

Running head: SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Shared neural representations of cognitive conflict and negative affect in the medial frontal cortex

Luc Vermeulen^{1*}, David Wisniewski¹, Carlos González-García¹, Vincent Hoofs¹, Wim Notebaert¹, Senne Braem^{1,2}

¹Department of Experimental Psychology, Ghent University, Belgium

²Department of Experimental and Applied Psychology, Vrije Universiteit Brussel, Belgium

*Correspondence

Luc Vermeulen

Department of Experimental Psychology

Henri Dunantlaan 2

B – 9000 Ghent

BELGIUM

Email: Luc.Vermeulen@ugent.be

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

20 **Abstract**

21 Influential theories of medial frontal cortex (MFC) function suggest that the MFC registers
22 cognitive conflict as an aversive signal, but no study directly tested this idea. Instead, recent
23 studies suggested that non-overlapping regions in the MFC process conflict and affect. In this
24 pre-registered human fMRI study, we used multivariate pattern analyses to identify which
25 regions respond similarly to conflict and aversive signals. The results reveal that, of all conflict-
26 and value-related regions, the ventral pre-supplementary motor area (often referred to as the
27 dorsal anterior cingulate cortex) showed a shared neural pattern response to different conflict and
28 affect tasks. These findings challenge recent conclusions that conflict and affect are processed
29 independently, and provide support for integrative views of MFC function.

30

31

32

33

34

35

36

37

38

39

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

40 **Introduction**

41 The medial frontal cortex (MFC), and dorsal anterior cingulate cortex (dACC) in particular, have
42 been implicated in various psychological processes such as cognitive control, somatic pain,
43 emotion regulation, reward learning and decision making (Ebitz & Hayden, 2016; Heilbronner &
44 Hayden, 2016; Shackman et al., 2011). In the domain of cognitive control, the dACC is
45 consistently activated by cognitive conflict, that is, the simultaneous activation of mutually
46 incompatible stimulus, task, or response representations (Botvinick et al., 2001). However, other
47 studies also showed dACC's involvement during the evaluation of negative outcomes, such as
48 reductions in reward (Gehring & Willoughby, 2002), negative feedback (Nieuwenhuis et al.,
49 2004) and pain (Rainville, 2002). Therefore, more integrative accounts of the dACC started to
50 redescribe its role in conflict monitoring as detecting a domain-general aversive learning signal
51 that can bias behavior away from the source of conflict (Botvinick, 2007; Shackman et al., 2011;
52 Shenhav et al., 2013, 2016). In other words, the dACC is thought to register “cognitive” conflict
53 as an aversive event. In accordance with this idea, recent studies have supported the presence of
54 a behavioral bias to avoid conflict, and have also shown that humans automatically evaluate
55 conflict as negative (Dignath et al., 2020; Dreisbach & Fischer, 2015; Inzlicht et al., 2015).

56 If conflict is registered as an aversive event in the dACC, one intriguing possibility is that
57 conflict and negative affect are encoded similarly in dACC (“shared or overlapping
58 representations”, Kragel et al., 2018). This conjecture also relates to a broader debate on the
59 functional organization of the MFC in regard to the psychological domains of cognitive control,
60 negative affect and pain (Inzlicht et al., 2015; Lieberman & Eisenberger, 2015; Shackman et al.,
61 2011). Towards the end of the last century, cognitive neuroscientists argued for the functional
62 segregation of affect (ventral part of the ACC) and cognitive control (dorsal part of the ACC) in

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

63 the MFC (Bush et al., 2000; Devinsky et al., 1995). In contrast, an influential meta-analysis by
64 Shackman and colleagues (2011) seemed to contradict this idea by showing substantial overlap
65 in the dorsal part of the ACC for the three-way conjunction of negative affect, conflict and pain
66 (Shackman et al., 2011). Similarly, one previous study tried to investigate the overlap in
67 activation between cognitive conflict and negative affect by using a repetition suppression
68 procedure, and found that dACC showed an attenuated response to negative affect following
69 cognitive conflict (Braem et al., 2017).

70 In more recent years, however, other studies failed to provide evidence for such functional
71 integration. For example, a number of recent studies and meta-analyses demonstrated that
72 distinct rather than overlapping parts of the MFC are associated with cognitive conflict and pain
73 processing (De La Vega et al., 2016; Jahn et al., 2016; Lieberman et al., 2016). Similarly, one
74 recent mega-analysis study reported a multivariate pattern analysis (MVPA) on full activation
75 maps from 18 studies to assess the similarity of patterns evoked by different domains. This study
76 also did not observe overlap between the activation patterns evoked by cognitive control, pain,
77 and negative emotion in the MFC (Kragel et al., 2018). Taken together, current studies seem to
78 be at odds with integrative views of MFC (Botvinick, 2007; Brown & Alexander, 2017; Calhoun
79 & Hayden, 2015; Heilbronner & Hayden, 2016; Shackman et al., 2011; Shenhav et al., 2013,
80 2016), which aim to explain the various responses of dACC by one underlying process (e.g.,
81 avoidance learning, value estimation, surprise processing). However, these studies relied on
82 group-averaged (often univariate) activation differences originating from distinct paradigms. For
83 example, the mega-analysis by Kragel and colleagues (2018) used a context-insensitive
84 “cognitive control” signal (i.e., across working memory, response inhibition, and conflict
85 processing tasks) considering paradigms from multiple studies that differ in experimental

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

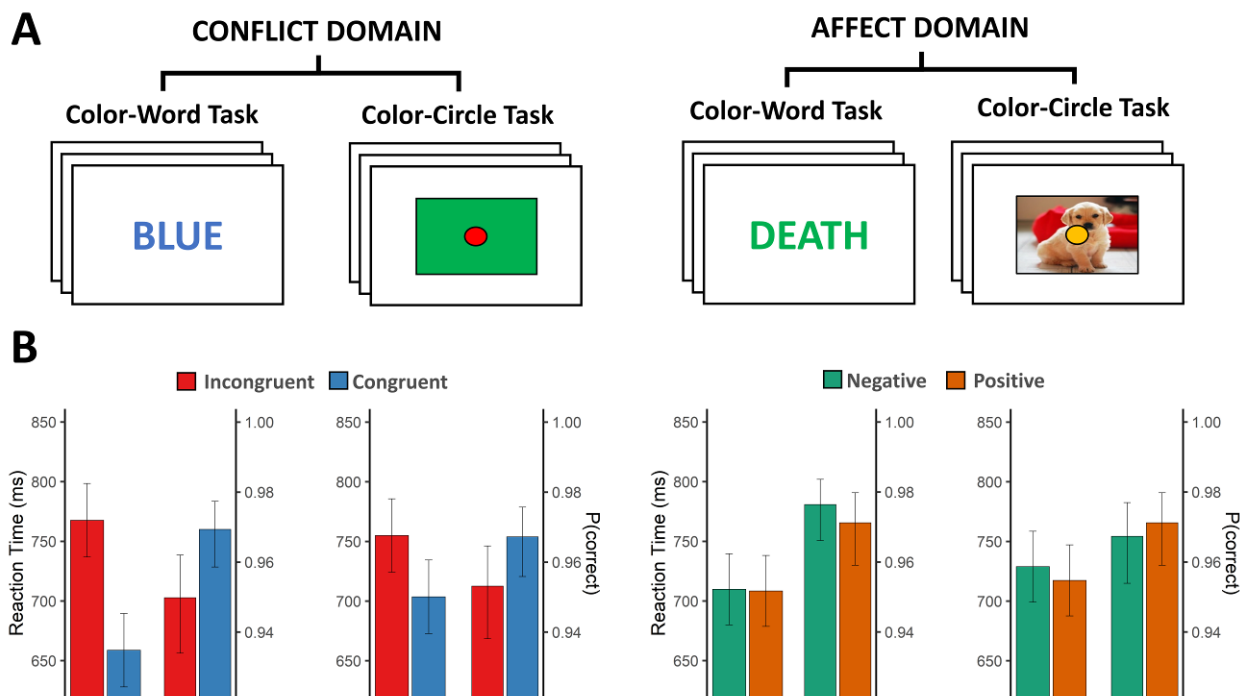
86 control. Moreover, many of these previous studies made use of intense pain responses that could
87 mask similarities with the arguably subtler affective evaluation of cognitive conflict.

