# In amygdala we trust: different contributions of the basolateral and central amygdala in learning whom to trust

## Amygdala subnuclei orchestrate trust learning

Ronald Sladky [+,1,*] & Federica Riva [+,1], Lisa Rosenberger [1], Jack van Honk [2,3], Claus Lamm [1,*]

[1] Social, Cognitive and Affective Neuroscience Unit, Department of Cognition, Emotion, and Methods in Psychology, Faculty of Psychology, University of Vienna, Vienna, Austria

[2] Department of Psychology, Utrecht University, 3584 CS Utrecht, the Netherlands

[3] Department of Psychiatry and Mental Health, MRC Unit on Risk & Resilience in Mental Disorders, University of Cape Town, Observatory, 7925 Cape Town, South Africa

+ **Shared authorship.** RS and FR contributed equally to this manuscript.

* **Corresponding authors.** Ronald Sladky and Claus Lamm.
Social, Cognitive and Affective Neuroscience Unit
Department of Cognition, Emotion, and Methods in Psychology
Faculty of Psychology, University of Vienna
Address. Liebiggasse 5, 1010 Vienna, Austria
*Phone.* +43 1 4277 47130, *E-mail.* ronald.sladky@univie.ac.at, claus.lamm@univie.ac.at

**Declaration of interests.** The authors declare no competing financial interests.

25 pages, 4 figures, 0 tables
abstract: 250 words, introduction: 521 words, discussion: 1,577 words

## ABSTRACT

Human societies are built on cooperation and mutual trust, but not everybody is trustworthy. Research on rodents suggests an essential role of the basolateral amygdala (BLA) in learning from social experiences (Hernandez-Lallement J et al., 2016), which was also confirmed in human subjects with selective bilateral BLA damage as they failed to adapt their trust behavior towards trustworthy vs. untrustworthy interaction partners (Rosenberger LA et al., 2019). However, neuroimaging in neurotypical populations did not consistently report involvement of the amygdala in trust behavior. This might be explained by the difficulty of differentiating between amygdala's structurally and functionally different subnuclei, i.e., the BLA and central amygdala (CeA), which have even antagonistic features particularly in trust behavior (van Honk J et al., 2013). Here, we used fMRI of the amygdala subnuclei of neurotypical adults (n=31f/31m) engaging in the repeated trust game. Our data show that both the BLA and the CeA play a role and indeed differentially: While the BLA was most active when obtaining feedback on whether invested trust had been reciprocated or not, the CeA was most active when subjects were preparing their next trust decision. In the latter phase, improved learning was associated with higher activation differences in response to untrustworthy vs. trustworthy trustees, in both BLA and CeA. Our data not only translate to rodent models and support our earlier findings in BLA-damaged subjects, but also show the specific contributions of other brain structures in the amygdala-centered network in learning whom to trust, and better not to trust.

## SIGNIFICANCE STATEMENT

In this fMRI study, the central amygdala was found active during trust behavior planning, while the basolateral amygdala was active during outcome evaluation. When planning trust behavior, central and basolateral amygdala activation differences between the players was related to whether participants learned to differentiate the players' trustworthiness. Nucleus accumbens tracked whether trust was reciprocated but was not related to learning. This suggests learning whom to trust is not related to reward processing in the nucleus

51 accumbens but rather to engagement of the basolateral amygdala. This study overcomes

52 major empirical gaps between animal models and human neuroimaging and shows how

53 different amygdala subnuclei and connected areas orchestrate learning to form different

54 subjective trustworthiness beliefs about others and guide trust choice behavior.

## INTRODUCTION

56 Human societies are built on cooperation and mutual trust. On the individual level, trusting

57 another person entails potential rewards, but also risks if the other person is abusing our

58 trust to our own disadvantage. Thus, learning to distinguish the trustworthiness of an

59 interaction partner is important for successful social interactions. Research on rodents

60 suggests an essential role of the basolateral amygdala (BLA) in learning from social

61 experiences (Hernandez-Lallement J,van Wingerden M,Schäble S and Kalenscher T, 2016). In

62 line with this, we showed in a previous study that human participants with selective

63 bilateral BLA damage failed to adapt their trust behavior towards trustworthy vs.

64 untrustworthy interaction partners in a repeated trust game (Rosenberger LA,Eisenegger

65 C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J, 2019). However, functional

66 reorganization after developmental brain damage might confine the generalizability of these

67 findings to neurotypical populations. Neuroimaging in neurotypical populations indeed did

68 not consistently report involvement of the amygdala in trust behavior. This might be

69 explained by difficulties in differentiating between the amygdala's structurally and

70 functionally different subnuclei, i.e., the BLA and central amygdala (CeA), which have even

71 antagonistic features particularly in trust behavior (van Honk J,Eisenegger C,Terburg

72 D,Stein DJ and Morgan B, 2013).

73 The amygdala is widely regarded as paramount for social cognition (Adolphs R, 2010), but it

74 has been investigated as a uniform structure in the majority of human neuroimaging studies

75 (Gupta R et al., 2011). While this approach may be due to the limited spatial specificity of

76 functional MRI particularly in the ventral brain (Sladky R et al., 2013;Sladky R et al., 2018), it

77 ignores the structural and functional heterogeneity of this brain area and its subnuclei

78    (Balleine BW and Killcross S, 2006). Here, we overcame the limitations of previous research

79    by using an acquisition protocol optimized for imaging ventral brain areas (Robinson S et

80    al., 2004) in combination with a multiband EPI sequence with high spatial and temporal

81    resolution (Moeller S et al., 2010), allowing for a time-resolved analysis of amygdalar

82    subnuclei.

83    Our recent research in participants with basolateral amygdala lesions (Rosenberger

84    LA,Eisenegger C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J, 2019) proposed that a

85    network centered around the basolateral amygdala adaptively subserves learning to trust

86    and distrust others. Importantly, this novel insight was based on a trust game task in which

87    the participants repeatedly interacted with a trustworthy and an untrustworthy interaction

88    partner. The task thus allowed us to investigate the dynamics of trust formation, as well as

89    the role that different decision-making processes play in that. Here, using functional MRI in

90    a healthy neurotypical population we employ the exact same behavioral paradigm to

91    confirm and extend these findings to the specific functions of the separate subnuclei of the

92    amygdala and the networks they are a part of. Our main aims were to derive what role the

93    different subnuclei of the amygdala play for different aspects relevant in learning whom to

94    trust, and to link them to neural activation in other sub-cortical regions that are highly

95    connected with the amygdala (i.e., the bed nucleus of the stria terminalis, the nucleus

96    accumbens, and the substantia nigra/VTA) (Janak PH and Tye KM, 2015).

97  # MATERIALS AND METHODS

98  **PARTICIPANTS**

99    62 heathy, neurotypical volunteers (age=23.83±3.15 years, f/m=31/31), mostly

100    undergraduate students from Vienna, Austria were recruited. Exclusion criteria were

101    standard MRI exclusion criteria (e.g.: pregnancy, claustrophobia, and MRI-incompatible

102    implants, clinically significant somatic diseases), a history of psychiatric or neurological

103    disorders, substance abuse, psychopharmacological medication, less than nine years of

104    education, as well as not being task-naive (e.g., having already participated in a similar

105    study or being a psychology student). All participants provided written informed consent in

106    accordance with the Declaration of Helsinki and were compensated for their participation.

