 0.95 | 4.65 | 1.50 |
| Positive (Filler) | This woman is letting a refugee stay with her for free. | 6.30 | 1.22 | 4.35 | 1.66 |
| Positive (Filler) | This woman is protecting endangered species in the jungle from wild predators. | 6.30 | 0.73 | 4.10 | 1.80 |
| Positive (Filler) | This man carried an injured climber down into the valley. | 6.40 | 0.75 | 5.05 | 1.50 |
| Positive (Filler) | This woman donated a kidney to her sick sister. | 6.53 | 0.77 | 4.65 | 1.50 |
| Positive (Filler) | This man rescued an injured woman from her crashed car. | 6.60 | 0.68 | 4.70 | 1.59 |

Table S2

*LMM Statistics for Trustworthiness and Facial Expression Ratings*

| Variable | Trustworthiness Rating | | | | | Expression Rating | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *df* | *t* | *p* | *b* | SE | *df* | *t* | *p* |
| Intercept | 4.01 | 0.09 | 31.67 | 47.49 | **<.001** | 4.00 | 0.09 | 38.44 | 44.74 | **<.001** |
| Phase (Post - Pre) | -0.37 | 0.05 | 1378.90 | -6.73 | **<.001** | -0.09 | 0.05 | 1415.21 | -1.84 | .066 |
| **Pre**/Appearance | 1.00 | 0.12 | 41.08 | 8.47 | **<.001** | 0.92 | 0.14 | 42.88 | 6.43 | **<.001** |
| **Post**/Appearance | 0.70 | 0.12 | 41.08 | 5.93 | **<.001** | 0.80 | 0.14 | 42.88 | 5.64 | **<.001** |
| **Pre**/Knowledge (Neg – Neu) | -0.06 | 0.12 | 51.33 | -0.54 | .593 | -0.07 | 0.09 | 62.36 | -0.75 | .457 |
| **Post**/Knowledge (Neg – Neu) | -1.00 | 0.12 | 51.33 | -8.68 | **<.001** | -0.41 | 0.09 | 62.36 | -4.63 | **<.001** |
| **Pre**/(Appearance: Knowledge) | 0.07 | 0.14 | 56.64 | 0.51 | .610 | 0.02 | 0.12 | 173.09 | 0.13 | .895 |
| **Post**/(Appearance: Knowledge) | 0.13 | 0.14 | 56.64 | 0.96 | .340 | 0.15 | 0.12 | 173.09 | 1.25 | .215 |
| Random effects | Var. | SD | | | | Var. | SD | | | |
| Participants (Intercept) | 0.20 | 0.45 | | | | 0.08 | 0.28 | | | |
| Knowledge | 0.21 | 0.46 | | | | 0.10 | 0.32 | | | |
| Appearance | 0.32 | 0.57 | | | | 0.18 | 0.42 | | | |
| Knowledge: Appearance | 0.04 | 0.20 | | | | 0.04 | 0.20 | | | |
| Items (Intercept) | 0.003 | 0.05 | | | | 0.12 | 0.34 | | | |
| Knowledge | 0.02 | 0.12 | | | | 0.001 | 0.04 | | | |
| Residual | 1.14 | 1.07 | | | | 0.90 | 0.95 | | | |

*Note.* Neg = negative, Neu = neutral, Pre = before learning, Post = after learning; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors, "/" indicates nesting of fixed factors.

- 3 -

Table S3

*LMM Statistics for Prediction of Mean ROI Amplitude (250 – 300 ms After Face Onset, Comprising Electrodes P4, P6, P8, POz, PO4, PO8, PO10, Oz, and O2) in the Rating Phase After Learning*

| Variable | *b* | SE | *df* | *t* | *p* |
|---|---|---|---|---|---|
| Intercept | 3.14 | 0.52 | 31.31 | 6.07 | **<.001** |
| Knowledge (Neg-Neu) | -0.43 | 0.17 | 20.50 | -2.45 | **.024** |
| Appearance | -0.19 | 0.17 | 22.38 | -1.12 | .276 |
| Knowledge:Appearance | -0.21 | 0.28 | 21.58 | -0.76 | .456 |
| Random effects | Variance | SD | | | |
| Participants (Intercept) | 8.30 | 2.88 | | | |
| Appearance | 0.25 | 0.50 | | | |
| Items (Intercept) | 0.05 | 0.22 | | | |
| Knowledge | 0.14 | 0.37 | | | |
| Residual | 25.62 | 5.06 | | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

