



UNIVERSITY OF LEEDS

This is a repository copy of *Many moral buttons or just one? Evidence from emotional facial expressions*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/135327/>

Version: Accepted Version

Article:

Franchin, L, Geipel, J, Hadjichristidis, C orcid.org/0000-0002-9441-6650 et al. (1 more author) (2019) Many moral buttons or just one? Evidence from emotional facial expressions. *Cognition and Emotion*, 33 (5). pp. 943-958. ISSN 0269-9931

<https://doi.org/10.1080/02699931.2018.1520078>

© 2018 Informa UK Limited, trading as Taylor & Francis Group. This is an author produced version of a paper published in *Cognition and Emotion* . Uploaded in accordance with the publisher's self-archiving policy.

Reuse

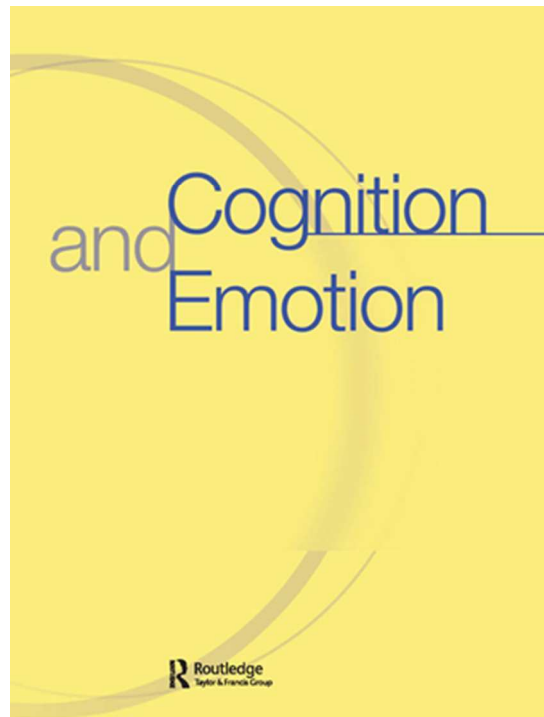
Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>



Many moral buttons or just one? Evidence from emotional facial expressions

Journal:	<i>Cognition and Emotion</i>
Manuscript ID	CEM-FA 444.17.R3
Manuscript Type:	Full Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Franchin, Laura; University of Trento Geipel, Janet; University of Chicago Hadjichristidis, Constantinos; University of Trento; University of Leeds, UK Surian, Luca; University of Trento
Keywords:	Facial expression, Moral judgment, Moral Foundation Theory, Harm, Purity

SCHOLARONE™
Manuscripts

1
2
3
4
5
6
7
8 **Many moral buttons or just one? Evidence from emotional facial expressions**
9

10
11
12
13
14
15 Laura Franchin¹, Janet Geipel², Constantinos Hadjichristidis^{3,4}, and Luca Surian¹
16
17

18
19
20 ¹Department of Psychology and Cognitive Sciences, University of Trento, Italy
21

22 ²Department of Psychology, University of Chicago, USA
23

24 ³Department of Economics and Management, University of Trento, Italy
25

26 ⁴Centre for Decision Research, University of Leeds, UK
27
28
29
30
31
32
33
34
35
36

37 *Acknowledgments.* We thank students who participated in our study and Ilaria Perrucci for her
38 help in data collection and in coding facial expressions. The authors declare no conflict of interest.
39
40

41
42 *Corresponding Author:* Laura Franchin, Department of Psychology and Cognitive Sciences,
43 University of Trento, Corso Bettini, 31, 38068 Rovereto (Trento), Italy; e-mail:
44 laura.franchin@unitn.it, telephone: +390464808633.
45
46
47
48
49
50
51

52 Word count = 9087
53
54
55
56
57
58
59
60

Abstract

We investigated whether moral violations involving harm selectively elicit anger, whereas purity violations selectively elicit disgust, as predicted by the Moral Foundations Theory (MFT). We analyzed participants' spontaneous facial expressions as they listened to scenarios depicting moral violations of harm and purity. As predicted by MFT, anger reactions were elicited more frequently by harmful than by impure actions. However, violations of purity elicited more smiling reactions and expressions of anger than of disgust. This effect was found both in a classic set of scenarios and in a new set in which the different kinds of violations were matched on weirdness. Overall, these findings are at odds with predictions derived from MFT and provide support for 'monist' accounts that posit harm at the basis of all moral violations. However, we found that smiles were differentially linked to purity violations, which leaves open the possibility of distinct moral modules.

Keywords: Facial expression; Moral judgment; Moral Foundations Theory; Harm; Purity

1
2
3 Many moral buttons or just one? Evidence from emotional facial expressions
4

5 One of the central goals in moral psychology is to uncover the mechanisms that underpin
6 moral judgment. To understand these mechanisms, we need to determine the extent to which their
7 input, output and characteristic computations are domain-specific, that is, the extent to which they
8 compute only particular types of information and produce certain kinds of outputs. The theoretical
9 positions given to this question fall into two camps, sometimes called 'pluralists' and 'monists'.
10 Pluralists claim that moral cognition is composed of many domain-specific mechanisms that are
11 selectively triggered by different kinds of violations. According to a famous proposal, the Moral
12 Foundations Theory (MFT), there are at least five specialized mechanisms that deal with
13 violations in the domains of harm, fairness, loyalty, authority and purity (e.g., Graham et al., 2013).
14 They are specialized in their input requirements—they are triggered by different input
15 conditions—but also in their computations and outputs, for example, they elicit distinct emotions.
16 Alternative 'monist' accounts propose that morality is about harm and that all moral disapprovals
17 involve, to a different degree, a concern for harm. Monist accounts may endorse a constructionist
18 perspective on how information is integrated during processing (for a review, see Cameron,
19 Lindquist, & Gray, 2015) or a modular breakdown of processing stages (e.g., Cushman, 2013;
20 Mikhail, 2007). In both cases, harm violations may occur in many different contexts and numerous
21 factors, such as directness and intentionality, can constraint moral judgments.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

41 In the present work, we focus on what is considered to be a crucial test for tackling the
42 issue of domain-specificity: Do different types of moral violations elicit different types of
43 emotions? A positive answer to this question has been one of the initial empirical claims for
44 multiple modules theories. Support for it came from studies using a variety of tasks and paradigms.
45 In their seminal paper, Haidt, Koller, and Dias (1993) presented children and adults with stories
46 purported to include instances of harmless actions that were either disrespecting, and therefore
47 violated the 'Community ethics', or were “unconventional food and sexual practices, designed to
48 trigger disgust” (p. 618), like eating a dog, kissing a sibling in the mouth, or masturbating using a
49
50
51
52
53
54
55
56
57
58
59
60

MANY MORAL BUTTONS OR JUST ONE?

4

1
2
3 dead chicken. They found that people in the United States and Brazil judged these acts as immoral,
4
5 often justifying their judgments by stating that the acts were disgusting. Note that such actions
6
7 were performed privately, to avoid the criticism that public disrespect and taboo violations may
8
9 cause harm on witnessing people, because, if they do, then they fall into the domain of harm-based
10
11 morality (Turiel, 1989; Turiel, Killen, & Helwig, 1987). In another influential study (Rozin,
12
13 Lowery, Imada, & Haidt, 1999), people were asked to pair different types of violations with
14
15 different facial emotional expressions or emotion words. Rozin and colleagues (1999)
16
17 demonstrated that violations of loyalty and group solidarity norms (the 'Community code') were
18
19 primarily associated with the emotion of contempt, violations involving harmful actions (the
20
21 'Autonomy code') were associated with anger, and violations concerning purity (the 'Divinity
22
23 code') were associated with disgust (hence the name of the model, CAD – Community-contempt,
24
25 Autonomy-anger, Divinity-disgust).
26
27

28
29 People's judgments of wrongness and disgust are affected by both content and
30
31 intentionality and such factors interact in a way consistent with MFT. Young and Saxe (2011)
32
33 found that intentional purity violations, such as incest, are judged as more wrong, but not more
34
35 disgusting, than accidental purity violations; and that purity violations are judged as more
36
37 disgusting, but not more wrong, than harm violations. Converging evidence comes also from
38
39 studies on spontaneous facial expressions, which reveal links between harm violations and
40
41 expressions of anger, and between purity violations and expressions of disgust (Cannon, Schnell,
42
43 & White, 2011). Some functional brain imaging studies also found that moral violations, such as
44
45 incest and pathogen-related actions, activate brain regions that have been associated with disgust
46
47 (e.g., Parkinson et al., 2011; Schaich Borg, Kahn, Sinnott-Armstrong, Kurzban, Robinson, &
48
49 Kiehl, 2013). Rottman, Kelemen, and Young (2014) found that the condemnation of suicide was
50
51 significantly associated with trait disgust and disgust reported in reading obituaries, regardless of
52
53 political orientation. Perceived harm did not predict judgments of suicide, suggesting that suicide
54
55
56
57
58
59
60

MANY MORAL BUTTONS OR JUST ONE?

5

1
2
3 is considered immoral not because of its harming consequences on close relatives, but primarily
4
5 on grounds of norms related to purity.

6
7 Despite the empirical support for MFT, the existence of specialized moral modules and of
8
9 non-harmful moral violations remains controversial. First, numerous results run against the
10
11 predictions of MFT. For example, Gutierrez and Giner-Sorolla (2007) reported that harm is
12
13 frequently associated with certain purity violations such as consensual incest. However, they
14
15 proposed that such concerns may reflect post hoc rationalizations and therefore may be consistent
16
17 with the social-intuitionist model (Haidt, 2001).

18
19
20 Gray, Schein, and Ward (2014) challenged this conclusion by reporting evidence
21
22 suggesting that harm is automatically activated by purity violations. More specifically, they
23
24 presented participants with instances of purity violations (e.g., covering a bible with feces) and
25
26 non-moral negative actions (e.g., a student failing an exam). Following each action, they presented
27
28 participants with a Chinese character and asked them to indicate whether its meaning was harmful,
29
30 sad, or wrong. Relative to neutral actions, both negative non-moral actions and moral violations
31
32 increased the sadness scores, but only moral violations increased also harmfulness and wrongness
33
34 scores.

35
36
37 Another challenging evidence for MFT comes from Royzman, Atanasov, Landy, Parks,
38
39 and Gepty (2016). These authors found that pathogen-free violations (related to the 'Divinity
40
41 code') were more strongly linked to anger than to disgust. Indeed, according to a recent review on
42
43 emotion-moral content studies (Cameron et al., 2015) the available evidence provides little
44
45 support for purity-disgust and harm-anger correspondences. Cameron and colleagues demonstrate
46
47 that most of the published studies do not report exclusive emotion-moral content links and the 25
48
49 studies that report such links fail short of ruling out that they are not due to global factors such as
50
51 core affect, typicality or severity.

52
53
54 Can such domain-general factors account for effects previously attributed to moral
55
56 content? Gray and Keeney (2015) claimed that scenarios commonly used to instantiate purity
57
58

MANY MORAL BUTTONS OR JUST ONE?

6

1
2
3 violations portray actions that are highly atypical (e.g., having a tail attached by plastic surgery)
4
5 and are, on average, less severe violations than those used to instantiate harm violations (e.g.,
6
7 armed assault). While this may be a common aspect of purity violations, especially in cultures
8
9 dominated by harm-based ethics (Graham, 2015), this raises the possibility that domain-general
10
11 aspects, rather than moral content per se, are responsible for effects previously thought to support
12
13 MFT. Gray and Keeney (2015) found that when controlling for such variables, the differences in
14
15 responses to purity and harm violations vanished. In their study, the focus was on an effect
16
17 originally found by Uhlmann and Zhu (2014): Harm violations are rated as more immoral than
18
19 purity violations, but purity violations are judged as more indicative of an agent's immoral
20
21 character than harm violations (see also Giner-Sorolla & Chapman, 2017). Gray and Keeney
22
23 orthogonally manipulated moral content, severity, and typicality, and found that the effect due to
24
25 moral content disappeared in a new set of violations – which they called ‘naturalistic’ – that had
26
27 been generated by lay people rather than by the experimenters.
28
29
30

31 The aim of the present study was to investigate the association between type of moral
32
33 violation and type of emotional reaction by focusing on the facial expressions elicited by harm and
34
35 purity violations. Two types of scenarios were employed: classic MFT scenarios that have been
36
37 extensively used in previous studies, and the new set of more 'naturalistic' scenarios which has
38
39 been developed by Gray and Keeney (2015). In the naturalistic scenarios, the harm and purity
40
41 violations were matched for weirdness. We employed both MFT and naturalistic scenarios to
42
43 examine whether they would yield similar patterns of results. According to MFT, purity and harm
44
45 scenarios should evoke differential affective responses because specific moral concerns are
46
47 associated with specific emotional states. More specifically, anger should be the expression most
48
49 frequently triggered by harm scenarios whereas disgust the expression most frequently triggered
50
51 by scenarios depicting violations of purity (e.g., Cannon et al., 2011; Rozin et al., 1999). We name
52
53 this the strong MFT prediction. We also considered a weaker MFT prediction according to which,
54
55
56
57
58
59
60

MANY MORAL BUTTONS OR JUST ONE?