107    The study was approved by the ethics committee of the Medical University of Vienna (EK-

108    Nr. 1489/2015).

109    **PROCEDURE AND TASK**

110    This study was part of a bigger project including two additional tasks and a sample of older

111    adults, which are not reported in the current article. Participants were first invited to a

112    screening session where they performed some cognitive tasks and filled in some self-

113    reported measures of psychological traits. The main session was usually conducted within

114    two weeks from the screening session. Participants were welcomed to the MRI facility

115    (University of Vienna MR Center) together with two other participants, who were in fact

116    two confederates of the experimenter invited to play the trustees' role. After having signed

117    the consent form and filled in the MR safety questionnaire, participants and confederates

118    were introduced to the protocol of the whole session. Afterwards, they went through the

119    training of the three tasks, including the trust game. At the end of the training, participants

120    were required to answer some questions in order to make sure they understood the task.

121    Participants were finally placed into the MR scanner, while the confederates were putatively

122    playing the task in the computer room next to the scanner room.

123    The repeated trust game was adapted from our previous study (Rosenberger LA,Eisenegger

124    C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J, 2019) and programmed in z-Tree

125    (version 3.3.7; (Fischbacher U, 2007)). The script of this trust game is deposited online

126    (Rosenberger LA,Eisenegger C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J, 2019). In

127    short, two players per round, an investor and a trustee, exchange monetary units with the

128    aim to maximize their monetary outcome. In total, 40 rounds were played and the

129    participant always played the role of the investor, while the trustees were allegedly played

130    by the two confederates in an alternate randomized order. In reality, the actions taken by the

131    two trustees were preprogrammed in a way that one of the confederates was behaving in a

132    trustworthy and the other one in an untrustworthy way. Confederates/trustees were of

133    similar age and same gender as the participant. At the beginning of each round (i.e., 20 per

134    trustworthy condition and 20 per untrustworthy condition) both players received an

135    endowment of 10 monetary units. Then each round encompasses four phases. In the

136    *preparation* phase, participants are presented with the picture of the trustee's face they are

137    playing with in the current round. In the *investment* phase, participants invest (part of) their

138    endowment (at least 1 unit) and the investment is tripled and then transferred to the trustee.

139    During the *waiting* phase, the trustees ostensibly perform their back-transfers. Finally,

140    during the *outcome* phase, participants are presented with the back-transfer outcome. In the

141    first two rounds, both the trustworthy and untrustworthy trustees back-transferred the same

142    amount of the money invested to the participants. In the following rounds, the trustworthy

143    trustee always back-transferred as much or more than the money invested by the player,

144    whereas the untrustworthy trustee always back-transferred less than or as much as the

145    money invested by the investor. The sums invested by the participants were considered as a

146    measure of trust given to the two trustees by the participants and used as the main variable

147    of interest. Points earned throughout the task were transformed to Euros and added to the

148    participants' compensation.

149    At the end of the task, participants were presented with the trustees' picture and were asked

150    to rate them on four adjectives: trustworthiness, fairness, attractiveness, and intelligence

151    (original German: *Wie vertrauenswürdig/attraktiv/intelligent/fair haben Sie den/die Teilnehmer/in*

152    *wahrgenommen?*). Ratings were provided on visual analogue scales and transformed off-line

153    to a numerical range between -10 and +10.

154    **FUNCTIONAL MRI DATA ACQUISITION AND PROCESSING**

155    MRI acquisitions were performed on a Skyra 3 Tesla MRI scanner (Siemens Healthineers,

156    Erlangen, Germany) using the manufacturer's 32 channel head coil at the MR Center of the

157    University of Vienna. In a single session, one run of the repeated trust game was performed

158    by the participant while we performed functional MRI using a gradient echo T2*-weighted

159    echo planar image sequence with the following parameters: MB-EPI factor=4, TR/TE =

160    704/34 ms, $2.2 \times 2.2 \times 3.5$ mm$^3$, $96 \times 92 \times 32$ voxels, flip angle=50°, n<2400 volumes.

161    Data processing and analyses of the functional MRI data were performed in SPM (SPM12,

162    http://www.fil.ion.ucl.ac.uk/spm/software/spm12/) and the Python projects nipype

163    (http://nipy.org/nipype) and nilearn (http://nilearn.github.io). Preprocessing comprised

164    slice-timing correction (Sladky R et al., 2011), realignment, non-linear normalization of the

165    EPI images to MNI space (final resolution = $1.5 \times 1.5 \times 1.5$ mm$^3$) using ANTs (Avants BB et

166    al., 2011), and spatial smoothing with a 6 mm FWHM Gaussian kernel.

167    **EXPERIMENTAL DESIGN AND STATISTICAL ANALYSES**

168    **Behavioral data analysis**

169    It is commonly understood that participants' investment behavior is a behavioral expression

170    of how they judged the trustees' trustworthiness and changes reflect the extent to which

171    they updated their beliefs (Bellucci G et al., 2017;Chang LJ et al., 2010;Rosenberger

172    LA,Eisenegger C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J, 2019). This *objective*

173    measure of trust was used to distinguish between learners and non-learners (using the

174    median as cut-off value) and for a Spearman correlation analysis between the subjective

175    ratings (trustworthiness, fairness, attractiveness, and intelligence) and the BOLD response in

176    the amygdala.

177    **Functional MRI data analysis**

178    First-level analyses of the data were implemented using nipype and performed using

179    SPM12's GLM approach. The GLM design matrix encompassed individual regressors for

180    each of the 4 task phases (i.e., preparation, investment, waiting, and outcome) and each of

181    the 2 interaction partners (trustworthy and untrustworthy, resulting in 8 effects of interest.

182    Additionally, 6 realignment parameters were added as nuisance regressors to account for

183    residual head motion effects. Second-level analyses of the data were implemented using

184    nipype and performed using SPM12's group-level approach for visual inspection of the

185    whole brain results.

186    Volume of interest analyses were performed on the mean timeseries extracted using

187    nilearn's fit_transform from anatomical masks from the BLA, CeA (Tyszka JM and Pauli

188     WM, 2016), NAc (AAL Atlas), BNST (Torrisi S et al., 2015), and SN/VTA (Talairach atlas

189     transformed to MNI space). To investigate phase-dependent activation, timeseries analyses

190     were conducted using custom python scripts that reproduced SPM's default GLM analysis,

191     using SPM's canonical HRF to convolve the regressors and a high-pass filter with the default

192     $f=1/128$ Hz cut-off frequency to account for signal drifts. Comparisons between learners and

193     non-learners were performed using two-sampled $t$-tests and based on their Spearman

194     correlations.

195     To verify that sensitivity of the fMRI dataset was sufficient to distinguish between BLA and

196     CeA activation, a functional connectivity analysis was conducted. Task fMRI data were

197     corrected for white matter and CSF signal and task effects (Ganger S et al., 2015) using

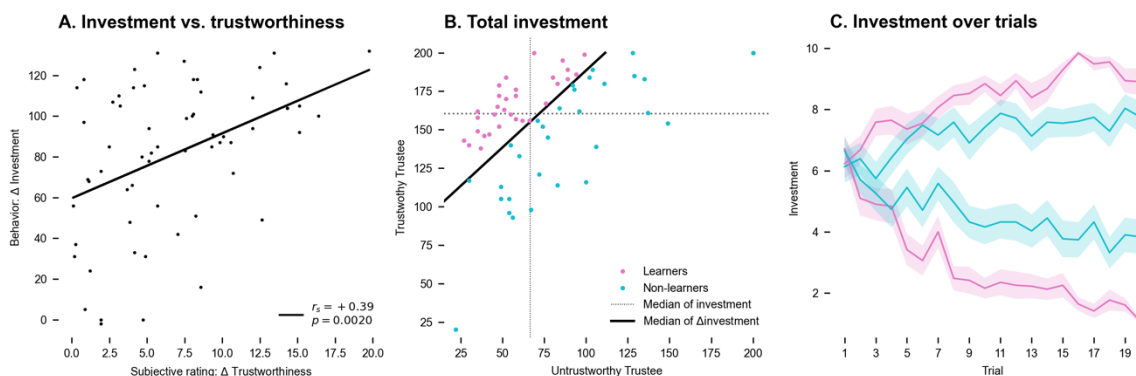198     regression before estimation of the functional connectivity maps of the BLA and CeA seeds.