- 4 -

Table S4

*LMM Statistics for Prediction of Mean ROI Amplitude (200 – 250 ms After Face Onset, Comprising Electrodes PO9, P5, TP7, CP5, P7, P7, and TP9) in the Rating Phase After Learning*

| Variable | *b* | SE | *df* | *t* | *p* |
|---|---|---|---|---|---|
| Intercept | -1.35 | 0.37 | 31.41 | -3.61 | **.001** |
| Knowledge (Neg-Neu) | -0.07 | 0.14 | 172.70 | -0.51 | .608 |
| Appearance | 0.32 | 0.13 | 20.62 | 2.46 | **.023** |
| Knowledge:Appearance | -0.14 | 0.23 | 70.43 | -0.60 | .550 |
| Random effects | Variance | SD | | | |
| Participants (Intercept) | 4.27 | 2.07 | | | |
| Knowledge | 0.03 | 0.19 | | | |
| Appearance | 0.05 | 0.22 | | | |
| Knowledge:Appearance | 0.16 | 0.40 | | | |
| Items (Intercept) | 0.03 | 0.18 | | | |
| Knowledge | 0.004 | 0.06 | | | |
| Residual | 20.62 | 4.54 | | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

Table S5

*LMM Statistics for Prediction of Mean ROI Amplitude in the LPP Time Range (500 – 600 ms, Comprising Electrodes Pz, Cz, C1, C2, CP1, and CP2) in the Rating Phase After Learning*

| Variable | *b* | SE | *df* | *t* | *p* |
|---|---|---|---|---|---|
| Intercept | 3.67 | 0.40 | 31.69 | 9.09 | **<.001** |
| Knowledge (Neg-Neu) | .019 | 0.16 | 77.80 | 1.21 | .232 |
| Appearance | 0.10 | 0.15 | 21.67 | 0.62 | .544 |
| Knowledge:Appearance | -0.29 | 0.25 | 56.77 | -1.16 | .250 |
| Random effects | Variance | SD | | | |
| Participants (Intercept) | 4.98 | 2.23 | | | |
| Knowledge | 0.05 | 0.22 | | | |
| Appearance | 0.14 | 0.38 | | | |
| Knowledge:Appearance | 0.13 | 0.36 | | | |
| Items (Intercept) | 0.06 | 0.24 | | | |
| Knowledge | 0.02 | 0.14 | | | |
| Residual | 23.23 | 4.82 | | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

Table S6

*GLMM Statistics for Analysis of the Attentional Blink Effect*

| Variable | Liberal hit criterion | | | | Strict hit criterion | | | |
|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *z* | *p* | *b* | SE | *z* | *p* |
| Intercept | 1.46 | 0.18 | 7.94 | **<.001** | -0.75 | 0.39 | -1.93 | .054 |
| Lag (Long-Short) | 0.95 | 0.17 | 5.77 | **<.001** | 0.76 | 0.17 | 4.52 | **<.001** |
| Random effects | Var. | SD | | | Var. | SD | | |
| Participants (Intercept) | 0.95 | 0.98 | | | 4.45 | 2.11 | | |
| Lag | 0.78 | 0.88 | | | 0.73 | 0.86 | | |
| Items (Intercept) | 0.08 | 0.29 | | | 0.22 | 0.47 | | |
| Lag | 0.01 | 0.12 | | | 0.01 | 0.11 | | |