88 Here, we took a different approach and developed a more targeted and well-controlled within-
89 subjects test of shared neural representations of conflict and negative affect. Namely, by using
90 multivariate cross-classification analyses, we assessed whether and where a classifier algorithm
91 trained to discern conflict (incongruent vs congruent events) can successfully predict affect
92 (negative vs positive events), and vice versa. Successful classification (i.e., classification above
93 chance) would be indicative of a similarity between the neural pattern response, and thus a
94 shared representational code between these two domains (Kaplan et al., 2015; Wisniewski,
95 2018). To test our predictions, 38 human subjects performed a color Stroop (Stroop, 1935) and
96 flanker task (Eriksen & Eriksen, 1974) in the conflict domain, and two closely matched tasks in
97 the affective domain (Fig. 1A). Importantly, we used two tasks in each domain in order to first
98 demonstrate an abstract or generalizable representation of conflict (and affect), that is
99 independent of conflict type (and affect source) (Jiang & Egner, 2014). Conflict and affect-
100 related brain signals were modeled using run-wise beta images retrieved from single subject
101 (first-level) analysis and were then used to perform a five-fold leave-one-run-out cross-
102 classification analysis using a linear Support Vector Machine (see Methods). We performed
103 preregistered Region of Interest (ROI, Fig. 2D) and whole brain searchlight analyses
104 (Supplementary Table 1), and report accuracy-minus-chance values for each ROI or searchlight
105 sphere (ROIs: Amygdala, Anterior Cingulate Cortex [ACC], dACC/pre-SMA, Anterior Insula
106 [AI], Posterior Cingulate Cortex [PCC], Ventral Striatum [VS], and the ventromedial Prefrontal
107 Cortex [vmPFC]). To construct the ROIs, we entered the name of the anatomical area as a search

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

108 terms in Neurosynth (Yarkoni et al., 2011) and constructed 10mm spheres around the peak z-
109 value of the corresponding meta-analytical activation map (see Methods).

110



111

112 **Figure 1.** Task Design and Behavioral Data. **(A)** Task design. Subjects either judged the color of
113 words or the color of circles. In the conflict domain, the color either matched or mismatched with
114 word meaning or background color creating congruent or incongruent conditions, respectively.
115 In the affective domain, positive or negative words and pictures were used to create the
116 respective conditions. Note that regardless of the domain, subjects always had to judge the color
117 of the word or the circle. These four task contexts were presented block-wise (with order
118 counterbalanced) in each run. **(B)** Reaction time and accuracy for the corresponding four task
119 contexts. Error bars denote the 95% CI.

120

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

121 **Results**

122 We first analyzed the behavioral data to check whether the conflict task showed the typical
123 congruency effects, and subjects processed the affective stimuli in the affective task. The
124 behavioral data from the conflict tasks (Fig. 1B, left panels and Supplementary Table 3) showed
125 typical congruency effects in reaction time (RT, $F_{(1,37)}=149.81$, $P<.001$, $BF_{10}>100$) and
126 accuracy ($F_{(1,37)}=11.72$, $P=.002$, $BF=52$), which were larger in the color-word task ($M=109$ ms,
127 $T_{(68)}=13.39$, $P<.001$) relative to the color-circle task ($M=51$ ms, $T_{(68)}=6.31$, $P<.001$) for the RT
128 measure (interaction between task and congruency effect: $F_{(1,37)}=35.55$, $P<.001$, $BF>100$). In the
129 affective tasks, we found slower RTs with negative relative to positive background stimuli
130 ($F_{(1,37)}=5.74$, $P=.025$, $BF=0.38$), but this was only the case for the color-circle task ($M=12$ ms,
131 $T_{(73)}=3.12$, $P=.003$) and not the color-word task ($M=1$ ms, $T_{(73)}=0.37$, $P=.710$) (interaction
132 between task and valence effect: $F_{(1,37)}=4.27$, $P=.046$, $BF=0.37$) (Fig. 1B, right panels and
133 Supplementary Table 4). No effects on accuracy were found in the affective tasks ($F_s<3.12$,
134 $P_s>.085$). In the conflict tasks, catch trials were used to draw attention to the conflicting nature
135 of the stimuli. Here, subjects had to judge the irrelevant rather than the relevant dimension (see
136 Method). We observed above-chance catch trial performance (chance level = 25 %), which did
137 not differ between the two conflict tasks, ($\chi^2_{(1)}=0.10$, $P=.755$, $BF=0.12$; color-circle task: 90.7 %;
138 color-word task 89.9 %). In the affective tasks, catch trials (where subjects had to make a
139 valence judgement instead of a color judgement) and a post-experiment incidental memory test
140 were used to inform processing of the (task-irrelevant) affective stimuli. We observed above-
141 chance catch trial performance (chance level = 50%), which did not differ between the two affect
142 tasks, ($\chi^2_{(1)}=0.19$, $P=.664$, $BF=0.12$; color-circle task: 86.4 %; color-word task 87.8 %), and

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

143 successful post-experiment recognition of the affective stimuli (Supplementary Figure 5),
144 ensuring that subjects processed the affective pictures.

145 In a first set of multivariate pattern analyses, we trained and tested a classifier within-task (within
146 the Stroop or flanker task; Fig. 2A, left panels; which regions respond to conflict within tasks?),
147 in each of our preregistered ROIs (for analysis details, see Method). Within-task ROI analyses in
148 the conflict domain (congruent vs. incongruent) revealed evidence for above chance-level
149 decoding in the ACC (*Wilcoxon* $V=388$, $P=.003$, $BF=15.61$), dACC/pre-SMA ($V=327$, $P=.009$,
150 $BF=8.48$) and PCC ($V=346$, $P=.009$, $BF=4.57$) but not in in any of the other regions (all $P>.145$,
151 $BF<0.61$) (Fig. 2A, right panel). The decoding accuracies did not differ by task (ACC:
152 $F_{(1,37)}=0.01$, $P=.915$, $BF=0.18$, dACC/preSMA: $F_{(1,37)}=0.72$, $P=.400$, $BF=0.24$, PCC:
153 $F_{(1,37)}=0.51$, $P=.476$, $BF=0.22$).

154 A second set of multivariate pattern analyses evaluated whether we could also cross-classify
155 conflict signals cross-task (train and test on different tasks; which regions respond similarly to
156 conflict independent of specific task features? Fig. 2B, left panels). To our knowledge, no study
157 has shown a voxel pattern response to conflict that is independent of conflict task in the ACC or
158 pre-SMA (e.g., Jiang & Egnér, 2014). Here, our within-conflict cross-task ROI analyses revealed
159 above chance level conflict decoding across tasks in the dACC/pre-SMA ($V=283$, $P=.012$,
160 $BF=5.57$) and PCC ($V=328$, $P=.023$, $BF=3.67$) (Fig. 2B, right panel). Decoding accuracy did not
161 differ between cross-task combination for the dACC/pre-SMA ($F_{(1,37)}=0.89$, $P=.352$, $BF=0.26$),
162 but was larger in the Stroop to flanker decoding relative to the flanker to Stroop decoding for the
163 PCC ($F_{(1,37)}=4.35$, $P=.043$, $BF=1.21$). These results were also replicated in an overall decoding
164 approach where the classifier was trained and tested in the whole domain regardless of task
165 (resulting in more samples to train the classifier; Supplementary Fig. 1A). Within the affective

