1    CONFIDENCE AS AN EXPRESSION OF COMMITMENT: WHY MISPLACED

2    EXPRESSIONS OF CONFIDENCE BACKFIRE

3

4    NON COPYEDITED VERSION, PLEASE DO NOT QUOTE

5    Vullioud, C.

6    Université de Neuchâtel

7    Centre de Sciences Cognitives

8    Espace Louis-Agassiz 1

9    Neuchâtel 2000

10    Switzerland

11    Clément, F.

12    Université de Neuchâtel

13    Centre de Sciences Cognitives

14    Espace Louis-Agassiz 1

15    Neuchâtel 2000

16    Switzerland

17    Scott-Phillips, T.

18    Evolutionary Anthropology Research Group

19    Department of Anthropology

20    Durham University

21    Dawson Building, South Road, Durham DH1 3LE, UK

22    Mercier, H.

23    Université de Neuchâtel

24    Centre de Sciences Cognitives

25    Espace Louis-Agassiz 1

26    Neuchâtel 2000

27    Switzerland

28    hugo.mercier@gmail.com

29

30

**Abstract**

Because communication can be abused by senders, it is not inherently stable. One way of stabilizing communication is for senders to commit to their messages. If a sender is committed to a message, she is willing to incur a cost (direct or reputational) if the message is found to be unreliable. This cost provides a reason for receivers to accept messages to which senders are committed. We suggest that expressions of confidence can be used as commitment signals: messages expressed more confidently commit their senders more. On this basis, we make three predictions: that confidently expressed messages are more persuasive (H1', already well established), that senders whose messages were accepted due to the senders' confidence but were then found to be unreliable should incur costs (H2'), and that if a message is accepted for reasons other than confidence, when it is found to be unreliable the sender should incur lower reputational costs than if the message had been accepted on the basis of the sender's confidence (H3'). A review of the literature revealed broadly supportive but still ambiguous evidence for H2' and no tests of H3'. In Experiments 1, 2, and 3 (testing H2') participants received the same advice from two senders, one being confident and the other unconfident. Participants were more likely to follow the advice of the confident sender, but once the advice was revealed to have been misguided, participants adjusted their trust so that they trusted the initially unconfident sender more than the confident sender. In Experiments 3 and 4 (testing H3') participants chose between either two senders differing in confidence or two senders differing in competence. Participants followed the advice of the confident sender and of the competent sender. When it was revealed that the advice was misguided, the confident sender suffered from a larger drop in trust than the competent sender. These results are relevant for communicative theories of overconfidence.

54

55          *Keywords*: Confidence, commitment, communication, overconfidence.

56                              *Word count*: 9089

57

58

59    Confidence as an Expression of Commitment: Why Misplaced Expressions of Confidence

60                                    Backfire

61

62            Communication between agents whose interests do not perfectly overlap is not

63    inherently stable. Even if both could benefit from communication, the danger is always

64    present that one would abuse communication for its own advantage. This observation holds at

65    the proximal level and at the ultimate level. At the proximal level, economists and other social

66    scientists have puzzled over the weight of 'cheap talk' (Farrell & Rabin, 1996): how can mere

67    words influence others when lying is not inherently costly? At the ultimate level, evolutionary

68    biologists have pointed out that communication can only be evolutionarily stable if it benefits

69    both senders and receivers (Dawkins & Krebs, 1978; Krebs & Dawkins, 1984; Maynard

70    Smith & Harper, 2003; Scott-Phillips, 2008). If senders do not benefit from communication,

71    they stop sending; if receivers do not benefit from communication, they stop receiving. But

72    what stops senders from sending signals that benefit only them, thereby threatening the

73    stability of communication?

74            Several mechanisms can stabilize communication (Maynard Smith & Harper, 2003).

75    For instance, some signals are inherently reliable because they cannot be faked—someone

76    who says "I am not a mute" cannot be lying, a Red Deer stag can only emit some types of

77    roars if it is large enough (see Maynard Smith & Harper, 2003). In humans, however, very

78    few signals are of this sort, so that we need to resort to other mechanisms to ensure the

79    stability of communication (Sperber et al., 2010). Here we focus on one of these mechanisms:

80    commitment. We suggest that in human communication, senders commit to various degrees to

81    their messages. A message to which the sender commits has, everything else equal, more

82    influence on the receiver. One way to express commitment is confidence: an assertion uttered

83    with more confidence commits its speaker more. We lay out and evaluate—through a

84    literature review and four experiments—consequences of this view of expressions of

85    confidence as commitment signals. In conclusion, we relate this view to theories that seek to

86    explain overconfidence through its communicative effects.

87

88    Commitment and communication

89

90          Commitment can take many forms. Some consider that commitment can be purely

91    internal. Such 'subjective commitment' (Fessler & Quintelier, 2013) consists in maintaining a

92    course of action not because of its instrumental value, but because of its intrinsic qualities.

93    Fessler and Quintelier (2013, p. 459) provide the example of a suicide bomber who follows

94    through on his plans because this course of action reflects his moral outrage towards the

95    targets of the bombing. In such a case, if the suicide bomber were to change his course of

96    action, he would suffer no external costs, but psychic costs such as feeling he has betrayed a

97    just cause. By contrast, objective commitment involves an actual cost attached to changing

98    one's course of action (Fessler & Quintelier, 2013). Opening a retirement account which

99    carries a heavy fee for withdrawals constitutes an objective commitment to saving for one's

100   retirement. In this example, the costs are purely personal but many instances of objective

101   commitment involve social costs. For instance, an individual who breaks a promise—which is

102   a typical form of commitment—often only incurs reputational costs.

103          The risks an individual takes in committing—i.e. the chances of having to pay some

104   costs if she fails to stay true to her commitment—should have a benefit, otherwise it is not

105   clear why anybody would commit to anything. These benefits can take many forms—for

106   instance, making sure that one isn't too poor upon retirement. In the context of

107   communication, the benefit of commitment is typically increased credibility, and the ability to

108   influence others credibility provides. When a receiver knows that a sender would incur some

109    costs if her communication proved unreliable, this provides him with a reason to believe her.

110    The role of commitment in communication can be more precisely laid out with the following

111    hypothesizes:

112

113        H1. Increased commitment should result in increased chances that a message is

114        accepted, or increased weight granted to the message.[1]

115        H2. If a message is found to have been unreliable (false, harmful), and the receiver had

116        accepted the message on the basis of the sender's commitment, then the sender should

117        suffer reputational costs.[2]

118

119    For commitment to play its hypothesized communicative role, it must be the case not only

120    that a sender of unreliable signals suffers some costs (per H2), but also that these costs be

121    higher than they would have been if she had not been committed. It is the cost added by

122    commitment that allows commitment to play its role. We can thus add the following

123    hypothesis:

124

125        H3. If a message is accepted on another basis than commitment, and if the message is

126        found to have been unreliable, then the sender's reputation should suffer less than if

127        the message had mostly been accepted on the basis of commitment. This would

128        happen for instance when a message is accepted because the receiver had deemed the

129        sender competent.

130

---

[1] Some caveats, which are not explored here, should be added to this hypothesis. The increased trust that results from increased commitment should be seen as multiplying the a priori trust in the sender rather than adding to it, so that completely mistrusted senders cannot rely on commitment to get their messages across. Moreover, expressed degrees of confidence that are implausibly high (e.g. "I am 100% sure my lottery number will come out") should also be dismissed.
[2] Reputational losses can affect either the perceived benevolence or the perceived competence of the sender (see Sperber *et al.*, 2010). In theory the loses due to failed commitments should mostly bear on the sender's perceived benevolence, but given that this prediction was not tested here, we do not elaborate further on this point.

