

1 **Neural dynamics between anterior insular cortex and right supramarginal gyrus dissociate genuine**
2 **affect sharing from automatic responses to pretended pain**

3 Yili Zhao¹, Lei Zhang¹, Markus Rütgen^{1,2}, Ronald Sladky¹, Claus Lamm^{1,2*}

4 ¹ Social, Cognitive and Affective Neuroscience Unit, Department of Cognition, Emotion, and Methods
5 in Psychology, Faculty of Psychology, University of Vienna, Liebiggasse 5, 1010 Vienna, Austria

6 ²Vienna Cognitive Science Hub, University of Vienna, Liebiggasse 5, 1010 Vienna, Austria

7

8 **Abstract**

9 Empathy for pain engages both shared affective responses and self-other distinction. In this study,
10 we addressed the highly debated question of whether neural responses previously linked to affect
11 sharing could result from the perception of salient affective displays. Moreover, we investigated how
12 affect sharing and self-other distinction interact to determine our response to a pain that is either
13 perceived as genuine or pretended (while in fact both were acted for reasons of experimental
14 control). We found stronger activations in regions associated with affect sharing (anterior insula,
15 aIns, and anterior mid-cingulate cortex, aMCC) as well as with affective self-other distinction (right
16 supramarginal gyrus, rSMG), in participants watching video clips of genuine vs. pretended facial
17 expressions of pain. Using dynamic causal modeling (DCM), we then assessed the neural dynamics
18 between the right aIns and rSMG in these two conditions. This revealed a reduced inhibitory effect
19 on the aIns to rSMG connection for genuine compared to pretended pain. For genuine pain only,
20 brain-to-behavior regression analyses highlighted a linkage between this inhibitory effect on the one
21 hand, and pain ratings as well as empathic traits on the other. These findings imply that if the pain of
22 others is genuine and thus calls for an appropriate empathic response, neural responses in the aIns
23 indeed seem related to affect sharing and self-other distinction is engaged to avoid empathic over-
24 arousal. In contrast, if others merely pretend to be in pain, the perceptual salience of their painful

25 expression results in neural responses that are down-regulated to avoid inappropriate affect sharing

26 and social support.

27

28 **Introduction**

29 As social beings, our own affective states are influenced by other people's feelings and affective
30 states. The facial expression of pain by others acts as a distinctive cue to signal their pain to others,
31 and thus results in sizeable affective responses in the observer. Certifying such responses as
32 evidence for empathy, however, requires successful self-other distinction, the ability to distinguish
33 the affective response experienced by ourselves from the affect experienced by the other person.

34 Studies using a wide variety of methods convergently have shown that observing others in pain
35 engages neural responses aligning with those coding for the affective component of self-experienced
36 pain, with the anterior insula (aIns) and the anterior mid-cingulate cortex (aMCC) being two key
37 areas in which such an alignment has been detected (Lamm et al., 2011; Rütgen et al., 2015;
38 Jauniaux et al., 2019; Xiong et al., 2019; Zhou et al., 2020; Fallon et al., 2020, for meta-analyses).

39 However, there is consistent debate on whether activity observed in these areas should indeed be
40 related to the sharing of pain affect, or whether it may not rather result from automatic responses
41 to salient perceptual cues - with pain vividly expressed on the face being one particularly prominent
42 example (Zaki et al., 2016, for review). It was thus one major aim of our study to address this
43 question. In this respect, contextual factors, individuals' appraisals, and attentional processes would
44 all impact their exact response to the affective states of others (Gu & Han, 2007; Hein & Singer,
45 2008, for review; Lamm et al., 2010; Forbes & Hamilton, 2020; Zhao et al., 2021). Recently, Coll et al.
46 (2017) have thus proposed a framework that attempts to capture these influences on affect sharing
47 and empathic responses. This model posits that individuals who see identical negative facial
48 expressions of others may have different empathic responses due to distinct contextual information,
49 and that this may depend on identification of the underlying affective state displayed by the other.

50 In the current functional magnetic resonance imaging (fMRI) study, we therefore created a situation
51 where we varied the genuineness of the pain affect felt by participants while keeping the perceptual
52 saliency (i.e., the quality and strength of pain expressions) identical. To this end, participants were

53 shown video clips of other persons who supposedly displayed genuine pain on their face vs. merely
54 pretended to be in pain. Note that for reasons of experimental control, all painful expressions on the
55 videos had been acted. This enabled us to interpret possible differences between conditions to the
56 observers' appraisal of the situation rather than to putative visual and expressive differences. This
57 way, we sought to identify the extent to which responses in affective nodes (such as the aIns and the
58 aMCC) genuinely track the pain of others, rather than resulting predominantly from the salient facial
59 expressions associated with the pain.

60 Another major aim of our study was to assess how self-other distinction allowed individuals to
61 distinguish between the sharing of actual pain vs. regulating an inappropriate and potentially
62 misleading "sharing" of what in reality is only a pretended affective state. We focused on the right
63 supramarginal gyrus (rSMG), which has been suggested to act as a major hub selectively engaged in
64 *affective* self-other distinction (Silani et al., 2013; Steinbeis et al., 2015; Hoffmann et al., 2016;
65 Bukowski et al., 2020). Though previous studies have indicated that rSMG is functionally connected
66 with areas associated with affect processing (Mars et al., 2011; Bukowski et al., 2020), we lack more
67 nuanced insights into how exactly rSMG interacts with these areas, and thus how it supports
68 accurate empathic responses. Hence, we used dynamic causal modeling (DCM) to investigate the
69 hypothesized distinct interactions between affective responses and self-other distinction for the
70 genuine and pretended pain situations, focusing on the aIns, aMCC, and their interaction with rSMG.
71 Furthermore, we investigated the relationship between neural activity and behavioral responses as
72 well as empathic traits. In line with the literature reviewed above, we expected that, on the
73 behavioral level, genuine pain would result in – alongside the obvious other-oriented higher pain
74 ratings – higher self-oriented unpleasantness ratings. On the neural level, we predicted aIns and
75 aMCC to show a stronger response to the genuine expressions of pain, but that these areas would
76 also respond to the pretended pain, but to a lower extent. Differences in rSMG engagement and
77 distinct patterns of this area's effective connectivity with aIns and aMCC were expected to relate to

78 self-other distinction, and thus to explain the different empathic responses to genuine vs. pretended
79 pain.

80 **Results**

81 **Behavioral results**

82 Three repeated-measures ANOVAs were performed with the factors *genuineness* (genuine vs.
83 pretended and *pain* (pain vs. no pain), for each of the three behavioral ratings. For ratings of painful
84 expressions in others (Figure 1C, left), there was a main effect of the factor genuineness: participants
85 showed higher ratings for the genuine vs. pretended conditions, $F_{\text{genuineness}}(1, 42) = 8.816, p = 0.005,$
86 $\eta^2 = 0.173$. There was also a main effect of pain: participants showed higher ratings for the pain vs.
87 no pain conditions, $F_{\text{pain}}(1,42) = 1718.645, p < 0.001, \eta^2 = 0.976$. The interaction term was significant
88 as well, $F_{\text{interaction}}(1, 42) = 7.443, p = 0.009, \eta^2 = 0.151$, and this was related to higher ratings of
89 painful expressions in others for the genuine pain compared to the pretended pain condition. For
90 ratings of painful feelings in others (Figure 1C, middle), there was a main effect of genuineness:
91 participants showed higher ratings for the genuine vs. pretended conditions, $F_{\text{genuineness}}(1, 42) =$
92 $770.140, p < 0.001, \eta^2 = 0.948$. There was also a main effect of pain, as participants showed higher
93 ratings for the pain vs. no pain conditions, $F_{\text{pain}}(1,42) = 1544.762, p < 0.001, \eta^2 = 0.974$. The
94 interaction for painful feelings ratings was significant as well, $F_{\text{interaction}}(1, 42) = 752.618, p < 0.001, \eta^2$
95 $= 0.947$, and this was related to higher ratings of painful feelings in others for the genuine pain
96 compared to the pretended pain condition. For ratings of unpleasantness in self (Figure 1C, right),
97 there was a main effect of genuineness: participants showed higher ratings for the genuine vs.
98 pretended conditions, $F_{\text{genuineness}}(1, 42) = 74.989, p < 0.001, \eta^2 = 0.641$. There was also a main effect
99 of pain: participants showed higher ratings for the pain vs. no pain conditions, $F_{\text{pain}}(1,42) = 254.709,$
100 $p < 0.001, \eta^2 = 0.858$. The interaction for unpleasantness ratings was significant as well, $F_{\text{interaction}}(1,$
101 $42) = 73.620, p < 0.001, \eta^2 = 0.637$, and this was related to higher ratings of unpleasantness in self
102 for the genuine pain compared to the pretended pain condition. In sum, the behavioral data

103 indicated higher ratings and large effect sizes of painful feelings in others and unpleasantness in self
104 for the genuine compared to the pretended pain condition. Ratings of pain expressions also differed
105 in terms of genuineness, at comparably low effect size, though they were expected to not show a
106 difference by way of our experimental design and the pilot study.

107 We also found a significant correlation between behavioral ratings of painful feelings in others and
108 unpleasantness in self in the genuine pain condition, $r = 0.691$, $p < 0.001$; while in the pretended
109 pain condition, the correlation was not significant, $r = 0.249$, $p = 0.107$ (Figure 1D). A bootstrapping
110 comparison showed a significant difference between the two correlation coefficients, $p = 0.002$, 95%
111 Confidence Interval (CI) = [0.230, 1.060].

112 [Insert Figure 1 here]

113 **fMRI results: mass-univariate analyses**

114 Three contrasts were computed: 1) genuine: pain – no pain, 2) pretended: pain – no pain, and 3)
115 genuine (pain – no pain) – pretended (pain – no pain). Across all three contrasts, we found
116 activations as hypothesized in bilateral aIns, aMCC, and rSMG (Figure 2A and Table 1).