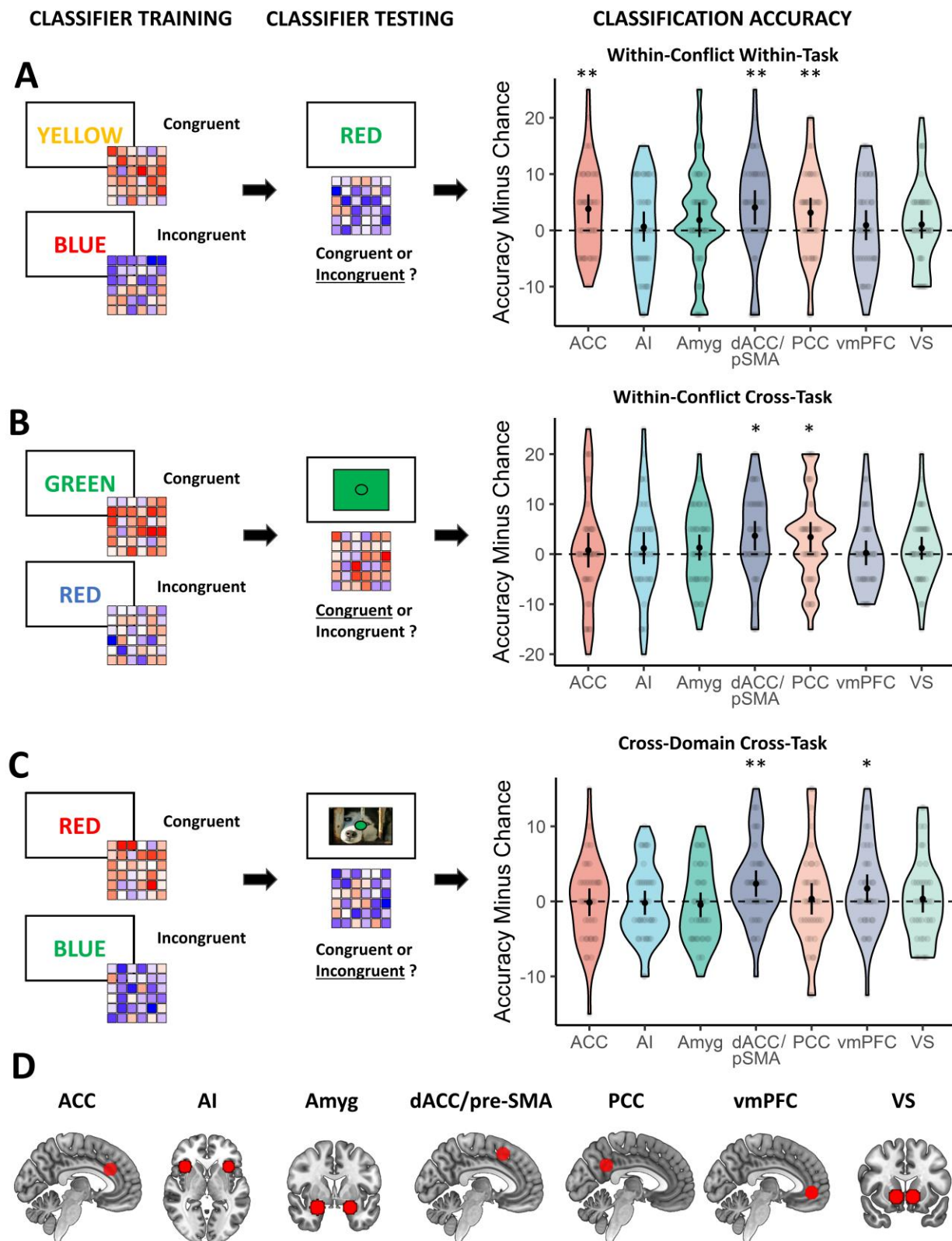
SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

166 domain (positive vs. negative), we also performed similar within- and cross-task decoding
167 analyses. However, while these analyses showed evidence for affect information in the AI and
168 vmPFC, they did not show evidence for decoding in the ACC, dACC/pre-SMA or PCC
169 (Supplementary Fig. 2; but see follow-up analyses below).

170 Finally, we investigated our main hypothesis by training a classifier on discerning conflict
171 (incongruent vs congruent) and testing its performance on discerning affect (negative vs
172 positive), and vice versa. For this analysis, we focused on the cross-domain cross-task decoding
173 (train and test in different domains on different tasks) as this analysis controls for low-level
174 features shared between the two tasks (Fig. 2C, left panel). The cross-domain cross-task ROI
175 decoding revealed evidence for cross-classification in the dACC/pre-SMA ($V=330$, $P=.007$,
176 $BF=8.43$; Fig. 2C, right panel) and vmPFC ($V=277$, $P=.045$, $BF=1.56$), which did not differ by
177 cross-task combination (dACC/pre-SMA: $F_{(1,37)}=0.36$, $P=.551$, $BF=0.20$, vmPFC: $F_{(1,37)}=0.02$,
178 $P=.880$, $BF=0.18$). None of the other ROIs reached significance (all $P_s > .444$, $BF_s < 0.24$). Using
179 the overall decoding approach (training on both tasks in one domain and testing on both tasks in
180 the other domain), we were only able to replicate successful cross-domain decoding in
181 dACC/pre-SMA ($V=449$, $P=.021$, $BF=4.65$; Supplementary Fig. 1C).

182

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT



SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

184 **Figure 2.** Main Results. **(A)** Training and testing the classifier within the conflict domain, within
185 the same task. **(B)** Training the classifier on one conflict task and testing its performance on
186 another conflict task. **(C)** Training the classifier to discern conflict and testing its performance on
187 classifying affect in another task (and vice versa). **(D)** ROIs: Anterior Cingulate Cortex (ACC),
188 Anterior Insula (AI), Amygdala (Amyg), dorsal Anterior Cingulate Cortex/pre-Supplementary
189 Motor Area (dACC/pre-SMA or dACC/pSMA), Posterior Cingulate Cortex (PCC), ventromedial
190 Prefrontal Cortex (vmPFC), Ventral Striatum (VS). * $P < .05$; ** $P < .01$; black dots and error bars
191 represent mean and ± 95 CI respectively; transparent dots represent individual data points; the
192 shape of the violin shows the distribution of the data.

193

194 In what follows, we report three additional sets of control analyses. A first set was designed to
195 evaluate the extent to which our analysis choices influenced our findings. A second set was
196 geared towards ruling out alternative hypotheses. Finally, a third set further zoomed in on the
197 above observation that we could not decode affect within- or across-tasks in the dACC/pre-SMA,
198 but did observe cross-domain decoding.

199 First, a number of control analyses using different analysis choices further confirmed our main
200 finding. For example, we replicated this result using different smoothing parameters
201 (Supplementary Fig. 3), or when using Harvard-Oxford atlas ROIs (Supplementary Fig. 4)
202 instead of ROIs retrieved from NeuroSynth. Moreover, also when using a set of functionally
203 (rather than anatomically) defined conflict-sensitive ROIs based on a recent meta-analysis (from
204 Chen et al., 2018, see their Table 2 for MNI coordinates), we again observed evidence for cross-
205 domain cross-task classification in the dACC/pre-SMA ($V=450$, $P=.013$, $BF=3.75$) but not for
206 other conflict-sensitive ROIs (left MOG, right AI, left AI, left IFG, left IPL, right IPL, left

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

207 MFG), except for the left AI ($V=425$, $P=.005$, $BF=8.61$). The result again replicated when using
208 the overall decoding approach in the dACC/pre-SMA ($V=449$, $p=.001$, $BF=41.06$), but not in the
209 left AI ($V=335$, $P=.260$, $BF=0.34$) (Supplementary Fig. 1, panel D).

210 Second, to further study whether the main effect was specific to the congruency identity and
211 valence of our experimental conditions, we also tested whether task difficulty differences or
212 differences in arousal could not explain our results. Notably, there was a small performance
213 difference in task performance on negative versus positive affect trials. Therefore, it is possible
214 that our cross-decoding result reflects the decoding of reaction time, rather than a shared conflict
215 and affect signal. However, the RT effects for the affective domain were rather small ($BF<1$) and
216 only present for one of the two affective tasks (also, note that there were no effects on accuracy).
217 Importantly, as reported above, the cross-domain decoding accuracy did not depend on which
218 task was used for testing. Moreover, cross-domain decoding accuracy was not higher when
219 splitting the sample for the subsample that did show the RT effect (RT effect >0 , $N=28$, $M=1.96$,
220 $P=.050$, $BF=1.66$) relative to a subsample that did not show the RT effect (RT effect <0 , $N=10$, M
221 $=3.50$, $P=.017$, $BF=4.58$) ($F_{(1,36)}=0.59$, $P=.444$). Nevertheless, in order to control for any
222 potential confounding effects of RT, we also ran another control analysis regressing out RT-
223 related effects from the neural data (via parametric modulation). This analysis led to the same
224 conclusions as above (cross-domain cross-task dACC/pre-SMA: $V=380$, $p=.036$, $BF=1.91$,
225 overall cross-domain dACC/pre-SMA: $V=348$, $p=.025$, $BF=2.61$). Note that this does make the
226 effect slightly smaller, which should not be surprising as this procedure also removes variance of
227 interest (as the effectiveness of the congruency manipulation is often defined by RT differences).
228 Next, we also found that the cross-domain decoding was unlikely to be due to arousal. First, we
229 were only able to decode arousal (high vs. low arousal; based on a median-split of arousal

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

230 ratings, orthogonal to valence) within-tasks in the ACC ($V=260$, $P=.039$, $BF=1.43$) and vmPFC
231 ($V=240$, $P=.048$, $BF=1.38$), but not in the dACC/pre-SMA ($V=230$, $P=.082$, $BF=0.83$) or other
232 ROIs ($P>.217$, $BF<0.39$). Second, training a classifier on conflict and testing on arousal (and
233 vice versa) did not show any evidence for cross-decoding in the dACC/pre-SMA ($V=294$,
234 $P=.404$, $BF=0.26$), nor any of the other ROIs ($P_s>.139$, $BF<0.51$), except for the vmPFC
235 ($V=401$, $P=.015$, $BF=3.25$).