7

1
2
3 relative to harm scenarios, purity scenarios trigger relatively fewer expressions of anger and
4
5 relatively more expressions of disgust.

6
7 Apart from anger and disgust, we also investigated the presence of contempt, smiling and
8
9 surprise expressions. Contempt is considered to be the most subtle and coldest of the triad of
10
11 hostility emotions (Izard, 1977; Rozin et al., 1999). Ekman (1994) considered contempt as
12
13 ‘disapproving of’ and ‘feeling morally superior to’ someone. It is usually claimed to involve a
14
15 negative evaluation of others and their actions and is associated with the violations of loyalty and
16
17 group solidarity norms, not with violations of individual rights or purity (Rozin et al., 1999). MFT
18
19 makes no clear prediction about contempt reactions in the two types of violations, apart that
20
21 contempt should not be the predominant reaction in either case. The main reason we decided to
22
23 examine contempt is that it is a crucial emotional response in the CAD model, associated
24
25 selectively with violations of the 'community morality' norms (Rozin et al., 1999).

26
27
28 Turning to smiling expressions, Cannon et al. (2011) found that some highly negative
29
30 reactions to purity violations were associated with an increased activity of the zygomaticus muscle
31
32 (which is associated with smiling). The authors explained this effect as the result of cross-talk
33
34 activity with the levator muscle that is activated in extreme disgust facial expressions (e.g., Vrana,
35
36 1993). However, they also discussed an alternative possibility: “participants may have found some
37
38 of the more extreme purity behaviors amusing as well as disgusting” (Cannon et al., 2011, p. 330).
39
40 This is plausible. According to the benign-violation hypothesis (McGraw & Warren, 2010), to
41
42 elicit humor a situation must be appraised as a violation and, simultaneously, must be perceived as
43
44 benign. According to McGraw and Warren (2010), situations that elicit humor range from
45
46 apparent physical threats to violations of personal dignity (e.g., physical deformities), linguistic
47
48 norms, social norms (e.g., strange behavior), and—related to the present purposes—moral norms
49
50 (e.g., bestiality or disrespectful behaviors). Because purity violations have been found to trigger
51
52 smiles, we also examined the presence of smiling expressions.
53
54
55
56
57
58
59
60

1
2
3 In addition, given that the traditional purity scenarios are regarded as highly unusual or
4 atypical (Gray & Keeney, 2015), it seems plausible that they might evoke surprise. Thus, we also
5 investigated the presence of surprise expressions. As it was the case with contempt, MFT does not
6 make a prediction about smiling or surprise expressions other than neither should be the most
7 prominent expression in response to either type of violation.
8
9

10
11
12
13 Finally, the reason we decided to study facial expressions, instead of relying on emotional
14 scale ratings, is twofold. First, participants may use and understand emotion words in different
15 ways than the one intended by the experimenter. For example, people may use the word 'disgust'
16 to refer to a state of irritation or anger (Nabi, 2002), or to achieve a rhetorical effect (Bloom, 2004).
17 Therefore, facial expressions provide a less ambiguous indicator of emotion than ratings on
18 emotion scales. Second, evidence from facial expressions is less open to interpretations such that a
19 particular response may be a product of a post hoc rationalization (see e.g., Gutierrez & Giner-
20 Sorolla, 2007). This is especially true in the present research in which we recorded the first
21 spontaneous facial expressions in response to vignettes depicting moral violations—note that
22 participants were simply asked to listen to the scenarios, they were not asked to morally evaluate
23 the depicted actions.
24
25
26
27
28
29
30
31
32
33
34
35
36

37 Method

38 Participants

39
40
41 Thirty-three Italian native speakers (17 females, 16 males, $M_{age} = 24.79$, age range =
42 18–38) participated in the study. Participants were recruited via an announcement through the
43 University of Trento mailing list. Two additional participants were excluded from the analysis
44 due to the low quality of the video recordings (i.e., participants were not centered on the camera
45 focus and therefore their facial expressions could not be recorded clearly). The study protocol
46 was conducted in accordance to the principles expressed in the Declaration of Helsinki.
47
48
49
50
51
52
53

54 We conducted a power analysis for the within-groups comparison (i.e., a t-test for one
55 sample case) using G-Power (Faul, Erdfelder, Lang, & Buchner, 2007). To detect a medium
56
57
58
59

MANY MORAL BUTTONS OR JUST ONE?

1
2
3 effect size (i.e., Cohen's $d = 0.61$, based on Cannon et al., 2011) with alpha set at .05 and a
4
5 power of .80, a minimum sample size of 24 participants was required. With our final sample
6
7 size of 33 participants, we therefore expect a power of .92.

Materials and Procedure

11 The experiment was conducted in a quiet room. Participants were tested individually
12
13 after being informed about the general aims of the study and having signed a consent form.
14
15 Each participant was asked to complete two tasks, which were separated by a short break. In
16
17 the first task, participants were asked to listen to 20 statements describing negative moral
18
19 actions while their spontaneous facial expressions in response to these statements were video-
20
21 recorded. In the second task, participants were asked to listen once again to the 20 statements
22
23 and to answer to a series of questions that were presented in a booklet.

26 *Task 1.* Participants were told that the aim of the study was to examine cognitive
27
28 processes while listening to statements about other individuals, but were not informed about
29
30 the specific hypothesis. Participants were informed that a video camera (Sony Handycam
31
32 HDR-SR5) recorded their upper part of the body. The camera was placed just below the
33
34 monitor, at a distance of 50 cm from the participants. Facial expressions were recorded during
35
36 the entire session. Participants simply had to look on the screen, read some instructions on it,
37
38 and listen to all the statements.

41 Specifically, participants listened to 20 statements: 5 depicting commonly used Moral
42
43 Foundations Theory (MFT) harm violations, 5 MFT purity violations (see Graham, Haidt, &
44
45 Nosek, 2009), 5 naturalistic harm violations, and 5 naturalistic purity violations (see Gray &
46
47 Keeney, 2015). The statements were presented in a pseudo-randomized order (see Table 1, for
48
49 the full list of the moral violations).

52 The presentation of the statements was computer-based and was created using the
53
54 open-source program OpenSesame (Mathôt, Schreijf, & Theeuwes, 2012). After a fixation cross
55
56 that appeared for 2 s, a blank screen appeared and a statement concerning a moral violation
57
58

MANY MORAL BUTTONS OR JUST ONE?

10

1
2
3 was presented orally (male adult voice). The presentation of moral violations lasted between 4
4
5 and 15 s. At the end of each moral violation, the written statement “Please think about this
6
7 action” appeared on the computer screen for 4 s. Each trial ended with a blank screen
8
9 (presented for 250 ms), which was followed by the instruction 'Please wait!' for 4 s between
10
11 trials (see Figure 1, for a schematic representation). During this task, participants were not
12
13 aware that they would later be asked to make judgments about these statements.
14

15 --INSERT FIGURE 1 ABOUT HERE--
16

17
18 *Task 2.* After a short break, the statements used in Task 1 were presented to the same
19
20 participants again but in a different pseudo-randomized order. After listening to each statement,
21
22 participants were asked to respond, in a booklet, to the following questions (Gray & Keeney,
23
24 2015): (1) “How morally wrong is this action?” (1 = *not at all wrong*, to 7 = *very wrong*), (2)
25
26 “How severe is this action?” (1 = *not at all severe*, to 7 = *very severe*), (3) “How atypical
27
28 (bizarre, weird, odd) is this action?” (1 = *not at all atypical*, 7 = *very atypical*), (4) “How
29
30 harmful (this implies physical and/or emotional) is this action? (1 = *not at all harmful*, to 7 =
31
32 *very harmful*), and (5) “How impure is this action?” (1 = *not at all impure*, to 7 = *very impure*).
33
34

35 At the end of the experimental session, participants were asked to guess what was the
36
37 purpose of the experiment. No one guessed the real purpose of the study. The experiment lasted
38
39 about 20 minutes.
40

41 **Coding of facial expressions**

42
43 We relied on a well-established objective coding system developed for the analysis of
44
45 facial expressions, the Facial Action Coding System (FACS; Ekman, Friesen, & Hager, 2002a).
46
47 This system allows the analysis of minimal units of facial activity, known as action units (AUs),
48
49 which are anatomically separate and visually distinguishable. We coded the presence of AUs
50
51 that are diagnostic of disgust, anger, contempt, smiling and surprise expressions, but also certain
52
53 additional AUs that could help disambiguate between the emotions of interest and other
54
55 additional AUs that could help disambiguate between the emotions of interest and other
56
57 emotions, namely sadness and fear. Specifically, we coded the following action units: AU1
58
59

MANY MORAL BUTTONS OR JUST ONE?

11

1
2
3 (inner brow raiser), AU2 (outer brow raiser), AU4 (brow lowered), AU5 (upper lid raiser), AU6
4 (cheek raiser), AU7 (lids tight), AU9 (nose wrinkle), AU10 (upper lip raiser), AU11 (nasolabial
5 furrow deepener), AU12 (lip corner puller), AU13 (sharp lip puller); AU14 (dimpler), AU15
6 (lip corner depressor), AU16 (lower lip depress), AU17 (chin raiser), AU18 (lip pucker), AU20
7 (lip stretch), 22 (lip funneler), AU23 (lip tightener), AU24 (lip presser), AU28 (lip suck), AU25
8 (lips part), AU26 (jaw drop), AU27 (mouth stretch). Below we present the scheme that we used
9 to classify the AUs into emotional expressions.
10
11
12
13
14
15
16

17
18 The facial expression of *disgust* was classified by the presence of AU9 (nose wrinkle),
19 AU10 (upper lip raise), or the combination of AU9 + AU10 (Ekman, Friesen, & Ancoli, 1980;
20 Smith, 1989). We recognize that AU10 is an ambiguous emotional indicator for disgust:
21
22 bilateral AU10 could also indicate anger, while unilateral AU10 could also indicate contempt
23 (Ekman & Friesen, 1986; Rozin et al., 1999). However, since our main aim was to assess the
24 MFT theory, we followed Rozin et al. (1999, p. 584) and counted both types of AU10 (either in
25 isolation or combined with AU9) as indicative of moral disgust. We also classified disgust by
26 the presence of AU25 (lips part) but only if it appeared in combination with AU9 and/or AU10
27 (Langner, Dotsch, Bijlstra, Wigboldus, Hawk, & van Knippenberg, 2010; Rozin, Lowery, &
28 Ebert, 1994).
29
30
31
32
33
34
35
36
37
38

39 The prototypical facial expression of *anger* involves the combination of frowning (AU4),
40 lid tightening (AU7), and lip tightening/lip pressing (AU23/AU24). However, several additional
41 variations of this expression have been reported (Durán, Reisenzein, & Fernández-Dols, 2017).
42 Following previous research, we classified anger by the following AU combinations: AU4 +
43 AU7 + AU23/AU24, AU4 + AU7, AU4 + AU5, AU7 + AU5, AU7 + AU23, AU7 + AU24, AU7
44 + AU17 (e.g., Durán et al., 2017; Ekman, Friesen, & Hager, 2002b; Langner et al., 2010;
45 Matsumoto, Keltner, Shiota, O'Sullivan, & Frank, 2008; Rosenberg, Ekman, & Blumenthal,
46 1998; Rozin et al., 1999; Sayette & Hufford, 1995; Smith, 1989; Wiggers, 1982). Furthermore,
47 we also classified anger by the presence of the following AUs in isolation: AU4 (brow lowered;
48
49
50
51
52
53
54
55
56
57
58
59

MANY MORAL BUTTONS OR JUST ONE?

12

1
2
3 Cannon et al., 2011; Smith, 1989; Wiggers, 1982), AU7 (lips tight; Smith, 1989), AU23 (lip
4 tightened; Alvarado & Jameson, 2002; Ekman et al., 2002b; Wiggers, 1982), and AU24 (lip
5 presser; Alvarado & Jameson, 2002; Rozin et al., 1999, p. 579).