117 To identify whether or which brain activity was specifically related to the behavioral ratings
118 described above, we performed a multiple regression analysis where we explored the relationship of
119 activation in the contrast genuine pain – pretended pain with the three behavioral ratings. We found
120 significant clusters in bilateral aIns, visual cortex, and cerebellum (Figure 2B); notably, when
121 statistically accounting for ratings of painful expressions in others and painful feelings in others, all
122 three clusters were exclusively explained by the ratings of self-unpleasantness.

123 [Insert Figure 2 here]

124 [Insert Table 1 here]

125 **DCM results**

126 We performed DCM analysis to specifically examine the modulatory effect of genuineness on the
127 effective connectivity between the right alns and rSMG. More specifically, we sought to assess
128 whether the experimental manipulation of genuine pain vs. pretended pain tuned the bidirectional
129 neural dynamics from alns to rSMG and *vice versa*, in terms of both directionality (sign of the DCM
130 parameter) and intensity (magnitude of the DCM parameter). If the experimental manipulation
131 modulated the effective connectivity, we would observe a strong posterior probability ($p_p > 0.95$) of
132 the modulatory effect. Our original analysis plan was to include aMCC in the DCM analyses, but
133 based on the fact that aMCC did not show as strong evidence (in terms of the multiple regression
134 analysis) as the alns of being involved in our task, we decided to use a more parsimonious DCM
135 model without the aMCC.

136 We found strong evidence of inhibitory effects on the alns to rSMG connection both in the genuine
137 pain condition and in the pretended pain condition (Figure 3A, 3B and 3C). Comparing the strength
138 of these modulatory effects on the alns to rSMG connection revealed a reduced inhibitory effect for
139 genuine pain as opposed to pretended pain, $t_{41} = 2.671$, $p = 0.011$ (Mean_{genuine pain} = -0.821, 95% CI =
140 [-0.878, -0.712]; Mean_{pretended pain} = -0.934, 95% CI = [-1.076, -0.822]; Figure 3C). There was no
141 evidence of a modulatory effect on the rSMG to alns connection.

142 **Individual associations between modulatory effects, behavioral ratings and questionnaires**

143 To examine how the modulatory effects from the DCM were related to the behavioral ratings, we
144 computed two stepwise linear regression models for each condition. The regression model was
145 significant for the genuine pain condition ($F_{\text{model}(1,41)} = 4.639$, $p = 0.037$, $R^2 = 0.104$), when painful
146 feelings in others were added to the model and the other two ratings were excluded ($B = 0.079$, β
147 = 0.322, $p = 0.037$). However, the model was not significant for the pretended pain condition (Figure
148 3D). The variance inflation factors (VIFs) for three ratings in both models were calculated to diagnose
149 collinearity, showing no severe collinearity problem (all VIFs < 5; the smallest VIF = 1.132 and the
150 largest VIF = 4.387).

151 In addition, we tested two stepwise linear regression models to investigate whether subscales of all
152 three questionnaires could explain modulatory effects for genuine pain and pretended pain. In the
153 genuine pain condition, we found that the modulatory effect was significantly explained by scores of
154 two subscales, i.e., affective ability and affective reactivity of the ECQ: $F_{\text{model}}(1,39) = 6.829$, $p =$
155 0.003 , $R^2 = 0.270$; $B_{\text{affective ability}} = 0.052$, $\beta = 0.497$, $p = 0.002$; $B_{\text{affective reactivity}} = -0.040$, $\beta = -0.421$,
156 $p = 0.008$. No significant predictor was found with the other questionnaires (i.e., IRI and TAS). In the
157 pretended pain condition, none of the three questionnaires significantly predicted variations of the
158 modulatory effect. No severe collinearity problem was detected for either regression model (all $VIFs$
159 < 2 ; the smallest $VIF = 1.011$ and the largest $VIF = 1.600$).

160 [Insert Figure 3 here]

161 Discussion

162 In this study, we developed and used a novel experimental paradigm in which participants watched
163 video clips of persons who supposedly either genuinely experienced or merely pretended to be in
164 strong pain. Combining mass-univariate analysis with effective connectivity (DCM) analyses, our
165 study provides evidence on the distinct neural dynamics between regions suggestive of affect
166 processing (i.e., aIns and aMCC) and self-other distinction (i.e., rSMG) for genuinely sharing vs.
167 responding to pretended, non-genuine pain. With this, we aimed to clarify two main questions: First,
168 whether neural responses in areas such as the aIns and aMCC to the pain of others are indeed
169 related to a veridical sharing of affect, as opposed to simply tracking automatic responses to salient
170 affective displays. And second, how processes related to self-other distinction, implemented in the
171 rSMG, enable appropriate empathic responses to genuine vs. merely pretended affective states.

172 The mass-univariate analyses suggest that the increased activity in aIns for genuine pain as opposed
173 to pretended pain properly reflects affect sharing. As aforementioned, the network of affective
174 sharing and certain domain-general processes (e.g., salience detection and automatic emotion
175 processing) overlap in aIns and aMCC (Zaki et al., 2016, for review). This indicates that indeed, part

176 of the activation in these areas could be related to perceptual salience, which is why it has been
177 widely debated as a potential confound of empathy and affect sharing models (Zaki et al., 2016;
178 Lamm et al., 2019, for review). However, when comparing genuine pain versus pretended pain,
179 activity in these areas was not only found to be stronger in response to genuine pain, but the
180 increased activation in aIns was also selectively correlated with ratings of self-oriented
181 unpleasantness (i.e., after statistically accounting for painful expressions and painful feelings in
182 others). That only aIns and not also aMCC shows such correlation may be explained by previous
183 studies, according to which aIns is more specifically associated with affective representations, while
184 the role of aMCC rather seems to evaluate and regulate emotions that arise due to empathy (Fan et
185 al., 2011; Lamm et al., 2011; Jauniaux et al., 2019). Taken together, the activation and brain-behavior
186 findings provide evidence that responses in aIns (and to a lesser extent also the aMCC) are not
187 simply automatic responses triggered by perceptually salient events. Rather, they seem to track the
188 actual affective states of the other person, and thus the shared neural representation of that
189 response (see Zhou et al., 2020, for similar recent conclusions based on multi-voxel pattern
190 analyses). Our findings are also in line with the proposed model of Coll et al. (2017), which suggests
191 that affect sharing is the consequence of emotion identification. More specifically, while part of the
192 activation in the aIns and aMCC is indeed related to an (presumably earlier) automatic response, the
193 added engagement of these areas once they have identified the pain as genuine shows that only in
194 this condition, they then also engage in proper affect sharing. Ideally, one should be able to discern
195 these processes in time, but neither the temporal resolution of our fMRI measurements nor the
196 paradigm in which we always announced the conditions beforehand would have been sensitive
197 enough to do so. Thus, future studies including complementary methods such as EEG and MEG, and
198 tailored experimental designs are needed to pinpoint the exact sequence of processes engaged in
199 automatic affective responses vs. proper affect sharing.

200 Beyond higher activation in affective nodes supporting (pain) empathy, increased activation was also
201 found in rSMG. This area was shown to be engaged in action observation and imitating emotions

202 (Bach et al., 2010; Pokorny et al., 2015; Gola et al., 2017; Hawco et al., 2017), and a specific role in
203 affective rather than cognitive self-other distinction has been identified for rSMG (Silani et al., 2013;
204 Steinbeis et al., 2015; Bukowski et al., 2020). Based on such findings, it has been proposed that the
205 rSMG allows for a rapid switching between or the integration of self- and other-related
206 representations, as two processes that may underpin the functional basis of successful self-other
207 distinction (Lamm et al., 2016, for review). Concerning the current findings, we thus propose that
208 the higher rSMG engagement in the genuine pain condition reflects an increasing demand for self-
209 other distinction imposed by the stronger shared negative affect experienced in this condition.
210 Theoretical models of empathy and related socio-affective responses suggest that such regulation is
211 especially important to avoid so-called empathic over-arousal, which would shift the focus away
212 from empathy and the other's needs, towards taking care of one's own personal distress (Batson et
213 al., 1987; Decety & Lamm, 2011, for review).

214 Beyond these differences in the magnitude of rSMG activation, the DCM analysis demonstrated less
215 inhibition on the aIns-to-rSMG connection for genuine pain compared to pretended pain. Various
216 theoretical accounts suggest that areas such as the aIns and rSMG may play a key role in comparing
217 self-related information with the sensory evidence (Decety & Lamm, 2007; Seth, 2013, for review).
218 According to recent theories on predictive processing (Clark, 2013, for review) and active inference
219 (Friston, 2010, for review), the brain can be regarded as a "prediction machine", in which the top-
220 down signals pass over predictions and the bottom-up signals convey prediction errors across
221 different levels of cortical hierarchies (Chen et al., 2009; Friston, 2010, for review; Bastos et al.,
222 2015). It is suggested that these top-down predictions are mediated by inhibitory neural connections
223 (Zhang et al., 2008; Bastos et al., 2015; Miska et al., 2018). Our findings align with such views, by
224 suggesting that the inhibitory connection from aIns to rSMG can be explained as the predictive
225 mismatch between the top-down predictions of self-related information (e.g., personal affect) and
226 sensory inputs (e.g., pain facial expressions). This suppression of neural activity leads to an
227 *explaining away* of incoming bottom-up prediction error. This is reflected by the absence of any

228 condition-dependent modulatory effects on the rSMG to alns connection, suggesting that the
229 influence of the task conditions is sufficiently modeled by the predictions from alns to rSMG.
230 Therefore, the stronger inhibition for pretended pain, compared to genuine pain, could indicate a
231 higher demand to overcome the mismatch between the visual inputs and the agent's prior beliefs
232 and contextual information about the situation (i.e., "this person looks like in pain, but I know
233 he/she does not actually feel it"). The reduced inhibition in the genuine pain condition could
234 moreover be a mechanism that explains the higher rSMG activation in this condition.