236 Third, and finally, we wanted to follow up on the unexpected result that we did not observe
237 within and cross-task decoding of affect in the dACC/pre-SMA. One potential explanation is that
238 the SNR was lower for our affect differences than the congruency differences. In line with this,
239 while the univariate affect contrast (negative > positive) showed no cluster-corrected activation
240 in the MFC (Supplementary Table 2), a large cluster of dmPFC activity does become visible,
241 when we lower the threshold ($p < .005$, uncorrected). Similarly, when restricting the whole-brain
242 decoding analysis to the MFC using the same mask as Kragel and colleagues (2018), within-
243 affect decoding did reveal a large cluster in the MFC, overlapping with our main dACC/pre-
244 SMA ROI, as well as the within-conflict decoding and cross-domain decoding results (see
245 Supplementary Fig. 6). A second reason for observing weaker affect decoding signals, could be
246 because the affect signals weakened or habituated throughout the experiment, as the same
247 affective words and pictures were repeated in each run. Therefore, we also investigated whether
248 affect might have been easier to decode in the beginning of the experiment, by studying cross-
249 task affect for the first half (run 1 and 2) and second half (run 4 and 5), separately. Indeed, there
250 was significant cross-task affect decoding in main dACC/pre-SMA ROI in the first half of the
251 experiment (run 1 and 2, $V=237$, $P=.019$, $BF=2.38$), but not in the second half of the experiment
252 (run 4 and 5, $V=138$, $P=.646$, $BF=0.10$). Finally, it is worth noting that our ROIs were primarily

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

253 selected to detect conflict-based differences (see also our literature-based ROIs on conflict
254 decoding, Supplementary Fig. 1, panel D), which again potentially lowered our chances to
255 observe affect decoding. When using ROIs based on an affect processing meta-analysis (from
256 Lindquist et al., 2016, their Table 2, cf. our conflict-sensitive ROIs), a pre-SMA ROI close to our
257 original dACC/pre-SMA ROI (but a bit more dorsal) did show significant affect ($V=314$, $P=.001$,
258 $BF=39.70$), conflict ($V=308$, $P=.002$, $BF=27.19$) and cross-domain decoding ($V=375$, $P=.042$,
259 $BF=1.52$). Together, these results do suggest that there was affect decoding in, or at least around
260 our main dACC/pre-SMA ROI.

261 **Discussion**

262 Together, our results reveal that the dACC/pre-SMA shows a similar voxel pattern response to
263 conflict and negative affect. Moreover, to the best of our knowledge, our study is also the first to
264 show decoding of conflict across conflict tasks in the dACC/pre-SMA, suggesting a shared
265 component in the detection of conflict across the Stroop and flanker task (Jiang & Egner, 2014).

266 These findings fit with integrative accounts of MFC function which propose that multiple
267 domains (e.g., cognitive conflict, negative emotion and pain) activate a single underlying process
268 in the MFC (e.g., adaptive control, avoidance learning, cost-benefit value estimation), and thus
269 predict similar activation patterns across domains (Botvinick, 2007; Brown & Alexander, 2017;
270 Calhoun & Hayden, 2015; Heilbronner & Hayden, 2016; Shackman et al., 2011; Shenhav et al.,
271 2013, 2016). Since our study focused on the domains of cognitive conflict and negative affect,
272 our finding is especially relevant for theories that have tried to explain dACC's sensitivity to
273 conflict and negative outcomes by a single unifying process (Botvinick, 2007; Shenhav 2013,
274 2016). These theories have proposed the idea that conflict is in itself an aversive outcome
275 (Botvinick, 2007; Shackman et al., 2011; Shenhav et al., 2013), and that the main mechanism of

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

276 the dACC is to bias behavior away from any costly, demanding or suboptimal outcome
277 (“domain-general avoidance learning”, Botvinick, 2007) or to allocate control based on the
278 expected benefits discounted by the expected costs of controlling behavior (“cost-benefit value
279 estimation” Shenhav et al., 2013).

280 Importantly, our study was not set up to disentangle the multiple candidate processes as to why
281 conflict and negative affect would elicit a shared neural pattern response in the MFC, so while
282 our main finding (i.e., similar voxel pattern response to conflict and negative affect) follows
283 naturally from above-mentioned models that predict similar patterns of brain activity for
284 different domains, it does not favor a specific mechanism (e.g. domain-general avoidance
285 learning). Still, our finding does lend credence to the idea that that cognitive control can be
286 understood as an emotional process (Inzlicht et al., 2015), and that conflict can be registered as
287 an aversive event (i.e., Botvinick, 2007; Dignath et al., 2020; Dreisbach & Fischer, 2015), at
288 least in the dACC/pre-SMA. Similarly, additional analyses further show that our effect is
289 unlikely to be driven by general differences in reaction time or task difficulty, or differences in
290 arousal.

291 Other recent studies failed to find similarities, and suggest that theories of MFC function should
292 not look for a unitary neural implementation, but rather focus on a unified computational
293 mechanism (Kragel et al., 2018). Our study does not argue against this idea, as shared neural
294 “representations” in fMRI could still be driven by different local neighboring neural populations.
295 However, we believe our findings do license consideration of a unitary neural implementation,
296 and offer support for the idea of a unified computational mechanism. Still, one might wonder
297 why we did, and others did not, find shared neural representations of conflict and affect. First, we
298 employed a multivariate approach rather than a univariate approach which is more sensitive and

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

299 more informative to assess similarity in brain activation (without necessarily being orthogonal to
300 a univariate approach; Davis et al., 2014). Second, our design was specifically set up to evaluate
301 similarity between two specific domains (conflict and negative affect) rather than very general
302 domains (cognitive control, negative emotion, pain). Third, related to this, we used a within-
303 subjects design (rather than the between-subject nature of meta- or mega-analyses), which
304 allowed for more sensitive analyses and to make the different task contexts as similar as possible
305 (high level of experimental control).

306 One unexpected finding relates to the fact that we did not observe a similar (significant) above-
307 chance decoding of affect in the dACC/pre-SMA, but did observe cross-domain decoding.
308 Although this finding was surprising to us at first, it most likely suggests differences in signal to
309 noise ratio (SNR) between the two domains and does not invalidate the cross-domain decoding
310 result (for a review and discussion of similar findings, see van den Hurk & de Beeck, 2019). This
311 idea was also further corroborated by follow-up analyses which showed affect decoding in the
312 dACC/pre-SMA when using other analysis techniques, or focusing on the first half of the
313 experiment only. Moreover, a lower SNR in the affect domain can also be explained by the fact
314 that affect was not relevant for the main task (which allowed us to keep the affective tasks as
315 similar as possible to the conflict tasks).

316 Finally, while our analyses were based on theories and a NeuroSynth ROI of the dACC,
317 investigating the (uncorrected) whole-brain searchlight decoding maps did reveal that the
318 conflict, affect and cross-domain decoding all trigger a dorsomedial frontal area, that might be
319 more accurately referred to as pre-SMA rather than dACC. This area roughly corresponds to our
320 main dACC/pre-SMA ROI which was retrieved by searching “dACC” in NeuroSynth (Yarkoni
321 et al., 2011), and therefore, we decided to refer to this ROI as dACC/pre-SMA. Nevertheless,

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

322 this illustrates the consequent mislabeling of the dACC in the literature and we are not the first to
323 report pre-SMA activation for domains such as cognitive control, pain and negative emotion
324 (Brown & Alexander, 2017; De La Vega et al., 2016; Jahn et al., 2016; Kolling et al., 2016;
325 Kragel et al., 2018; Lieberman & Eisenberger, 2015; Lindquist et al., 2016).