131     Expressions of confidence as commitment signals

132

133          At least since Schelling's foundational work (Schelling, 1960), the communicative

134     benefits of commitment have received much attention (in an evolutionary perspective, see,

135     e.g. Fessler & Quintelier, 2013; Nesse, 2001). This attention has mostly focused on explicit

136     commitments, such as promises (e.g. Schelling, 2001). However, other speech acts also

137     commit their sender. In particular, assertions commit their sender to the truth of the

138     proposition expressed (e.g. Searle, 1969). This suggests that a sender whose assertions are

139     found to be false would suffer reputational costs. In practice, the distinction between speech

140     acts is often blurred (e.g. Astington, 1988), and what matters is not simply whether one's

141     speech act is, say, a promise or an assertion, but the degree of commitment that the sender

142     expresses.

143          Human languages possess a variety of devices that enable senders to modulate their

144     degree of commitment (Moeschler, 2013; Morency, Oswald, & de Saussure, 2008). For

145     instance, a sender is more committed to the propositional content of her utterances than to

146     their implicatures (Moeschler, 2013). Expressions of confidence also likely affect the degree

147     to which the sender is understood by receivers to be committed to her statements. Expressions

148     of confidence are ubiquitous in human communication, be they verbal ("I'm sure," "I guess,"

149     etc.) or non-verbal (gestures, tones, facial expressions). Indeed, the mechanisms which allow

150     senders to gauge their level of confidence might have evolved for the purpose of

151     communication (Shea et al., 2014). If expressions of confidence play the role of commitment

152     signals, then the hypotheses formulated above about commitment in general should apply to

153     expressions of confidence:

154

155     H1'. Increased confidence should result in increased chances that a message is

156     accepted, or increased weight granted to the message (the same caveats as above

157     apply).

158     H2'. If a message is found to have been unreliable (false, harmful), and the receiver

159     had accepted the message on the basis of the sender's confidence, then the sender

160     should suffer reputational costs.

161     H3'. If a message is accepted on another basis than sender confidence, and if the

162     message is found to have been unreliable, then the sender's reputation should suffer

163     less than if the message had mostly been accepted on the basis of confidence. This

164     would happen for instance when a message is accepted because the receiver had

165     deemed the sender competent.

166

167   H1' and H2' are similar to the hypotheses laid out about calibration in Tenney *et al.* (2008, p.

168   1369). In support of H1', many experiments have revealed that confidence tends to increase

169   the influence messages have on receivers (see, e.g., Price & Stone, 2004; Tenney, Small,

170   Kondrad, Jaswal, & Spellman, 2011; Yaniv, 1997; and, for children, Brosseau-Liard, Cassels,

171   & Birch, 2014). The goal of this article is to review the evidence relevant to H2', to further

172   test H2', and to offer the first—to the best of our knowledge—tests of H3'.

173

174   Do receivers punish senders who were confident but wrong?

175

176     Experiments relevant to evaluating H2' have yielded contradictory results. A first

177   series of experiments unambiguously supports H2'. Tenney and her colleagues (Tenney,

178   MacCoun, Spellman, & Hastie, 2007; Tenney et al., 2011, 2008) confronted participants with

179   the testimony of two senders whose confidence calibration was manipulated. For instance, in

180    Experiment 1 of Tenney *et al.* (2008), the participants had to evaluate the testimony of two

181    witnesses on the basis of the accuracy of two collateral statements (i.e. statements unrelated to

182    the case used to evaluate the reliability of a witness' testimony). One witness was confident

183    for both statements, while the other was confident for one statement and unconfident for the

184    other. At first, the participants did not know whether the statements were accurate; they were

185    then more likely to trust the confident witness. It was then revealed that each witness had been

186    mistaken about one statement. As a result, the confident witness was poorly calibrated, having

187    held confidently an inaccurate statement. In one condition, the less confident witness was well

188    calibrated since she had been wrong on the uncertain statement. In this condition, after the

189    accuracy feedback the participants found the less confident but better calibrated witness to be

190    more credible than the more confident but less well calibrated witness, and they were more

191    likely to believe her testimony. This experiment offers strong support for H2'. The

192    participants initially accepted a piece of testimony because its sender was confident. When the

193    confidence of the sender was revealed to have been unwarranted, the participants chose to

194    trust a sender who had been less confident but who had not expressed unwarranted

195    confidence.

196         Other experiments have found ambiguous support for H2'. Sah *et al.* (2013) asked

197    participants to gauge the weight of individuals on the basis of a picture of these individuals

198    and someone else's opinion (the senders' opinion). The senders had either high or low

199    confidence, and they were either very accurate or very inaccurate. Inaccurate and confident

200    senders were deemed, after the task had been completed, to be less credible than inaccurate

201    and unconfident senders. However, the opinions of inaccurate but confident senders were not

202    taken into account less than that of the inaccurate and unconfident senders (in spite of the

203    absence of floor effects). In another type of experiment, participants had to evaluate two

204    candidates: one who was very confident in his abilities, and one who was more cautious

205 (Tenney & Spellman, 2011). At first, confidence paid off, with better ratings for the confident

206 candidate. Once it was revealed that both candidates had in fact the same qualities, they were

207 both rated equally well. Thus, although confidence had no positive affect after it was revealed

208 to have been mistaken, it had no negative effects either (at least in the short term, see below).

209       Another set of studies, using a very different methodology, reached similar

210 conclusions (Kennedy, Anderson, & Moore, 2013). Participants first completed half of a task

211 in a small group. They were then asked their perception of the status and the competence of

212 each group member, including themselves. The participants' actual performance was then

213 revealed to all. The groups reconvened and completed the second half of the task before

214 answering the same status and competence questions. Participants were considered

215 overconfident if they estimated their status to be higher than warranted by their actual

216 performance. As in Anderson et al. (2012), before the participants had received performance

217 feedback, those who were overconfident were seen as being more competent and as having a

218 higher status (supporting H1'). After the performance feedback, the positive effects of

219 overconfidence disappeared, but there were no negative effects (i.e. the participants who were

220 initially overconfident were not perceived less well than those who had been initially well

221 calibrated).

222       These studies (Kennedy et al., 2013; Sah et al., 2013; Tenney & Spellman, 2011) seem

223 to suggest that, contrary to H2', overconfident senders do not see their reputation suffer much.

224 After their inaccuracy has been revealed, overconfident senders are not trusted less (or not

225 much else in the case of Sah *et al.* 2013) than unconfident senders. However, these studies can

226 be interpreted in a way that is compatible with H2'. They report a drop in trust or in status

227 once someone is revealed to have been overconfident. That this drop does not compensate for

228 the initial benefits of overconfidence might only reflect the scope of the experiments. If an

229 individual had kept being overconfident, and this overconfidence had kept causing drops in

230   trust and status, then that individual would have become less trustworthy, and would have

231   been attributed lower status, than her better calibrated peers. Thus these studies do not flatly

232   contradict H2'. Instead they suggest that for mistaken confidence to become costly, in some

233   cases, it has to be large enough, or clear enough, or repeated enough times. That this is the

234   case is suggested not only by the studies of Tenney *et al*. cited above—in which mistaken

235   confidence might have been particularly salient—but also by the results of Paulhus (1998). In

236   one of these experiments, participants met repeatedly over the course of several weeks. At

237   first, self-enhancers—individuals who tend to be overconfident in their abilities—were

238   perceived positively. After seven weeks, however, they were rated negatively on a variety of

239   traits.

240        Overall, the evidence regarding H2' is thus ambiguous, although we surmise that if the

241   experiments cited above that do not directly support H2' had been extended, the costs of

242   being confident but wrong would have become clearer, and thus their support for H2' clearer

243   as well.