235 We also found the strength of the inhibitory effect in the genuine pain condition to correlate with
236 ratings of painful feelings in others, but not with the ratings of pain expression in others or
237 unpleasantness in self. For the pretended pain condition none of the ratings showed a correlation.
238 The latter could in principle be due to a lack of variation in the ratings (which by way of the design
239 were mostly close to zero or one). We deem it more plausible, though, that the correlation findings
240 provide further evidence that the modulation of alns to rSMG is implicated in encoding others'
241 emotional states when participants engaged in genuine affect sharing. It is also interesting to note
242 that the found correlation relates to cognitive evaluations of the other's pain rather than to own
243 affect, as tracked by the unpleasantness in self-ratings. This would to some extent be in line with
244 DCM findings by Kanske et al. (2016). These authors found that the inhibition of the temporoparietal
245 junction (TPJ) by the alns was linked to interactions between Theory of Mind (ToM) and empathic
246 distress, i.e., the interaction of "cognitive" vs. "affective" processes engaged in understanding
247 others' cognitive and affective states. Note that the right TPJ is an overarching area involved in self-
248 other distinction of which rSMG is considered a part or at least closely connected to (Decety &
249 Lamm, 2007, for review).

250 The correlations between the DCM inhibitory effect and empathic traits assessed via questionnaires
251 provide further refinements for the relevance of rSMG in implementing self-other distinction to
252 allow for an appropriate empathic response. When participants shared genuine affect, the inhibitory

253 effect on the aIns to rSMG connection was positively correlated with affective ability and negatively
254 correlated with affective reactivity. Affective ability reflects the capacity to subjectively share
255 emotions with others, while affective reactivity plays a role in the susceptibility to vicarious distress
256 and thus to more automatic responses to another's emotion (Batchelder et al., 2017). Again, as for
257 the correlations with the three rating scales, we did not find correlations of empathic traits for the
258 pretended pain condition. Taken together, the DCM results and their qualification by the correlation
259 findings suggest that in the genuine pain condition, which requires an accurate sharing of pain, rSMG
260 interacts with aIns to achieve "affective-to-affective" self-other distinction – i.e., disambiguating
261 affective signals originating in the self from those attributable to the other person. The aIns to rSMG
262 connection in the pretended pain condition may reflect a related, yet slightly distinct mechanism.
263 Here, it seems that "cognitive-to-affective" self-other distinction is at play, which helps resolve
264 conflicting information between the top-down contextual information (i.e., that the demonstrator is
265 not actually in pain) from what seems an unavoidable affective response to the highly salient
266 perceptual cue of the facial expression of pain. Given our behavioral and trait data did not allow us
267 to distinguish more precisely between these different types of self-other distinction, this however
268 remains an interpretation and a hypothesis that will require further investigation.

269 One potential limitation of the study could be the slightly higher ratings of other-oriented pain
270 expressions for genuine pain, which were hypothesized to have no difference, as compared to
271 pretended pain. As we found the enhanced aIns activation in the genuine pain condition mainly
272 tracked personal unpleasantness rather than perceptually domain-general processes, and because
273 the effect size of the pain expression difference was much smaller than for the affect ratings, we
274 consider this difference did not fundamentally influence the interpretation of our findings.

275 In conclusion, the current study advances our understanding of two main aspects of empathy. First,
276 we provide evidence that empathy-related responses in the aIns can indeed be linked to affective
277 sharing, rather than attributing them to responses triggered only by perceptual saliency. Second, we

278 show how aIns and rSMG are orchestrated to track what another person really feels, thus enabling
279 us to appropriately respond to their actual needs. Beyond these basic research insights, our study
280 provides novel avenues for clinical application, and the investigation of contextual and interpersonal
281 factors in the accurate diagnosis of pain and its expression.

282 **Materials and Methods**

283 **Participants**

284 Forty-eight participants took part in the study. Five of them were excluded because of excessive
285 head motion (> 15% scans with the frame-wise displacement over 0.5 mm in one session). Data of
286 the remaining 43 participants (21 females; age: Mean = 26.72 years, S.D. = 4.47) were entered into
287 analyses. This sample size was determined on a priori power analysis in Gpower 3.1 (Faul et al.,
288 2007). We assumed a medium effect size of *Cohen's d* = 0.5. After calculation, the minimum sample
289 size statistically required for this study was 34 ($\alpha = 0.05$, two-tailed, $1-\beta = 0.80$). Participants were
290 pre-screened by an MRI safety-check questionnaire, assuring normal or corrected to normal vision
291 and no presence or history of neurologic, psychiatric, or major medical disorders. All participants
292 were being right-handed (self-reported) and provided written consent including post-disclosure of
293 any potential deception. The study was approved by the ethics committee of the Medical University
294 of Vienna and was conducted in line with the latest version of the Declaration of Helsinki (2013).

295 **Manipulation of facial expressions**

296 As part of our study we developed a novel experimental design and corresponding stimuli, which
297 consisted of video clips showing different demonstrators ostensibly in four different situations: 1)
298 Genuine pain: the demonstrator's right cheek was penetrated by a hypodermic needle attached to a
299 syringe, and the demonstrator's facial expression changed from neutral to a strongly painful facial
300 expression. 2) Genuine no pain: the demonstrator maintained a neutral facial expression when a Q-
301 tip fixed on the backend of the same syringe touched their right cheek. 3) Pretended pain: the

302 demonstrator's right cheek was approached by the same syringe and the hypodermic needle, with
303 the latter covered by a protective cap; upon touch by the cap, the demonstrator's facial expression
304 changed from neutral to a strongly painful facial expression. 4) Pretended no pain: the demonstrator
305 maintained a neutral facial expression when a Q-tip fixed on the backend of the same syringe
306 touched their right cheek.

307 To create these stimuli, we recruited 20 demonstrators (10 females), with experience in acting, and
308 filmed them in front of a dark blue background. An experimenter who stood on the right side of the
309 demonstrators, but of whom only the right hand holding the syringe could be seen, administered the
310 injections and touches. Unbeknownst to the participants, all painful expressions were acted, as the
311 needle was a telescopic needle (i.e., a needle that seemed to enter the cheek upon contact, but in
312 reality, was invisibly retracting into the syringe). The reason for using a protective cap in the
313 pretended pain condition was to match the perceptual situation that an aversive object was
314 approaching a body part in both pain conditions. In all situations, the demonstrator was instructed
315 to look naturally towards the camera 1.5 m in front of them. As soon as the needle or the cap
316 touched the demonstrator's cheek, the demonstrator made a painful facial expression, as naturally
317 and vividly as possible. In the neutral control conditions, demonstrators maintained a neutral facial
318 expression when a Q-tip fixed at the backend of the syringe touched their cheek. Again, a syringe
319 with a needle attached to the other end was used to perceptually control for the presence of an
320 aversive object in all four conditions. Note that in another set of conditions, demonstrators showed
321 disgusted or neutral expressions. Data from these conditions will be reported elsewhere. All
322 demonstrators signed an agreement that their video clips and static images could be used for
323 scientific purposes.

324 **Stimulus validation and pilot study**

325 To validate the stimuli, we performed an online validation study with N = 110 participants, who were
326 asked to rate a total of 120 video clips of 2 s duration of the two conditions (60 of each condition)

327 showing painful expressions (i.e., the genuine and the pretended pain conditions). The main aim of
328 the validation study was to identify a set of demonstrators that expressed pain with comparable
329 intensity and quality, and whose pain expressions in the genuine and pretended conditions were
330 comparable. After each video clip, participants rated three questions on a visual analog scale with 9
331 tick-marks and the two end-points marked as “almost not at all” to “unbearable”: 1) How much pain
332 did the person *express* on his/her face? 2) How much pain did the person *actually* feel? 3) How
333 *unpleasant* did you feel to watch the person in this situation? The order of these three questions
334 was pseudo-randomized. Moreover, eight catch trials randomly interspersed across the validation
335 study to test whether participants maintained attention to the stimuli. Here, participants were asked
336 to correctly select the demonstrator they had seen in the last video, between two static images of
337 the correct and a distractor demonstrator displayed side by side, both showing neutral facial
338 expressions.

339 The validation study was implemented within the online survey platform SoSci Survey
340 (<https://www.soscisurvey.de>), with a study participation invite published on Amazon Mechanical
341 Turk (<https://www.mturk.com/>), a globally commercial platform allowing for online testing. Survey
342 data of 62 out of 110 participants (34 females; age: Mean = 28.71 years, S.D. =10.11) were entered
343 into analysis (inclusion criteria: false rate for the test questions < 2/8, survey duration > 20 min and <
344 150 min, and the maximum number of continuous identical ratings < 5). Based on this validation
345 step, we had to exclude videos of 6 demonstrators (3 females) for which participants showed a
346 significant difference in painful expressions in others between the genuine pain and the pretended
347 pain conditions. As a result of this validation, videos of 14 demonstrators (7 females), which showed
348 no difference in the pain *expression* rating between genuine and pretended conditions, and which
349 overall showed comparable mean ratings in all three ratings, were selected for the subsequent pilot
350 study.