326 In sum, by using a well-controlled within-subject design and multivariate analysis techniques, we
327 show that conflict and negative affect evoke a similar voxel pattern response in the MFC. Our
328 finding brings important nuance to other recent studies suggesting different neural activations or
329 representations, and helps to further constrain theories of integrative MFC function.

330 **Materials and Methods**

331 *Participants*

332 The study was pre-registered with the pre-registration template from AsPredicted.org on the
333 Open Science Framework (<https://osf.io/p5frq/>). As pre-registered, 40 participants participated in
334 our study. Two participants were excluded (one due to excessive head motion [$>2.5\text{mm}$
335 translation] and one aborted the scanning session). The average age of the remaining 38
336 participants (13 male) was 23.71 years ($SD=3.53$, $\text{min}=18$, $\text{max}=33$). Thirty-six participants
337 were right-handed, one was left-handed and one was ambidextrous (as assessed by the Edinburgh
338 Handedness Inventory (Oldfield, 1971)). Every participant had normal or corrected to normal
339 vision and reported no current or history of neurological, psychiatric or major medical disorder.
340 Every participant gave their informed written consent before the experiment, and was paid 35
341 euros for participating afterwards. The study was approved by the local ethics committee
342 (University Hospital Ghent University, Belgium).

343 *Experimental Paradigm*

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

344 The experiment was implemented using Psychopy 2 version 1.85.2 (Peirce, 2007). On each trial,
345 participants had to judge the color of a target stimulus in the center of the screen, using two MR-
346 compatible response boxes (each box had two buttons) to indicate one out of four possible
347 response options (red, blue, green and yellow). The key-to-color mapping was counterbalanced
348 between participants. The exact features of the target stimulus varied block-wise, depending on
349 one of four different task-contexts. Specifically, participants either had to respond to the color of
350 words (“color-word naming task”) or respond to the color of circles (“color-circle naming task”),
351 which both had a conflict and affective version.

352 The conflict-version of the color-word naming task was a Stroop task (Stroop, 1935), where the
353 meaning of the words could either be congruent or incongruent with the actual color of the word.
354 For example, participants could see the words “BLUE”, “RED”, “GREEN” or “YELLOW”
355 (Dutch: “ROOD”, “BLAUW”, “GROEN” or “GEEL”) presented in a blue, red, green or yellow
356 font. The conflict version of the color-circle naming task was essentially a color-based variant on
357 the Eriksen flanker task (Eriksen & Eriksen, 1974), where the irrelevant feature consisted of a
358 colored background square which could either be congruent or incongruent with the color of the
359 circle. Here, participants could see blue, red, green or yellow circles presented on a blue, red,
360 green or yellow background square. In both tasks, half of the trials were congruent (e.g., “RED”
361 in a red font; a red circle presented on a red square background) while the other half of the trials
362 were incongruent (e.g., “RED” in a blue font; a red circle on a blue square background).

363 The affect-versions of the color-word naming and color-circle naming tasks made use of
364 irrelevant affective words or pictures, respectively. In the color-word naming task, 16 positive
365 and 16 negative words were presented (Moors et al., 2013) that were matched on arousal, power,
366 age of acquisition, Dutch word frequency (Keuleers et al., 2010), word length and grammatical

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

367 category (Noun, Adjective and Verbs). The affective picture distractors in the background of the
368 color-circle naming task were retrieved from the OASIS database (Kurdi et al., 2017). Sixteen
369 positive and 16 negative pictures were presented that were matched on semantic category
370 (Animals, Objects, People, Scenery) and arousal. This resulted in a total of eight conditions:
371 congruent, incongruent, positive or negative trials, that either involved words or pictures/colored
372 backgrounds.

373 Each trial started with a fixation sign (“+”) that was presented for 3 to 6.5 seconds (in steps of
374 0.5 s; $M=3.5$ s; drawn from an exponential distribution). Next, the target stimulus was presented
375 for 1.5 seconds (fixed presentation time regardless of RT). In order to increase the saliency of the
376 irrelevant dimension (conflict and affect), the onset of the affective word or picture preceded the
377 presentation of the target feature by 200 ms during which the color of the target feature (word or
378 circle) was white.

379 Participants performed five scanning runs and during each run the subjects performed each of the
380 four task contexts in separate blocks. The order of the four blocked task contexts was fixed
381 within participant but counterbalanced between participants. Each block hosted 32 trials (16
382 congruent/positive and 16 incongruent/negative) which were presented in a pseudo-random
383 fashion with the following restriction: neither relevant nor irrelevant features of the target
384 stimulus could be repeated. This restriction was used to investigate confound-free congruency
385 sequence effects (see Braem et al., 2019; Schmidt, 2019; but this was not the aim of the current
386 study and will not be discussed further). In total, each participant made 640 trials (i.e., five runs
387 of four blocks of 32 trials).

388 In each task context (block), we also included one catch trial (at random, but not in the first two
389 or last two trials of each block). In these catch trials, the presentation of the task-irrelevant word,

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

390 picture, or colored square would not be followed by the presentation of the target color, and
391 remain on screen for three seconds. Participants were instructed that during these catch trials,
392 when no color information was present in the relevant dimension, their goal was to judge the
393 irrelevant dimension depending on the cognitive domain. In the conflict domain, participants had
394 to respond to the meaning of the word (“RED”, “BLUE”, “GREEN” or “YELLOW”) or to the
395 color of the background square (red, blue, green or yellow) by using the respective key that
396 would be used to judge the relevant dimension. In the affective domain, participants had to judge
397 the affective word or background picture as either positive or negative by pressing all keys once
398 or twice (response mapping for positive and negative stimuli counterbalanced between
399 participants). The purpose of these catch trials was to increase the saliency of the irrelevant
400 dimension.

401 Before the scanning session, participants were welcomed and instructed to read the informed
402 consent after which they started practicing the experimental paradigm. After the scanning
403 sessions, participants performed an unannounced recognition memory test on old and new
404 affective words and pictures. Here, participants had to indicate whether they had previously seen
405 the word or picture in the experiment (old/new judgement). The new words were matched with
406 the old words in terms of valence, arousal, power, age of acquisition, word length, frequency,
407 grammatical category. The new pictures were matched on valence, arousal and semantic
408 category. In both a behavioral ($n = 20$) and fMRI pilot ($n = 20$), we already established that
409 participants showed adequate performance on both the main task and the recognition memory
410 task. Finally, participants completed four questionnaires (Need for Cognition, Behavioral
411 Inhibition/Activation Scale, Positive and Negative Affect Schedule, Barret Impulsivity Scale)

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

412 and were thanked for their participation. No significant correlations between these questionnaire
413 scales and cross-classification accuracies were found, so we do not report these results.

414 *Behavioral Data Analysis*

415 Behavioral analyses were performed in R (RStudio version 1.1.463, www.rstudio.com). For the
416 reaction time (RT) analyses, we removed incorrect, premature (< 150 ms), and extreme
417 responses (RTs outside 3 SD from each condition mean for each participant). This resulted in an
418 average of 94.42 % of the trials left for the RT analyses ($SD=3.18$, min=84.22, max=98.28). We
419 conducted a repeated measures ANOVA on the reaction time and accuracy measure with the
420 within-subject factors Condition (conflict domain: congruent vs. incongruent, affective domain:
421 positive vs. negative) and Task (color-word naming vs. color-circle naming). We also assessed
422 post-scanning recognition memory of affective stimuli with a probit generalized linear mixed
423 effects model on the probability to say that the stimulus was ‘old’ with fixed effects for
424 Experience (old vs. new), Valence (positive vs. negative) and Task Type (word vs. picture) and
425 crossed random effects for Participant and Item. We also pre-registered some exclusion criteria
426 based on behavioral performance. Participants with a mean RT outside 3 SD from the sample
427 mean or a hit rate below 3 SD or 60 % (chance level=25 %) from the sample mean were
428 excluded. Participants that performed poorly on the post-scanning recognition memory test, i.e.,
429 hit rate or false alarm rate outside 3 SD of the sample mean were also excluded. In the end, no
430 exclusions based on task performance had to be made. While performance on catch trials was not
431 a pre-registered exclusion criterion, we found that two participants responded on chance level in
432 the catch trials of the affective domain (chance level=50 %, positive vs. negative judgement).
433 Excluding these participants did not change our conclusions.