244

245   The present experiments

246

247        The literature offers ambiguous support in favor of H2', and H3' has not been tested.

248   With the overarching goal of testing the role of commitment in the expression of confidence,

249   the present experiments further test H2' and offer the first tests of H3'. All the experiments

250   follow a similar template. Two senders provide advice to the participants, with varying

251   degrees of confidence (all experiments) or competence (Experiment 3 and 4). The participants

252   take this information into account. It is then revealed that at least one of the senders was

253   mistaken, and participants are asked to decide which of the senders they would rather punish

254   and which sender they would trust in the future.

255    In Experiments 1 and 2, both senders are equally wrong in their advice, only varying

256    in the degree of confidence with which the advice is expressed. H2' predicts that the

257    participants will inflict a higher cost on the more confident sender (through lower trust in

258    particular). In Experiments 3 and 4 a sender, who is more confident (in one condition) or

259    more competent (in the other condition) than the other sender, is proven wrong. H3' predicts

260    that the confident but wrong senders will see their reputation suffer more than that of the

261    equally wrong but competent sender.

262    The experiments were designed to be engaging for the participants. Experiment 1,

263    which was conducted in a classroom, used a simple, realistic situation and videos of the

264    senders. Experiments 2, 3, and 4, which were conducted online, used textual advice, but they

265    entailed a real, immediate cost for participants who accepted the wrong message. Finally,

266    Experiment 4 tested the evolution of the participants' trust in the senders by asking them to

267    make another potentially costly decision between advice provided by the same two senders.

268

269                        **Experiment 1 (a and b)**

270                               **Method**

271    **Participants**

272    Ninety undergraduate students (59 females; $M_{Age}$ 20.1; SD = 1.77) from a Swiss

273    University took part in Experiment 1a and 42 undergraduate students (27 females; $M_{Age}$ 21.6;

274    SD = 3.36), also from a Swiss University, took part in Experiment 1b. Both groups were

275    French speaking.

276    **Materials and procedure**

277    Experiments 1a and 1b took place in classrooms, before students attended a lecture.

278    The students had been asked to bring headphones and either a laptop or a smartphone, and

279  they answered the questions online on these devices. Experiments 1a and 1b were identical

280  except for the wording of one question that was found to have been unclear in 1a.

281      In order to set up the situation in which the advice would be given, the participants

282  were asked to imagine: "You have just started working as a middle manager in a big

283  company. You have to meet the Swiss manager for international coordination to organize an

284  important trip. As you don't know how to reach this manager, you rely on the coffee break to

285  ask two of your colleagues"[3]

286      A short movie then started in which two actors from the university theater group

287  played the colleagues. The first shot went from the corridor to the break room where the two

288  colleagues were standing. A screen appeared stating: "Hello, I'm trying to reach the Swiss

289  manager for international coordination. Do you know who is he and where I can find him?" In

290  the following shot, one of the colleagues answered: "Hi! International coordination, I know

291  him! It's Mr. Descloux, in building L, for Lausanne. You can believe me, I'm sure it's him."

292  His tone was confident (the surnames were not counterbalanced, but we see no plausible

293  reason to believe that this will have caused the effects observed).

294      In the final shot, which only showed the remaining colleague, he said: "Hi, hmm, I

295  don't know but I think that for the international coordination, it's Mr. Grandjean, in building

296  B, for Bern. But I'm really not sure". His tone was unconfident. Actors and presentation order

297  (i.e. confident first / unconfident first) were counterbalanced across participants, for a total of

298  four different films.

299      After they had watched the short movie, half of the participants were asked to choose

300  which of the two pieces of advice they wanted to follow. This was done to see whether

301  participants who had not explicitly stated that they would follow the advice of the confident

302  sender would still be more likely to punish him. All the participants were then told that both

---

[3] All texts are translated from French.

13

303  colleagues had been wrong, that the international coordination manager the participants were

304  looking for was in fact someone else in a different building. Thus the only difference between

305  the two colleagues was how confident they had been in their wrong answers.

306      The participants were then asked two questions (here for Experiment 1a). The first was

307  aimed at testing the participants' choice of which colleague they would like to directly punish:

308  "A few day later, your boss asked you to find somebody to put the 2000 invitations for the

309  collaborators' dinner in their envelope and to stamp them, during lunch break. You are the

310  team leader of the two colleagues seen during the coffee break. Whom do you give the task

311  to?" (punishment question). The second question bore on the reputation of the colleagues as a

312  sender in an unrelated area: "Since you are new in town, you are looking for a good

313  kindergarten for your kids. Whom do you ask advice from?" (trust question). For both

314  questions, participants had to pick one of the two senders.

315      As explained in the results section, the answers to the direct punishment question

316  proved surprising in light of the other results. To test whether participants had understood the

317  direct punishment question as intended, a few days after they had taken part in the

318  experiment, they were asked if they had understood the question as asking about a punishment

319  (forcing someone to do a boring task) or about a reward (trusting someone with a task). Fifty

320  percent of the participants had interpreted the question as being about a reward, thereby

321  invalidating the answers.

322      To fix this problem, in Experiment 1b we asked two questions instead of a single

323  punishment question, and made the wording unambiguous: "You are the team leader of the

324  two colleagues you saw during the coffee break. You are about to start two new projects. The

325  first project is a project that has no interest or importance. Taking part in this project is

326  demeaning and can be seen as a kind of punishment." Participants answered the first question,

327  then read and answered the second: "The second project is an important and interesting

328    project for a big client. Taking part in this project is gratifying and can be seen as a kind of

329    reward." The order of presentation of the questions was counterbalanced between participants

330    in both Experiment 1a and 1b.

331

332                                **Results**

333    *Experiment 1a.*

334            Eighty-five percent of the participants, who had been asked whose advice they wanted

335    to follow, decided to follow the advice of the confident sender (39/46; binomial $p$ = .001).[4]

336    For the remaining two questions, there were no differences between these participants and

337    those who had not been asked to specify which of the pieces of advice they wanted to follow

338    (Mann-Whitney; punishment question: $Z$= -.65, $p$ = .514; reputation question: $Z$= -.46, $p$ =

339    .644), and their results were aggregated. For the question intended to bear on punishment,

340    63% of these participants chose the unconfident sender (57/90; binomial $p$ = .015). However,

341    as mentioned above, the punishment question was problematic. For the trust question, 71% of

342    the participants trusted the unconfident sender (64/90; binomial $p$ < .001).

343    *Experiment 1b.*

344            Seventy-six percent of the participants, who had been asked whose advice they wanted

345    to follow, decided to follow the advice of the confident sender (16/21; binomial $p$ = .027).

346    Overall, 83% of the participants (35/42; binomial $p$ < .001) answered that they would punish

347    the highly confident sender while 69% (29/42; binomial $p$ = 0.02) answered that they would

348    reward the unconfident sender.

349            In Experiment 1, participants behaved in line with H1'—they were more likely to

350    believe a confident than an unconfident sender—and with H2'—they were more likely to

351    punish, and less likely to trust, on an unrelated matter, a sender who was confident but was

_____

[4] All data are available in the ESM.

352   then proven to have been wrong, than a sender who had been equally wrong but had been less

353   confident. Experiment 2 seeks to replicate the result regarding H2' using a different context

354   and different tools (online experiment).

355

356   **Experiment 2 (a and b)**

357   **Method**

358   **Participants**

359   Forty participants (17 females; $M_{Age}$ 33.30; SD = 11.01) took part in experiment 2a,

360   and 42 participants (16 females; $M_{Age}$ 34.05; SD = 10.94) in experiment 2b. The participants

361   were recruited through Amazon Mechanical Turk (MTurk).