Table S7

*GLMM Statistics for Prediction of Hits by Mean N2 Amplitude (220 – 300 ms, Comprising Electrodes TP9, TP10, P7, P8, PO9, PO10, O1, and O2) in Short Lag Attentional Blink Trials*

| Variable | Liberal hit criterion | | | | Strict hit criterion | | | |
|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *z* | *p* | *b* | SE | *z* | *p* |
| Intercept | 0.93 | 0.20 | 4.664 | **<.001** | -1.18 | 0.38 | -3.14 | **.002** |
| N2 | -0.09 | 0.01 | -13.66 | **<.001** | -0.09 | 0.01 | -12.81 | **<.001** |
| Random effects | Var. | SD | | | Var. | SD | | |
| Participants (Intercept) | 1.12 | 1.06 | | | 4.08 | 2.02 | | |
| Items (Intercept) | 0.09 | 0.31 | | | 0.25 | 0.51 | | |

Table S8

*GLMM Statistics for Prediction of Hits by Mean VAN Amplitude (150 – 400 ms After T2 Onset, Comprising Electrodes TP9, TP10, P7, P8, PO9, PO10, O1, and O2) in Short Lag Attentional Blink Trials*

| Variable | Liberal hit criterion | | | | Strict hit criterion | | | |
|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *z* | *p* | *b* | SE | *z* | *p* |
| Intercept | 0.92 | 0.21 | 4.47 | **<.001** | -1.17 | 0.38 | -3.10 | **.002** |
| VAN | -0.10 | 0.01 | -10.46 | **<.001** | -0.09 | 0.01 | -9.14 | **<.001** |
| Random effects | Var. | SD | | | Var. | SD | | |
| Participants (Intercept) | 1.19 | 1.09 | | | 4.12 | 2.03 | | |
| Items (Intercept) | 0.09 | 0.31 | | | 0.25 | 0.50 | | |

Table S9

*LMM Statistics for Prediction of Mean ROI Amplitude (200 – 250 ms After T2 Onset, Comprising Electrodes PO9, P5, TP7, CP5, P7, P7, and TP9) in Short Lag T2-Present Attentional Blink Trials*

| Variable | *b* | SE | *df* | *t* | *p* |
|---|---|---|---|---|---|
| Intercept | -1.42 | 0.28 | 31.87 | -5.04 | **< .001** |
| Knowledge (Neg-Neu) | -0.06 | 0.12 | 20.58 | -0.51 | .613 |
| Appearance | 0.20 | 0.11 | 23.14 | 1.78 | .089 |
| Knowledge:Appearance | 0.19 | 0.19 | 21.35 | 0.99 | .334 |
| Random effects | Variance | SD | | | |
| Participants (Intercept) | 2.42 | 1.56 | | | |
| Appearance | 0.12 | 0.34 | | | |
| Items (Intercept) | 0.03 | 0.17 | | | |
| Knowledge | 0.09 | 0.30 | | | |
| Residual | 20.45 | 4.52 | | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

Table S10

*LMM Statistics for Prediction of Mean ROI Amplitude (250 – 300 ms After T2 Onset, Comprising Electrodes P4, P6, P8, POz, PO4, PO8, PO10, Oz, and O2) in Short Lag T2-Present Attentional Blink Trials*

| Variable | $b$ | SE | $df$ | $t$ | $p$ |
|---|---|---|---|---|---|
| Intercept | -4.03 | 0.32 | 31.97 | -12.74 | **< .001** |
| Knowledge (Neg-Neu) | 0.14 | 0.14 | 38.21 | 1.06 | .298 |
| Appearance | -0.14 | 0.11 | 22.28 | -1.22 | .234 |
| Knowledge:Appearance | -0.09 | 0.18 | 410.87 | -0.52 | .601 |
| Random effects | Variance | SD | | | |
| Participants (Intercept) | 3.04 | 1.74 | | | |
|     Knowledge | 0.20 | 0.45 | | | |
|     Appearance | 0.04 | 0.19 | | | |
|     Knowledge:Appearance | 0.02 | 0.16 | | | |
| Items (Intercept) | 0.04 | 0.21 | | | |
|     Knowledge | 0.0005 | 0.02 | | | |
|     Residual | 24.71 | 4.97 | | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