434 *fMRI data acquisition*

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

435 fMRI data was collected using a 3T Magnetom Prisma MRI scanner system (Siemens Medical
436 Systems, Erlangen, Germany), with a sixty-four-channel radio-frequency head coil. A 3D high-
437 resolution anatomical image of the whole brain was acquired for co-registration and
438 normalization of the functional images, using a T1-weighted MPRAGE sequence (TR=2250 ms,
439 TE=4.18 ms, TI=900 ms, acquisition matrix=256 × 256, FOV=256 mm, flip angle=9°, voxel
440 size=1 × 1 × 1 mm). Furthermore, a field map was acquired for each participant, in order to
441 correct for magnetic field inhomogeneities (TR=520 ms, TE1=4.92 ms, TE2=7.38 ms, image
442 matrix=70 x 70, FOV=210 mm, flip angle=60°, slice thickness=3 mm, voxel size=3 x 3 x 2.5
443 mm, distance factor=0%, 50 slices). Whole brain functional images were collected using a T2*-
444 weighted EPI sequence (TR=1730 ms, TE=30 ms, image matrix=84 × 84, FOV=210 mm, flip
445 angle=66°, slice thickness=2.5 mm, voxel size=2.5 x 2.5 x 2.5 mm, distance factor=0%, 50
446 slices) with slice acceleration factor 2 (Simultaneous Multi-Slice acquisition). Slices were
447 orientated along the AC-PC line for each subject.

448 *fMRI data analysis*

449 fMRI data analysis was performed using Matlab (version R2016b 9.1.0, MathWorks) and
450 SPM12 (www.fil.ion.ucl.ac.uk/spm/software/spm12/). Raw data was imported according to
451 BIDS standards (<http://bids.neuroimaging.io/>) and functional data was subsequently realigned,
452 slice-time corrected, normalized (resampled voxel size 2 mm³) and smoothed (full-width at half
453 maximum of 8 mm). The preprocessed data was then entered into a first-level general linear
454 model analysis (GLM), and subsequently into a multivariate pattern analysis (MVPA(Cox &
455 Savoy, 2003; Haxby, 2012; Haynes, 2015; Kriegeskorte et al., 2006). Results were analyzed
456 using a mass-univariate approach. Although we pre-registered that we would not normalize and
457 smooth the data for our classification analyses, we found that Signal-to-Noise Ratio (SNR) was

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

458 significantly improved with these additional preprocessing steps (Supplementary Fig. 3A). In
459 addition, an independent classification analysis (classifying left vs. right responses in primary
460 motor cortex) showed that decoding accuracies were significantly higher with these additional
461 preprocessing steps (Supplementary Fig. 3B). Knowing that decoding information in the PFC is
462 notoriously difficult as decoding accuracies are close to chance (relative to decoding in
463 occipitotemporal cortex (Bhandari et al., 2018), and the finding that smoothing can and does
464 often improve SNR and decoding performance (Hendriks et al., 2017; Kamitani & Sawahata,
465 2010; Op de Beeck, 2010), we decided to optimize our MVPA analyses by decoding on
466 normalized and smoothed data. For completeness, however, we also depict the results from our
467 main cross-classification analysis for different levels of smoothing (FWHM 0, 4 and 8 mm; see
468 Supplementary Fig. 3C).

469 First-level GLM analyses consisted of 5 identically modeled sessions (i.e., the five runs). Each
470 session consists of eight regressors of interest (for the eight conditions, see above), four block
471 regressors (to account for the blocked presentation of each combination of word versus picture
472 versions of the conflict versus affect tasks), two nuisance regressors (that model performance
473 errors and catch trials) and six movement regressors. The regressors were convolved with the
474 canonical HRF. The modeled duration of the regressors of interest (the eight conditions) and
475 nuisance regressors (errors, catch trials) was zero, while the modeled duration of the block
476 regressors was equal to the length of the blocks.

477 Next, the beta images from the first-level GLM were submitted to leave-one-run-out decoding
478 scheme with ‘The Decoding Toolbox’ (Hebart et al., 2015) using a linear support-vector
479 classification algorithm ($C=1$). We performed whole-brain searchlight decoding (sphere radius: 3
480 voxels; Supplementary Table 1) as well as ROI decoding (see below for ROI methods). Cross-

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

481 validation decoding was conducted within the affective (positive vs. negative) and conflict
482 (congruent vs. incongruent) domain for each task separately (“within-domain within-task
483 classification”). To assess the generalizability of the classifier within the domain, we also
484 conducted cross-classification analyses where we trained the classifier on one task and tested its
485 performance on the other task for each task type combination (from color-circle naming to color-
486 word naming and vice versa) separately (“within-domain cross-task classification”). To
487 investigate the generalizability of these classifiers across the domain (our main hypothesis), we
488 trained the classifier in the conflict domain and tested its performance in the affective domain,
489 and vice versa. We conducted these analyses cross task type combinations (i.e., from color-circle
490 naming to color-word naming, or from color-word naming to color-circle naming) to further
491 control for low-level task features, following the same reasoning as the within-domain cross-task
492 classification analyses (“cross-domain cross-task classification”). For each of these three
493 decoding analyses, we also ran ANOVAs to evaluate whether the result differed depending on
494 the task (e.g., color-circle naming versus color-word naming) or task-to-task direction (i.e., from
495 color-circle naming to color-word naming, or from color-word naming to color-circle naming).
496 Finally, we also report an “overall decoding” analysis, where the classifier was trained across the
497 two task types at once, thereby ignoring whether the event featured words or pictures/colored
498 backgrounds.

499 Each classification analysis resulted in ‘accuracy-minus-chance’ decoding maps for each subject.
500 These maps were then entered into a group second-level GLM analysis in SPM12. Here, a one-
501 sample t-test determined which voxels show significant accuracy above chance level.

502 Next to MVPA, we also conducted classic univariate analyses. Here, we constructed a set of
503 contrasts subtracting (A) positive from negative conditions and (B) congruent from incongruent

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

504 conditions for (1) each task separately as well as across both tasks. These contrast images were
505 then entered into a second-level analysis in which a one-sample t-test determined which voxels
506 show significant activation for each contrast. We applied a statistical threshold of $p < 0.001$
507 (uncorrected) at the voxel level, and $p < 0.05$ (family-wise error corrected) at the cluster level on
508 all analyses (Supplementary Table 2).

509 *ROI analyses*

510 As part of our pre-registered main analysis plan, we conducted ROI decoding analyses. We set
511 out to study the Amygdala, Anterior Cingulate Cortex (ACC), dorsal Anterior Cingulate
512 Cortex/pre-SMA (dACC/pre-SMA), Anterior Insula (AI), Parietal Cingulate Cortex (PCC),
513 Ventral Striatum (VS), and the ventromedial PFC (vmPFC). Our preregistration noted that all
514 ROIs would be obtained from the Harvard-Oxford cortical and subcortical structural atlases,
515 thresholded at 25%. However, as the dACC ROI was not defined in the Harvard-Oxford atlas,
516 we decided to retrieve this ROI from Neurosynth (Yarkoni et al., 2011) by entering “dacc” as
517 search term (returning 162 studies reporting 4547 activations). Although this ROI was based on
518 the “dacc” search term, the peak effect of studies reporting dACC activity actually lies more
519 dorsally than the cingulate gyrus, overlapping with the pre-SMA (Lieberman & Eisenberger,
520 2015). Therefore, we refer to this ROI as the dACC/pre-SMA. Next, we built a 10 mm sphere
521 around the peak activation point in this activation map (association map). Because the dACC
522 ROI was spherical (in contrast to the other six atlas ROIs), we decided to retrieve all ROIs from
523 NeuroSynth for comparability. However, because our preregistration mentioned that the ROIs
524 would be obtained from the Harvard-Oxford cortical and subcortical structural atlases, we also
525 report the analyses where all non-dACC regions were evaluated using these ROIs, which
526 returned similar results and did not change our main conclusions (see Supplementary Figure 4).