362

363   **Materials and procedure**

364   Experiments 2a and 2b were conducted online. Experiment 2a is a conceptual

365   replication of Experiment 1 in which the two senders provide the same advice, only varying in

366   their degree of confidence. Experiment 2b is a control experiment in which the senders give

367   correct pieces of advice, designed to insure that participants do not have a general bias against

368   confident senders. We predict that in Experiment 2b, by contrast to Experiment 2a,

369   participants will punish less and trust more, after feedback, the more confident sender.

370   Each experiment comprised two tasks. In the first task, participants took on the role of

371   an adviser. They were told that another participant would have to type a text, and that they

372   had to advise them about which text they should type in order to make the typing easier and

373   faster. The participants were presented, for a short time (5s) with two texts, one of which

374   contained many difficult words which made it longer to type than the other (see ESM). The

375   participants then had to say which text they would advise another participant to choose, and to

376   write a short statement qualifying their answer. This first task had two goals. First, to make it

377   more believable that the advice the participants received in the second task could have been

378   given by another participant. Second, to make participants believe that the task of the advisor

379   was not trivially easy, so that bad advice could be attributed to an honest error.

380        In the second task, the participants took on the role of the advisee. They had to choose

381   one text among two to type, and were thus motivated to choose the text that would take less

382   time to type. The only indication they had as to which text would take less time to type took

383   the form of two pieces of advice provided by senders described as previous participants in the

384   experiment (in fact we created the pieces of advice ourselves). Both senders advised to select

385   the same text, but they offered different statements in support of their advice. The confident

386   sender's advice was accompanied by this statement "I'm 100% sure this text is the easiest of

387   the two," while the unconfident sender's advice was accompanied by "It was very quick, I

388   couldn't see well, so I'm not so sure." Both senders were either wrong (Experiment 2a) or

389   right (Experiment 2b), only differing in their degree of confidence. After participants had

390   chosen a text to type, they were told that they had chosen either the longer text (Experiment

391   2a) or the shorter text (Experiment 2b). The participants then typed the texts, which took

392   approximately one minute for the long text—in Experiment 2a—and 30 seconds for the short

393   text—Experiment 2b). The texts were provided to them in a picture format, so that they could

394   not cut and paste their content, and the participants could not move on to the next screen until

395   they had typed the exact text provided.

396        After they had typed the texts, participants were asked two forced-choice questions

397   similar to the questions asked in the first experiment. The first was aimed at testing the

398   participants' choice of which senders they would like to directly punish: "If you could stop

399   one of the two participants whose advice you saw from receiving a bonus, who would you

400   pick?" (punishment question). The second question bore on the reputation of the senders: "If

401   you had to do the experiment again, and you could only receive advice from one of these two

402    participants, who would you pick?" (trust question). Question order as well as order of

403    presentation of the senders (confident vs. unconfident) were counterbalanced between

404    participants. The detail of all the information provided to the participants, screen by screen, is

405    available in the ESM.

406

407                                                    **Results**

408    *Experiment 2a.*

409            Eighty-eight percent of the participants (35/40; binomial $p < .001$) picked the text

410    advised by the two senders. Participants who had not followed the recommendation of the

411    senders were excluded of further analysis. All of the participants who had followed the advice

412    preferred to punish the confident sender (35/35, binomial $p < .001$) and 91% (32/35; binomial

413    $p < .001$) indicated that they would trust the unconfident sender.

414    *Experiment 2b and comparison with 2a.*

415            Seventy-six percent of the participants (32/42; binomial $p = .001$) followed the advice

416    given by the two senders. Participants who had not followed the recommendation of the

417    senders were excluded of further analysis. Twenty-two percent of the participants who had

418    followed the advice (7/32; binomial $p = .002$) preferred to punish the confident sender and 9%

419    (3/32; binomial $p < .001$) indicated that they would trust the unconfident sender to complete

420    the task again. Compared to Experiment 2a, in Experiment 2b participants were more likely to

421    punish the unconfident sender (Mann-Whitney $Z= -6.6$, $p < .001$), and more likely to trust the

422    confident sender as a better sender (Mann-Whitney $Z= -6.7$, $p < .001$).

423            Experiment 2a supports H2': between two senders who were equally wrong,

424    participants tend to punish more, and to trust less, the more confident sender. Experiment 2b

425    shows that these results do not stem from a general bias against confident senders. When both

426    senders are equally right, participants tend to punish more, and to trust less, the less confident

427  sender. Experiment 3 uses the same procedure as Experiment 2 in order to test H3' (and

428  incidentally to replicate Experiment 2a).

429

430                                **Experiment 3**

431                                  **Method**

432  **Participants**

433       Ninety-nine participants (34 females; $M_{Age}$ 32.58; SD = 8.84) were recruited through

434  Mturk to participate in experiment 3.

435

436  **Materials and procedure**

437       Experiment 3 was designed to test H3' through a between-participant design with two

438  conditions: a Competence Condition and a Confidence Condition. The Competence Condition

439  was broadly similar to Experiment 2a with three crucial differences. As in Experiment 2a, the

440  participants were confronted with two senders. The first departure from Experiment 2a is that

441  the difference in confidence between the sender was removed: both used neutral expressions

442  to accompany their advice ("Text one looked like the shorter one to type" and "Seemed to be

443  overall the easier of the two" respectively). Second, a difference of competence between the

444  two senders was introduced. Participants were told that the two senders had different track

445  records of success at the task of picking the easier: "in previous experiments, he or she [i.e.

446  the sender] correctly chose the easiest text 12 out of 13 times [respectively 3 out of 13

447  times]." Third, the two senders advised to select different texts—while they advised to select

448  the same text in Experiment 2a. We introduced this change in order to test whether competent

449  senders were initially believed more or less than confident senders (relative to incompetent

450  and unconfident senders respectively).

451    The Confidence Condition was identical to Experiment 2a with one exception: as in

452    the Competence Condition, the two senders provided different advice about which texts to

453    select. Since the test of H3' consists in a comparison across conditions, what matters is not

454    that the senders within each condition are equally wrong, but that the confident sender and the

455    competent sender are equally wrong, as is the case (since they are both wrong while the other

456    sender is right). Order of presentation of the senders, question order, and, in the Competence

457    Condition, matching between the neutral statements and the senders, were counterbalanced.

458    The questions were the same as in Experiment 2.

459    To summarize, participants in the Competence Condition were exposed to two senders

460    of different initial competence, while participants in the Confidence Condition were exposed

461    to two senders of different confidence. In both conditions, the sender who we expected to be

462    initially believed—the competent sender in the Competence Condition and the confident

463    sender in the Confidence Condition—was proven wrong. We then asked participants

464    questions about which sender they would like to inflict costs on—either directly or by trusting

465    them less.

466

467                                                    **Results**

468    Participants were as likely to select the text advised by the confident speaker in the

469    Confidence Condition (76%, 37/49; binomial $p < .001$) than they were to select the text

470    advised by the competent speaker in the Competence Condition (78%, 39/50; binomial $p <$

471    $.001$) (Mann-Whitney $Z= -.29$, $p = .770$). Participants who had not followed the advice of the

472    competent or the confident senders were excluded from further analyses.

473    In the Confidence Condition, 86% of the participants (32/37; binomial $p < .001$)

474    preferred to punish the confident sender and 89% (33/37; binomial $p < .001$) indicated that

475    they would trust the unconfident sender to complete the task again. In the Competence

476     Condition, 77% of the participants preferred to punish the competent sender (30/39; binomial

477     $p = .001$) but only 56% (22/39; binomial $p = .522$) indicated that they would trust the

478     incompetent sender to complete the task again. There was no difference between the

479     conditions in the answers to the punishment question (Mann-Whitney $Z = -1.1$, $p = .286$), but

480     the participants who had accepted the advice of the competent senders were more likely to

481     trust them in the future than the participants who had accepted the advice of the confident

482     senders were to trust the confident senders in the future (Mann-Whitney $Z = -3.7$, $p = .002$).