## RESULTS

200     Participants played the repeated trust game inside the MRI scanner with a trustworthy and

201     an untrustworthy trustee, both simulated (2×20 rounds). In general, participants were able

202     to adapt their trust behavior, i.e., investments in the trust game, to the trustworthy and the

203     untrustworthy trustee. However, there was a marked variability within our study sample,

204     which allowed for a partition into a *learner* and *non-learner* sub-group (FIGURE 1). The task

205     consisted of four different task phases (i.e., the *preparation*, *investment*, *waiting*, and *outcome*

206     phase). A detailed time-resolved analysis of the BLA and CeA revealed that activation

207     changed over the course of the different task phases. We found maximum BLA activation in

208     the *outcome* evaluation phase and maximum CeA activation in the *preparation* phase. Yet,

209     there was no overall BLA and CeA activation difference between the trustworthy or

210     untrustworthy trustee in any of the task phases (FIGURE 2). However, when differentiating

211     between learners and non-learners, we observed more activation in the BLA and the CeA for

212     the untrustworthy trustee during the *introduction* phase of a trust game round (FIGURE 3).

213     Additionally, while nucleus accumbens (NAc), substantia nigra and ventral tegmental area

214     (SN/VTA), and bed nucleus of the stria terminalis (BST) activity was increased for the

215     trustworthy trustee during *outcome* evaluation, there was no group difference between

216     learners and non-learners (**FIGURE 4**).

## BEHAVIORAL RESULTS

218     Marked trust differences emerged across the whole sample in the investment behavior

219     towards the trustworthy as opposed to the untrustworthy trustee, with participants

220     generally investing more in the trustworthy trustee on average, and increasingly so over the

221     course of the repeated rounds of the task (**FIGURE 1B & C**). Morever, we find that individual

222     differences in behavioral trust($\Delta$ *investment = investment*$_{trustworthy}$ - *investment*$_{untrustworthy}$) showed

223     a positive correlation with subjective trustworthiness ratings ($\Delta$ *trustworthiness =*

224     *trustworthiness*$_{trustworthy}$ - *trustworthiness*$_{untrustworthy}$), $r_s = +0.39$, $p=0.002$ (**FIGURE 1A**). On the

225     subjective level, the trustworthy trustee was rated as significantly more trustworthy, fair,

226     and intelligent than the untrustworthy trustee (all $p<0.05$, Bonferroni corrected), but not as

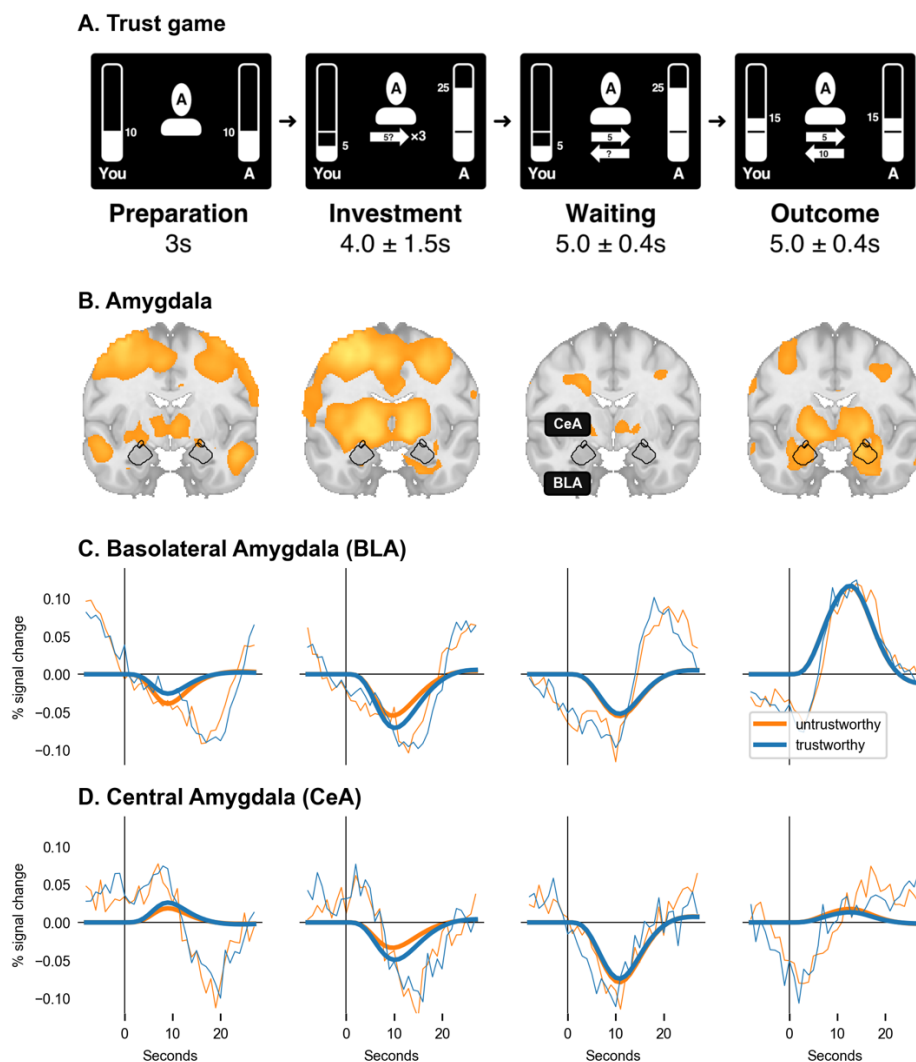227     more attractive (n.s., after Bonferroni correction).

228



FIGURE 1 | **A. Investment vs. trustworthiness.** Behavioral trust ($\Delta$ investment) correlates with subjective ratings ($\Delta$ trustworthiness rating), $r_s$=+0.39, $p$=0.002. **B. Participants' investment behavior.** In total, participants invested more in the trustworthy trustee. The difference between the investment into the trustworthy and untrustworthy trustee ($\Delta$ investment) was used to median-split the population into a subgroup that learned to differentiate (learners, magenta color) and those who did not (non-learners, cyan color). **C. Participants' investment behavior over time.** After a few trials, learners adapted their investment behavior to favor the trustworthy trustee. This differentiation was reduced in non-learners. Plot displays mean and SEM.

## NEUROIMAGING RESULTS

239    We find that different subnuclei of the amygdala were engaged in the trust game show

240    increased activation during different phases of the task paradigm. This suggests that they

241    are supposedly related to different aspects and processes required by the formation of trust.