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

527 In addition to the pre-registered ROI analyses which were based on anatomically determined
528 ROIs, we also ran another set of ROI analyses with functionally informed ROIs. Namely, we
529 created 10 mm sphere ROIs for all conflict-sensitive regions based on the most recent and
530 inclusive meta-analysis we could find on cognitive conflict (Chen et al., 2018, Table 2).

531 Each ROI decoding analysis returned one accuracy-minus-chance value per ROI and participant.
532 We tested whether these values were significantly higher than zero (one-tailed) with the non-
533 parametric Wilcoxon signed-rank test and a Bayesian t-test (using the default priors from the
534 BayesFactor package in R; Cauchy prior width: $r=.707$). We report the Bayes Factor (BF) that
535 quantifies the evidence for the alternative hypothesis (i.e., decoding accuracy is higher than
536 zero). Our pre-registered stopping criterion was if the main finding was $BF>6$ (i.e., or if we had
537 reached 40 subjects, for financial reasons), but we would like to note that, if so, this result was
538 typically also $p<.00714$, which is the Bonferroni-corrected alpha for the main set of 7 ROIs.

539 Finally, we investigated whether the significant cross-task cross-domain classification accuracy
540 correlated with the following behavioral indices: post-scanning affective recognition memory (d-
541 prime), congruency sequence effects in reaction time and error rate and congruency sequence
542 effects in reaction time and error rates (p-values of reported correlations are Holm-corrected for
543 five tests) (see Supplementary Figure 5).

544 **Acknowledgements**

545 We would like to thank Tobias Egner for valuable comments on a previous draft of the
546 manuscript. W.N., S.B. (G.0660.17N) and L.V. (11H5619N) were supported by the FWO –
547 Research Foundation Flanders. C.G.G. was supported by the Special Research Fund of Ghent
548 University (BOF.GOA.2017.0002.03). D.W. was supported by the FWO
549 (FWO.KAN.2019.0023.01), and the European Union’s Horizon 2020 research and innovation

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

550 program under the Marie Skłodowska-Curie grant agreement No 665501. All procedures applied
551 in the present experiment were carried out with adequate understanding and written consent of
552 the subjects and are in accordance with the Declaration of Helsinki.

553 **Author Contributions**

554 S.B. and W.N. developed the study concept. S.B., W.N. and L.V. contributed to the study design.
555 Data collection was performed by L.V. and V.H.. Data analysis was performed by L.V. under the
556 supervision of S.B., D.W. and C.G.C.. The manuscript was drafted by L.V. in cooperation with
557 S.B., W.N., D.W. and C.G.C.. All authors approved the final version of the manuscript for
558 submission.

559 **Declaration of Interests**

560 The authors have no competing interests to declare.

561

562

563

564

565

566

567

568

569

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

570 **References**

- 571 Bhandari, A., Gagne, C., & Badre, D. (2018). Just above Chance: Is It Harder to Decode
572 Information from Prefrontal Cortex Hemodynamic Activity Patterns? *Journal of*
573 *Cognitive Neuroscience*, 30(10), 1473–1498. https://doi.org/10.1162/jocn_a_01291
- 574 Botvinick, M. M. (2007). Conflict monitoring and decision making: Reconciling two
575 perspectives on anterior cingulate function. *Cognitive, Affective, & Behavioral*
576 *Neuroscience*, 7(4), 356–366. <https://doi.org/10.3758/CABN.7.4.356>
- 577 Botvinick, M. M., Braver, T. S., Link to external site, this link will open in a new window,
578 Barch, D. M., Carter, C. S., & Cohen, J. D. (2001). Conflict monitoring and cognitive
579 control. *Psychological Review*, 108(3), 624–652. [http://dx.doi.org/10.1037/0033-](http://dx.doi.org/10.1037/0033-295X.108.3.624)
580 295X.108.3.624
- 581 Braem, S., Bugg, J. M., Schmidt, J. R., Crump, M. J. C., Weissman, D. H., Notebaert, W., &
582 Egner, T. (2019). Measuring Adaptive Control in Conflict Tasks. *Trends in Cognitive*
583 *Sciences*, 23(9), 769–783. <https://doi.org/10.1016/j.tics.2019.07.002>
- 584 Braem, S., King, J. A., Korb, F. M., Krebs, R. M., Notebaert, W., & Egner, T. (2017). The role
585 of anterior cingulate cortex in the affective evaluation of conflict. *Journal of Cognitive*
586 *Neuroscience*, 29(1), 137–149.
- 587 Brown, J. W., & Alexander, W. H. (2017). Foraging value, risk avoidance, and multiple control
588 signals: How the anterior cingulate cortex controls value-based decision-making. *Journal*
589 *of Cognitive Neuroscience*, 29(10), 1656–1673.
- 590 Bush, G., Luu, P., & Posner, M. I. (2000). Cognitive and emotional influences in anterior
591 cingulate cortex. *Trends in Cognitive Sciences*, 4(6), 215–222.

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

- 592 Calhoun, A. J., & Hayden, B. Y. (2015). The foraging brain. *Current Opinion in Behavioral*
593 *Sciences*, 5, 24–31.
- 594 Chen, T., Becker, B., Camilleri, J., Wang, L., Yu, S., Eickhoff, S. B., & Feng, C. (2018). A
595 domain-general brain network underlying emotional and cognitive interference
596 processing: Evidence from coordinate-based and functional connectivity meta-analyses.
597 *Brain Structure and Function*, 223(8), 3813–3840.
- 598 Cox, D. D., & Savoy, R. L. (2003). Functional magnetic resonance imaging (fMRI) “brain
599 reading”: Detecting and classifying distributed patterns of fMRI activity in human visual
600 cortex. *NeuroImage*, 19(2), 261–270. [https://doi.org/10.1016/S1053-8119\(03\)00049-1](https://doi.org/10.1016/S1053-8119(03)00049-1)
- 601 Davis, T., LaRocque, K. F., Mumford, J. A., Norman, K. A., Wagner, A. D., & Poldrack, R. A.
602 (2014). What do differences between multi-voxel and univariate analysis mean? How
603 subject-, voxel-, and trial-level variance impact fMRI analysis. *NeuroImage*, 97, 271–
604 283. <https://doi.org/10.1016/j.neuroimage.2014.04.037>
- 605 De La Vega, A., Chang, L. J., Banich, M. T., Wager, T. D., & Yarkoni, T. (2016). Large-scale
606 meta-analysis of human medial frontal cortex reveals tripartite functional organization.
607 *Journal of Neuroscience*, 36(24), 6553–6562.
- 608 Devinsky, O., Morrell, M. J., & Vogt, B. A. (1995). Contributions of anterior cingulate cortex to
609 behaviour. *Brain*, 118(1), 279–306.
- 610 Dignath, D., Eder, A. B., Steinhauser, M., & Kiesel, A. (2020). Conflict monitoring and the
611 affective-signaling hypothesis—An integrative review. *Psychonomic Bulletin & Review*,
612 1–24.
- 613 Dreisbach, G., & Fischer, R. (2015). Conflicts as aversive signals for control adaptation. *Current*
614 *Directions in Psychological Science*, 24(4), 255–260.