483            By contrast with the other experiments, in Experiment 3 one sender was right while

484     the other was wrong. We could hardly have expected participants to punish the sender who

485     actually gave them sound advice over the one who gave them unreliable advice, even if the

486     latter was more competent. As a result, the answers to the punishment question are not as

487     relevant here as in the other experiments: they cannot properly test H3'. By contrast, the trust

488     question can adequately test H3', since a participant can trust someone who has been wrong

489     once over someone who has been right once, if other factors make up for this difference. The

490     results of the trust question support H3'. Even though initial trust was equally high in the

491     confident sender and the competent sender, and that both senders proved to be equally wrong,

492     final trust was higher in the competent sender than in the confident sender. This shows that,

493     when other factors are controlled for, the reputation of the confident sender suffered more

494     than that of the competent sender. Experiment 4 replicates Experiment 3 while increasing its

495     validity by introducing stakes in the final trust question.

496

497                                              **Experiment 4**

498                                                **Method**

499     **Participants**

500    Seventy-nine participants (37 females; $M_{Age}$ 32.49; SD = 9.99) were recruited through

501    MTurk to take part in the experiment.

502

503    **Materials and procedure**

504    Experiment 4 was similar to Experiment 3 with two differences. First, instead of

505    asking participants who they would pick if they had to complete the task again, participants

506    had to actually perform the same task, receiving advice from the same two senders. In this

507    second task, both senders advised to select different texts so that we could measure which

508    sender was trusted more. In the Competence Condition, the success rates of the two senders

509    were provided again, having been updated to account for their failure in the first task. In the

510    Confidence Condition, two new statements were adapted from those previously written by

511    participants to express confidence ("I am absolutely sure in my decision") and lack of

512    confidence ("looked like easier to type but I'm not really sure"). Second, given that we had

513    established in Experiment 3 that confidence and competence had the same influence on

514    participants' choices in the first task, in this first task both senders in each condition gave the

515    same advice. This makes the results of the punishment question more interesting.

516

517                                              **Results**

518    In the Confidence Condition, 93% of the participants (37/40; binomial $p < .001$)

519    selected the text advised by both senders in the first task; in the Competence Condition, 95%

520    of the participants did so (37/39; binomial $p < .001$). Participants not following the advice of

521    the two participants were excluded from further analyses.

522    In the Confidence Condition, 85% of the participants (31/37; binomial $p < .001$)

523    preferred to punish the confident sender and 65% (24/37; binomial $p = .099$) trusted the

524    advice of the unconfident sender in the second task. In the Competence Condition, 49% of the

525    participants (18/37; binomial $p$ = 1.00) preferred to punish the competent sender and 30%

526    (9/37; binomial $p$ = .020) trusted advice of the incompetent sender in the second task.

527    Participants were more likely to punish the confident sender than the competent sender

528    (Mann-Whitney $Z$= -3.2, $p$ = .002) and they were more likely to trust the competent sender

529    than the confident sender (Mann-Whitney $Z$= -3.0, $p$ = .003).

530          These results offer strong support for H3'. Even though the confident sender's

531    message and the competent sender's message were revealed to have been equally wrong, the

532    confident sender was subsequently punished more, and trusted less than the competent sender.

533

534                                    **Replications**

535          To ensure the reliability of our findings, we replicated the results from all online

536    experiments (Experiments 2, 3, and 4). A total of 413 participants were recruited through

537    MTurk in three sessions:  Experiment 2a and 2b (8 excluded, final Ns: 2a=37, 2b=38, 27

538    females; $M_{Age}$ 34.61; SD = 9.57), experiment 3 (11 excluded, final N=189; 83 females; $M_{Age}$

539    35.92; SD = 11.94), experiment 4  (12 excluded, final N=149; 78 females; $M_{Age}$ 33.42; SD =

540    9.20). For Experiments 2a and 2b, which were simple conceptual replications of previous

541    experiments, we used the same Ns as in the first version of the present experiments. For

542    Experiments 3 and 4, which were more novel—therefore potentially more contentious—and

543    which included a comparison across conditions, we doubled the number of participants

544    recruited in the first version of these experiments.

545          The 31 participants who were excluded had said they had already taken part in the

546    same experiment previously or were not sure that they had not.

547    **Experiment 2a**

548          Eighty-four percent of the participants (31/37; binomial $p$ < .001) selected the text

549    advised by the two senders. Participants who had not followed the recommendation of the

550    senders were excluded of further analysis. Seventy-one percent of the participants who had

551    followed the advice preferred to punish the confident sender (22/31, binomial $p = .029$) and

552    74% (23/31; binomial $p < .011$) indicated that they would trust the unconfident sender. These

553    results thus closely replicate those obtained previously.

**Experiment 2b and comparison with 2a**

555    Ninety-two percent of the participants (35/38; binomial $p < .001$) followed the advice

556    given by the two senders. Participants who had not followed the recommendation of the

557    senders were excluded of further analysis. Thirty-one percent of the participants who had

558    followed the advice (11/35; binomial $p = .041$) preferred to punish the confident sender and

559    6% (2/35; binomial $p < .001$) indicated that they would trust the unconfident sender to

560    complete the task again. Compared to Experiment 2a, in Experiment 2b participants were

561    more likely to punish the unconfident sender (Mann-Whitney $Z = -3.2$, $p = .001$), and more

562    likely to trust the confident sender as a better sender (Mann-Whitney $Z = -5.7$, $p < .001$). These

563    results thus closely replicate those obtained previously.

**Experiment 3.**

565    Participants were as likely to select the text advised by the confident speaker in the

566    Confidence Condition (66%, 61/92; binomial $p = .002$) than they were to select the text

567    advised by the competent speaker in the Competence Condition (71%, 71/97; binomial $p <$

568    $.001$) (Mann-Whitney $Z = -1.03$, $p = .303$). Participants who had not followed the advice of the

569    competent or the confident senders were excluded from further analyses.

570    In the Confidence Condition, 80% of the participants (49/61; binomial $p < .001$)

571    preferred to punish the confident sender and 72% (44/61; binomial $p = .001$) indicated that

572    they would trust the unconfident sender to complete the task again. In the Competence

573    Condition, 54% of the participants preferred to punish the competent sender (38/71; binomial

574    $p = .635$) but only 45% (32/71; binomial $p = .477$) indicated that they would trust the

575    incompetent sender to complete the task again. Participants in the Confidence Condition were

576 more likely to punish the confident sender than the participants in the Competence Condition

577 were to punish the competent sender  (Mann-Whitney $Z$= -3.2, $p$ = .001), and the participants

578 who had accepted the advice of the competent senders were more likely to trust them in the

579 future than the participants who had accepted the advice of the confident senders were to trust

580 the confident senders in the future (Mann-Whitney $Z$= -3.1, $p$ = .002). These results thus

581 closely replicate those obtained previously. The only potential difference was that participants

582 were significantly more likely to punish the confident sender than the competent sender,

583 whereas this difference was not significant in the original experiment. The results of the

584 replications are thus, if anything, even more in line with H3'.

585 **Experiment 4**

586       In the Confidence Condition, 88% of the participants (68/77; binomial $p$ < .001)

587 selected the text advised by both senders in the first task; in the Competence Condition, 83%

588 of the participants did so (60/72; binomial $p$ < .001). Participants not following the advice of

589 the two participants were excluded from further analyses.