Table S11

*GLMM Statistics for Analysis of Short Lag Attentional Blink Trials Including Mean Amplitude During the EPN Time Range as a Predictor*

| Variable | Liberal hit criterion | | | | Strict hit criterion | | | |
|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *z* | *p* | *b* | SE | *z* | *p* |
| Intercept | 0.99 | 0.21 | 4.84 | **<.001** | -1.12 | 0.38 | -2.97 | **.003** |
| Knowledge (Neg-Neu) | -0.0002 | 0.07 | -0.004 | .997 | 0.17 | 0.10 | 1.81 | .071 |
| Appearance | -0.21 | 0.12 | -1.83 | .067 | -0.34 | 0.17 | -2.04 | **.042** |
| EPN | -0.05 | 0.01 | -7.71 | **<.001** | -0.05 | 0.01 | -7.29 | **<.001** |
| Knowledge:Appearance | 0.09 | 0.11 | 0.84 | .402 | 0.23 | 0.15 | 1.55 | .122 |
| Knowledge:EPN | -0.05 | 0.01 | -3.65 | **<.001** | -0.03 | 0.01 | -2.01 | **.044** |
| Appearance:EPN | 0.0004 | 0.01 | 0.04 | .969 | 0.01 | 0.01 | 0.94 | .347 |
| Knowledge: Appearance:EPN | -0.01 | 0.02 | -0.38 | .701 | -0.01 | 0.02 | -0.51 | .614 |
| Random effects | Var. | SD | | | Var. | SD | | |
| Participants (Intercept) | 1.21 | 1.10 | | | 4.16 | 2.04 | | |
| Knowledge | 0.02 | 0.14 | | | 0.05 | 0.22 | | |
| Appearance | 0.10 | 0.31 | | | 0.10 | 0.32 | | |
| Knowledge: Appearance | 0.01 | 0.11 | | | 0.09 | 0.30 | | |
| Items (Intercept) | 0.08 | 0.28 | | | 0.20 | 0.45 | | |
| Knowledge | 0.02 | 0.15 | | | 0.05 | 0.22 | | |

*Note.* Neg = negative, Neu = neutral; EPN = mean amplitude during time range of 180 to 240 ms at electrodes : P7, P5, P3, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, PO10, O1, Oz, and O2, based on CBPT results; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors, "/" indicates nesting of fixed factors.

Additional notes on Table S11:

When including the additional predictor *EPN* in the GLMM model, the previously observed main effect of knowledge (for the *strict* hit criterion) does not reach statistical significance anymore. This might be due to the fact that the continuous predictor *EPN* represents crucial differences associated with knowledge more precisely than the factor *knowledge* with two levels. The main effect of *appearance*, which was observable only as a trend in the previous model, now passed the threshold of statistical significance (for the *strict* criterion). This might be due to the better control of variance through the inclusion of mean EPN amplitudes. However, no connection of the effect with neural activity was observed.

Table S12

*GLMM Statistics for Analysis of Short Lag Attentional Blink Trials Including Mean Amplitude During the EPN Time Range as a Predictor (Appearance and Mean Amplitude Nested Within Knowledge Condition).*