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

- 615 Ebitz, R. B., & Hayden, B. Y. (2016). Dorsal anterior cingulate: A Rorschach test for cognitive
616 neuroscience. *Nature Neuroscience*, *19*(10), 1278.
- 617 Eriksen, B. A., & Eriksen, C. W. (1974). Effects of noise letters upon the identification of a
618 target letter in a nonsearch task. *Perception & Psychophysics*, *16*(1), 143–149.
- 619 Gehring, W. J., & Willoughby, A. R. (2002). The medial frontal cortex and the rapid processing
620 of monetary gains and losses. *Science*, *295*(5563), 2279–2282.
- 621 Haxby, J. V. (2012). Multivariate pattern analysis of fMRI: The early beginnings. *NeuroImage*,
622 *62*(2), 852–855. <https://doi.org/10.1016/j.neuroimage.2012.03.016>
- 623 Haynes, J.-D. (2015). A Primer on Pattern-Based Approaches to fMRI: Principles, Pitfalls, and
624 Perspectives. *Neuron*, *87*(2), 257–270. <https://doi.org/10.1016/j.neuron.2015.05.025>
- 625 Hebart, M. N., Görden, K., & Haynes, J.-D. (2015). The Decoding Toolbox (TDT): A versatile
626 software package for multivariate analyses of functional imaging data. *Frontiers in*
627 *Neuroinformatics*, *8*. <https://doi.org/10.3389/fninf.2014.00088>
- 628 Heilbronner, S. R., & Hayden, B. Y. (2016). Dorsal anterior cingulate cortex: A bottom-up view.
629 *Annual Review of Neuroscience*, *39*, 149–170.
- 630 Hendriks, M. H. A., Daniels, N., Pegado, F., & Op de Beeck, H. P. (2017). The Effect of Spatial
631 Smoothing on Representational Similarity in a Simple Motor Paradigm. *Frontiers in*
632 *Neurology*, *8*. <https://doi.org/10.3389/fneur.2017.00222>
- 633 Inzlicht, M., Bartholow, B. D., & Hirsh, J. B. (2015). Emotional foundations of cognitive
634 control. *Trends in Cognitive Sciences*, *19*(3), 126–132.
635 <https://doi.org/10.1016/j.tics.2015.01.004>

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

- 636 Jahn, A., Nee, D. E., Alexander, W. H., & Brown, J. W. (2016). Distinct regions within medial
637 prefrontal cortex process pain and cognition. *Journal of Neuroscience*, *36*(49), 12385–
638 12392.
- 639 Jiang, J., & Egner, T. (2014). Using neural pattern classifiers to quantify the modularity of
640 conflict–control mechanisms in the human brain. *Cerebral Cortex*, *24*(7), 1793–1805.
- 641 Kamitani, Y., & Sawahata, Y. (2010). Spatial smoothing hurts localization but not information:
642 Pitfalls for brain mappers. *NeuroImage*, *49*(3), 1949–1952.
643 <https://doi.org/10.1016/j.neuroimage.2009.06.040>
- 644 Kaplan, J. T., Man, K., & Greening, S. G. (2015). Multivariate cross-classification: Applying
645 machine learning techniques to characterize abstraction in neural representations.
646 *Frontiers in Human Neuroscience*, *9*, 151.
- 647 Keuleers, E., Brysbaert, M., & New, B. (2010). SUBTLEX-NL: A new measure for Dutch word
648 frequency based on film subtitles. *Behavior Research Methods*, *42*(3), 643–650.
649 <https://doi.org/10.3758/BRM.42.3.643>
- 650 Kolling, N., Wittmann, M. K., Behrens, T. E., Boorman, E. D., Mars, R. B., & Rushworth, M. F.
651 (2016). Value, search, persistence and model updating in anterior cingulate cortex.
652 *Nature Neuroscience*, *19*(10), 1280.
- 653 Kragel, P. A., Kano, M., Van Oudenhove, L., Ly, H. G., Dupont, P., Rubio, A., Delon-Martin,
654 C., Bonaz, B. L., Manuck, S. B., & Gianaros, P. J. (2018). Generalizable representations
655 of pain, cognitive control, and negative emotion in medial frontal cortex. *Nature*
656 *Neuroscience*, *21*(2), 283.

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

- 657 Kriegeskorte, N., Goebel, R., & Bandettini, P. (2006). Information-based functional brain
658 mapping. *Proceedings of the National Academy of Sciences*, *103*(10), 3863–3868.
659 <https://doi.org/10.1073/pnas.0600244103>
- 660 Kurdi, B., Lozano, S., & Banaji, M. R. (2017). Introducing the Open Affective Standardized
661 Image Set (OASIS). *Behavior Research Methods*, *49*(2), 457–470.
662 <https://doi.org/10.3758/s13428-016-0715-3>
- 663 Lieberman, M. D., Burns, S. M., Torre, J. B., & Eisenberger, N. I. (2016). Reply to Wager et al.:
664 Pain and the dACC: The importance of hit rate-adjusted effects and posterior
665 probabilities with fair priors. *Proceedings of the National Academy of Sciences*, *113*(18),
666 E2476–E2479.
- 667 Lieberman, M. D., & Eisenberger, N. I. (2015). The dorsal anterior cingulate cortex is selective
668 for pain: Results from large-scale reverse inference. *Proceedings of the National*
669 *Academy of Sciences*, *112*(49), 15250–15255.
- 670 Lindquist, K. A., Satpute, A. B., Wager, T. D., Weber, J., & Barrett, L. F. (2016). The brain basis
671 of positive and negative affect: Evidence from a meta-analysis of the human
672 neuroimaging literature. *Cerebral Cortex*, *26*(5), 1910–1922.
- 673 Moors, A., De Houwer, J., Hermans, D., Wanmaker, S., van Schie, K., Van Harmelen, A.-L., De
674 Schryver, M., De Winne, J., & Brysbaert, M. (2013). Norms of valence, arousal,
675 dominance, and age of acquisition for 4,300 Dutch words. *Behavior Research Methods*,
676 *45*(1), 169–177. <https://doi.org/10.3758/s13428-012-0243-8>
- 677 Nieuwenhuis, S., Holroyd, C. B., Mol, N., & Coles, M. G. (2004). Reinforcement-related brain
678 potentials from medial frontal cortex: Origins and functional significance. *Neuroscience*
679 *& Biobehavioral Reviews*, *28*(4), 441–448.

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

- 680 Oldfield, R. C. (1971). The assessment and analysis of handedness: The Edinburgh inventory.
681 *Neuropsychologia*, 9(1), 97–113. [https://doi.org/10.1016/0028-3932\(71\)90067-4](https://doi.org/10.1016/0028-3932(71)90067-4)
- 682 Op de Beeck, H. P. (2010). Against hyperacuity in brain reading: Spatial smoothing does not
683 hurt multivariate fMRI analyses? *NeuroImage*, 49(3), 1943–1948.
684 <https://doi.org/10.1016/j.neuroimage.2009.02.047>
- 685 Peirce, J. W. (2007). PsychoPy—Psychophysics software in Python. *Journal of Neuroscience*
686 *Methods*, 162(1), 8–13. <https://doi.org/10.1016/j.jneumeth.2006.11.017>
- 687 Rainville, P. (2002). Brain mechanisms of pain affect and pain modulation. *Current Opinion in*
688 *Neurobiology*, 12(2), 195–204.
- 689 Schmidt, J. R. (2019). Evidence against conflict monitoring and adaptation: An updated review.
690 *Psychonomic Bulletin & Review*, 26(3), 753–771.
- 691 Shackman, A. J., Salomons, T. V., Slagter, H. A., Fox, A. S., Winter, J. J., & Davidson, R. J.
692 (2011). The Integration of Negative Affect, Pain and Cognitive Control in the Cingulate
693 Cortex. *Nature Reviews Neuroscience*, 12(3), 154–167. <https://doi.org/10.1038/nrn2994>
- 694 Shenhav, A., Botvinick, M. M., & Cohen, J. D. (2013). The expected value of control: An
695 integrative theory of anterior cingulate cortex function. *Neuron*, 79(2), 217–240.
- 696 Shenhav, A., Cohen, J. D., & Botvinick, M. M. (2016). Dorsal anterior cingulate cortex and the
697 value of control. *Nature Neuroscience*, 19(10), 1286.
- 698 Stroop, J. R. (1935). Studies of interference in serial verbal reactions. *Journal of Experimental*
699 *Psychology*, 18(6), 643.
- 700 van den Hurk, J., & de Beeck, H. P. O. (2019). Generalization asymmetry in multivariate cross-
701 classification: When representation A generalizes better to representation B than B to A.
702 *BioRxiv*, 592410.

SHARED REPRESENTATIONS OF CONFLICT AND NEGATIVE AFFECT

- 703 Wisniewski, D. (2018). Context-dependence and context-invariance in the neural coding of
704 intentional action. *Frontiers in Psychology, 9*.
- 705 Yarkoni, T., Poldrack, R. A., Nichols, T. E., Essen, D. C. V., & Wager, T. D. (2011). Large-scale
706 automated synthesis of human functional neuroimaging data. *Nature Methods, 8*(8), 665–
707 670. <https://doi.org/10.1038/nmeth.1635>
- 708