590       In the Confidence Condition, 66% of the participants (45/68; binomial $p$ = .010)

591 preferred to punish the confident sender and 50% (34/68; binomial $p$ = 1.000) trusted the

592 advice of the unconfident sender in the second task. In the Competence Condition, 37% of the

593 participants (22/60; binomial $p$ = .052) preferred to punish the competent sender and 28%

594 (17/60; binomial $p$ = .001) trusted advice of the incompetent sender in the second task.

595 Participants were more likely to punish the confident sender than the competent sender

596 (Mann-Whitney $Z$= -3.2, $p$ = .001) and they were more likely to trust the competent sender

597 than the confident sender (Mann-Whitney $Z$= -2.5, $p$ = .013). These results thus closely

598 replicate those obtained previously.

599

600 **General discussion**

601     The goal of this series of experiments was to test two hypotheses: H2', that senders

602 whose messages are accepted because they are confident suffer a reputation loss when their

603 messages are found to have been misleading; and H3', that this reputation loss is greater than

604 that incurred by senders whose messages were accepted for other reasons (here, competence).

605 Incidentally, all experiments also found support for H1' (confidence increases message

606 acceptance). Experiments 1, 2, and 4 support H2'. In these experiment, participants receive

607 the same advice from two senders, one being confident and the other unconfident. At first,

608 participants are more likely to follow the advice of the confident sender. However, once the

609 advice is revealed to have been misguided, participants adjust their trust so that they trust the

610 initially unconfident sender more.

611     Experiments 3 and 4 support H3'. Participants choose between either two senders

612 differing in confidence or two senders differing in competence. At first, participants follow

613 the advice of the confident sender and of the competent sender—and they do so equally

614 strongly. When it is revealed that the advice is misguided, participants are more likely to trust

615 the initially unconfident sender. By comparison, the drop in trust incurred by the competent

616 but wrong sender is less severe, since after the feedback, the participants either do not trust

617 the competent sender less than the incompetent sender (Experiment 3), or they keep trusting

618 the competent sender more (Experiment 4). Experiments 2, 3, and 4 were successfully

619 replicated using the same population, demonstrating the robustness of their results. The results

620 from all the experiments are summarized in Table 1.

621

26

622

| Experiment | 1a | 1b | 2a | 3 | | 4 | |
|---|---|---|---|---|---|---|---|
| N | 90 | 42 | 40 *37* | 49 *92* | 50 *97* | 40 *77* | 39 *72* |
| Which sender incorrect? | Both senders incorrect | | Both senders incorrect | Only confident or competent sender incorrect | | Both senders incorrect | |
| Trait of the favored sender | Confident | | Confident | Confident | Competent | Confident | Competent |
| Initially trust the favored sender | 85% | 76% | 88% *84%* | 76% *66%* | 78% *71%* | 93% *88%* | 95% *83%* |
| After feedback: Punish the favored sender | Question unclear | 83%*** | 100% *** *71%\** | 86% *80%* ** | 77% *54%* | 85% ** *66%* ** | 49% *37%* |
| After feedback: Trust the favored sender | 29%*** | 31%** | 9% *** *26%\** | 11% ** *28%* ** | 44% *55%* | 35% ** *50%* * | 70% *72%* |

623

624  Table 1: Summary of results for Experiments 1 to 4. In all the experiments the favored sender

625  (confident or competent) was wrong. The results of the replications of experiments 2 to 4 are in italics.

626  For the two 'After feedback' lines, the percentages are computed on the basis of the participants who

627  trusted the favored sender (confident or competent). Stars denote the level of significance (* < 0.05, **

628  < 0.01, *** < 0.001). For Experiments 1 and 2, stars represent differences from chance performance.

629  For Experiments 3 and 4, stars represent differences between conditions.

630

631      Besides offering support for H3', which had not been previously tested, the present

632  experiments extend the literature related to H2' in different ways. Given the ambiguity in the

633  existing literature regarding H2', the simple adjunction of more evidence in support of this

634  hypothesis is pertinent. Moreover, the experiments extend previous results supporting H2' (in

635  particular the experiments of Tenney *et al.* 2007, 2008, 2011) in at least two ways. First, the

636    experiments reveal that the drop in reputation incurred by overconfident senders extends

637    beyond the domain in which they have been found to be overconfident: In Experiment 1a, the

638    senders were proven to be wrong on a work matter, and they were then less trusted on a

639    family matter. This suggests that experiments which only test for a drop of trust following

640    confident but unreliable messages in the same domain as that of the message might be

641    underestimating the costs of mistaken confidence. Second, three of the experiments

642    introduced costs for following the misguided advice (having to type a much longer text,

643    Experiments 2, 3, and 4) and one incentivized the choice of which sender to trust after the

644    feedback (following the best advice might lead participants to type a shorter text, Experiment

645    4).

646                                      **Conclusion**

647    One of the mechanisms senders rely on to get receivers to accept their messages is

648    commitment. By committing to their messages, they accept to incur a cost if the messages are

649    found to be unreliable (H2), a cost that has to be greater than the cost they would have

650    incurred if their unreliable messages had been accepted for reasons other than commitment

651    (H3). Knowing of this cost, receivers have a reason to accept the messages senders commit to

652    (H1). We suggested that expressions of confidence could play the role of commitment signals,

653    leading to the formulation of the equivalent hypotheses for confidence instead of commitment

654    more generally (H1', H2', and H3', see General Discussion above). Four experiments

655    provided incidental support for H1' (which was already solidly supported), new support for

656    H2' (which was supported, but only ambiguously), and some initial support for H3' (which

657    had never been tested).

658    Considered with the existing evidence reviewed above, we now believe there is strong

659    support for H2': it seems that mistaken confidence backfires and hurts senders. Even though

660    what we have developed here is a (partial) theory of expressed confidence, and not a theory of

661    overconfidence more generally, the hypotheses we examined (H2' in particular) are relevant

662    for some theories of overconfidence. H2' provides another reason why overconfidence should

663    be costly: not only can overconfidence lead to personally damaging decisions (e.g., Barber &

664    Odean, 2001), but, if expressed, it can hurt one's reputation. This extra cost makes it even

665    more puzzling that overconfidence seems to be such a common phenomenon (e.g.,

666    Kahneman, 2011).

667        Some theories of overconfidence posit that overconfidence yields benefits that

668    outweigh its costs. For instance, a model suggests that overconfidence allows agents to

669    compete more effectively over resources (Johnson & Fowler, 2011), and another that

670    overconfidence leads to a better mental health (Taylor & Brown, 1988). By contrast, other

671    theories have suggested that some forms of overconfidence exist because the expression of

672    overconfidence yields benefits that are conferred by others (social benefits). In particular,

673    according to the status-enhancement theory of overconfidence, overconfidence confers social

674    benefits because "overly positive self-views help individuals convince others that they are

675    more capable than they actually are" (Anderson, Brion, Moore, & Kennedy, 2012, p. 718; see

676    also, e.g. Trivers, 2011).

677        The status-enhancement theory of overconfidence predicts that individuals who

678    express overconfidence should get social benefits, and that these benefits should be higher

679    than the costs they might incur if their overconfidence were revealed (Kennedy et al., 2013).

680    In this theory, overconfidence is not necessarily attached to specific statements (as in the

681    present experiments), but rather with one's general abilities. However, we believe that such

682    confidence might still constitute a form of commitment: people would commit not to a

683    specific statement, but to the strength of their general abilities. If this were the case, then

684    overconfidence should be punished, in that individuals who are consistently confident beyond

685    their abilities should be seen as less reliable than individuals whose confidence matches their

686    abilities.