351 In the pilot study, 47 participants (24 females; age: Mean = 26.28 years, S.D. = 8.80) were recruited
352 for a behavioral experiment in the behavioral laboratory. The aim was to verify the experimental
353 effects and the feasibility of the experimental procedures that we intended to use in the main fMRI
354 experiment, as well as to identify video stimuli that may not yield the predicted responses. Thus, all
355 four conditions described above were presented to the participants. Participants were explicitly
356 instructed that they would watch other persons' genuine painful expressions in some blocks, while
357 in other blocks, they would see other persons acting out painful expressions (recall that in reality, all
358 demonstrators had been actors, and the information about this type of necessary deception was
359 conveyed to participants at the debriefing stage). They would see all demonstrators' neutral
360 expressions as well. Participants were instructed to rate the three questions mentioned above. Upon
361 screening for video clips that showed aberrant responses, we excluded videos of two demonstrators
362 (1 female), for whom the pain *expression* rating difference between the pretended vs. genuine
363 expressions was large. 48 videos of 12 demonstrators entered the following analyses. Three separate
364 repeated-measures ANOVAs were respectively performed for the three rating questions. For the
365 main effect of *genuineness* (genuine vs. pretended), it was not significant and low in effect size for
366 painful expressions in others ($F_{\text{genuineness}}(1, 46) = 2.939, p = 0.093, \eta^2 = 0.060$), but was significant
367 with high effect size for the painful feelings in others ($F_{\text{genuineness}}(1, 46) = 280.112, p < 0.001, \eta^2 =$
368 0.859) as well as the unpleasantness in self ($F_{\text{genuineness}}(1, 46) = 43.143, p < 0.001, \eta^2 = 0.484$). The
369 main effects of *pain* (pain vs. no pain) for all three questions were found significant with high effect
370 size (the smallest effect size was for the rating of unpleasantness in self, $F_{\text{pain}}(1, 46) = 82.199, p <$
371 $0.001, \eta^2 = 0.641$). Our pilot study thus a) provided assuring evidence that the novel experimental
372 paradigm worked as expected, and b) made it possible to select video clips that we could match for
373 the two conditions (i.e., genuine pain and pretended pain). More specifically, as expected and
374 required for the main study, participants rated the painfulness of the demonstrators to be
375 substantially higher when it was genuine as compared to those that were pretended, and this also
376 resulted in much higher unpleasantness experienced in the self. It is worth noting that, the two

377 conditions did not differ with respect to the ratings of the painful facial expressions, implying that
378 putative differences in ratings as well as the subsequent brain imaging data could only be attributed
379 to the contextual appraisal of the demonstrators' actual painful states, rather than the differences in
380 facial pain perception. Based on this pilot study, we thus decided on video clips of 12 demonstrators
381 (6 females) in the main fMRI experiment.

382 **Experimental design and procedure of the fMRI study**

383 The experiment was implemented using Cogent 2000 (version 1.33;
384 http://www.vislab.ucl.ac.uk/cogent_2000.php). MRI scanning took place at the University of Vienna
385 MRI Center. Once participants arrived at the scanner site, an experimenter instructed them that they
386 would watch videos from the four conditions outlined above. Participants were explicitly instructed
387 to recreate the feelings of the demonstrators shown in the videos as vividly and intensely as
388 possible. Based on the validation and pilot study, the painful *expressions* for the genuine and
389 pretended conditions were matched. We also counterbalanced the demonstrators appearing in the
390 genuine and pretended conditions across participants, thus controlling for differences in behavioral
391 and brain response that could be explained by differences between the stimulus sets. Note that, all
392 video clips were validated and piloted multiple times to ensure the experimental effect (details can
393 be found in the section above).

394 The participant performed the fMRI experiment in two runs (Figure 1A and 1B). Each run was
395 composed of two blocks showing genuine pain and two blocks showing pretended pain. In each
396 block, the participant watched nine video clips containing both painful and neutral videos. To remind
397 participants' the condition of the upcoming block, a label of 4 s duration appeared at the beginning
398 of each block, showing either "genuine" or "pretended" (in German). Each trial started with a
399 fixation cross (+) presented for 4 – 7 s (in steps of 1.5 s, Mean = 5.5 s). After that, the video (duration
400 = 2 s) was played. A short jitter was inserted after the video for 0.5 – 1.0 s (in steps of 0.05 s, Mean =
401 0.75 s). After the jitter, the following three questions were displayed (in German) one after the other

402 in a pseudo-randomized order: 1) How much pain did the person *express* on his/her face? 2) How
403 much pain did the person *actually feel*? 3) How *unpleasant did you feel* to watch the person in this
404 situation? Beneath each question, a visual analog scale ranging from 0 (not at all) to 8 (unbearable)
405 with 9 tick-marks was positioned. The participant moved the marker along the scale by pressing the
406 left or right keys on the button box, and they pressed the middle key to confirm their answer. The
407 marker initially was always located at the midpoint (“4”) of the scale. When the confirmed key was
408 pressed, the marker turned from black to red. All ratings lasted for 4 s even when the participant
409 pressed the confirmed key before the end of this period. Between the two runs, the participant had
410 a short break (1-2 min).

411 Before entering the scanner, participants conducted practice trials on the computer to get
412 familiarized with the button box and the experimental interface. After that, participants were moved
413 into the scanner and performed the task. Following the functional imaging runs, a 6.5 min structural
414 scanning was employed. When participants finished the scanning session, they were scheduled for a
415 date to complete three questionnaires in the lab: the Empathy Components Questionnaire (ECQ)
416 (Batchelder, 2015; Batchelder et al., 2017), the Interpersonal Reactivity Index (IRI) (Davis, 1980), and
417 the Toronto Alexithymia Scale (TAS) (Bagby et al., 1994). For the ECQ, there are 27 items in total to
418 be categorized into five subscales: cognitive ability, cognitive drive, affective ability, affective drive,
419 and affective reactivity, using a 4-point Likert scale ranging from 1 (“strongly disagree”) to 4
420 (“strongly agree”) (Batchelder, 2015; Batchelder et al., 2017). For the IRI, there are 28 items divided
421 into four subscales: perspective taking, fantasy, empathic concern, and personal distress, using a 5-
422 point Likert scale ranging from 0 (“does not describe me well”) to 4 (“describes me very well”)
423 (Davis, 1980). For the TAS, there are 20 items and three subscales - difficulty describing feelings,
424 difficulty identifying feelings, and externally oriented thinking, using a 5-point Likert scale ranging
425 from 1 (“strongly disagree”) to 5 (“strongly agree”) (Bagby et al., 1994). The average interval
426 between the scanning session and the lab survey was one week. The participant was debriefed after
427 completing the whole study.

428 **Behavioral data analysis**

429 We applied repeated-measures ANOVAs to investigate the main effects and the interaction of the
430 two factors genuine vs. pretended and pain vs. no pain, using SPSS (version 26.0; IBM). Furthermore,
431 we conducted Pearson correlations to examine whether ratings of painful feelings in others were
432 correlated with unpleasantness in self for the genuine pain and the pretended pain. The correlation
433 coefficients were further compared using a bootstrap approach with the R package bootcorci
434 (<https://github.com/GRousselet/bootcorci>).

435 **fMRI data acquisition**

436 fMRI data were collected using a Siemens Magnetom Skyra MRI scanner (Siemens, Erlangen,
437 Germany) with a 32-channel head coil. Functional whole-brain scans were collected using a
438 multiband-accelerated T2*-weighted echoplanar imaging (EPI) sequence (multiband acceleration
439 factor = 4, interleaved ascending acquisition in multi-slice mode, 52 slices co-planar to the
440 connecting line between anterior and posterior commissure, TR = 1200 ms, TE = 34 ms, acquisition
441 matrix = 96×96 voxels, FOV = 210×210 mm², flip angle = 66°, inter-slice gap = 0.4 mm, voxel size =
442 $2.2 \times 2.2 \times 2$ mm³). Two functional imaging runs, each lasting around 16 min (~800 images per run),
443 were performed. Structural images were acquired with a magnetization-prepared rapid gradient-
444 echo (MPRAGE) sequence (TE/TR = 2.43/2300 ms, flip angle = 8°, ascending acquisition, single-shot
445 multi-slice mode, FOV= 240×240 mm², voxel size = $0.8 \times 0.8 \times 0.8$ mm³, 208 sagittal slices, slice
446 thickness = 0.8 mm).

447 **fMRI data processing and mass-univariate functional segregation analyses**

448 Imaging data were preprocessed with a combination of Nipype (Gorgolewski et al., 2011) and
449 MATLAB (version R2018b 9.5.0; MathWorks) with Statistical Parametric Mapping (SPM12;
450 <https://www.fil.ion.ucl.ac.uk/spm/software/spm12/>). Raw data were imported into BIDS format
451 (<http://bids.neuroimaging.io/>). Functional data were subsequently preprocessed using slice timing

452 correction to the middle slice (Sladky et al., 2011), realignment to the first image of each session, co-
453 registration to the T1 image, segmentation between grey matter, white matter and cerebrospinal
454 fluid (CSF), normalization to MNI template space using Diffeomorphic Anatomical Registration
455 Through Exponentiated Lie Algebra (DARTEL) toolbox (Ashburner, 2007), and smoothing with a 6
456 mm full width at half-maximum (FWHM) three-dimensional Gaussian kernel.

457 To improve data quality, we performed data scrubbing of the functional scans for those whose
458 frame-wise displacements (FD) were over 0.5 mm (Power et al., 2012; Power et al., 2014). In other
459 words, we identified individual outlier scans and flagged the volume indices as nuisance regressors
460 in the general linear model (GLM) for the first-level analysis.

461 In order to perform mass-univariate functional segregation analyses, a first-level GLM design matrix
462 was created and composed of two identically modeled runs for each participant. Seven regressors of
463 interest were entered in each model: stimulation phase of the four conditions (i.e., genuine pain,
464 genuine no pain, pretended pain, pretended no pain; 2000 ms), rating phase of the three questions
465 (i.e., painful expressions in others, painful feelings in others, and unpleasantness in self; 12000 ms).
466 Six head motion parameters and the scrubbing regressors (FD > 0.5 mm; if applicable) were
467 additionally entered as nuisance regressors. Individual contrasts of the four conditions and the three
468 ratings (all across the two runs) against implicit baseline were respectively created.

469 On the second level, a flexible factorial design was employed to perform the group-level analysis.
470 The design included three factors: a between-subject factor (i.e., subject) that was specified
471 independent and with equal variance, a within-subject factor (i.e., genuine or pretended) that was
472 specified dependent and with equal variance, and a second within-subject factor (i.e., pain or no
473 pain) that was specified dependent and with equal variance (Gläscher & Gitelman, 2008). Three
474 contrasts were computed: (1) main effect of genuine: pain – no pain, (2) main effect of pretended:
475 pain – no pain, and (3) interaction: genuine (pain – no pain) – pretended (pain – no pain). We
476 applied an initial threshold of $p < 0.001$ (uncorrected) at the voxel level and a family-wise error

477 (FWE) correction ($p < 0.05$) at the cluster level. The cluster extent threshold was determined by the
478 SPM extension “cp_cluster_Pthresh.m” (<https://goo.gl/kjVydz>).