242    The two subnuclei that played the most specific role (**FIGURE 2B**) were the basolateral (BLA)

243    and the central amygdala (CeA). Notably, the activation differences in these subnuclei and

244    the validity of our analysis approach is supported by differences in their functional

245    connectivity profiles, determined in our data. While the BLA connected to sensory

246    integration areas and lateral PFC, the CeA connected to the ventral striatum, including the

247    nucleus accumbens, and areas in the medial PFC (**SUPPLEMENTARY FIGURE 1**). The role of

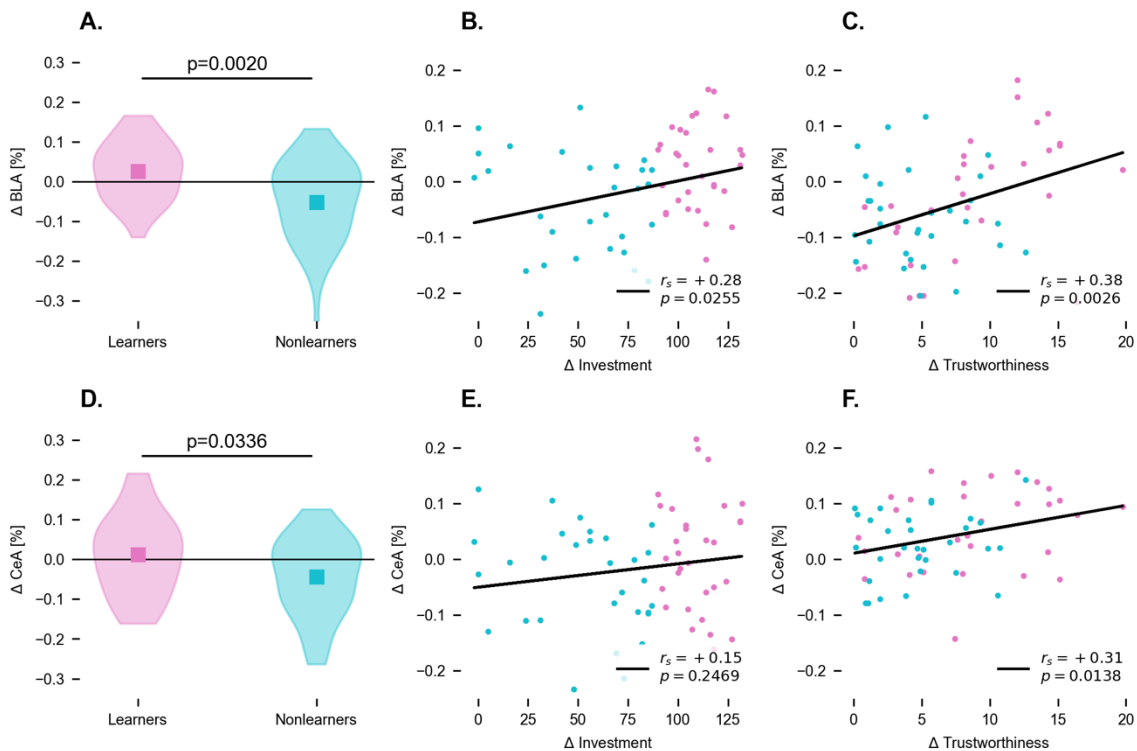248    these subnuclei in the different task phases is as follows.



249

10

250 **FIGURE 2 | A. fMRI implementation of the trust game.** Inside the MRI scanner,
251 participants played the repeated trust game alternating with a (simulated) trustworthy
252 and an untrustworthy trustee (2×20 rounds). *Preparation Phase.* Participants were
253 presented with the face of the trustee they played with in this round. Both received an
254 endowment of 10 points at the outset of each round. *Investment Phase.* Participants
255 were asked to select an amount of 1 to 10 points to invest in the present trustee. The
256 amount invested was tripled and added to the trustee's account. *Waiting Phase.* While
257 the trustees made their decision, the participant needed to wait. *Outcome Phase.*
258 Finally, the trustee transferred back points to the participant, resulting in a non-negative
259 outcome for the trustworthy (as shown in the example) and a non-positive outcome for
260 the untrustworthy trustee. **B. Statistical parametric maps (SPMs) and outline of the**
261 **anatomically defined Volumes of Interest (VOIs) of BLA and CeA.** SPMs show contrast for
262 both trustees combined vs. baseline and are thresholded at p<0.001 for display
263 purposes. **C & D. Time course of BLA and CeA BOLD responses.** CeA but not BLA was
264 activated during the *preparation* phase, while BLA but not CeA was activated during the
265 *outcome* phase. There were no activation differences between the trustworthy trustee
266 (blue) and the untrustworthy trustee (orange). Thick lines represent the estimated BOLD
267 model and fine lines represent the actual data (average VOI time courses).

268 In the *preparation* phase, activity in the BLA was reduced ($T=-8.9$, $p<0.0001$) and in the CeA

269 increased ($T=9.9$, $p<0.0001$), compared to the fixation baseline. Both BLA and CeA activity

270 were reduced during *investment* ($T=-15.4$, $p<0.0001$ and $T=-9.5$, $p<0.0001$) and *waiting* phase

271 ($T=-9.7$, $p<0.0001$ and $T=-13.0$, $p<0.0001$). During the *outcome* evaluation phase, activity in the

272 BLA was increased ($T=14.3$, $p<0.0001$), while it was reduced in the CeA ($T=-3.1$, $p=0.002232$)

273 (**FIGURE 2C & 2D**). All reported *p*-values survive Bonferroni correction for multiple

274 comparisons (at p<0.05 Bonferroni FWE-corrected). Note that these response patterns were

275 irrespective of whether a participant played with a trustworthy or untrustworthy trustee, as

276 there were no significant differences between these two conditions. These findings thus

277 relate to the general role of the amygdala subnuclei in the different parts of the task, and the

278 overall processes and subfunctions engaged by the trust decision.

279 As a next step, we aimed to pinpoint how the engagement of the amygdala was related to

280 differential evaluations of trustworthiness, and the resulting trust behavior toward the two

281 trustees. Individual difference analyses showed a relationship between BLA and the CeA

282 activation in the *preparation* phase and subjective trustworthiness and behavioral trust

283 measures. More specifically, we first used a median split of *Δ investment* (**FIGURE 1B**) to

284 distinguish learners from non-learners (i.e., those who adjusted their investment behavior

285    less to the trustworthiness of the trustee), and then assessed how they differed in their

286    amygdala activations. A Mann-Whitney *U*-test showed that during the *preparation phase*, the

287    activation difference between untrustworthy–trustworthy trustee was significantly larger in

288    learners than in non-learners in the BLA ($p$=0.0020, $u$=275.0, **FIGURE 3A**) and in the CeA

289    ($p$=0.0336, $u$=350.0, **FIGURE 3D**). Moreover, the BLA activation differences between

290    untrustworthy vs. trustworthy trustee in this phase correlated positively with behavioral

291    trust (*Δ investment)*, $r_s$=+0.28, $p$=0.0255 (**FIGURE 3B**) and subjective trustworthiness (*Δ*

292    *trustworthiness)*, $r_s$=+0.38, $p$=0.0026 (**FIGURE 3C**). CeA activation differences correlated with

293    subjective trustworthiness ratings, *Δ trustworthiness*, $r_s$=+0.31, p=0.0138 (**FIGURE 3F**), but not

294    with behavioral trust, *Δ investment* (**FIGURE 3E**). While considering whether to trust or

295    distrust a trustee, the CeA in learners thus seems primarily linked to evaluations of

296    trustworthiness, whereas the BLA is additionally relevant for the actual behavioral outcome

297    as well as whether someone efficiently learns to adapt behavior to the actually reciprocated

298    trust or not. Moreover, these relationships are driven by stronger engagement for rounds

299    with the untrustworthy (compared to the trustworthy) trustee, suggesting that what is