6
7
8
9 Unlike disgust, the coding of anger is problematic because certain single AUs or AU
10 combinations that are associated with anger are also associated with sadness and fear (e.g., AU4
11 is associated with all three emotions, while AU4+AU5 is associated with both anger and fear).
12
13 For this reason, we coded certain additional AUs—AU1, AU2, AU15 and AU20—whose
14 presence may indicate sadness or fear (Durán et al., 2017; Matsumoto et al., 2008). Specifically,
15 the facial expression of sadness was classified by the presence of AU1+AU4+AU15, whereas
16 the facial expression of fear by the presence of AU1/2+AU4+AU5+AU20 (Durán et al., 2017).

17
18
19
20
21
22
23
24 The facial expression of *contempt* was classified by the presence of the unilateral AU14
25 (dimple) (Ekman & Friesen, 1986; Matsumoto, 1992; Rozin et al., 1999). The *smiling* facial
26 expression was classified by the presence of AU12 (lip corner pull) or AU6 + AU12 (cheek
27 raise with lip corner pull) (Ekman et al., 1980; Ekman, Davidson, & Friesen, 1990; Sayette &
28 Hufford, 1995; Smith, 1989). Finally, the facial expression of *surprise* comprises three
29 components: eyebrow raising (AU1, AU2), eye widening (AU5), and mouth opening (AU25)
30 (Durán et al., 2017; Reisenzein, 2000). We classified this expression by the following
31 combination of AUs—AU1+AU2+AU5+AU25. Note that the expression of surprise shares
32 several AUs (i.e., AU1, AU2, AU5) with the expression of fear.

33
34
35
36
37
38
39
40
41
42
43
44 The discussion of the anger and surprise expressions indicates a potential problem of
45 classification ambiguity as these facial expressions share action units with sadness and fear. To
46 deal with such ambiguities, we employed a similarity-based rule.¹ Whenever the facial
47 expression displayed by a subject in response to a particular item fully matched one emotion
48 (i.e., it contained all AUs associated with a variant of that emotion) while only partially matched
49 another, we classified it as an expression of the former emotion. For example, the combination

50
51
52
53
54
55
56
57
58
59
60

¹ We thank an anonymous reviewer for this suggestion.

1
2
3 AU4+AU7+AU20+AU24 was classified as an expression of anger as it included all the features
4 required in one variant of anger (AU4+AU7+AU24) while only partially matched the
5 expression of fear (it matched two of its components [AU4, AU20] but also missed two [AU1/2,
6 AU5]). Whenever a facial expression partially matched two or more emotions then we counted
7 it as an expression of the emotion that it matched more closely. For example, the combination
8 AU1+AU15 was classified as an expression of sadness (as it misses one component for the
9 facial expression of sadness, while it misses three components for the facial expressions of
10 surprise or fear). In cases where some AU/AUs that are not core components of an emotion
11 were present in isolation (e.g., AU18; AU26+AU28), we did not classify them as an expression
12 of a particular emotion.

13
14
15
16
17
18
19
20
21
22
23
24 A certified FACS rater watched all video-recordings. Forty-two percent of the data were
25 also comparison coded by an independent certified FACS coder. Both coders viewed videotapes
26 in slow motion using VLC media player on a laptop and listed the presence of every single AU
27 of interest on a coding sheet. According to MFT, violations automatically activate distinct
28 emotions. Therefore, we focused on the first facial expression, and consequently the first AUs,
29 among those of our interest, that a participant displayed immediately after listening to a
30 violation. Our window of analysis started with the oral presentation of a moral violation and
31 continued until the instruction “Please think about this action” had disappeared from the screen
32 (see Figure 1).

33
34
35
36
37
38
39
40
41
42
43
44 The interrater reliability (Cohen’s Kappa) was based on the classification of the AUs,
45 and it was calculated on 42.4% of all data-points (280 out of 660 items). This analysis showed
46 an excellent consistency between the two independent coders for the classification of AU1, AU2,
47 AU4, AU6, AU10, AU15, AU17, AU18, AU20, AU23, AU24, AU26, AU28 ($\kappa > .80$), and a
48 good consistency for the classification of AU5, AU7, AU9, AU12, AU14, AU25 ($.65 < \kappa < .74$)
49 (e.g., Fleiss, 1981). Inconsistencies were resolved through discussion between the coders.

50 51 52 53 54 55 56 57 58 59 60 **Results**

MANY MORAL BUTTONS OR JUST ONE?

14

1
2
3 *Facial expressions (Task 1)*. We first examined the emotional expressions elicited by
4 listening to the scenarios. MFT predicts that anger would be the most frequent reaction to
5 harm violations, and disgust the most frequent reaction to purity violations. Consistent with
6 this prediction, in MFT harm scenarios anger was the most frequent expression. It was more
7 frequent than disgust ($t(32) = 7.21, p < .001$), smiling expressions ($t(32) = 5.79, p < .001$),
8 and contempt ($t(32) = 4.80, p < .001$). However, in MFT purity scenarios disgust was not the
9 most frequent expression. It was more frequent than contempt ($t(32) = 2.70, p = .011$), but
10 less frequent than anger ($t(32) = 2.90, p = .007$) and smiling expressions ($t(32) = 3.82, p$
11 = .001). Expressions of surprise were rare, and therefore excluded from the analyses.²

12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
Analyses of the naturalistic scenarios produced similar results. As predicted by MFT,
anger was the most frequent expression in naturalistic harm scenarios. It was more frequent
than smiling ($t(32) = 5.23, p < .001$), disgust ($t(32) = 4.99, p < .001$), and contempt ($t(32) =$
2.81, $p = .008$). However, the disgust expression was the least frequent expression in
response to purity scenarios. There were fewer expressions of disgust than of anger, smiling
or contempt (respectively, $t(32) = 5.66, p < .001$; $t(32) = 6.49, p < .001$; $t(32) = 2.73, p$
= .010).

54
55
56
57
58
59
60
We then examined the weaker version of MFT (i.e., the hypothesis that in
comparison to harm scenarios, purity scenarios trigger less anger and more disgust) by
comparing the presence of anger and disgust expressions across harm and purity scenarios.
As before, we first focused on MFT and then on naturalistic scenarios. As predicted by weak
MFT, the anger expression was more frequently exhibited in response to MFT harm scenarios
than MFT purity scenarios ($t(32) = 2.06, p = .048$), while the disgust expression was more
frequently exhibited in response to MFT purity scenarios than MFT harm scenarios, $t(32) =$

²The expressions classified as surprise were rare. However, there were several cases where one or two components of the surprise expression was present (e.g., AU1, AU1+AU2) (see Table 2). Because such cases are difficult to classify in terms of our similarity-based scheme (as the match with surprise was not higher than that with sadness or fear), we decided to leave them unclassified.

MANY MORAL BUTTONS OR JUST ONE?

15

1
2
3 3.08, $p = .004$ (see Figure 2a). In additional analyses, we found that contempt showed the
4 same trend as the anger expression, that is, it was more frequently exhibited in MFT harm
5 scenarios than MFT purity scenarios ($t(32) = 3.436, p = .002$). In contrast, smiling
6 expressions followed the opposite pattern: they were more frequent in response to MFT
7 purity scenarios than MFT harm scenarios ($t(32) = 5.50, p < .001$, see Figure 2a).
8
9
10
11
12

13 However, the analyses of naturalistic scenarios produced less favorable results for
14 weak MFT. Contrary to its predictions, no difference was found in the frequency of anger and
15 disgust (or contempt) expressions across harm and purity scenarios (see Figure 2b). The only
16 significant difference between naturalistic harm and naturalistic purity scenarios was in smile
17 expressions: smiles were more frequent in response to purity scenarios than in response to
18 harm scenarios, $t(32) = 5.75, p < .001$ (see Figure 2b).
19
20
21
22
23
24
25

26 --INSERT FIGURE 2 ABOUT HERE--
27

28 Tables 2 and 3 show the frequency of action units and facial expressions of interest,
29 respectively, elicited by specific items. As can be seen in Table 3, among the MFT harm
30 scenarios, items HM1, HM2, HM5 elicited selectively the expression of anger. Among MFT
31 purity scenarios, only item PM2 most frequently elicited the expression of disgust. Item PM5
32 elicited mostly smiles, PM4 elicited anger, while PM1 and PM3 elicited equally smiles and
33 anger. So, contrary to strong MFT, only one of the purity scenarios selectively elicited disgust.
34
35
36
37
38
39
40
41

42 Turning to naturalistic harm scenarios, items H3, H4 elicited selectively the
43 expression of anger. None of the naturalistic purity violations predominantly elicited disgust
44 expressions. Items P2, P3, and P5 elicited mostly smiles, P4 elicited anger and P1 elicited
45 equally smiles, anger and contempt. In sum, in both MFT and naturalistic purity scenarios the
46 expressions of disgust were rare: instead, there were frequent smiles and expressions of anger.
47
48
49
50
51

52 The use of single AUs as indicators of anger is potentially problematic because it
53 makes our similarity-based rule to favor anger classifications. In an attempt to address this
54 issue, in Table 4 we distinguish between anger expressions classified on the basis of a
55
56
57
58
59

MANY MORAL BUTTONS OR JUST ONE?

16

1
2
3 combination of AUs (and thus are less ambiguous) from those classified on the basis of a
4
5 single AU (and thus are more ambiguous). Note that the pattern of results remains essentially
6
7 unchanged if we replace the 'Anger' column of Table 3 with the 'Combination' column of
8
9 Table 4. Specifically, the emotions that originally emerged as predominant in each of the four
10
11 types of scenarios, also emerge as predominant even if we restrict anger expressions to cases
12
13 where a combination of AUs was present. This is also the case if we consider instead
14
15 individual items.
16

17
18 --INSERT TABLE 1, TABLE 2, TABLE 3 AND TABLE 4 ABOUT HERE--
19

20 It is noteworthy to mention that Ekman and Friesen (2003, Chapter 7) have proposed
21
22 an even more restrictive classification of the facial expression of anger, which they call 'the
23
24 full anger expression'. According to this classification, a simple combination of AUs is not
25
26 sufficient to classify an expression as anger; in addition, the facial signals must be present in
27
28 all three facial areas ("the brow/forehead; the eyes/lids and root of the nose; and the lower
29
30 face, including the cheeks, mouth, most of the nose, and chin", p. 28). The AU combinations
31
32 that we observed with this characteristic were: AU4+AU7+AU23/AU24/AU17;
33
34 AU4+AU5+AU17. Even if we follow this highly restrictive classification scheme, the pattern
35
36 of results with respect to the hypotheses of interest is similar to that observed from AU
37
38 combinations (Table 4). In relation to the MFT scenarios, there were more anger expressions
39
40 for the harm versus the purity items (full facial expression: 10 vs. 3; all AU combinations: 41
41
42 vs. 32), whereas in relation to the Naturalistic scenarios, we observed an equal number of
43
44 anger expressions across the harm and purity items (full facial expression: 2 vs. 3; all AU
45
46 combinations: 21 vs. 20).
47
48
49

50 *Moral judgments (Task 2)*. We expected to replicate the findings of Gray and Keeney
51
52 (2015). Specifically, we hypothesized that, compared to MFT harm scenarios, MFT purity
53
54 scenarios would be rated as less severe and weirder. We further expected that MFT harm
55
56 scenarios would be perceived as more harmful and, strangely, more impure than MFT purity
57
58
59

MANY MORAL BUTTONS OR JUST ONE?

scenarios. We examined these hypotheses via an analysis of variance (ANOVA) with Content (Harm vs. Purity) as a within-subjects factor on four dependent variables: Severity, Weirdness, Harm, and Impurity.

We found that, with respect to MFT purity scenarios, MFT harm scenarios received higher severity ratings, $F(1, 32) = 50.18, p < .001, \eta_p^2 = .61$, lower weirdness ratings, $F(1, 32) = 44.99, p < .001, \eta_p^2 = .58$, and higher harmfulness ratings, $F(1, 32) = 66.81, p < .001, \eta_p^2 = .68$.

Strangely, but in line with Gray and Keeney (2015, Study 1), MFT harm scenarios also received higher impurity ratings than MFT purity scenarios, $F(1, 32) = 18.18, p < .001, \eta_p^2 = .36$. Harm and impurity ratings were highly correlated in the scenarios, $r(8) = .82, p = .004$. The findings are fully in line with Gray and Keeney's (2015). MFT purity scenarios were perceived as less severe, less harmful, less impure, and weirder than MFT harm scenarios (see Figure 3).