479 **Brain-behavior relationships**

480 A multiple regression model was built on the group level to investigate the relationship between
481 specific brain activations and behavioral ratings. In this model, the contrast genuine pain –
482 pretended pain was set as the dependent variable, and three behavioral ratings were specified as
483 independent variables. All covariates were mean-centered. The model aimed to test which brain
484 activations of the contrast could be explained by an independent variable after accounting for the
485 other two. Note that, we performed the regression model with the contrast genuine pain –
486 pretended pain instead of the more exhaustive contrast genuine (pain - no pain) - pretended (pain –
487 no pain), and this was because the genuine and the pretended pain conditions were the main focus
488 of our work. Moreover, the pain contrast showed more robust (in terms of statistical effect size) and
489 widespread activations across the brain, making it more likely to pick up possible brain-behavior
490 relationships. The same threshold as above was applied in this analysis.

491 We aimed to assess these brain-behavior relationships for the following regions of interest (ROI): 1)
492 alns and aMCC, i.e., two regions associated with affective processes and specifically with empathy
493 for pain, 2) rSMG, an area implicated in affective self-other distinction. The ROI masks were defined
494 as the conjunction of the averaging contrast between genuine and pretended: pain – no pain
495 (threshold: voxel-wise FWE correction, $p < 0.05$) and the anatomical masks created by the Wake
496 Forest University (WFU) Pick Atlas SPM toolbox (<http://fmri.wfubmc.edu>) with the automated
497 anatomical atlas (AAL). The ROI masks were created with Marsbar ROI Toolbox implemented in
498 SPM12 (Brett et al., 2002). Note that we specifically selected the ROIs this way, such that they were
499 orthogonal (i.e., independent) to the subsequent analyses of interest. As exploratory analyses found
500 significant correlations mainly in alns, rather than in aMCC, we will focus in the results section on

501 two ROIs: the right aIns and the rSMG. Focusing on the right aIns instead of the left one was because
502 the right aIns is on the ipsilateral hemisphere as rSMG.

503 **Analyses using dynamic causal modeling (DCM)**

504 To investigate the functional network involved in affective processes and self-other distinction and
505 how it was modulated by our experimental manipulations (i.e., genuine pain and pretended pain),
506 we used DCM to estimate the effective connectivity between the ROIs based on the tasked-related
507 brain responses (Stephan & Friston, 2010, for review). The DCM analyses were conducted with
508 DCM12.5 implemented in SPM12 (v. 7771). Firstly, we extracted individual time series separately for
509 each ROI. To ensure the selected voxels engaged in a task-relevant activity but not random signal
510 fluctuations, we determined the voxels both on a group-level threshold and an individual-level
511 threshold (Holmes et al., 2020). An initial threshold was set as $p < 0.05$, uncorrected. The significant
512 voxels in the main effect of genuine pain and pretended pain were further selected by an individual
513 threshold. For each participant, an individual peak coordinate within the ROI mask was searched and
514 an individual mask was consequently defined using a sphere of the 6 mm radius around the peak. As
515 a result, the individual time series for each ROI was extracted from the significant voxels of the
516 individual mask and summarized by the first eigenvariate. One participant was excluded as no voxels
517 survived significance testing. Secondly, we specified three regressors of interest: genuine pain,
518 pretended pain, and the video input condition (the combination of genuine pain and pretended
519 pain). That we did not specify no-pain conditions was because 1) the pain conditions were our main
520 focus, and 2) adding no-interest conditions would inevitably increase the model complexity. Then, a
521 fully connected DCM model for each participant was created. Three parameters were specified: 1)
522 bidirectional connections between regions and self-connections (matrix A), 2) modulatory effects
523 (i.e., genuine pain and pretended pain) on the between-region connections (matrix B), and 3) driving
524 inputs (i.e., the video input condition) into the model on both regions (matrix C) (Zeidman et al.,
525 2019a). To remain parsimonious, we did not set modulatory effects on the self-connections in Matrix

526 A. Then the full DCM model was individually estimated. Finally, group-level DCM inference was
527 performed using parametric empirical Bayes (Zeidman et al., 2019b). We conducted an automatic
528 search over the entire model space (max. $n = 256$) using Bayesian model reduction (BMR) and
529 random-effects Bayesian model averaging (BMA), resulting in a final group model that takes
530 accuracy, complexity, and uncertainty into account (Zeidman et al., 2019b). The threshold of the
531 Bayesian posterior probability was set to $p_p > 0.95$ (i.e., *strong evidence*) but we reported all
532 parameters above $p_p > 0.75$ (i.e., *positive evidence*) for full transparency of the DCM results. Finally, a
533 paired sample *t*-test was performed to compare modulatory effects between the genuine pain and
534 the pretended pain conditions.

535 To probe whether task-related modulatory effects were associated with behavioral measurements,
536 we performed stepwise linear regression analyses of modulatory parameters with, 1) the three
537 behavioral ratings, and 2) the empathy-related questionnaires (i.e., IRI, ECQ, and TAS). We set up
538 two regression models for the genuine pain condition and the pretended pain condition,
539 respectively, in which the DCM parameters of modulatory effects were determined as dependent
540 variables and the ratings of painful expressions in others, painful feelings in others, and
541 unpleasantness in self as independent variables. Accordingly, we performed additional two
542 regression models for both conditions in which DCM modulatory effects were set as dependent
543 variables and scores of each subscale of all questionnaires were set as independent variables,
544 respectively. As two participants did not complete all three questionnaires, we excluded their data
545 from the regression analyses. The statistical significance of the regression analysis was set to $p <$
546 0.05 . The multicollinearity for independent variables was diagnosed using the variance inflation
547 factor (VIF) that measures the correlation among independent variables, in the R package *car*
548 (<https://cran.r-project.org/web/packages/car/index.html>). Here we used a rather conservative
549 threshold of $VIF < 5$ as a sign of no severe multicollinearity (Menard, 2002; James et al., 2013).

550 **Acknowledgements**

551 This work was supported by Chinese Scholarship Council (CSC) Grant (201604910515) and Vienna
552 Doctoral School in Cognition, Behavior and Neuroscience (VDS CoBeNe) completion grant fellowship
553 to Y.Z.; the Vienna Science and Technology Fund (WWTF VRG13-007) to C.L., and the Austrian
554 Science Fund (FWF P 32686) to C.L. and M.R.. We thank Michael Schnödt, Lukas Repnik, Elisa
555 Warmuth, Betty Geidel, Phan Ri, Sven Sander, Gvantsa Gogisvanidze, Robert Meyka, and Anja Tritt
556 for help with data acquisition.

557 **Conflicts of interest**

558 The authors declare no competing financial interests.

559 **References**

- 560 Ashburner, J. (2007). A fast diffeomorphic image registration algorithm. *NeuroImage*, *38*(1), 95-113.
561 doi: <https://doi.org/10.1016/j.neuroimage.2007.07.007>
- 562 Bach, P., Peelen, M. V., & Tipper, S. P. (2010). On the Role of Object Information in Action
563 Observation: An fMRI Study. *Cerebral Cortex*, *20*(12), 2798-2809. doi:
564 <http://doi.org/10.1093/cercor/bhq026>
- 565 Bagby, R. M., Taylor, G. J., & Parker, J. D. A. (1994). The twenty-item Toronto Alexithymia Scale: II.
566 Convergent, discriminant, and concurrent validity. *Journal of Psychosomatic Research*, *38*(1),
567 33-40. doi: [http://doi.org/10.1016/0022-3999\(94\)90006-X](http://doi.org/10.1016/0022-3999(94)90006-X)
- 568 Bastos, A. M., Litvak, V., Moran, R., Bosman, C. A., Fries, P., & Friston, K. J. (2015). A DCM study of
569 spectral asymmetries in feedforward and feedback connections between visual areas V1 and
570 V4 in the monkey. *NeuroImage*, *108*, 460-475. doi:
571 <https://doi.org/10.1016/j.neuroimage.2014.12.081>
- 572 Batchelder, L. (2015). *Characterising the components of empathy: implications for models of autism*.
573 University of Bath.
- 574 Batchelder, L., Brosnan, M., & Ashwin, C. (2017). The Development and Validation of the Empathy
575 Components Questionnaire (ECQ). *PLOS ONE*, *12*(1), e0169185. doi:
576 <http://doi.org/10.1371/journal.pone.0169185>
- 577 Batson, C. D., Fultz, J., & Schoenrade, P. A. (1987). Distress and Empathy: Two Qualitatively Distinct
578 Vicarious Emotions with Different Motivational Consequences. *Journal of Personality*, *55*(1),
579 19-39. doi: <https://doi.org/10.1111/j.1467-6494.1987.tb00426.x>
- 580 Brett, M., Anton, J.-L., Valabregue, R., & Poline, J.-B. (2002). *Region of interest analysis using an SPM*
581 *toolbox*. Paper presented at the 8th international conference on functional mapping of the
582 human brain.
- 583 Bukowski, H., Tik, M., Silani, G., Ruff, C. C., Windischberger, C., & Lamm, C. (2020). When differences
584 matter: rTMS/fMRI reveals how differences in dispositional empathy translate to distinct