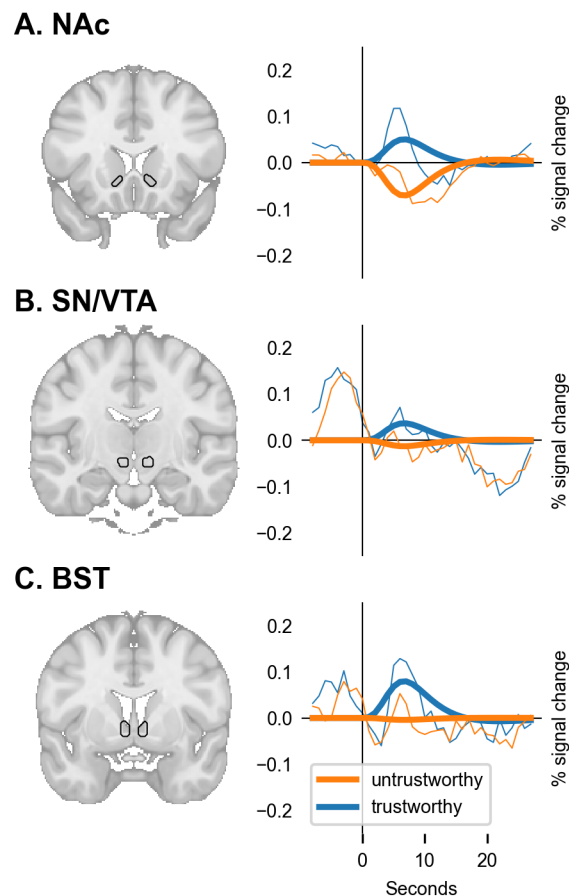300    coded is rather the absence than the presence of trust.

**FIGURE 3 | Activation differences between untrustworthy and trustworthy trustee in the *preparation* phase.** BLA activation differences (contrast: untrustworthy - trustworthy) were higher for learners (magenta) vs. non-learners (cyan) (**A**), correlated with investment differences (**B**) and post-experiment subjective trustworthiness rating differences (**C**). The same relationship was found for CeA (**D & F**), except the correlation with investment differences was not significant (**E**).

The neural responses in the *preparation* phase mainly provide insights into how the acquired information about a trustee's trustworthiness drives the decisions of participants. The activation in the *outcome* evaluation phase, on the other hand, tells us about how this information is acquired and possibly updated. As outlined above, we observed overall activation in the BLA during *outcome* evaluation phase (**FIGURE 2**), and this may be linked to reward processing (Lüthi A and Lüscher C, 2014). Surprisingly, though, we did not find differences between the trustworthy and untrustworthy trustee in the BLA or CeA in the outcome phase, and neither did we find correlations with trust behavior and trustworthiness rating. We thus extended our analyses to subcortical regions with particularly strong anatomical and functional connections to the amygdala. These were the nucleus accumbens (NAc), as well as the dopaminergic midbrain, comprising substantia

319    nigra and the ventral tegmental area (SN/VTA), relevant for encoding reward, and the bed

320    nucleus of the stria terminalis (BST), relevant for encoding threat (Avery SN et al.,

321    2016;Clauss JA et al., 2019;Siminski N et al., 2020).

322    When the *outcome* of the trustee decision was presented, higher activation in the NAc,

323    SN/VTA, and BST were observed for the trustworthy compared to the untrustworthy

324    trustee (NAc $t$=+7.21, $p$<0.0001, SN/VTA: $t$=+3.31, $p$<0.0010, and BST $t$=+4.38, $p$<0.0001)

325    (FIGURE 4). Moreover, the gain or loss (i.e., *back-transfer - investment amount*) correlated with

326    NAc ($r_s$=+0.19, $p$<0.0001) and BST ($r_s$=+0.10, $p$<0.0001), but this was irrespective of the

327    activation difference between trustworthy and untrustworthy trustee.



328

329    FIGURE 4 | More activity for the trustworthy (blue) vs. untrustworthy trustee (orange)
330    during the *outcome* event (A) in the nucleus accumbens (NAc), T=+7.21, p<0.0001, (B)
331    the substantia nigra (SN) and ventral tegmental area (VTA), T=+3.31, p<0.0010, and (C)

14

332    the bed nucleus of the stria terminalis (BST), T=+4.38, p<0.0001. Thick lines represent
333    the estimated BOLD model.

## DISCUSSION

335    Our previous study in BLA-damaged participants highlighted that the BLA is indispensable

336    for learning to differentiate between trustworthy and untrustworthy trustees in the trust

337    game (Rosenberger LA,Eisenegger C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J,

338    2019). This has important implications for our understanding of social decision-making in

339    humans and, most likely, other mammals (O'Connell LA and Hofmann HA, 2012).

340    However, extending these findings to the neural networks connected to the amygdala in

341    healthy, neurotypical, human participants is of the essence. Here we confirm the relevance

342    of the BLA for distinguishing between trustworthy and untrustworthy trustees based on

343    previous experience and how, in conjunction with the CeA, it plays a role in the guiding of

344    trust behavior. Specifically, BLA activity was increased during the processing of the

345    outcome of the trustee's behavior but unselectively for trustworthy vs. untrustworthy

346    trustee. Instead, we found increased activation in the NAc, BST, and SN/VTA for the

347    trustworthy vs. untrustworthy trustee during outcome processing. Importantly, here we did

348    not observe an activation difference between learners and non-learners. This could indicate

349    that learners and non-learners processed the outcome in a similar fashion, suggesting that

350    their understanding of the task and motivation were comparable. This further highlights the

351    central role of the BLA for trust learning.

352    Indeed, we found the BLA to be most active during outcome evaluation, i.e., when

353    participants learned whether their trust was reciprocated or not, suggesting that it plays an

354    important role in acquiring beliefs about the trustworthiness of others. It appears, however,

355    that the BLA is not directly involved in building specific outcome expectations during the

356    *waiting* and *evaluation* phase. The BOLD response in the BLA was not modulated by the

357    trustworthiness or the trustees' back-transfer amount, unlike activity in the NAc, SN, and

358    BST. This highlights that the BLA, although indispensable for learning whom to trust

359    (Rosenberger LA,Eisenegger C,Naef M,Terburg D,Fourie J,Stein DJ and van Honk J, 2019),

15

360    as indicated by our previous research, is only a component of a complex brain network for

361    reward processing and social evaluation.

362    In addition, we found that while participants prepared for their next investment, the BLA

363    together with the CeA exhibits increased activation for the untrustworthy trustee.

364    Importantly, this activation difference was only found in those participants who learned to

365    differentiate between the trustees, indicating its role in (1) *guiding trust behavior* as BLA

366    activation differences directly precede the participant's investment behavior and also (2) in

367    *trustworthiness evaluation*, as BLA and CeA BOLD responses correlated with the subjective

368    rating after the experiment.

369    Nowadays, it is a well-established finding that a sub-population of BLA's neurons

370    selectively responds to reward, whereas other sub-populations either only respond to

371    aversive stimuli (Pryce CR, 2018), or selectively increase their firing rate when the rewarding

372    or aversive stimulus was unexpected, i.e., not predicted (Belova MA et al., 2007) (which

373    means that something novel has to be learned about the environment). In the context of our

374    findings, this view supports the notion that the BLA is relevant for encoding both the

375    rewarding behavior of the trustworthy trustee and the aversive behavior of the

376    untrustworthy trustee. Additionally, we can speculate that optimal performance in the trust

377    game does not only rely on reward learning and threat detection, but also on predicting

378    affective consequences based on abstract information. Supporting evidence for this theory

379    can be found in a recent study in a patient with acquired complete bilateral amygdala

380    lesions (patient SM, 49 years old, female), who showed impairments in making good

381    predictions about what kind of written statements will induce fear (Cardinale EM et al.,

382    2021).