| Variable | Liberal hit criterion | | | | Strict hit criterion | | | |
|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *z* | *p* | *b* | SE | *z* | *p* |
| Intercept | 0.99 | 0.21 | 4.84 | **<.001** | -1.12 | 0.38 | -2.97 | **.003** |
| Knowledge (Neg-Neu) | -0.0002 | 0.07 | 0.00 | .997 | 0.17 | 0.10 | 1.81 | .071 |
| **Neu/EPN** | -0.03 | 0.01 | -3.03 | **.002** | -0.04 | 0.01 | -3.83 | **<.001** |
| **Neg/EPN** | -0.08 | 0.01 | -8.09 | **<.001** | -0.07 | 0.01 | -6.75 | **<.001** |
| **Neu/Appearance** | -0.26 | 0.13 | -1.96 | .050 | -0.46 | 0.17 | -2.74 | **.006** |
| **Neg/Appearance** | -0.17 | 0.12 | -1.35 | .178 | -0.22 | 0.20 | -1.13 | .259 |
| **Neu/(EPN:Appearance)** | 0.004 | 0.01 | 0.30 | .761 | 0.02 | 0.02 | 0.99 | .324 |
| **Neg/(EPN:Appearance)** | -0.004 | 0.02 | -0.24 | .813 | 0.005 | 0.01 | 0.33 | .738 |
| Random effects | Var. | SD | | | Var. | SD | | |
| Participants (Intercept) | 1.21 | 1.10 | | | 4.16 | 2.04 | | |
| Knowledge | 0.02 | 0.14 | | | 0.05 | 0.22 | | |
| Appearance | 0.10 | 0.31 | | | 0.10 | 0.32 | | |
| Knowledge:Appearance | 0.01 | 0.11 | | | 0.09 | 0.30 | | |
| Items (Intercept) | 0.08 | 0.28 | | | 0.20 | 0.45 | | |
| Knowledge | 0.02 | 0.15 | | | 0.05 | 0.22 | | |

*Note.* Neg = negative, Neu = neutral; EPN = mean amplitude during time range of 180 to 240 ms at electrodes : P7, P5, P3, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, PO10, O1, Oz, and O2, based on CBPT results; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors, "/" indicates nesting of fixed factors.

Table S13

*LMM Statistics for Prediction of Mean ROI Amplitude (180 – 240 ms After T2 Onset, Comprising Electrodes P7, P5, P3, P4, P6, P8, PO9, PO7, PO3, POz, PO4, PO8, PO10, O1, Oz, and O2) in Short Lag T2-Hit Attentional Blink Trials*

| Variable | Liberal hit criterion | | | | | Strict hit criterion | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | *b* | SE | *df* | *t* | *p* | *b* | SE | *df* | *t* | *p* |
| Intercept | -0.26 | 0.30 | 32.29 | -0.62 | .395 | -0.56 | 0.33 | 29.00 | -1.69 | .102 |
| Knowledge (Neg-Neu) | -0.19 | 0.14 | 31.36 | -1.29 | .207 | -0.19 | 0.17 | 23.25 | -1.13 | .271 |
| Appearance | -0.12 | 0.12 | 22.70 | -1.07 | .296 | -0.10 | 0.13 | 24.40 | -0.81 | .428 |
| Knowledge:Appearance | -0.08 | 0.18 | 127.44 | -0.47 | .639 | -0.20 | 0.27 | 26.52 | -0.76 | .453 |
| Random effects | Var. | SD | | | | Var. | SD | | | |
| Participants (Intercept) | 2.70 | 1.64 | | | | 2.89 | 1.70 | | | |
| Knowledge | 0.25 | 0.50 | | | | 0.04 | 0.19 | | | |
| Appearance | 0.02 | 0.15 | | | | 0.01 | 0.11 | | | |
| Knowledge:Appearance | 0.02 | 0.16 | | | | 0.21 | 0.46 | | | |
| Items (Intercept) | 0.04 | 0.19 | | | | 0.02 | 0.13 | | | |
| Knowledge | 0.01 | 0.10 | | | | 0.08 | 0.27 | | | |
| Residual | 16.62 | 4.08 | | | | 16.07 | 4.01 | | | |

*Note.* Neg = negative, Neu = neutral; higher *Appearance* value corresponds to more trustworthy looking appearance; ":" indicates interactions between fixed factors.