687          As noted above, some experimental results suggested that overconfident individuals

688    were not punished in this way (Kennedy et al., 2013). However, in these experiments

689    participants still decreased the trust they granted overconfident individuals when their

690    overconfidence was revealed. If we extrapolate from this trend, then an individual who would

691    remain overconfident, or who would be too overconfident from the start, would end up being

692    trusted less than a better calibrated individual. Indeed, as suggested in the introduction, this is

693    what the rest of the literature (to which we can now add the current results) suggests.

694          This does not mean that the status-enhancement theory of overconfidence cannot

695    apply in some cases. Individuals could be mistaken about the risks of overconfidence.

696    Individuals could also find themselves in situations in which overconfidence has low costs,

697    either because the senders' relative lack of competence is unlikely to be revealed (e.g. experts

698    who make vague predictions), or because the senders are mostly engaged in one shot

699    interactions (e.g. car dealers).

700          In spite of these potential exceptions, the idea that expressing overconfidence is not

701    generally a successful strategy fits well with many results suggesting that some forms of

702    overconfidence are not as robust as was once thought. Overconfidence can take at least the

703    three following forms (Moore & Healy, 2008). Overplacement is saying that we are better

704    than others when we are not (e.g. when most people believe they are smarter than the median

705    individual). Overestimation is saying we are better than we are (e.g. when people say they can

706    solve problems they can't solve). Overprecision is making statements that are more precise

707    than warranted (e.g. when people say they are 95% sure the value of a stock will increase

708    when in fact it has only 75% chances of increasing).

709    Overplacement and overestimation are not robust. Many studies that were supposed to

710    demonstrate overplacement and overestimation have been contested on statistical grounds

711    (Benoît & Dubra, 2011; Harris & Hahn, 2011). The amount of overplacement and

712    overestimation varies widely as a function of different factors: the relative difficulty of the

713    questions on which participants have to estimate their performance (Lichtenstein & Fischhoff,

714    1977), the participants' culture (Heine & Lehman, 1995), the ease with which overplacement

715    can be justified (Dunning, Meyerowitz, & Holzberg, 1989), the amount of feedback provided

716    to the participants (Rose & Windschitl, 2008), and so forth (e.g. Galesic, Olsson, &

717    Rieskamp, 2012). The amount of variation in overplacement and overestimation is such that

718    reversals are common. For instance, participants tend to underestimate their performance on

719    easy problems (Lichtenstein & Fischhoff, 1977), and they believe they are below average

720    when it comes to uncommon abilities (Moore, 2007). Note that in these experiments

721    confidence is usually not measured behaviorally (e.g. by testing which tasks the participants

722    are willing to engage in), but by asking participants to express their degree of confidence. As

723    a result, the current hypotheses should apply. Considerations of the potential social costs

724    caused by unwarranted expressions of confidence might help explain the pattern of data.

725    By contrast with overplacement and overestimation, overprecision is much more

726    robust (Moore, Tenney, & Haran, in press). Moreover, overprecision is the form of

727    overconfidence which is closest to the overconfidence displayed by the senders in our

728    experiments. We thus seem to face the following puzzle: being overprecise is costly yet

729    common. We suggest that the conversational norm theory of overprecision can solve this

730    puzzle (Yaniv & Foster, 1995). According to this theory, if people are overprecise, it is

731    because they favor informativeness in the tradeoff between informativeness and accuracy.

732    Since more precise statements are less likely to be accurate, overprecision tends to decrease

733    accuracy. However, more precise statements are more informative. To take an extreme

734  example, if you ask a realtor to estimate the value of your house and she says "between $10

735  and $100,000,000," she is bound to be right, but her statement is also so vague as to be

736  irrelevant (see, Sperber & Wilson, 1995).

737      The conversational norm theory of overconfidence is supported by data both on the

738  sender's side and on the receiver's side. On the sender's side, individuals appropriately tailor

739  the preciseness of their statements to the context—for instance by providing more precise

740  time when the individual who is asking is going to catch a train (Van der Henst, Carles, &

741  Sperber, 2002). On the receiver's side, participants prefer a precise estimate (e.g. between 140

742  and 150 for the number of countries belonging to the U.N.) to a vague one (50 to 300), even

743  after the second is revealed to be more accurate (the correct answer was 159) (Yaniv &

744  Foster, 1995). It thus seems that by making statements more relevant, overprecision yields

745  some benefits for receivers.

746      Crucially, it is also possible that overprecision doesn't entail any costs for receivers.

747  For overprecise statements to be costly, they have to be taken at face value. For instance,

748  when a participant discovers that another participant thought the number of countries

749  belonging to the U.N. was between 140 and 150, she might not take that to mean that the

750  participant is certain of this estimate, only that this is her best guess that would still be

751  relevant enough to be useful. That this is how receivers understand messages is suggested by

752  the fact that, everything else equal, receivers tend to heavily discount senders' opinions when

753  these opinions contradict their own views (see, e.g., Yaniv, 2004). Receivers would often be

754  better off taking the senders' opinion into account more, not less, so that even if the senders

755  have been overprecise, this overprecision is likely to have played a positive role (i.e. by

756  making receivers take the message into account more, even if not quite enough).

757      Moreover, overprecision does not seem to get in the way of the effective

758  communication of degrees of confidence. When participants have to complete a perceptual

759    task in dyads, they are able to determine which dyad member is more confident, so that the

760    dyad is able to select the answer favored by the more confident participant, which is generally

761    the correct answer (Bahrami et al., 2010). On balance, it thus seems that overprecision is

762    beneficial rather than costly for receivers. It is thus not surprising that it should not be

763    punished and that it should be so common and robust.

764         These considerations suggest that the prevalence of overprecision can be reconciled

765    with the current hypotheses, in particular H2' which posits that expressed mistaken

766    confidence should be costly. Our hypotheses about commitment and confidence do not bear

767    on the literal meaning of the statements, but on the meaning attributed to the sender. For

768    instance, a sender would obviously not commit to the literal meaning of an ironic or a

769    metaphorical statement. Similarly, some contexts call for modesty, others for bluster. To the

770    extent that senders manage to get their thoughts across effectively, then they should be

771    understood to be committed to the appropriate degree. If overprecision does not, on average,

772    lead receivers to attribute to senders a misleadingly high degree of confidence, then it should

773    not be routinely punished.

774         An important challenge for future research is to integrate these insights from the

775    experimental literature with general theories of pragmatics, and human communicative

776    behavior more generally. For instance, some approaches emphasize the importance of 'face'

777    (see Brown & Levinson, 1987). From that perspective expressions of different degrees of

778    confidence—lower confidence in particular—allow the speaker to preserve face even when

779    what they communicate might be false. In any case, there is much scope for further

780    experimental research: existing experiments are rudimentary in contrast to the complexity of

781    the expressions of commitment found in everyday dialogue. The current framework can create

782    a useful bridge between linguistic, psychological, and evolutionary theories pertaining to the

783    expression of confidence and commitment more generally.