We next compared the two types of scenarios, MFT versus naturalistic, with respect to weirdness ratings, moral judgment, severity ratings, impurity ratings, and harmfulness ratings. We found that MFT scenarios were judged as more weird than naturalistic scenarios, $t(32) = 7.06, p < .001$. In particular, MFT purity scenarios were judged as more weird than their counterpart naturalistic scenarios, $t(32) = 9.34, p < .001$ (see Figure 3). These findings are in line with Gray and Keeney's (2015) results. MFT scenarios were also judged as less severe than naturalistic scenarios, $t(32) = 2.21, p = .034$. However, analyzing separately harm and purity scenarios, no significant differences were found. Thus, neither MFT harm nor MFT purity scenarios were judged as less severe than their naturalistic scenario counterparts (see Figure 3). In this case, only the general t test output is in line with Gray and Keeney's (2015) results. Finally, MFT scenarios were judged less morally wrong, impure, and harmful than naturalistic scenarios, $t(32) = 2.91, p = .006$; $t(32) = 3.55, p = .001$; $t(32) = 2.78, p = .009$, respectively. Moreover, when we conducted paired-sample t tests for harm and purity, naturalistic and MFT scenarios, we found that only MFT purity scenarios were judged as less morally wrong, impure

1
2 and harmful than their naturalistic counterparts, $t(32) = 3.71, p = .011$; $t(32) = 2.61, p = .014$;
3
4
5 $t(32) = 3.58, p = .001$ (Figure 3). In sum, the results of the moral judgment task provide a useful
6
7 cross-cultural replication of Gray and Keeney (2015, Study 2).
8

9 --INSERT FIGURE 3 ABOUT HERE--
10

11 **Discussion**

12
13 We tested a strong and a weak claim of the moral foundation theory. According to the
14 strong claim, harm scenarios predominantly elicit anger reactions whereas purity scenarios
15 predominantly elicit disgust reactions. According to the weaker claim, with respect to harm
16 scenarios, purity scenarios elicit fewer anger reactions and more disgust reactions. Participants
17 listened to scenarios depicting harm or purity violations, and their spontaneous facial expressions
18 were analyzed. The scenarios included classic MFT scenarios (e.g., Haidt et al., 1993), and more
19 naturalistic scenarios, in which the harmful and impure actions were matched for typicality
20 ('weirdness', Gray & Keeney, 2015). The results—at least a straightforward interpretation of
21 them—do not support either MFT prediction.
22
23

24 Specifically, we found that the expression of anger was elicited equally often by harm and
25 purity violations, and this was particularly true for naturalistic scenarios. Disgust reactions were
26 rare. In line with weak MFT, disgust reactions were more frequent in response to purity than to
27 harm violations, but only for MFT scenarios. But even for MFT purity scenarios, contrary to
28 strong MFT, disgust expressions were less frequent than anger or smiling expressions. For
29 naturalistic scenarios, disgust expressions were extremely rare for both harm and purity violations.
30 The higher frequency of disgust expressions in response to MFT purity violations could be due to
31 the fact that these violations included physically disgusting elements. The only significant
32 difference between harm and purity violations was in the frequency of smiling expressions: purity
33 violations triggered significantly more smiles than harm violations. Intriguingly, this was the case
34 both for MFT and naturalistic scenarios. Had this been the case only for MFT purity scenarios,
35 one could attribute it to their weirdness.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Why did participants smile more in response to purity than harm violations? The versatile
4 nature of the human smile was highlighted in several studies (e.g., Ekman, 1992; Keltner, 1995;
5 Rychlowska, Jack, Garrod, Schyns, Martin, & Niedenthal, 2017). In our case, one possibility is
6 that smiles indicate covert disgust. Under this interpretation, the present results could still be
7 consistent with both the strong and weak MFT. Notice that for purity scenarios (both MFT and
8 naturalistic) the sum of smiling and disgust expressions was higher than that of anger expressions
9 (inspect Figure 2). Another, more plausible, explanation is that participants found purity violations
10 more amusing and less threatening than harm violations (see McGraw & Warren, 2010). This
11 explanation appears to be consistent with a domain general proposal that considers purity
12 violations as involving personal harm, but in a more subtle and indirect way than harm violations.
13 Notice, however, that there were also more smiles in response to naturalistic purity versus
14 naturalistic harm scenarios, which were balanced in terms of perceived weirdness. Thus, a
15 domain-general proposal cannot easily explain the observed association between smiles and purity
16 violations.

17
18 Smiles might be associated with pleasure, relief, or amusement (enjoyment smiles), but
19 non-enjoyment smiles also exist. Theorists have identified various types of non-enjoyment smiles
20 such as ‘embarrassed smiles’, ‘masking smiles’ (while experiencing negative emotions, the smile
21 at least partially covers muscular movements associated with another emotion), ‘false smiles’
22 (appearing as if enjoyment is felt when it is not), ‘anticipatory smiles’ or ‘miserable smiles’
23 (representing the willingness to endure unpleasant circumstances) (Ekman, 1985; Ekman &
24 Friesen, 1982). These different forms of non-enjoyment smiles could be associated with
25 compliance, embarrassment, shame, grin-and-bear-it (Ekman, 1992; Keltner, 1995; Keltner &
26 Buswell, 1996).

27
28 To which types of smile do the smiles we found belong to? Although there are no
29 standardized procedures for classifying different types of non-enjoyment smiles, there is a wide
30 agreement in the literature that when the lip corner raising is accompanied by the cheek raising

MANY MORAL BUTTONS OR JUST ONE?

20

1
2
3 (commonly known as a Duchenne smile) the smile denotes enjoyment (Ekman et al., 1990). In
4
5 view of that, we reanalyzed our data on purity scenarios classifying smiles as either Duchenne
6
7 smiles or non-Duchenne smiles (or simple smiles). We counted about twice as many Duchenne
8
9 smiles than simple smiles in response to both MFT scenarios (31 Duchenne, 18 simple) and
10
11 naturalistic scenarios (24 Duchenne, 14 simple). Thus, consistent with McGraw and Warren
12
13 (2010), these results suggest that participants found some purity violations amusing.
14

15
16 The discussion on smiles highlights two limitations of the present study. The first is that
17
18 emotional responses could be complex, that is, they could involve a combination of different
19
20 emotions (e.g., Gross & Levenson, 1995; Kuppens, Tuerlinckx, Russell, & Barrett, 2013). To
21
22 tackle this critical point, future research could investigate whether the present findings hold also
23
24 with implicit physiological measures and other behavioral responses. Moreover, future research
25
26 could use experimental paradigms involving ecological social contexts, which would prompt
27
28 people to explicitly display spontaneous and stronger facial expressions in order to communicate a
29
30 social message to others.
31

32
33 A second limitation of the present study concerns the use of facial expressions as an
34
35 indicator of emotions. Emotions are not always accompanied by a facial expression and this is
36
37 especially true for anger and disgust—the main emotions of interest in the present study—but less
38
39 so for amusement (see e.g., Durán et al., 2017). Moreover, the fact that the participants knew that
40
41 they were videotaped could have reinforced this tendency (Durán et al., 2017). Thus, a defender of
42
43 MFT could argue that the purity scenarios did elicit mostly disgust, but people did not displayed it
44
45 in their facial expressions. However, for this argument to go through, the defender of MFT must
46
47 also assume that expressions of disgust were more strongly inhibited than expressions of anger.
48
49 Therefore, this objection is not really strong.³
50
51
52
53
54
55

56
57 ³We thank an anonymous reviewer for pointing this limitation to us and for discussing its
58 implications for the MFT hypotheses.
59

1
2
3 However, as this is a critical issue, we decided to address it directly in a follow up study
4 (for details, see Supplementary materials). We presented a new sample of 34 participants with the
5 20 scenarios of the main study, but instead of focusing on facial expressions of emotions, we
6 gathered self-ratings. Following each scenario, we asked participants to select the first emotion
7 they felt while reading it (“While you were reading about this situation, what was your first
8 emotion?”) and to rate the intensity of that emotion (“Please express a judgment on the intensity of
9 this emotion”). Participants could respond to the first question by selecting one of the following
10 emotions: enjoyment, grin-and-bear-it, embarrassment, shame, disgust, anger, contempt, surprise.
11 The reason we included enjoyment, grin-and-bear-it, embarrassment and shame was to help
12 disambiguate the meaning of smiles in our main study (for a more detailed discussion, see
13 Supplementary materials).
14
15
16
17
18
19
20
21
22
23
24
25

26 Notwithstanding the limitations of measuring emotions through self-reports (see
27 Introduction), with respect to the MFT hypotheses, the results are largely convergent with those
28 of the main study (see Figure S1).⁴ For both MFT and naturalistic harm scenarios anger was the
29 most frequent response, whereas for purity scenarios there was no clear prevalent emotion. In
30 the MFT purity scenarios the most frequent response was disgust, but it was closely followed by
31 surprise, enjoyment, and grin-and-bear-it. In the naturalistic purity scenarios the most frequent
32 response was contempt, closely followed by grin-and-bear-it and embarrassment. Thus, the
33 results do not support either MFT hypotheses. Strong MFT can explain the results for harm
34 scenarios (anger was the predominant response) but not for purity scenarios (no emotion
35 prevailed). Weak MFT can explain the responses to the MFT scenarios (harm scenarios elicited
36 more anger than disgust, whereas purity scenarios elicited more disgust than anger), but not to
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

52 ⁴In the Introduction, we raised two issues with respect to measuring emotion through self-reports:
53 (a) certain emotional words might be ambiguous and (b) the results could be influenced by post
54 hoc rationalizations. With regards to (a) we found that five participants expressed doubts about the
55 meaning of certain emotional words (e.g., ‘*sdegno*’, the Italian word for contempt). With regards
56 to (b) since we asked directly about emotional judgments—these judgments were not preceded by
57 moral judgments—the impact of post hoc rationalizations should be limited.
58
59

MANY MORAL BUTTONS OR JUST ONE?

22

1
2
3 naturalistic scenarios (harm and purity scenarios elicited equal amounts of disgust). The low
4
5 frequency of disgust expression in the new study further questions the possibility that
6
7 participants in the main study experienced disgust, but they suppressed the facial expression of
8
9 it because they were video-recorded.

10
11 Returning to the main study on facial expressions, the present pattern of results is also
12
13 consistent with a semi-exclusive correspondences model, in which harm violations elicit anger,
14
15 and purity violations elicit not only disgust but also other emotions (see Cameron et al., 2015,
16
17 Figure 4 for an adaptation to our results of the schema proposed by Cameron et al., 2015). The
18
19 results of the follow up study provide further support for this assertion (see Figure S1). This
20
21 happens probably because harm violations instantiate clear and direct cases of personal harm,
22
23 while purity violations are more complex, which, in turn, is reflected in the different emotions
24
25 they elicit. For example, the purity items showing rape and adultery are instances of violations not
26
27 only in the purity domain, but also in other domains such as harm and loyalty (Graham, 2015).
28
29 Thus, it is of no surprise that, by using such ‘mixed’ items, one does not find the specific
30
31 association effects predicted by MFT.
32
33

34
35 --INSERT FIGURE 4 ABOUT HERE--
36

37
38 This argument rests on the assumption that pure violations of purity exist, which is exactly
39
40 what has been highly controversial from the beginning of this debate (Turiel, 1989). Domain
41
42 general views deny the existence of harmless moral violations of purity and claim that all purity
43
44 violations are perceived as potentially harmful (Royzman et al., 2016). The data reported by Gray
45
46 and Keeney (2015) show that people do associate purity items, both naturalistic and MFT ones, to
47
48 the harm domain, as claimed by the domain general view (on the lack of coherence and clarity in
49
50 defining the domain of purity, see Russell & Giner-Sorolla, 2013). So, defenders of pluralist
51
52 models should not assume, by simply relying on intuition, that some instances of purity violations
53
54 are fine while others are not, but should instead provide an objective way of identifying such cases.
55
56
57
58
59

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Meanwhile, the use of alternative sets of scenarios that allow control over relevant factors, such as typicality, seems a fruitful research strategy.

In the present study we analyzed the emotional reactions that were immediately and spontaneously expressed while listening to moral violations, but we cannot specify at what stage of the information processing such reactions occurred. They could have occurred as a response to a covert judgment, or they could have been the basis upon which the judgment took place, along the lines of classic works in the empiricist tradition (Haidt et al., 1993; Hume, 1751/1957). This ambiguity does not allow us to determine which particular aspect of domain specificity is challenged by the present results. The results of the follow up study similarly do not shed light on this matter. Future research may wish to tackle this issue.