383    The fact that we did not observe any habituation in any of the amygdala subregions

384    (SUPPLEMENTARY FIGURE 2) indicates that the BLA not only responds to novel stimuli but is

385    relevant for the continuous encoding and updating of information of social experiences. In

386    the light of the recent debate on amygdala BOLD signal habituation (Geissberger N et al.,

16

387  2020;Infantolino ZP et al., 2018;McDermott TJ et al., 2020;Plichta MM et al., 2012;Sladky R et

388  al., 2012) this finding could be important for the development of additional tasks that

389  robustly activate the amygdala.

390  While BLA's activation during outcome evaluation suggests its involvement in

391  discriminating and tracking outcome-specific effects, the CeA is involved in general

392  motivational aspects of reward-related events (Corbit LH and Balleine BW, 2005) and, thus,

393  might not play a role in the actual learning process in the *outcome* phase. Instead, we found it

394  active during the *preparation* phase, which immediately preceded the *investment* phase. This

395  could indicate that the CeA is regulated by the BLA output, which has been demonstrated

396  before for a different task in a cross-species model (Terburg D et al., 2018). As CeA activity

397  was increased before the participant's investment, it might play a role in controlling trust

398  behavior. More importantly, CeA activity during the preparation phase correlated with the

399  subjective rating of trustworthiness of the trustee, indicating that it could be relevant for

400  encoding the affective value attached to the trustee.

401  During *outcome* evaluation, we observed increased activation in the bed nucleus of the stria

402  terminalis (BST), which, together with the CeA, is considered the *extended amygdala* complex

403  (Alheid G and Heimer L, 1988;de Olmos JS and Heimer L, 1999). The BST has been

404  suggested to play a role in both reward processing and social cognition (O'Connell LA and

405  Hofmann HA, 2011) and exhibits strong connections to the NAc (Avery SN et al., 2014).

406  While the CeA is associated with fast fear responses (e.g., startle reflex), the BST is

407  responsible for slower affective learning processes (Gewirtz JC et al., 1998) and has been

408  linked to adaptive and maladaptive responses to sustained stress and threat (Avery

409  SN,Clauss JA and Blackford JU, 2016;Somerville LH et al., 2013). Of note, the BST plays a

410  particular role in dealing with unpredictable threat (Goode TD et al., 2019), which could be

411  the case in an uncertain social investment. However, these two views are still part of

412  ongoing debates (Pedersen WS et al., 2019;Shackman AJ and Fox AS, 2016). Most recently,

413  the BST was shown to be more involved in fear-related anticipation processes, whereas the

414  CeA was linked to threat confrontation (Siminski N,Böhme S,Zeller J,Becker M,Bruchmann

17

415 M,Hofmann D,Breuer F,Mühlberger A,Schiele M and Weber H, 2020). In this study we

416 found the BST to be involved in the *outcome* evaluation phase. Based on the literature, it

417 could be expected that the BST would show more activation for the aversive *untrustworthy*

418 trustee, which was not the case. Instead, we observed that the BOLD responses of BST and

419 NAc were both more activated by the trustworthy trustee. The NAc and other striatal areas

420 are known to be involved in evaluating the trustees trustworthiness based on their back-

421 transfer behavior (Baumgartner T et al., 2008;Delgado MR et al., 2005;King-Casas B et al.,

422 2005) and amygdala to NAc coactivation is relevant for social decision making (Haruno M et

423 al., 2014). Rodent research has shown that BLA to NAc connections mediate reward learning

424 (Namburi P et al., 2015;Sesack SR and Grace AA, 2010). Importantly, stimulus-evoked

425 excitation of NAc neurons depends on input from the BLA and is required for dopamine to

426 enhance the stimulus-evoked firing of NAc neurons, ultimately, leading to reward-seeking

427 behavior (Ambroggi F et al., 2008). This could mean that both regions might engage in a

428 synergetic fashion, where the NAc would be particularly relevant for tracking rewards. The

429 BST, on the other hand, could be responsible for increasing arousal as generous investments

430 in the trustworthy trustee also entail a potential threat of betrayal. These findings suggest a

431 functional dissociation between reward and risk evaluation based on the observed outcome

432 of one's behavior, which appeared to be comparable in non-learners, and the mechanisms of

433 trust learning.

434 In sum, we confirm that the BLA is indeed involved in learning whom to trust and that

435 observations from amygdala-lesioned participants can be translated to healthy neurotypical

436 participants. Additionally, our fine-grained, time-resolved analyses of the amygdala

437 subnuclei and the functionally-connected brain areas provide important insights into

438 different cognitive mechanisms involved in trust learning. We found that the BLA is

439 relevant for *discriminating* between trustworthy and untrustworthy trustees based on

440 previous experience and for *optimizing trust behavior*. Only in those participants who learned

441 to optimize their investments, we found selectively more activation in the BLA during the

442 planning of a new investment that required trust. The BLA was also active during outcome

443 evaluation suggesting its involvement in the process of *belief formation* based on the trustees'

444 back-transfer amount. As we did not observe a difference between the trustworthy or

445 untrustworthy trustee, we can assume that encoding of potential rewards and risks is

446 mediated by the NAc and BST, respectively, which showed a selectively increased activity

447 for the trustworthy trustee or an increased investment. Finally, the CeA is known to receive

448 inputs from the BLA and BST, and exhibited the largest BOLD response during the *planning*

449 phase. CeA activity did not correlate with the participant's trust behavior, however, there

450 was a correlation with the participant's subjective belief of the trustees' trustworthiness. This

451 suggests that the CeA could encode subjective value, possibly also indirectly affecting trust

452 behavior via the BLA. Taken together, our work suggests that there is a high demand for

453 translational work on the amygdala, its subnuclei, and connected brain regions. Based on

454 the present results, we propose that careful variations of the trust game in combination with

455 computational modeling may serve as an experimental model to further uncover the neural

456 mechanisms underlying human social cognition.