## 3. Figures

**Figure S1**

*Topographies of Grand Average ERPs Time-locked to T2 Face Onset Showing the Differences Between Hit and Miss Short Lag Trials in the Attentional Blink Task (Liberal Hit Criterion)*
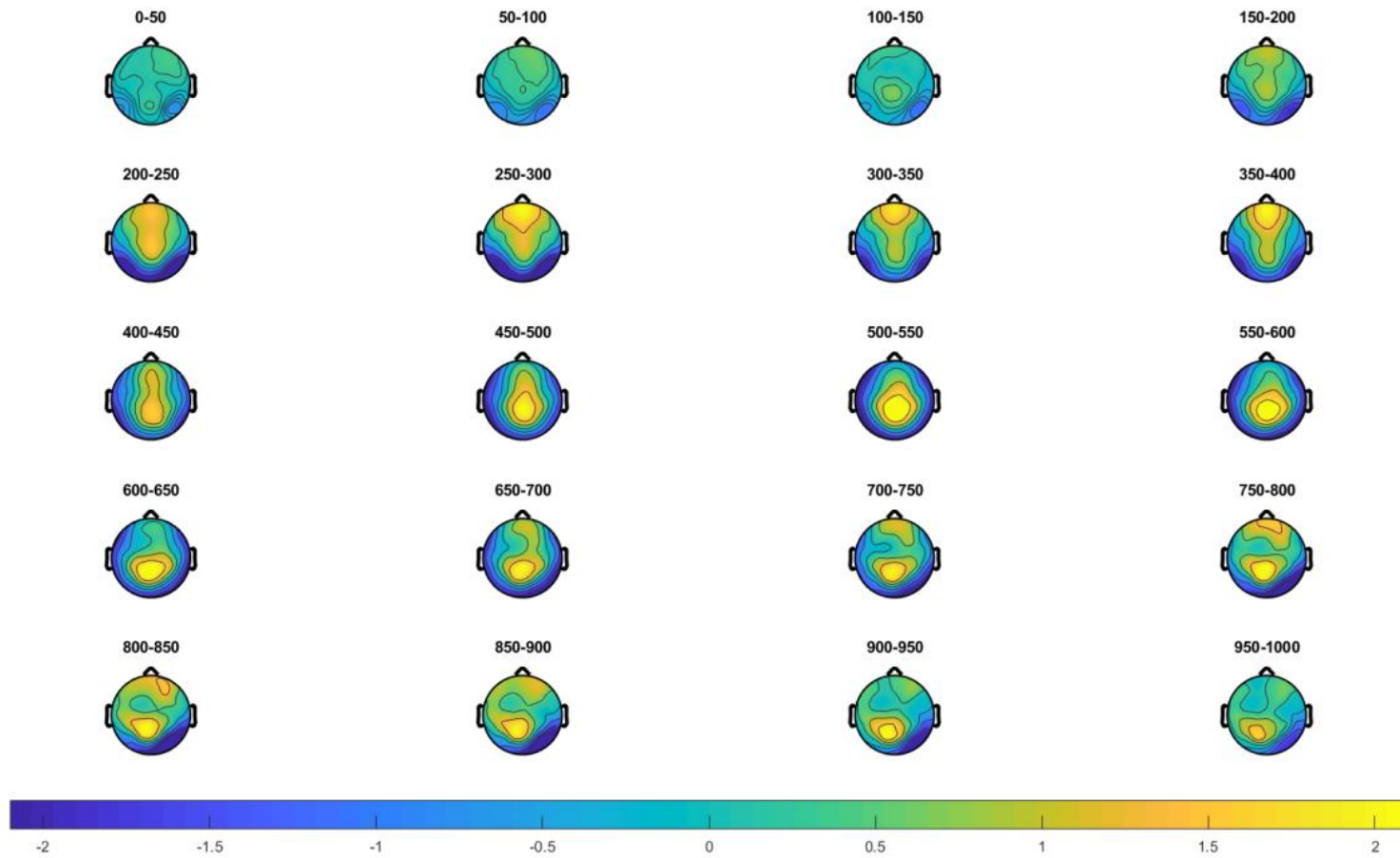
*Topographies of Grand Average ERPs Time-Locked to T2 Face Onset Showing the Differences Between Hit and Miss Short Lag Trials in the Attentional Blink Task (Strict Hit Criterion)*