In conclusion, the present results go against MFT theory. Although anger was the most frequent expression in response to harm violations, disgust was not the most frequent response to purity violations. The present results also go against a weaker version of MFT theory. In the main study, the naturalistic harm and purity scenarios—which were matched in weirdness—were virtually identical in the amounts of anger, disgust, and contempt expressions they triggered. This result fits particularly well with monist claims that eventual differences between harm and purity scenarios are due to lack of control of several important factors such as weirdness. In the follow up study with self-ratings a slightly different pattern emerged: naturalistic harm and purity scenarios elicited equal amounts of disgust and contempt, but harm scenarios elicited more anger. However, one result in the present research keeps alive the possibility of distinct moral foundations: Smiles were selectively associated with violations of purity. Therefore, the general theory may be right—different violations have distinct emotional footprints—but the original claims about violation-emotion pairs may be wrong or incomplete. Violations of purity, like benign norm violations, may be characterized by the elicitation of amusement, but also by other emotions that may underlie the smiling expression. This opens exciting avenues for new research.

References

- Alvarado, N., & Jameson, K. A. (2002). Varieties of anger: the relation between emotion terms and components of anger expressions. *Motivation and Emotion*, 26, 153-182.
doi:10.1023/A:1019815402873
- Bloom, P. (2004). *Descartes' baby: How the science of child development explains what makes us human*. New York: Basic.
- Cannon, P. R., Schnall, S., & White, M. (2011). Transgressions and expressions: Affective facial muscle activity predicts moral judgments. *Social Psychological and Personality Science*, 2, 325-331. doi:10.1177/1948550610390525
- Cameron C. D., Lindquist K. A., & Gray K. (2015). A constructionist review of morality and emotions: No evidence for specific links between moral content and discrete emotions. *Personality and Social Psychology Review*, 9, 371-394. doi:10.1177/1088868314566683
- Cushman, F. (2013). Action, outcome, and value: a dual-system framework for morality. *Personality and Social Psychology Review*, 17, 273-92. doi:10.1177/1088868313495594.
- Durán, J. I., Reisenzein, R., & Fernández-Dols, J.-M. (2017). Coherence between emotions and facial expressions. In J.-M. Fernández-Dols & J. A. Russell (Eds.), *The science of facial expression* (pp.107-129). Oxford University Press.
doi:10.1093/acprof:oso/9780190613501.001.0001
- Ekman, P. (1985). *Telling lies: Clues to deceit in the marketplace, marriage, and politics*. New York: W.W. Norton.
- Ekman, P. (1992). Facial expressions of emotion: New findings, new questions. *Psychological science*, 3, 34-38. doi:10.1111/j.1467-9280.1992.tb00253.x
- Ekman, P. (1994). Antecedent events and emotion metaphors. In P. Ekman & R. J. Davidson (Eds.), *The nature of emotion: Fundamental questions* (pp. 146-149). New York: Oxford University Press.
- Ekman, P., Davidson, R. J., & Friesen, W. V. (1990). The Duchenne smile: Emotional expression

- 1
2
3 and brain physiology: II. *Journal of Personality and Social Psychology*, 58, 2, 342-353.
4
5 Ekman, P., & Friesen, W. V. (1982). Felt, false, and miserable smiles. *Journal of Nonverbal*
6
7 *Behavior*, 6, 238-252. doi:10.1007/BF00987191
8
9 Ekman, P., & Friesen, W. V. (1986). A new pancultural facial expression of emotion. *Motivation*
10
11 *and Emotion*, 10, 159-168. doi:10.1007/BF00992253
12
13 Ekman, P., & Friesen, W. V. (2003). *Unmasking the face. A guide to recognizing emotions from*
14
15 *facial expressions*. Cambridge, MA: Malor Books.
16
17 Ekman, P., Friesen, W. V., & Ancoli, S. (1980). Facial signs of emotional experience. *Journal of*
18
19 *Personality and Social Psychology*, 39, 1125-1134. doi:10.1037/h0077722
20
21 Ekman, P., Friesen, W. V., & Hager, J. C. (2002a). *Facial Action Coding System (FACS)*. Salt
22
23 Lake City, UT: Research Nexus division of Network Information Research Corporation.
24
25 Ekman, P., Friesen, W. V., & Hager, J. C. (2002b). *Facial Action Coding System. Investigator's*
26
27 *guide*. Salt Lake City, UT: Research Nexus division of Network Information Research
28
29 Corporation.
30
31
32
33 Faul, F., Erdfelder, E., Lang, A. G., & Buchner, A. (2007). G*Power 3: A flexible statistical
34
35 power analysis program for the social, behavioral, and biomedical sciences. *Behavior*
36
37 *Research Methods*, 39, 175-191. doi:10.3758/BRM.41.4.1149
38
39 Giner-Sorolla, R., & Chapman, H. A. (2017). Beyond purity: Moral disgust toward bad character.
40
41 *Psychological Science*, 28, 80-91. doi:10.1177/0956797616673193
42
43
44 Graham, J. (2015). Explaining away differences in moral judgment: Comment on Gray and
45
46 Keeney (2015). *Social Psychological and Personality Science*, 6, 869-873.
47
48 doi:10.1177/1948550615592242
49
50 Graham, J., Haidt, J., Koleva, S., Motyl, M., Iyer, R., Wojcik, S., & Ditto, P. H. (2013). Moral
51
52 Foundations Theory: The pragmatic validity of moral pluralism. *Advances in Experimental*
53
54 *Social Psychology*, 47, 55-130. doi:10.1016/B978-0-12-407236-7.00002-4
55
56
57 Graham, J., Haidt, J., & Nosek, B. A. (2009). Liberals and conservatives rely on different sets of
58
59

MANY MORAL BUTTONS OR JUST ONE?

26

1
2
3 moral foundations. *Journal of Personality and Social Psychology*, 96, 1029–1046.

4
5 doi:10.1037/a0015141

6
7 Gray, K., & Keeney, J. E. (2015). Impure or just weird? Scenario sampling bias raises questions
8
9 about the foundation of morality. *Social Psychological and Personality Science*, 6, 859-868.

10
11 M doi:10.1177/1948550615592241

12
13 Gray, K., Schein, C., & Ward, A. F. (2014). The myth of harmless wrongs in moral cognition:

14
15 Automatic dyadic completion from sin to suffering. *Journal of Experimental Psychology:*

16
17 *General*, 143, 1600-1615. doi:10.1037/a0036149

18
19 Gross, J. J., & Levenson, R. W. (1995). Emotion elicitation using films. *Cognition and Emotion*, 9,

20
21 87-108. doi:10.1080/02699939508408966

22
23 Gutierrez, R., & Giner-Sorolla, R. (2007). Anger, disgust, and presumption of harm as reactions to

24
25 taboo-breaking behaviors. *Emotion*, 7, 853-868. doi:10.1037/1528-3542.7.4.853

26
27 Haidt, J. (2001). The emotional dog and its rational tail: a social intuitionist approach to moral

28
29 judgment. *Psychological Review*, 108, 814-834. doi:10.1037/0033-295X.108.4.814

30
31 Haidt, J., Koller, S. H., & Dias, M. G. (1993). Affect, culture, and morality, or is it wrong to eat

32
33 your dog? *Journal of Personality and Social Psychology*, 65, 613-628. doi:10.1037//0022-

34
35 3514.65.4.613

36
37 Hume, D. (1957) *An inquiry concerning the principle of morals* (Vol. 4). New York: Liberal Arts

38
39 Press. (Original work published 1751).

40
41 Izard, C. E. (1977). *Human emotions*. New York: Plenum Press.

42
43 Keltner, D. (1995). Signs of appeasement: Evidence for the distinct displays of embarrassment,

44
45 amusement, and shame. *Journal of Personality and Social Psychology*, 68, 441-454.

46
47 doi:10.1037/0022-3514.68.3.441

48
49 Keltner, D., & Buswell, B. N. (1996). Evidence for the distinctness of embarrassment, shame, and

50
51 guilt: A study of recalled antecedents and facial expressions of emotion. *Cognition and*

52
53 *Emotion*, 10, 155-171. doi:10.1080/026999396380312

MANY MORAL BUTTONS OR JUST ONE?

27

- 1
2
3 Kuppens, P., Tuerlinckx, F., Russell, J. A., & Barrett, L. F. (2013). The relation between valence
4 and arousal in subjective experience. *Psychological Bulletin*, *139*, 917-940.
5
6 doi:10.1037/a0030811
7
8
9 Langner, O., Dotsch, R., Bijlstra, G., Wigboldus, D. H. J., Hawk, S. T., & van Knippenberg, A.
10 (2010). Presentation and validation of the Radboud Faces Database. *Cognition and Emotion*,
11 *24*, 1377-1388. doi:10.1080/02699930903485076
12
13
14
15 Mathôt, S., Schreij, D., & Theeuwes, J. (2012). OpenSesame: An open-source, graphical
16 experiment builder for the social sciences. *Behavior Research Methods*, *44*, 314-324.
17
18 doi:10.3758/s13428-011-0168-7
19
20
21
22 Matsumoto, D. (1992). More evidence for the universality of a contempt expression. *Motivation*
23 *and Emotion*, *16*, 363-368. doi: 10.1007/BF00992972
24
25
26 Matsumoto, D., Keltner, D., Shiota, M. N., O'Sullivan, M., & Frank, M. (2008). Facial
27 expressions of emotion. In M. Lewis, J. M. Haviland-Jones, & L. Feldman Barrett (Eds.),
28 *Handbook of emotions* (pp. 211-234). New York: The Guilford Press.
29
30
31
32 McGraw, A. P., & Warren, C. (2010). Benign violations: Making immoral behavior funny.
33 *Psychological Science*, *21*, 1141-1149. doi:10.1177/0956797610376073.
34
35
36
37 Mikhail, J. (2007). Universal moral grammar: Theory, evidence and the future. *Trends in*
38 *Cognitive Sciences*, *11*, 143-152. doi:10.1016/j.tics.2006.12.007
39
40
41
42 Nabi, R. (2002). Anger, fear, uncertainty, and attitudes: a test of the cognitive-functional model.
43 *Communication Monographs*, *69*, 204-216. doi:10.1080/03637750216541
44
45
46 Parkinson, C., Sinnott-Armstrong, W., Koralus, P. E., Mendelovici, A., McGeer, V., & Wheatley,
47 T. (2011). Is morality unified? Evidence that distinct neural systems underlie moral
48 judgments of harm, dishonesty, and disgust. *Journal of Cognitive Neuroscience*, *23*, 3162-
49 3180. doi:10.1162/jocn_a_00017
50
51
52
53
54 Reisenzein, R. (2000). Exploring the strength of association between the components of emotion
55 syndromes: The case of surprise. *Cognition and Emotion*, *14*, 1-38.
56
57
58
59

MANY MORAL BUTTONS OR JUST ONE?

28

1
2
3 doi:10.1080/026999300378978

4 Rosenberg, E. L., Ekman, P., & Blumenthal, J. A. (1998). Facial expression and the affective
5 component of cynical hostility in male coronary heart disease patients. *Health Psychology*,
6 *17*, 376-380. doi:10.1037/0278-6133.17.4.376

7
8
9
10
11 Rottman, J., Kelemen, D., & Young, L. (2014). Tainting the soul: Purity concerns predict moral
12 judgments of suicide. *Cognition*, *130*, 217–226. doi:10.1016/j.cognition.2013.11.007

13
14
15 Royzman, E., Atanasov, P., Landy, J. F., Parks, A., & Gepty, A. (2016). CAD or MAD? Anger
16 (not disgust) as the predominant response to pathogen-free violations of the Divinity code.
17
18
19
20
21 *Emotion*, *14*, 892-907. doi:10.1037/a0036829

22
23
24
25
26
27 Rozin, P., Lowery, L., & Ebert, R. (1994). Varieties of disgust faces and the structure of disgust.
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
Journal of Personality and Social Psychology, *66*, 870-881. doi:10.1037/0022-
3514.66.5.870

Rozin, P., Lowery, L., Imada, S., & Haidt, J. (1999). The CAD triad hypothesis: A mapping
between three moral emotions (contempt, anger, disgust) and three moral codes (community,
autonomy, divinity). *Journal of Personality and Social Psychology*, *76*, 574-586.
doi:10.1037/0022-3514.76.4.574

Russell, P. S., & Giner-Sorolla, R. (2013). Bodily moral disgust: what it is, how it is different
from anger, and why it is an unreasoned emotion. *Psychological Bulletin*, *139*, 328-351.
doi:10.1037/a0029319

Rychlowska, M., Jack, R. E., Garrod, O. G. B., Schyns, P. G., Martin, J. D., & Niedenthal, P. M.
(2017). Functional smiles: Tools for love, sympathy, and war. *Psychological Science*, *28*,
1259-1270. doi: 10.1177/0956797617706082

Sayette, M. A., & Hufford, M. R. (1995). Urge and Affect: A Facial Coding Analysis of Smokers.
Experimental and Clinical Psychopharmacology, *3*, 417-423. doi:1064-1297/95/\$3.00

Schaich Borg, J., Kahn, R. E., Sinnott-Armstrong, W., Kurzban, R., Robinson, P. H., & Kiehl, K.
A. (2013). Subcomponents of psychopathy have opposing correlations with punishment

MANY MORAL BUTTONS OR JUST ONE?