- 585 neural underpinnings of self-other distinction in empathy. *Cortex*, 128, 143-161. doi:
- 586 <https://doi.org/10.1016/j.cortex.2020.03.009>
- 587 Chen, C. C., Henson, R. N., Stephan, K. E., Kilner, J. M., & Friston, K. J. (2009). Forward and backward
- 588 connections in the brain: A DCM study of functional asymmetries. *NeuroImage*, 45(2), 453-
- 589 462. doi: <https://doi.org/10.1016/j.neuroimage.2008.12.041>
- 590 Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive
- 591 science. *Behavioral and Brain Sciences*, 36(3), 181-204. doi:
- 592 <http://doi.org/10.1017/S0140525X12000477>
- 593 Coll, M.-P., Viding, E., Rütgen, M., Silani, G., Lamm, C., Catmur, C., & Bird, G. (2017). Are we really
- 594 measuring empathy? Proposal for a new measurement framework. *Neuroscience &*
- 595 *Biobehavioral Reviews*, 83, 132-139. doi: <https://doi.org/10.1016/j.neubiorev.2017.10.009>
- 596 Davis, M. H. (1980). A multidimensional approach to individual differences in empathy.
- 597 Decety, J., & Lamm, C. (2007). The Role of the Right Temporoparietal Junction in Social Interaction:
- 598 How Low-Level Computational Processes Contribute to Meta-Cognition. *The Neuroscientist*,
- 599 13(6), 580-593. doi: <http://doi.org/10.1177/1073858407304654>
- 600 Decety, J., & Lamm, C. (2011). Empathy versus Personal Distress: Recent Evidence from Social
- 601 Neuroscience. In J. Decety & W. Ickes (Eds.), *The social neuroscience of empathy* (pp. 199 -
- 602 213): MIT Press.
- 603 Fallon, N., Roberts, C., & Stancak, A. (2020). Shared and distinct functional networks for empathy
- 604 and pain processing: A systematic review and meta-analysis of fMRI studies. *Social cognitive*
- 605 *and affective neuroscience*. doi: <https://doi.org/10.1093/scan/nsaa090>
- 606 Fan, Y., Duncan, N. W., de Greck, M., & Northoff, G. (2011). Is there a core neural network in
- 607 empathy? An fMRI based quantitative meta-analysis. *Neuroscience & Biobehavioral Reviews*,
- 608 35(3), 903-911. doi: <https://doi.org/10.1016/j.neubiorev.2010.10.009>

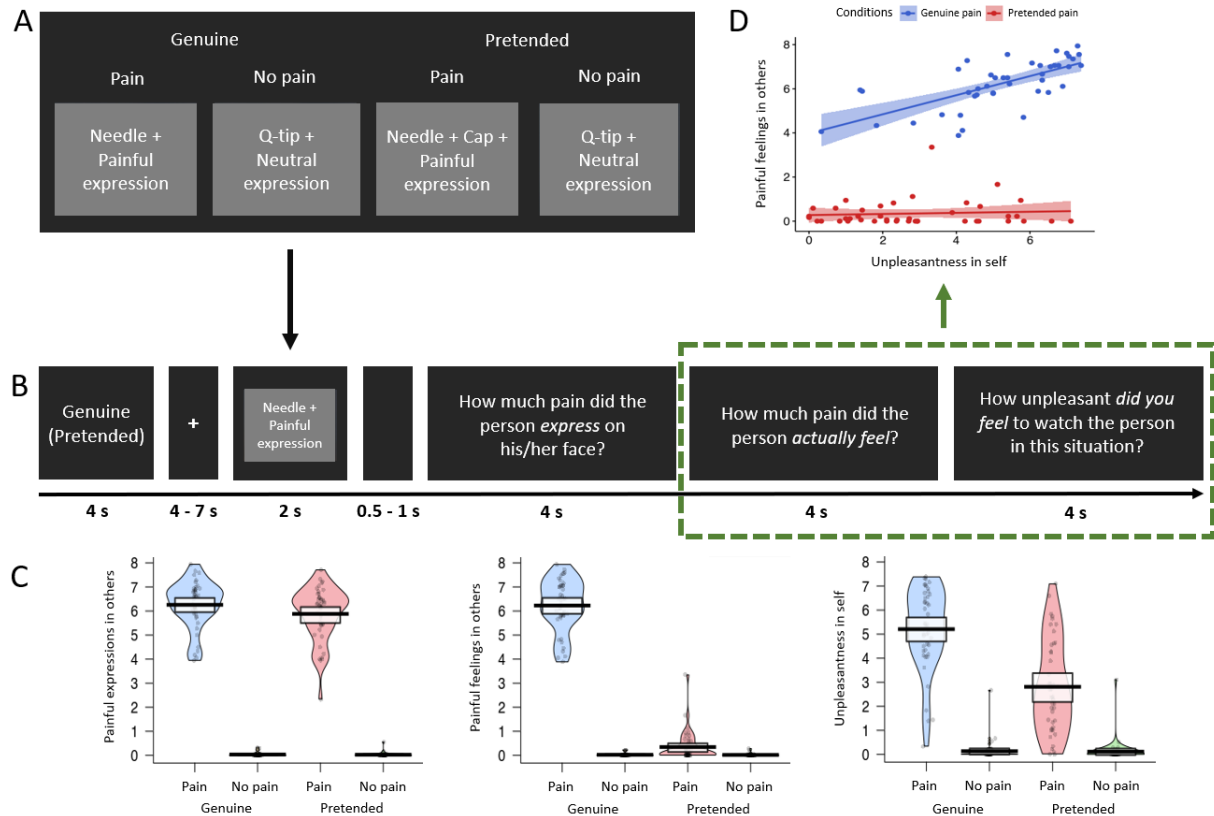
- 609 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G*Power 3: A flexible statistical power
610 analysis program for the social, behavioral, and biomedical sciences. *Behavior Research*
611 *Methods*, 39(2), 175-191. doi: <http://doi.org/10.3758/BF03193146>
- 612 Forbes, P. A. G., & Hamilton, A. F. d. C. (2020). Brief Report: Autistic Adults Assign Less Weight to
613 Affective Cues When Judging Others' Ambiguous Emotional States. *Journal of Autism and*
614 *Developmental Disorders*, 50(8), 3066-3070. doi: [http://doi.org/10.1007/s10803-020-04410-](http://doi.org/10.1007/s10803-020-04410-w)
615 [w](http://doi.org/10.1007/s10803-020-04410-w)
- 616 Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*,
617 11(2), 127-138. doi: <http://doi.org/10.1038/nrn2787>
- 618 Gläscher, J., & Gitelman, D. (2008). Contrast weights in flexible factorial design with multiple groups
619 of subjects. *SPM@ JISCMail. AC. UK) Sml, editor*, 1-12.
- 620 Gola, K. A., Shany-Ur, T., Pressman, P., Sulman, I., Galeana, E., Paulsen, H., Nguyen, L., Wu, T.,
621 Adhimoolam, B., Poorzand, P., Miller, B. L., & Rankin, K. P. (2017). A neural network
622 underlying intentional emotional facial expression in neurodegenerative disease.
623 *NeuroImage: Clinical*, 14, 672-678. doi: <https://doi.org/10.1016/j.nicl.2017.01.016>
- 624 Gorgolewski, K., Burns, C., Madison, C., Clark, D., Halchenko, Y., Waskom, M., & Ghosh, S. (2011).
625 Nipype: A Flexible, Lightweight and Extensible Neuroimaging Data Processing Framework in
626 Python. *Frontiers in Neuroinformatics*, 5(13). doi: <http://doi.org/10.3389/fninf.2011.00013>
- 627 Gu, X., & Han, S. (2007). Attention and reality constraints on the neural processes of empathy for
628 pain. *NeuroImage*, 36(1), 256-267. doi: <https://doi.org/10.1016/j.neuroimage.2007.02.025>
- 629 Hawco, C., Kovacevic, N., Malhotra, A. K., Buchanan, R. W., Viviano, J. D., Iacoboni, M., McIntosh, A.
630 R., & Voineskos, A. N. (2017). Neural Activity while Imitating Emotional Faces is Related to
631 Both Lower and Higher-Level Social Cognitive Performance. *Scientific Reports*, 7(1), 1244.
632 doi: <http://doi.org/10.1038/s41598-017-01316-z>

- 633 Hein, G., & Singer, T. (2008). I feel how you feel but not always: the empathic brain and its
634 modulation. *Current Opinion in Neurobiology*, *18*(2), 153-158. doi:
635 <https://doi.org/10.1016/j.conb.2008.07.012>
- 636 Hoffmann, F., Koehne, S., Steinbeis, N., Dziobek, I., & Singer, T. (2016). Preserved Self-other
637 Distinction During Empathy in Autism is Linked to Network Integrity of Right Supramarginal
638 Gyrus. *Journal of Autism and Developmental Disorders*, *46*(2), 637-648. doi:
639 <http://doi.org/10.1007/s10803-015-2609-0>
- 640 Holmes, E., Zeidman, P., Friston, K. J., & Griffiths, T. D. (2020). Difficulties with Speech-in-Noise
641 Perception Related to Fundamental Grouping Processes in Auditory Cortex. *Cerebral Cortex*,
642 *31*(3), 1582-1596. doi: <http://doi.org/10.1093/cercor/bhaa311>
- 643 James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning* (Vol.
644 112): Springer.
- 645 Jauniaux, J., Khatibi, A., Rainville, P., & Jackson, P. L. (2019). A meta-analysis of neuroimaging studies
646 on pain empathy: investigating the role of visual information and observers' perspective.
647 *Social cognitive and affective neuroscience*, *14*(8), 789-813. doi:
648 <https://doi.org/10.1093/scan/nsz055>
- 649 Kanske, P., Böckler, A., Trautwein, F.-M., Parianen Lesemann, F. H., & Singer, T. (2016). Are strong
650 empathizers better mentalizers? Evidence for independence and interaction between the
651 routes of social cognition. *Social cognitive and affective neuroscience*, *11*(9), 1383-1392. doi:
652 <http://doi.org/10.1093/scan/nsw052>
- 653 Lamm, C., Bukowski, H., & Silani, G. (2016). From shared to distinct self-other representations in
654 empathy: evidence from neurotypical function and socio-cognitive disorders. *Philosophical
655 transactions of the Royal Society of London. Series B, Biological sciences*, *371*(1686),
656 20150083. doi: <http://doi.org/10.1098/rstb.2015.0083>