## ACKNOWLEDGMENTS

## SUPPLEMENTARY INFORMATION

463 Supplemental Information can be found online at [tbc].
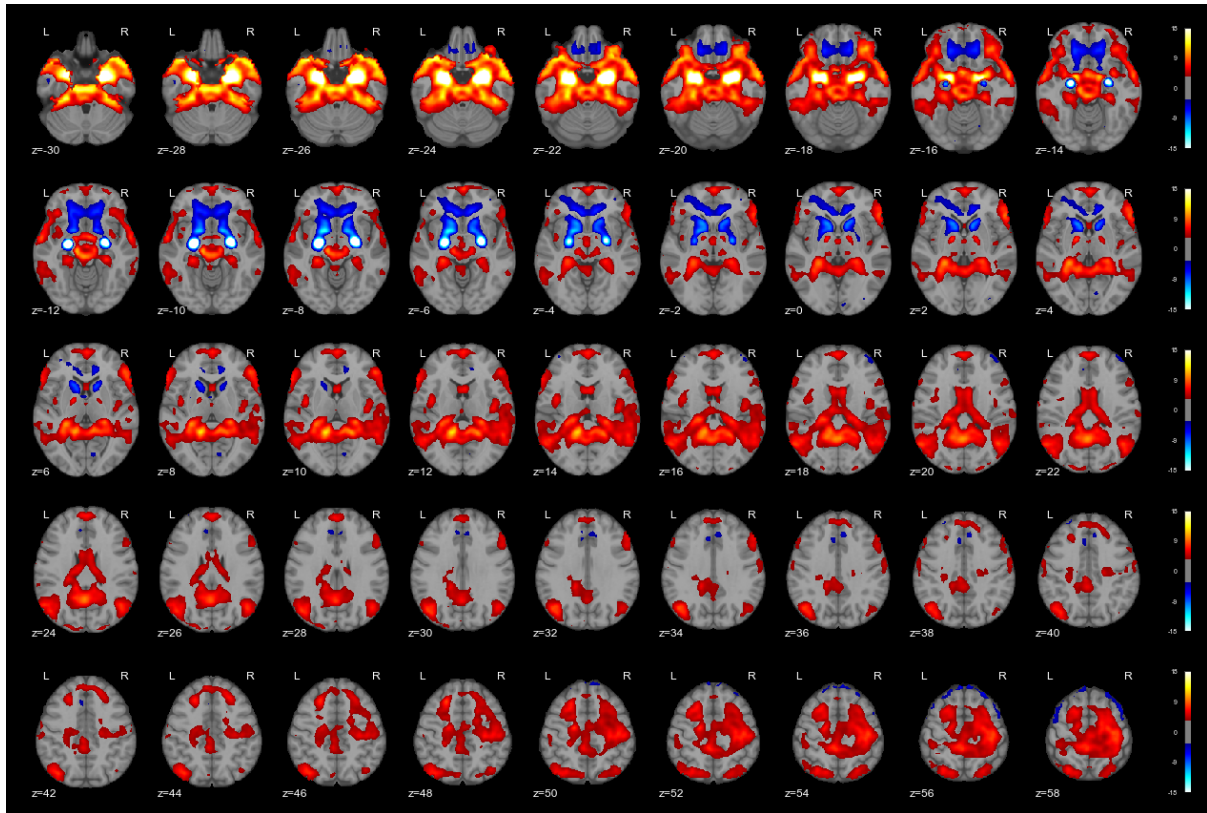
## AUTHOR CONTRIBUTIONS

465 Conceptualization and Methodology, R.S., F.R., L.R., J.v.H., C.L.; Investigation, F.R.; Formal

466 Analysis, R.S., F.R.; Writing – Original Draft, R.S., F.R., L.R., J.v.H., C.L.; Writing – Review &

467 Editing, R.S., F.R., L.R., J.v.H., C.L.; Funding Acquisition, C.L.

## DECLARATION OF INTERESTS

468

469 The authors declare no competing interests.
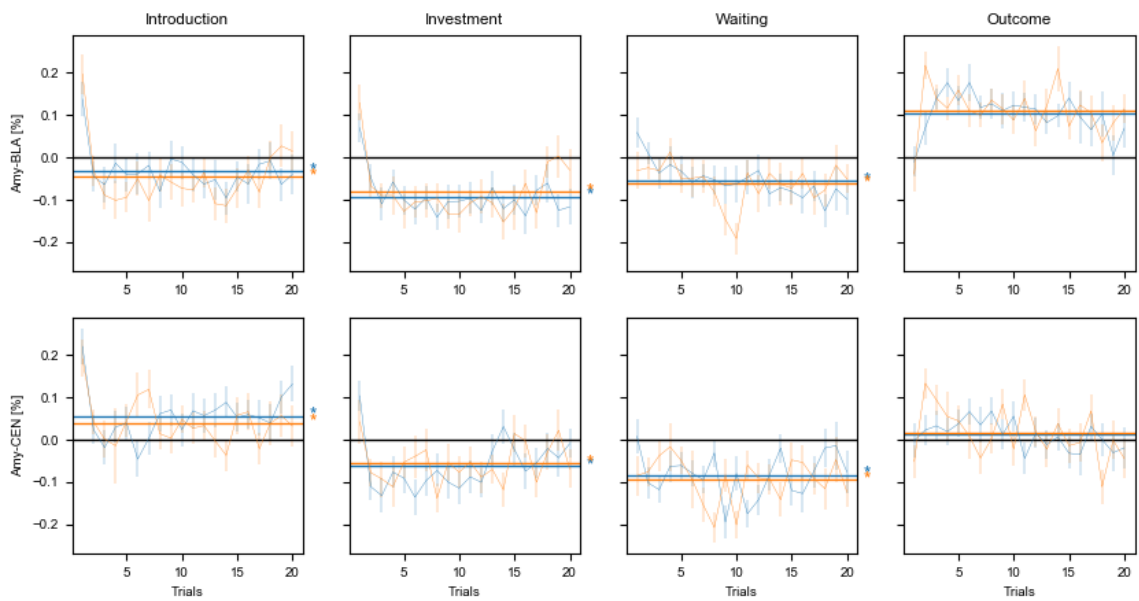
470

471 **SUPPLEMENT**

472



473 **Supplementary Figure S1.** Differences in functional connectivity of BLA>CeA (hot) and
474    CeA>BLA (cool).



475

476 **Supplementary Figure S2. No evidence for amygdala habituation.** Averaged percent
477 signal change for the different task phases for the trustworthy (blue) and untrustworthy
478 trustee (orange).

479

22

# REFERENCES

Adolphs R (2010), What does the amygdala contribute to social cognition? Annals of the New York Academy of Sciences 1191:42-61.

Alheid G, Heimer L (1988), New perspectives in basal forebrain organization of special relevance for neuropsychiatric disorders: the striatopallidal, amygdaloid, and corticopetal components of substantia innominata. Neuroscience 27:1-39.

Ambroggi F, Ishikawa A, Fields HL, Nicola SM (2008), Basolateral amygdala neurons facilitate reward-seeking behavior by exciting nucleus accumbens neurons. Neuron 59:648-661.

Avants BB, Tustison NJ, Song G, Cook PA, Klein A, Gee JC (2011), A reproducible evaluation of ANTs similarity metric performance in brain image registration. NeuroImage 54:2033-2044.

Avery SN, Clauss JA, Blackford JU (2016), The Human BNST: Functional Role in Anxiety and Addiction. Neuropsychopharmacology : official publication of the American College of Neuropsychopharmacology 41:126-141.

Avery SN, Clauss JA, Winder DG, Woodward N, Heckers S, Blackford JU (2014), BNST neurocircuitry in humans. NeuroImage 91:311-323.

Balleine BW, Killcross S (2006), Parallel incentive processing: an integrated view of amygdala function. Trends in neurosciences 29:272-279.

Baumgartner T, Heinrichs M, Vonlanthen A, Fischbacher U, Fehr E (2008), Oxytocin shapes the neural circuitry of trust and trust adaptation in humans. Neuron 58:639-650.

Bellucci G, Chernyak SV, Goodyear K, Eickhoff SB, Krueger F (2017), Neural signatures of trust in reciprocity: A coordinate-based meta-analysis. Human brain mapping 38:1233-1248.

Belova MA, Paton JJ, Morrison SE, Salzman CD (2007), Expectation modulates neural responses to pleasant and aversive stimuli in primate amygdala. Neuron 55:970-984.

Cardinale EM, Reber J, O'Connell K, Turkeltaub PE, Tranel D, Buchanan TW, Marsh AA (2021), Bilateral amygdala damage linked to impaired ability to predict others' fear but preserved moral judgements about causing others fear. Proceedings of the Royal Society B 288:20202651.

Chang LJ, Doll BB, van't Wout M, Frank MJ, Sanfey AG (2010), Seeing is believing: Trustworthiness as a dynamic belief. Cognitive psychology 61:87-105.

Clauss JA, Avery SN, Benningfield MM, Blackford JU (2019), Social anxiety is associated with BNST response to unpredictability. Depression and anxiety 36:666-675.

Corbit LH, Balleine BW (2005), Double dissociation of basolateral and central amygdala lesions on the general and outcome-specific forms of pavlovian-instrumental transfer. Journal of Neuroscience 25:962-970.

de Olmos JS, Heimer L (1999), The concepts of the ventral striatopallidal system and extended amygdala. Annals of the New York Academy of Sciences 877:1-32.