1
2
3 judgments. *Journal of Personality and Social Psychology*, 105, 667-687.

4
5 doi:10.1037/a0033485

6
7 Smith, C. A. (1989). Dimensions of appraisal and physiological response in emotion. *Journal of*
8
9 *Personality and Social Psychology*, 56, 339-353. doi:0022-3514/89/S00.75

10
11 Turiel, E. (1989). Domain-specific social judgments and domain ambiguities. *Merrill- Palmer*
12
13 *Quarterly*, 35, 89-114.

14
15 Turiel, E., Killen, M., & Helwig, C. C. (1987). Morality: Its structure, functions, and vagaries. In J.
16
17 Kagan & S. Lamb (Eds.), *The emergence of morality in young children* (pp. 155-243).
18
19 Chicago: University of Chicago Press.

20
21 Uhlmann, E. L., & Zhu, L. (2014). Acts, persons, and intuitions: Person-centered cues and gut
22
23 reactions to harmless transgressions. *Social Psychological and Personality Science*, 5, 279-
24
25 285. doi:10.1177/1948550613497238

26
27 Vrana, S. R. (1993). The psychophysiology of disgust: Differentiating negative emotional contexts
28
29 with facial EMG. *Psychophysiology*, 30, 279-286. doi:10.1111/j.1469-8986.1993.tb03354.x

30
31 Wiggers, M. (1982). Judgments of facial expressions of emotion predicted from facial behavior.
32
33 *Journal of Nonverbal Behavior*, 7, 101-116. doi:10.1007/BF00986872

34
35 Young, L., & Saxe, R. (2011). When ignorance is no excuse: Different roles for intent across
36
37 moral domains. *Cognition*, 2, 202-214. doi:10.1016/j.cognition.2011.04.005

Table 1

Moral Foundations Theory (MFT) and naturalistic scenarios (respectively, see Graham et al., 2009, and Gray & Keeney, 2015).

	HARM	PURITY
MFT scenarios	HM1. Someone kicks a dog in the head, hard	PM1. Someone signs a piece of paper that says: “ I hereby sell my soul after my death, to whoever has this piece of paper”
	HM2. Someone shoots and kills an animal that is a member of an endangered species	PM2. Someone cooks and eats their dog, after it dies of natural causes
	HM3. Someone makes cruel remarks to an overweight person about his or her appearance	PM3. Someone gets plastic surgery that adds a 2 inch tail onto the end of their spine
	HM4. Someone steps on an anthill, killing thousand of ants	PM4. Someone gets a blood transfusion of 1 liter of disease-free, compatible blood from a convicted child molester
	HM5. Someone sticks a pin into the palm of a child they don't know	PM5. Someone attends a performance art piece in which all participants (including that person) have to act like animals for 30 minutes, including crawling around naked and urinating on stage
Naturalistic scenarios	H1. Someone physically strikes another person	P1. Someone has an affair with another person, while they are married to someone else
	H2. Someone speaks to another person in a cruel and offensive manner	P2. Someone hires a prostitute for an evening of sex
	H3. Someone harasses and intimidates another person	P3. Someone has a sex on camera, making a pornographic film that will be distributed for profit
	H4. Someone steals a valuable item from another person	P4. Someone forces another person to have sexual intercourse with them, without that person's consent
	H5. Someone intentionally kills another person	P5. Someone strips nude on stage in a room full of strangers

Table 2

Frequencies of the Action Units codified in correspondence with the Moral Foundations Theory (MFT) and the naturalistic scenarios.

MFT Harm scenarios																			
Item	AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU12	AU14	AU15	AU17	AU18	AU20	AU23	AU24	AU25	AU26	AU28
HM1	1	1	12	1	1	7	0	0	4	1	1	1	0	0	3	0	1	0	0
HM2	3	3	10	0	0	6	0	0	0	5	0	4	1	1	3	1	0	0	0
HM3	1	1	6	1	0	2	0	0	1	9	0	0	0	0	0	0	0	2	2
HM4	4	4	7	1	3	3	1	0	6	0	1	2	0	0	2	0	0	0	0
HM5	2	2	23	1	0	18	1	2	0	2	1	2	0	1	3	1	0	0	0
Tot.	11	11	58	5	4	36	2	2	11	17	3	9	1	2	11	2	1	2	2
MFT Purity scenarios																			
Item	AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU12	AU14	AU15	AU17	AU18	AU20	AU23	AU24	AU25	AU26	AU28
PM1	4	3	7	0	7	3	0	0	11	0	0	1	0	0	0	0	1	0	0
PM2	3	3	7	4	1	7	2	12	2	1	0	0	0	0	0	0	0	0	0
PM3	4	3	10	1	8	5	2	2	13	0	1	2	0	0	0	0	2	0	0
PM4	2	2	12	5	0	9	1	2	0	3	1	2	0	0	0	0	0	0	0
PM5	3	3	3	0	15	2	1	2	23	0	0	1	0	0	0	0	6	0	0
Tot.	16	14	39	10	31	26	6	18	49	4	2	6	0	0	0	0	9	0	0

Naturalistic Harm scenarios

Item	AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU12	AU14	AU15	AU17	AU18	AU20	AU23	AU24	AU25	AU26	AU28
H1	1	1	4	0	2	2	0	0	3	0	0	0	0	1	1	0	0	0	0
H2	2	1	5	0	0	3	0	0	0	5	0	0	1	0	1	0	0	2	2
H3	1	1	10	1	0	3	0	2	0	2	1	0	0	1	2	0	0	0	0
H4	3	3	6	1	0	4	1	1	0	3	0	3	0	0	4	0	1	0	0
H5	3	3	5	3	0	3	0	0	0	5	0	1	0	0	1	1	0	0	1
Tot.	10	9	30	5	2	15	1	3	3	15	1	4	1	2	9	1	1	2	3

Naturalistic Purity scenarios

Item	AU1	AU2	AU4	AU5	AU6	AU7	AU9	AU10	AU12	AU14	AU15	AU17	AU18	AU20	AU23	AU24	AU25	AU26	AU28
P1	5	4	3	0	2	4	0	0	4	5	0	3	0	0	1	1	0	1	1
P2	4	3	3	0	9	0	0	1	11	3	2	1	0	1	1	0	1	1	0
P3	3	3	1	2	8	1	0	0	12	2	1	1	0	0	0	1	3	2	1
P4	1	1	13	1	0	6	0	0	0	0	0	1	2	0	3	1	0	0	0
P5	5	3	7	0	5	4	0	0	11	1	1	1	0	0	0	0	1	1	0
Tot.	18	14	27	3	24	15	0	1	38	11	4	7	2	1	5	3	4	5	2

Note. Each AU was coded as present if it was visible in the slightest degree and bilaterally (except for AU10 and AU14). All AU10s reported in the table were bilateral. Although AU10 is indicative of both disgust and anger, we followed Rozin et al. (1999) and coded it as disgust. In several cases some AUs appeared in isolation or in an ambiguous combination with other AUs. Because the

1 classification of these cases is problematic, we did not count them as indicative of an emotional expression, but we report them in the
2
3 table. Here is a list of these cases: AU1 appeared alone twice; AU2 appeared in combination with only AU1 eighteen times; AU5
4
5 appeared in combination with only AU1+AU2 four times; AU15 appeared alone six times, and twice in combination with only AU1
6
7 or AU1+AU2+AU17; AU20 appeared alone once; AU18 appeared alone four times; AU26 appeared with only AU28 six times.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

Table 3

Frequencies of the facial expressions of interest for the Moral Foundations Theory (MFT) and the naturalistic scenarios.

		Harm					Purity						
		Item	Angry	Disgusted	Contemptuous	Smiling	Surprised	Item	Angry	Disgusted	Contemptuous	Smiling	Surprised
MFT scenarios	HM1	14	0	1	4	0	PM1	8	0	0	11	0	
	HM2	12	0	5	0	0	PM2	9	12	1	2	0	
	HM3	6	0	9	1	1	PM3	10	2	0	13	0	
	HM4	7	1	0	6	1	PM4	13	2	3	0	0	
	HM5	24	3	2	0	1	PM5	3	2	0	23	0	
	Tot.	63	4	17	11	3	Tot.	43	18	4	49	0	
Naturalistic scenarios	H1	5	0	0	3	0	P1	7	0	5	4	0	
	H2	6	0	5	0	0	P2	3	1	3	11	0	
	H3	10	2	2	0	0	P3	3	0	2	12	1	
	H4	10	1	3	0	1	P4	16	0	0	0	0	
	H5	6	0	5	0	2	P5	7	0	1	11	0	
	Tot.	37	3	15	3	3	Tot.	36	1	11	38	1	

1
2
3
4 *Note.* Considering the high frequency of smiles elicited by purity scenarios, we reanalyzed the data distinguishing between two types of smile:
5
6 the simple smile (non-Duchenne), consisting of pulling the lip corners back (AU12, without any other AU activity), and the Duchenne smile,
7
8 defined by the combination of AU6 + AU12 (cheek raise with lip corner pull) (e.g., Ekman, 1992b). The Duchenne smile is thought to be
9
10 associated with enjoyment, while non-Duchenne smiles are not (Ekman et al., 1990). In MFT purity scenarios there were 18 simple and 31
11
12 Duchenne smiles, while in naturalistic purity scenarios there were 14 simple and 24 Duchenne smiles. In regards to the emotional expressions
13
14 which were not of our immediate interest, we counted 9 partial expressions of sadness (7 of these were based on the isolated presence of AU15)
15
16 and 1 of fear (based on the isolated presence of AU20).
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

Table 4

Additional information about the frequencies of the facial expression of anger. We differentiate between expressions of anger identified by a combination of AUs (less ambiguous) from those identified by the presence of a single AU (more ambiguous).

		Harm			Purity		
		Item	Combination	Single	Item	Combination	Single
MFT scenarios	HM1		7	7	PM1	4	4
	HM2		8 ^Δ	4	PM2	7	2
	HM3		2	4	PM3	6*	4
	HM4		5*	2	PM4	12*	1
	HM5		19* ^{Δ†}	5	PM5	3	0
	Tot.		41	22	Tot.	32	11
Naturalistic scenarios	H1		2	3	P1	5	2
	H2		3	3	P2	1 ^Δ	2
	H3		4 ^Δ	6	P3	2 [†]	1
	H4		7	3*	P4	8*	8
	H5		5	1*	P5	4	3*
	Tot.		21	16	Tot.	20	16

Note. We observed the following AU combinations for anger: AU4+AU7+AU17/23/AU24; AU4+AU5+AU7/AU17; AU4+AU17+AU23/24; AU4+AU5/AU7; AU7+AU5/AU17/AU23/AU24. The single AUs we observed were: AU4; AU7; AU23; AU24. * Indicates the existence of a single ambiguous case where, in addition to the AUs that are characteristic of anger, we also observed the presence of AU1/AU2, which could be associated with fear (e.g., AU1+AU2+ AU4 and/or AU7). ^ΔIndicates the existence of a single ambiguous case where we also observed the presence of AU20, which could be associated with fear (HM2: AU4+AU17+AU20+AU24; HM5: AU4+AU7+AU20+AU24; H3: AU4+AU20+AU23; P2:

1
2
3 AU4+AU17+AU20+AU23). †Indicates the existence of a single ambiguous case where we also
4
5 observed the presence of AU15, which could be associated with sadness (HM5:
6
7 AU4+AU7+AU15; P3: AU4+AU7+AU15+AU17).
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

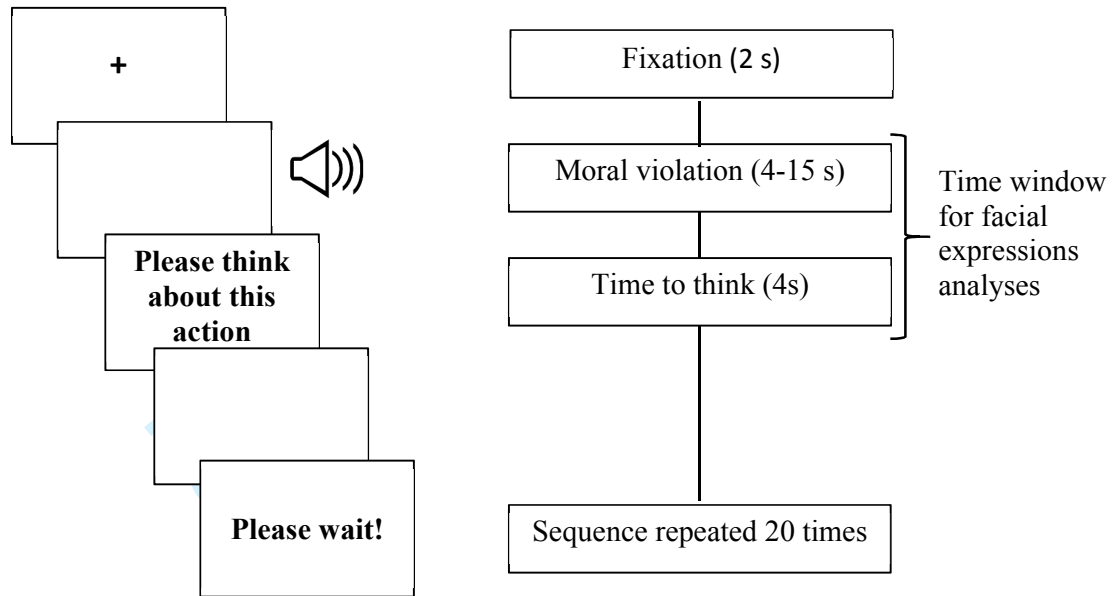


Figure 1. Schema of the experimental procedure of Task 1.