- 657 Lamm, C., Decety, J., & Singer, T. (2011). Meta-analytic evidence for common and distinct neural
658 networks associated with directly experienced pain and empathy for pain. *NeuroImage*,
659 54(3), 2492-2502. doi: <https://doi.org/10.1016/j.neuroimage.2010.10.014>
- 660 Lamm, C., Meltzoff, A. N., & Decety, J. (2010). How do we empathize with someone who is not like
661 us? A functional magnetic resonance imaging study. *J. Cognitive Neuroscience*, 22(2), 362–
662 376. doi: <http://doi.org/10.1162/jocn.2009.21186>
- 663 Lamm, C., Rütgen, M., & Wagner, I. C. (2019). Imaging empathy and prosocial emotions.
664 *Neuroscience Letters*, 693, 49-53. doi: <https://doi.org/10.1016/j.neulet.2017.06.054>
- 665 Mars, R. B., Sallet, J., Schüffegen, U., Jbabdi, S., Toni, I., & Rushworth, M. F. S. (2011). Connectivity-
666 Based Subdivisions of the Human Right “Temporoparietal Junction Area”: Evidence for
667 Different Areas Participating in Different Cortical Networks. *Cerebral Cortex*, 22(8), 1894-
668 1903. doi: <http://doi.org/10.1093/cercor/bhr268>
- 669 Menard, S. (2002). *Applied logistic regression analysis* (Vol. 106): Sage.
- 670 Miska, N. J., Richter, L. M. A., Cary, B. A., Gjorgjieva, J., & Turrigiano, G. G. (2018). Sensory experience
671 inversely regulates feedforward and feedback excitation-inhibition ratio in rodent visual
672 cortex. *eLife*, 7, e38846. doi: <http://doi.org/10.7554/eLife.38846>
- 673 Pokorny, J. J., Hatt, N. V., Colombi, C., Vivanti, G., Rogers, S. J., & Rivera, S. M. (2015). The Action
674 Observation System when Observing Hand Actions in Autism and Typical Development.
675 *Autism Research*, 8(3), 284-296. doi: <https://doi.org/10.1002/aur.1445>
- 676 Power, J. D., Barnes, K. A., Snyder, A. Z., Schlaggar, B. L., & Petersen, S. E. (2012). Spurious but
677 systematic correlations in functional connectivity MRI networks arise from subject motion.
678 *NeuroImage*, 59(3), 2142-2154. doi: <https://doi.org/10.1016/j.neuroimage.2011.10.018>
- 679 Power, J. D., Mitra, A., Laumann, T. O., Snyder, A. Z., Schlaggar, B. L., & Petersen, S. E. (2014).
680 Methods to detect, characterize, and remove motion artifact in resting state fMRI.
681 *NeuroImage*, 84, 320-341. doi: <https://doi.org/10.1016/j.neuroimage.2013.08.048>

- 682 Rütgen, M., Seidel, E. M., Silani, G., Riechansky, I., Hummer, A., Windischberger, C., Petrovic, P., &
683 Lamm, C. (2015). Placebo analgesia and its opioidergic regulation suggest that empathy for
684 pain is grounded in self pain. *Proceedings of the National Academy of Sciences*, *112*(41),
685 E5638-E5646. doi: <https://doi.org/10.1073/pnas.1511269112>
- 686 Seth, A. K. (2013). Interoceptive inference, emotion, and the embodied self. *Trends in Cognitive*
687 *Sciences*, *17*(11), 565-573. doi: <https://doi.org/10.1016/j.tics.2013.09.007>
- 688 Silani, G., Lamm, C., Ruff, C. C., & Singer, T. (2013). Right Supramarginal Gyrus Is Crucial to Overcome
689 Emotional Egocentricity Bias in Social Judgments. *The Journal of Neuroscience*, *33*(39),
690 15466-15476. doi: <http://doi.org/10.1523/jneurosci.1488-13.2013>
- 691 Sladky, R., Friston, K. J., Tröstl, J., Cunnington, R., Moser, E., & Windischberger, C. (2011). Slice-timing
692 effects and their correction in functional MRI. *NeuroImage*, *58*(2), 588-594. doi:
693 <https://doi.org/10.1016/j.neuroimage.2011.06.078>
- 694 Steinbeis, N., Bernhardt, B. C., & Singer, T. (2015). Age-related differences in function and structure
695 of rSMG and reduced functional connectivity with DLPFC explains heightened emotional
696 egocentricity bias in childhood. *Social cognitive and affective neuroscience*, *10*(2), 302-310.
697 doi: <https://doi.org/10.1093/scan/nsu057>
- 698 Stephan, K. E., & Friston, K. J. (2010). Analyzing effective connectivity with functional magnetic
699 resonance imaging. *WIREs Cognitive Science*, *1*(3), 446-459. doi:
700 <https://doi.org/10.1002/wcs.58>
- 701 Xiong, R.-C., Fu, X., Wu, L.-Z., Zhang, C.-H., Wu, H.-X., Shi, Y., & Wu, W. (2019). Brain pathways of
702 pain empathy activated by pained facial expressions: a meta-analysis of fMRI using the
703 activation likelihood estimation method. *Neural regeneration research*, *14*(1), 172-178. doi:
704 <http://doi.org/10.4103/1673-5374.243722>
- 705 Zaki, J., Wager, T. D., Singer, T., Keysers, C., & Gazzola, V. (2016). The Anatomy of Suffering:
706 Understanding the Relationship between Nociceptive and Empathic Pain. *Trends in Cognitive*
707 *Sciences*, *20*(4), 249-259. doi: <https://doi.org/10.1016/j.tics.2016.02.003>

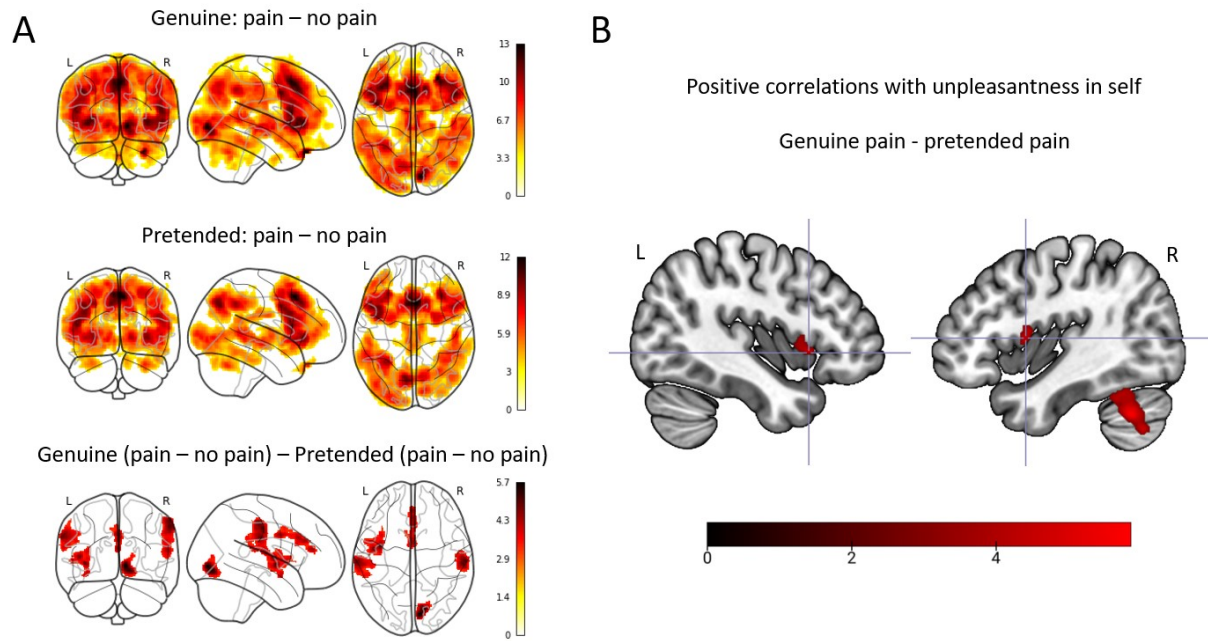
- 708 Zeidman, P., Jafarian, A., Corbin, N., Seghier, M. L., Razi, A., Price, C. J., & Friston, K. J. (2019a). A
709 guide to group effective connectivity analysis, part 1: First level analysis with DCM for fMRI.
710 *NeuroImage*, 200, 174-190. doi: <https://doi.org/10.1016/j.neuroimage.2019.06.031>
- 711 Zeidman, P., Jafarian, A., Seghier, M. L., Litvak, V., Cagnan, H., Price, C. J., & Friston, K. J. (2019b). A
712 guide to group effective connectivity analysis, part 2: Second level analysis with PEB.
713 *NeuroImage*, 200, 12-25. doi: <https://doi.org/10.1016/j.neuroimage.2019.06.032>
- 714 Zhang, M., Chung, S. H., Fang-Yen, C., Craig, C., Kerr, R. A., Suzuki, H., Samuel, A. D. T., Mazur, E., &
715 Schafer, W. R. (2008). A Self-Regulating Feed-Forward Circuit Controlling *C. elegans* Egg-
716 Laying Behavior. *Current Biology*, 18(19), 1445-1455. doi:
717 <https://doi.org/10.1016/j.cub.2008.08.047>
- 718 Zhao, Y., Rütgen, M., Zhang, L., & Lamm, C. (2021). Pharmacological fMRI provides evidence for
719 opioidergic modulation of discrimination of facial pain expressions. *Psychophysiology*, 58(2),
720 e13717. doi: <https://doi.org/10.1111/psyp.13717>
- 721 Zhou, F., Li, J., Zhao, W., Xu, L., Zheng, X., Fu, M., Yao, S., Kendrick, K. M., Wager, T. D., & Becker, B.
722 (2020). Empathic pain evoked by sensory and emotional-communicative cues share common
723 and process-specific neural representations. *eLife*, 9, e56929. doi:
724 <http://doi.org/10.7554/eLife.56929>
- 725
- 726



727 **Figure 1. fMRI experimental design and behavioral results.** (A) Overview of the experimental design
 728 with the four conditions genuine vs. pretended, pain vs. no pain. Examples show static images, while
 729 in the experiment participants were shown video clips. (B) Overview of experimental timeline. At the
 730 outset of each block, a reminder of “genuine” or “pretended” was shown (both terms are shown
 731 here for illustrative purposes, in the experiment either genuine or pretended was displayed). After a
 732 fixation cross, a video in the corresponding condition appeared on the screen. Followed by a short
 733 jitter, three questions about the video were separately presented and had to be rated on a visual
 734 analogue scale. These would then be followed by the next video clip and questions (not shown). (C)
 735 Violin plots of the three types of ratings for all conditions. Participants generally demonstrated
 736 higher ratings for painful expressions in others, painful feelings in others, and unpleasantness in self
 737 in the genuine pain condition than in the pretended pain condition. Ratings of all three questions
 738 were higher in the painful situation than in the neutral situation, regardless of whether in the
 739 genuine or pretended condition. The thick black lines illustrate mean values, and the white boxes
 740 indicate a 95% CI. The dots are individual data, and the “violin” outlines illustrate their estimated

741 density at different points of the scale. (D) Correlations of painful feelings in others and
742 unpleasantness in self for the genuine pain and the pretended pain (the relevant questions were
743 highlighted with a green rectangular). Results revealed a significant Pearson correlation between the
744 two questions in the genuine pain condition, but no correlation in the pretended pain condition. The
745 lines represent the fitted regression lines, bands indicate a 95% CI.