Delgado MR, Frank RH, Phelps EA (2005), Perceptions of moral character modulate the neural systems of reward during the trust game. Nature neuroscience 8:1611-1618.

Fischbacher U (2007), z-Tree: Zurich toolbox for ready-made economic experiments. Experimental economics 10:171-178.

Ganger S, Hahn A, Kublbock M, Kranz GS, Spies M, Vanicek T, Seiger R, Sladky R, et al. (2015), Comparison of continuously acquired resting state and extracted analogues from active tasks. Human brain mapping 36:4053-4063.

Geissberger N, Tik M, Sladky R, Woletz M, Schuler A-L, Willinger D, Windischberger C (2020), Reproducibility of amygdala activation in facial emotion processing at 7T. NeuroImage 211:116585.

Gewirtz JC, Mcnish KA, Davis M (1998), Lesions of the bed nucleus of the stria terminalis block sensitization of the acoustic startle reflex produced by repeated stress, but not fear-potentiated startle. Progress in Neuro-Psychopharmacology and Biological Psychiatry 22:625-648.

Goode TD, Ressler RL, Acca GM, Miles OW, Maren S (2019), Bed nucleus of the stria terminalis regulates fear to unpredictable threat signals. Elife 8:e46525.

Gupta R, Koscik TR, Bechara A, Tranel D (2011), The amygdala and decision-making. Neuropsychologia 49:760-766.

Haruno M, Kimura M, Frith CD (2014), Activity in the nucleus accumbens and amygdala underlies individual differences in prosocial and individualistic economic choices. Journal of Cognitive Neuroscience 26:1861-1870.

Hernandez-Lallement J, van Wingerden M, Schäble S, Kalenscher T (2016), Basolateral amygdala lesions abolish mutual reward preferences in rats. Neurobiology of Learning and Memory 127:1-9.

Infantolino ZP, Luking KR, Sauder CL, Curtin JJ, Hajcak G (2018), Robust is not necessarily reliable: From within-subjects fMRI contrasts to between-subjects comparisons. NeuroImage 173:146-152.

Janak PH, Tye KM (2015), From circuits to behaviour in the amygdala. Nature 517:284-292.

King-Casas B, Tomlin D, Anen C, Camerer CF, Quartz SR, Montague PR (2005), Getting to know you: reputation and trust in a two-person economic exchange. Science 308:78-83.

Lüthi A, Lüscher C (2014), Pathological circuit function underlying addiction and anxiety disorders. Nature neuroscience 17:1635-1643.

McDermott TJ, Kirlic N, Akeman E, Touthang J, Cosgrove KT, DeVille DC, Clausen AN, White EJ, et al. (2020), Visual cortical regions show sufficient test-retest reliability while salience regions are unreliable during emotional face processing. NeuroImage 220:117077.

Moeller S, Yacoub E, Olman CA, Auerbach E, Strupp J, Harel N, Ugurbil K (2010), Multiband Multislice GE-EPI at 7 Tesla, With 16-Fold Acceleration Using Partial Parallel Imaging With Application to High Spatial and Temporal Whole-Brain FMRI. Magnetic Resonance in Medicine 63:1144-1153.

Namburi P, Beyeler A, Yorozu S, Calhoon GG, Halbert SA, Wichmann R, Holden SS, Mertens KL, et al. (2015), A circuit mechanism for differentiating positive and negative associations. Nature 520:675-678.

O'Connell LA, Hofmann HA (2011), The vertebrate mesolimbic reward system and social behavior network: a comparative synthesis. Journal of Comparative Neurology 519:3599-3639.

O'Connell LA, Hofmann HA (2012), Evolution of a vertebrate social decision-making network. Science 336:1154-1157.

Pedersen WS, Muftuler LT, Larson CL (2019), A high-resolution fMRI investigation of BNST and centromedial amygdala activity as a function of affective stimulus predictability, anticipation, and duration. Social cognitive and affective neuroscience 14:1167-1177.

Plichta MM, Schwarz AJ, Grimm O, Morgen K, Mier D, Haddad L, Gerdes AB, Sauer C, et al. (2012), Test-retest reliability of evoked BOLD signals from a cognitive-emotive fMRI test battery. NeuroImage 60:1746-1758.

Pryce CR (2018), Comparative evidence for the importance of the amygdala in regulating reward salience. Current Opinion in Behavioral Sciences 22:76-81.

Robinson S, Windischberger C, Rauscher A, Moser E (2004), Optimized 3 T EPI of the amygdalae. NeuroImage 22:203-210.

Rosenberger LA, Eisenegger C, Naef M, Terburg D, Fourie J, Stein DJ, van Honk J (2019), The Human Basolateral Amygdala Is Indispensable for Social Experiential Learning. Curr Biol.

Sesack SR, Grace AA (2010), Cortico-basal ganglia reward network: microcircuitry. Neuropsychopharmacology : official publication of the American College of Neuropsychopharmacology 35:27-47.

Shackman AJ, Fox AS (2016), Contributions of the central extended amygdala to fear and anxietycontributions of the central extended amygdala to fear and anxiety. Journal of Neuroscience 36:8050-8063.

Siminski N, Böhme S, Zeller J, Becker M, Bruchmann M, Hofmann D, Breuer F, Mühlberger A, et al. (2020), BNST and amygdala activation to threat: effects of temporal predictability and threat mode. Behavioural Brain Research:112883.

Sladky R, Baldinger P, Kranz GS, Trostl J, Hoflich A, Lanzenberger R, Moser E, Windischberger C (2013), High-resolution functional MRI of the human amygdala at 7 T. Eur J Radiol 82:728-733.

Sladky R, Friston KJ, Trostl J, Cunnington R, Moser E, Windischberger C (2011), Slice-timing effects and their correction in functional MRI. NeuroImage 58:588-594.

Sladky R, Geissberger N, Pfabigan DM, Kraus C, Tik M, Woletz M, Paul K, Vanicek T, et al. (2018), Unsmoothed functional MRI of the human amygdala and bed nucleus of the stria terminalis during processing of emotional faces. NeuroImage 168:383-391.

Sladky R, Hoflich A, Atanelov J, Kraus C, Baldinger P, Moser E, Lanzenberger R, Windischberger C (2012), Increased neural habituation in the amygdala and orbitofrontal cortex in social anxiety disorder revealed by FMRI. PloS one 7:e50050.

Somerville LH, Wagner DD, Wig GS, Moran JM, Whalen PJ, Kelley WM (2013), Interactions between transient and sustained neural signals support the generation and regulation of anxious emotion. Cerebral cortex 23:49-60.

Terburg D, Scheggia D, Del Rio RT, Klumpers F, Ciobanu AC, Morgan B, Montoya ER, Bos PA, et al. (2018), The basolateral amygdala is essential for rapid escape: A human and rodent study. Cell 175:723-735. e716.

Torrisi S, O'Connell K, Davis A, Reynolds R, Balderston N, Fudge JL, Grillon C, Ernst M (2015), Resting State Connectivity of the Bed Nucleus of the Stria Terminalis at Ultra-High Field. Human brain mapping 36:4076-4088.

Tyszka JM, Pauli WM (2016), In vivo delineation of subdivisions of the human amygdaloid complex in a high-resolution group template. Human brain mapping 37:3979-3998.

van Honk J, Eisenegger C, Terburg D, Stein DJ, Morgan B (2013), Generous economic investments after basolateral amygdala damage. Proceedings of the National Academy of Sciences 110:2506-2510.