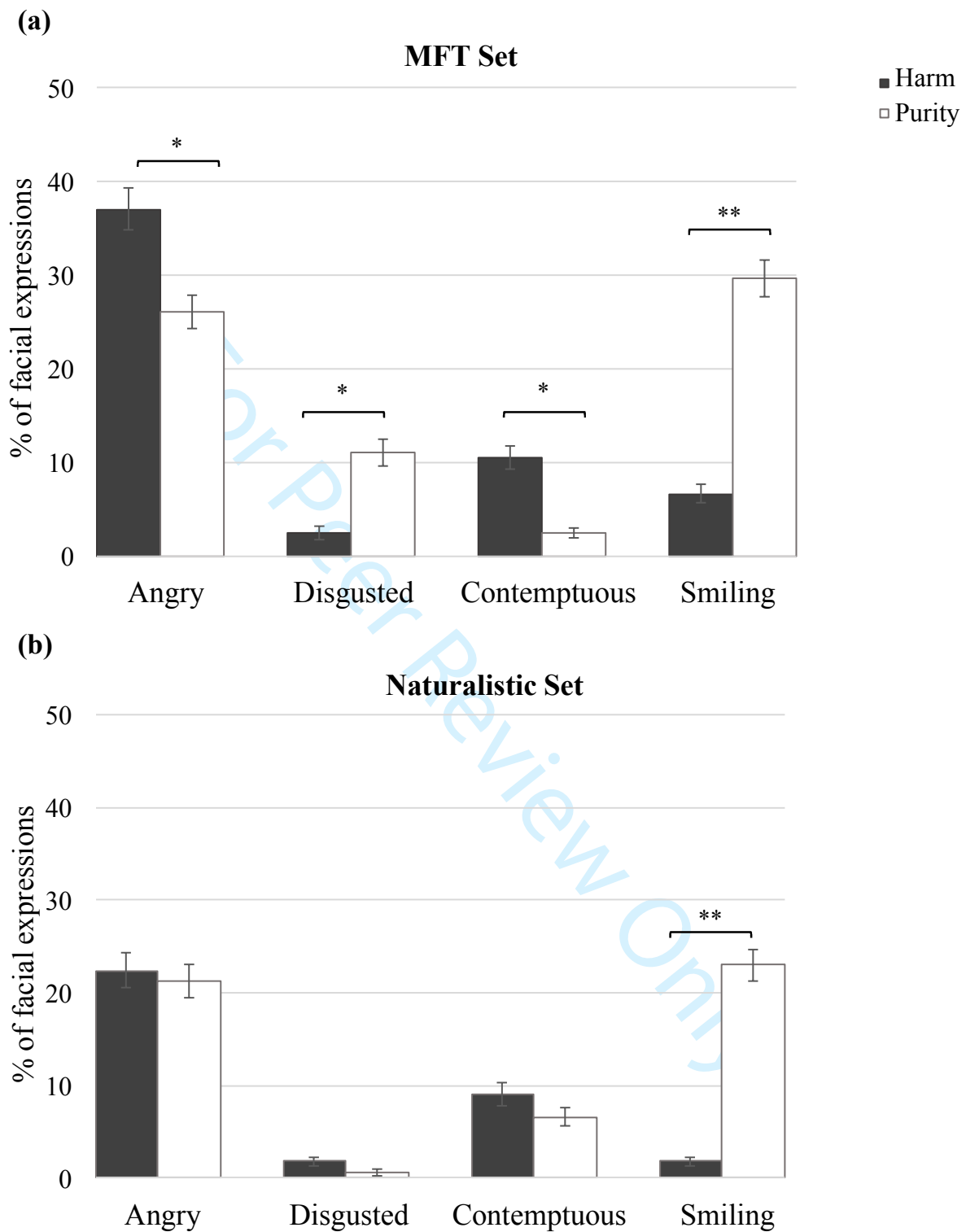


Figure 2. Mean percentage of facial expressions in response to the harm and purity scenarios, separately for the MFT set (a) and the naturalistic set (b). The surprise expression was not included due to its low frequency of exhibition. Error bars indicate standard error. * $p < .005$; ** $p \leq .001$

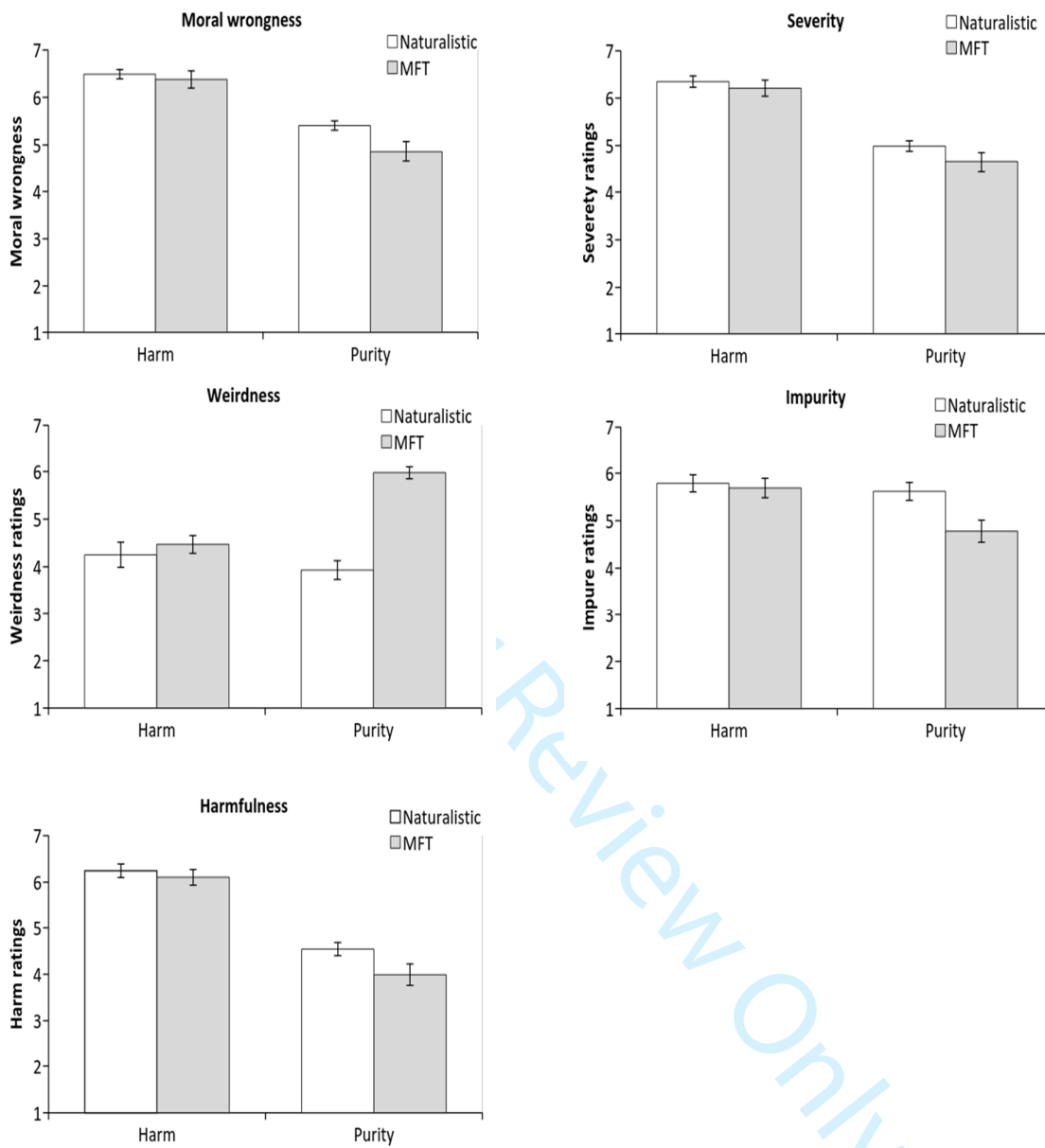


Figure 3. Moral wrongness, severity, weirdness, impurity, and harmfulness ratings by content (harm vs. purity) and source (MFT vs. naturalistic). Error bars indicate standard error of the means.

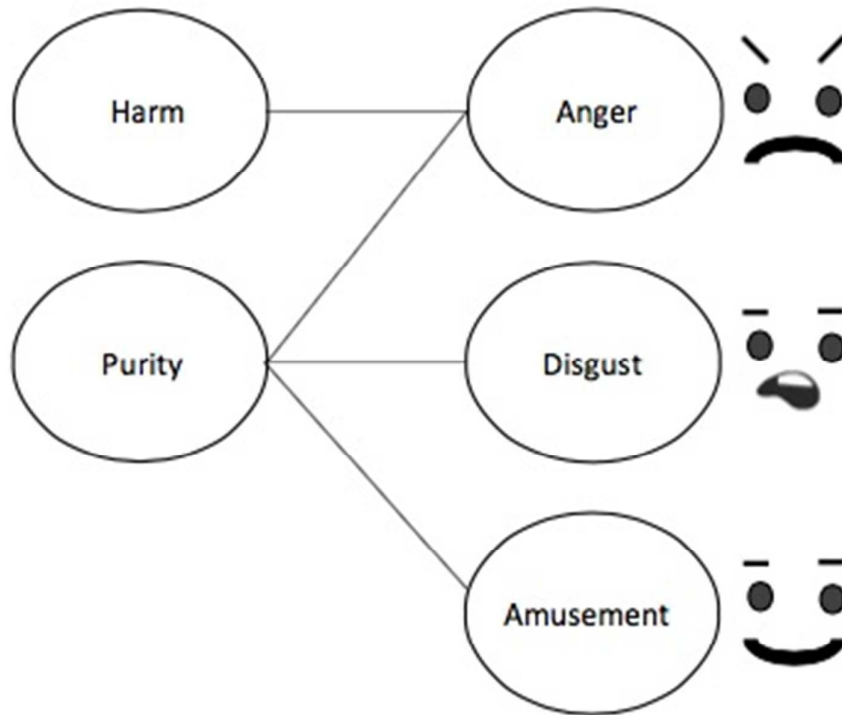


Figure 4. Semi-exclusive morality-emotion correspondences model (an adaptation of the model proposed by Cameron et al. 2015, p.8). The lack of lines between harm and disgust or amusement means non-significant relationships.

Supplementary Information

Study on self-reported emotions

This study used the same materials as the main study, but here we focused on self-reported emotions rather than on facial expressions. This study aimed to address two issues related to the main study. First, facial expressions are not a perfect proxy for emotions because emotions are not always accompanied by facial expressions. As stated in the paper, this is particularly true for the two emotions of interest: anger and disgust (e.g., Durán et al., 2017). Therefore, we wanted to investigate whether we would obtain similar results with regards to the strong and weak MFT hypotheses with self-reported emotions. We expected that neither of these hypotheses would be fully supported. Second, we wanted to disambiguate the meaning of smiles (enjoyment versus non-enjoyment), which were observed in response to scenarios depicting purity violations. To this end, we included various emotions that could underlie a smile: enjoyment, embarrassment, shame, and grin-and-bear-it.

Participants

Thirty-four native-Italian speakers (22 females, 12 males, $M_{age} = 22$, age range = 19–28) participated in the study. They were recruited via an announcement through the University of Trento mailing list.