746



747

748 **Figure 2. Neuroimaging results: Mass-univariate analyses.** (A) Activation maps of genuine: pain – no

749 pain (top), pretended: pain - no pain (middle), and genuine (pain – no pain) – pretended (pain – no

750 pain) (bottom). As expected, we found brain activations in the bilateral alns, aMCC, and rSMG in all

751 three contrasts (except for the bottom contrast, where the right alns is only close to the significance

752 threshold). (B) The multiple regression analysis demonstrated significant clusters in the left (peak: [-

753 42, 15, -2]) and right anterior insular cortex (peak: [45, 5, 8]) for the ratings of unpleasantness in self.

754 All activations are thresholded with cluster-level FWE correction, $p < 0.05$ ($p < 0.001$ uncorrected

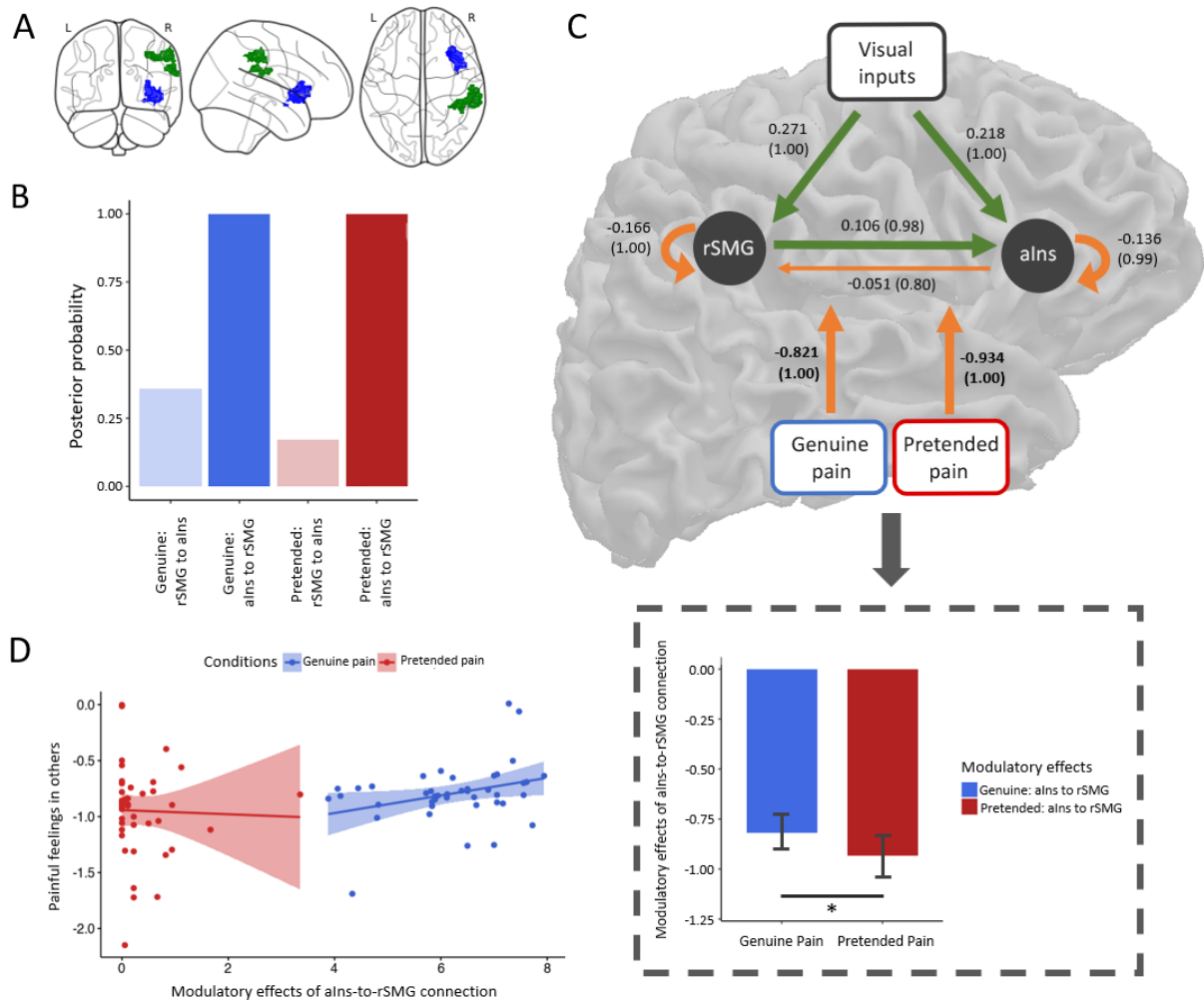
755 initial selection threshold).

756

757

758

759



760

761 **Figure 3. DCM results and brain-behavior analyses.** (A) ROIs included in the DCM: aIns (blue; peak:
 762 [33, 29, 2]) and rSMG (green; peak: [41, -39, 42]). (B) Posterior probability of modulatory effects for
 763 the genuine pain and the pretended pain. (C) The group-average DCM model. Green arrows indicate
 764 neural excitation, and orange arrows indicate neural inhibition. Importantly, we found strong
 765 evidence of inhibitory effects on the connection of aIns to rSMG for both the genuine pain condition
 766 and the pretended pain condition. Values without the bracket quantify the strength of connections
 767 and values in the bracket indicate the posterior probability of connections. All DCM parameters of
 768 the optimal model showed greater than a 95% posterior probability (i.e., strong evidence) except for
 769 the intrinsic connection of aIns to rSMG ($p_p = 0.80$). Paired sample t-test showed less inhibitory
 770 effects of the aIns-to-rSMG connection for the genuine pain than the pretended pain. This result is

771 highlighted with a grey rectangular. Data are mean \pm 95% CI. (D) The stepwise linear regression
772 model revealed a positive correlation between the inhibitory effect and painful feelings in others
773 (after accounting for the other two ratings) for genuine pain but no correlation for pretended pain.
774

775 **Table 1.** Results of mass-univariate functional segregation analyses in the MNI space. Region names
 776 were labeled with the AAL atlas, threshold $p < 0.05$ cluster-wise FWE correction (initial selection
 777 threshold $p < 0.001$, uncorrected). BA = Brodmann area, L = left hemisphere, R = right hemisphere.

778

Region label	BA	Cluster size	x	y	z	t-value
Genuine: pain - no pain						
Lingual_R	18	183732	11	-84	-3	13.38
Temporal_Pole_Sup_R	38		30	33	-33	13.31
Supp_Motor_Area_R	8		5	15	51	12.96
Supp_Motor_Area_R	8		3	17	50	12.92
Supp_Motor_Area_L	8		-5	17	48	12.56
Insula_L	45		-32	26	6	12.32
Insula_R	45		33	29	3	12.09
Frontal_Inf_Oper_R	44		51	14	15	12.01
Frontal_Inf_Oper_R	44		50	12	18	11.79
Precentral_L	6		-42	3	39	11.72
Fusiform_R	20	463	36	-5	-41	5.58
Pretended: pain - no pain						
Supp_Motor_Area_R	8	59665	5	20	48	11.80
Supp_Motor_Area_L	8		-6	18	50	11.14
Frontal_Inf_Oper_L	44		-50	15	15	10.39
Insula_R	45		33	29	0	9.81
Insula_L	45		-29	30	0	9.60
Frontal_Inf_Tri_R	44		47	15	26	9.21
Precuneus_L	7	35136	-9	-71	41	10.27
Parietal_Inf_L	39		-32	-51	41	9.39
Precuneus_R	7		9	-69	38	8.44
Temporal_Mid_L	21		-53	-47	5	7.67
Occipital_Mid_L	19		-44	-78	2	7.47
Parietal_Inf_R	39		39	-50	41	7.25
Temporal_Mid_R	22	12970	51	-20	-6	7.70
Lingual_R	17		12	-86	-2	7.40
Fusiform_R	37		47	-33	-27	5.32
Occipital_Mid_R	18		33	-86	3	5.23
Cingulum_Mid_R	23	1666	-3	-14	27	6.35
Cingulum_Mid_L	23		-3	-24	32	5.57

Temporal_Pole_Sup_R	47	589	32	35	-33	7.18
Frontal_Sup_Orb_R	11		17	41	-24	3.36
Genuine (pain – no pain) – pretended (pain – no pain)	19	18	24	-81	39	5.27
SupraMarginal_L	40	1877	-66	-21	32	4.94
Postcentral_L	1		-50	-21	26	3.75
SupraMarginal_R	40	1833	63	-20	42	5.09
Rolandic_Oper_R	40		59	-15	14	4.47
Insula_L	13	1299	-38	-3	-2	5.01
Rolandic_Oper_L	4		-45	-6	8	4.8
Cingulum_Ant_L	32	1138	0	41	17	4.54
Cingulum_Mid_R	32		2	24	32	4.45
Cingulum_Mid_L	24		0	2	35	4.43
Cingulum_Ant_R	8		2	32	27	4.42
Lingual_R	18	1003	9	-84	-3	5.72
Calcarine_R	17		18	-78	8	3.61
Insula_R	13	225	39	8	-3	3.91
Rolandic_Oper_R	13		41	0	11	3.77

779

780