Materials and Procedure

Participants were asked to complete a paper-and-pencil questionnaire. The questionnaire contained the 20 statements which were used in the main study: five involved MFT harm violations, five MFT purity violations, five naturalistic harm violations, and five naturalistic purity violations (for the full list of statements, see Table 1). The statements were presented in a pseudorandomized order. Following each statement, participants were asked to respond to two questions: (1) “While you were reading about the situation, what was your first emotion?” (*enjoyment, grin-and-bear-it, embarrassment, shame, disgust, anger, contempt, surprise*), and

1
2
3 then (2) “Please express a judgment on the intensity of this emotion” (1 = *very weak*, to 9 = *very*
4 *strong*). The entire experiment lasted about 10 minutes.

7 **Results and Discussion**

9 Figure S1 illustrates the percentage of times participants selected a specific emotion in
10 response to each class of items (harm versus purity) and set of scenarios (MFT versus
11 naturalistic). As it can be seen in Figure S1, for harm items, across both MFT and naturalistic
12 scenarios, anger was the most frequent response, followed by contempt. For purity items, no
13 single emotion prevailed. Specifically, for MFT purity items the most frequent response was
14 disgust, but it was closely followed by surprise, enjoyment, and grin-and-bear-it. For naturalistic
15 purity scenarios, the most frequent response was contempt, which was closely followed by grin-
16 and-bear-it and embarrassment. Turning to emotions that could underlie smiles, MFT purity
17 scenarios elicited an equal amount of positive and negative emotions (enjoyment and grin-and-
18 bear-it), whereas naturalistic purity scenarios predominantly elicited negative emotions (grin-and-
19 bear-it and embarrassment).
20
21
22
23
24
25
26
27
28
29
30
31
32

33 Table S1 shows the frequency and intensity of emotional responses to each of the 20
34 items. With respect to harm items, for both MFT and naturalistic scenarios, the predominant
35 response was anger. There was a single exception, item HM4, for which the prevalent response
36 was grin-and-bear-it. With respect to purity items, disgust was the prevalent response for just
37 two MFT items, PM2 and PM5. Other prevalent responses included enjoyment and surprise
38 (PM1, PM3), and grin-and-bear-it (PM4). For naturalistic items, prevalent responses were grin-
39 and-bear-it (P2, P3), contempt (P1), anger (P4), and embarrassment (P5).
40
41
42
43
44
45
46
47

48 In sum, just like in the main study, the results from the MFT scenarios provide support
49 for the weak MFT hypothesis. Specifically, anger was more prevalent than disgust in response
50 to harm scenarios, whereas disgust was more prevalent than anger in response to purity
51 scenarios. This result also replicates previous findings (Rozin et al., 1999). However, just like
52 in the main study, the results from the naturalistic scenarios conflict with weak MFT
53
54
55
56
57
58
59

1
2
3 hypothesis. Although anger was more prevalent than disgust in response to harm scenarios,
4
5 disgust was equally infrequent in response to both harm and purity scenarios. In sum, the
6
7 results do not offer support for the MFT hypothesis.
8

9
10 However, the results of the two studies also show some differences (compare Figure
11
12 S1 with Figure 2). For example, the differences in the frequency of anger between harm and
13
14 purity scenarios is much more pronounced in the follow up than in the main study. More
15
16 generally, the results of the purity scenarios are more evenly distributed among the eight
17
18 emotions. This finding is consistent with the hypothesis that purity violations elicit more than
19
20 one emotion. Specifically, it could be a consequence of people's difficulty to remember which
21
22 emotion came first or to correctly label that emotion. Indeed, three participants mentioned that
23
24 they felt a mix of emotions in response to certain scenarios, while two others mentioned that
25
26 they were uncertain about the difference between shame and contempt. In our opinion, the
27
28 results of the main study are 'cleaner' because participants were unaware about the purpose of
29
30 the study, and because the measure used was more direct (self-reports are more subject to
31
32 meta-judgments).
33
34

35
36 It is noteworthy to mention that the present analyses lend support to Gray and
37
38 Keeney's (2015) claim that the MFT purity items are weird (weirder than the MFT harm items,
39
40 and weirder than the naturalistic harm and purity items). First, surprise was a frequent response
41
42 to MFT purity items, almost as frequent as disgust. Indeed, for the naturalistic purity items in
43
44 which surprise responses were rare, disgust responses were also rare, and the results no longer
45
46 supported weak MFT. Second, MFT purity items elicited more enjoyment than their
47
48 naturalistic counterparts. Together, the increased surprise and enjoyment associated with MFT
49
50 purity items are consistent with the idea that these items are weird.
51
52
53
54
55
56
57
58
59
60

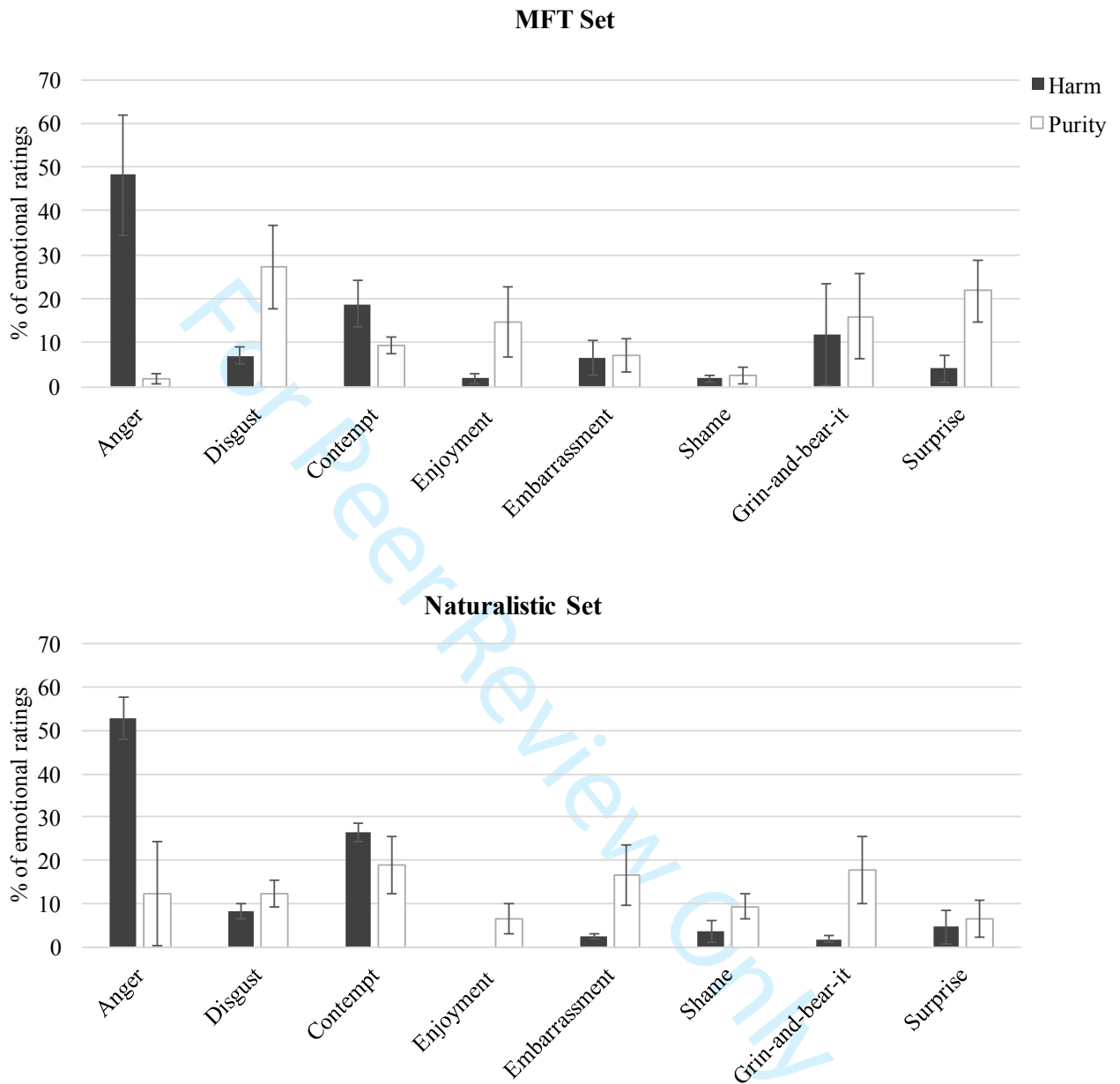


Figure S1. Mean percentage of times an emotion was chosen as a response response to the harm and purity scenarios of the MFT set (a) and the naturalistic set (b). Error bars indicate standard deviation.

Table S1

Frequency of selection of each emotion (in brackets, mean intensity rating) in correspondence with the Moral Foundations Theory (MFT) and the naturalistic scenarios.

		Harm							
	Item	Anger	Disgust	Contempt	Enjoyment	Embarrassment	Shame	Grin-and-bear-it	Surprise
MFT scenarios	HM1	25(7.8)	4(7.5)	4(6.5)	0	0	0	0	1(9)
	HM2	26(7.8)	1(8)	4(6.5)	0	1(4)	1(9)	1(4)	0
	HM3	16(7.4)	1(8)	8(7.5)	0	7(6.7)	1(9)	1(5)	0
	HM4	4(7.5)	3(6.3)	4(5.8)	2(5)	1(5)	1(2)	18(4.8)	1(1)
	HM5	11(7.5)	3(8.3)	12(6.4)	1(5)	2(4)	0	0	5(6.6)
	Mean	16.4(7.6)	2.4(7.6)	6.4(6.5)	0.6(5)	2.2(4.9)	0.6(6.7)	4(4.6)	1.4(5.5)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

	SD	9.3(0.2)	1.3(0.8)	3.6(0.6)	0.9 (0)	2.8(1.3)	0.5(4)	7.8(0.5)	2.1(4.1)
Naturalistic scenarios	H1	17(7.1)	2(6.5)	8(7.1)	0	1(5)	0	0	6(6.3)
	H2	15(6.9)	3(7.7)	10(6.6)	0	1(9)	4(6.5)	1(3)	0
	H3	21(7.3)	4(8.3)	8(7.1)	0	1(6)	0	0	0
	H4	15(6.6)	4(6.8)	11(6.4)	0	1(5)	2(5.5)	1(2)	0
	H5	22(7.8)	1(9)	8(8.1)	0	0	0	1(2)	2(7)
	Mean	18(7)	2.8(7.6)	9(7.1)	0(0)	0.8(6.3)	1.2(6)	0.6(2.3)	1.6(6.7)
	SD	3.3(0.5)	1.3(1)	1.4(0.7)	0(0)	0.4(1.9)	1.8(0.7)	0.5(0.6)	2.6(0.5)

		Purity								
		Item	Anger	Disgust	Contempt	Enjoyment	Embarrassment	Shame	Grin-and-bear-it	Surprise
MFT scenarios	PM1	0	4(5.7)	2(6.5)	12(5.9)	2(5.5)	0	1(5)	13(6.2)	
	PM2	1(9)	20(7.6)	4(6.3)	0	1(8)	1(9)	3(6.3)	4(7.8)	
	PM3	0	5(4.8)	2(7.5)	9(6.6)	2(5)	0	4(5.3)	12(6.1)	
	PM4	2(6.5)	7(6.6)	3(6.3)	0	0	0	17(5.5)	5(6.6)	
	PM5	0	10(7.6)	5(6.2)	4(7)	7(7.4)	3(8)	2(5.5)	3(7)	
	Mean	0.6(7.8)	9.2(6.5)	3.2(6.3)	5(6.5)	2.4(6.5)	0.8(8.5)	5.4(5.5)	7.4(6.7)	
	SD	0.9(1.8)	6.5(1.2)	1.3(0.1)	5.4(0.5)	2.7(1.5)	1.3(0.7)	6.6(0.5)	4.7(0.7)	
Naturalistic scenarios	P1	0	4(7.2)	13(7.7)	3(5)	5(5.2)	3(7.7)	4(3.5)	2(4.5)	
	P2	0	6(6.7)	8(6.4)	1(7)	4(5)	2(8.5)	13(4)	0	
	P3	1(8)	4(7)	4(5.8)	1(4)	6(5.5)	6(6)	10(5.9)	2(4)	
	P4	20(8.5)	6(7.8)	6(8.3)	0	0	1(8)	1(6)	0	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

P5	0	1(3)	1(5)	6(6.3)	13(7.3)	4(7)	2(6.5)	7(6.7)
Mean	4.2(8.2)	4.2(6.4)	6.4(6.6)	2.2(5.6)	5.6(5.8)	3.2(7.4)	6(5.2)	2.2(5.1)
SD	8.8(0.3)	2(1.9)	4.5(1.4)	2.4(1.3)	4.7(1.1)	1.9(1)	5.2(1.3)	2.8(1.4)

For Peer Review Only