784

785

788

789                                     References

790    Anderson, C., Brion, S., Moore, D. A., & Kennedy, J. A. (2012). A status-enhancement
791    account of overconfidence. *Journal of Personality and Social Psychology*, *103*(4), 718.
792    Astington, J. W. (1988). Children's understanding of the speech act of promising. *Journal of*
793    *Child Language*, *15*(1), 157–173.
794    Bahrami, B., Olsen, K., Latham, P. E., Roepstorff, A., Rees, G., & Frith, C. D. (2010).
795    Optimally interacting minds. *Science*, *329*(5995), 1081–1085.
796    Barber, B. M., & Odean, T. (2001). Boys will be boys: Gender, overconfidence, and common
797    stock investment. *Quarterly Journal of Economics*, 261–292.
798    Benoît, J.-P., & Dubra, J. (2011). Apparent overconfidence. *Econometrica*, *79*(5), 1591–1625.
799    Brosseau-Liard, P., Cassels, T., & Birch, S. (2014). You Seem Certain but You Were Wrong
800    Before: Developmental Change in Preschoolers' Relative Trust in Accurate versus Confident
801    Speakers. *PloS One*, *9*(9), e108308.
802    Brown, P., & Levinson, S. C. (1987). *Politeness: Some universals in language usage*.
803    Cambridge: Cambridge university press.
804    Dawkins, R., & Krebs, J. R. (1978). Animal signals: Information or manipulation? In J. R.
805    Krebs & N. B. Davies (Eds.), *Behavioural Ecology: An Evolutionary Approach* (pp. 282–
806    309). Oxford: Basil Blackwell Scientific Publications.
807    Dunning, D., Meyerowitz, J. A., & Holzberg, A. D. (1989). Ambiguity and self-evaluation:
808    the role of idiosyncratic trait definitions in self-serving assessments of ability. *Journal of*
809    *Personality and Social Psychology*, *57*(6), 1082–1090.
810    Farrell, J., & Rabin, M. (1996). Cheap talk. *Journal of Economic Perspectives*, *10*(110–118).
811    Fessler, D. M., & Quintelier, K. (2013). Suicide bombers, weddings, and prison tattoos: an
812    evolutionary perspective on subjective commitment and objective commitment. In K.
813    Sterelny, R. Joyce, B. Calcott, & B. Fraser (Eds.), *Cooperation and its Evolution* (pp. 459–
814    484). Cambridge: MIT Press.
815    Galesic, M., Olsson, H., & Rieskamp, J. (2012). Social sampling explains apparent biases in
816    judgments of social environments. *Psychological Science*, 956797612445313.
817    Harris, A. J., & Hahn, U. (2011). Unrealistic optimism about future life events: a cautionary
818    note. *Psychological Review*, *118*(1), 135.
819    Heine, S. J., & Lehman, D. R. (1995). Cultural variation in unrealistic optimism: Does the
820    West feel more vulnerable than the East? *Journal of Personality and Social Psychology*,
821    *68*(4), 595.
822    Johnson, D. D., & Fowler, J. H. (2011). The evolution of overconfidence. *Nature*, *477*(7364),
823    317–320.
824    Kahneman, D. (2011). *Thinking, fast and slow*. New York: Farrar Straus & Giroux.
825    Kennedy, J. A., Anderson, C., & Moore, D. A. (2013). When overconfidence is revealed to
826    others: Testing the status-enhancement theory of overconfidence. *Organizational Behavior*
827    *and Human Decision Processes*, *122*(2), 266–279.
828    Krebs, J. R., & Dawkins, R. (1984). Animal signals: Mind-reading and manipulation? In J. R.
829    Krebs & N. B. Davies (Eds.), *Behavioural Ecology: An Evolutionary Approach* (Vol. 2ème,
830    pp. 390–402). Oxford: Basil Blackwell Scientific Publications.
831    Lichtenstein, S., & Fischhoff, B. (1977). Do those who know more also know more about
832    how much they know? *Organizational Behavior and Human Performance*, *20*(2), 159–183.
833    Maynard Smith, J., & Harper, D. (2003). *Animal signals*. Oxford: Oxford University Press.
834    Moeschler, J. (2013). Is a speaker-based pragmatics possible? Or how can a hearer infer a
835    speaker's commitment? *Journal of Pragmatics*, *48*(1), 84–97.

836    Moore, D. A. (2007). Not so above average after all: When people believe they are worse than
837    average and its implications for theories of bias in social comparison. *Organizational*
838    *Behavior and Human Decision Processes*, *102*(1), 42–58.
839    Moore, D. A., & Healy, P. J. (2008). The trouble with overconfidence. *Psychological Review*,
840    *115*(2), 502.
841    Moore, D. A., Tenney, E. R., & Haran, U. (in press). Overprecision in judgment. In G. Wu &
842    G. Keren (Eds.), *Handbook of Judgment and Decision Making*. New York: Wiley.
843    Morency, P., Oswald, S., & de Saussure, L. (2008). Explicitness, implicitness and
844    commitment attribution: A cognitive pragmatic approach. *Belgian Journal of Linguistics*,
845    *22*(1), 197–219.
846    Nesse, R. (Ed.). (2001). *Evolution and the capacity for commitment*. New York: Russell Sage
847    Foundation.
848    Paulhus, D. L. (1998). Interpersonal and intrapsychic adaptiveness of trait self-enhancement:
849    A mixed blessing? *Journal of Personality and Social Psychology*, *74*(5), 1197.
850    Price, P. C., & Stone, E. R. (2004). Intuitive evaluation of likelihood judgment producers:
851    evidence for a confidence heuristic. *Journal of Behavioral Decision Making*, *17*, 39–57.
852    Rose, J. P., & Windschitl, P. D. (2008). How egocentrism and optimism change in response to
853    feedback in repeated competitions. *Organizational Behavior and Human Decision Processes*,
854    *105*(2), 201–220.
855    Sah, S., Moore, D. A., & MacCoun, R. J. (2013). Cheap talk and credibility: The
856    consequences of confidence and accuracy on advisor credibility and persuasiveness.
857    *Organizational Behavior and Human Decision Processes*, *121*(2), 246–255.
858    Schelling, T. C. (1960). *The strategy of conflict*. Cambridge, MA.: Harvard University Press.
859    Schelling, T. C. (2001). Commitment: Deliberate Versus Involuntary. In R. Nesse (Ed.),
860    *Evolution and the capacity for commitment* (pp. 48–56). New York: Russell Sage Foundation.
861    Scott-Phillips, T. C. (2008). Defining biological communication. *Journal of Evolutionary*
862    *Biology*, *21*(2), 387–395.
863    Searle, J. (1969). *Speech Acts: An Essay in the Philosophy of Language*. Cambridge:
864    Cambridge University Press.
865    Shea, N., Boldt, A., Bang, D., Yeung, N., Heyes, C., & Frith, C. D. (2014). Supra-personal
866    cognitive control and metacognition. *Trends in Cognitive Sciences*, *18*(4), 186–193.
867    Sperber, D., Clément, F., Heintz, C., Mascaro, O., Mercier, H., Origgi, G., & Wilson, D.
868    (2010). Epistemic vigilance. *Mind and Language*, *25*(4), 359–393.
869    Sperber, D., & Wilson, D. (1995). *Relevance: Communication and cognition*. New York:
870    Wiley-Blackwell.
871    Taylor, S. E., & Brown, J. D. (1988). Illusion and well-being: a social psychological
872    perspective on mental health. *Psychological Bulletin*, *103*(2), 193.
873    Tenney, E. R., MacCoun, R. J., Spellman, B. A., & Hastie, R. (2007). Calibration trumps
874    confidence as a basis for witness credibility. *Psychological Science*, *18*(1), 46–50.
875    Tenney, E. R., Small, J. E., Kondrad, R. L., Jaswal, V. K., & Spellman, B. A. (2011).
876    Accuracy, confidence, and calibration: how young children and adults assess credibility.
877    *Developmental Psychology*, *47*(4), 1065.
878    Tenney, E. R., & Spellman, B. A. (2011). Complex social consequences of self-knowledge.
879    *Social Psychological and Personality Science*, *2*(4), 343–350.
880    Tenney, E. R., Spellman, B. A., & MacCoun, R. J. (2008). The benefits of knowing what you
881    know (and what you don't): How calibration affects credibility. *Journal of Experimental*
882    *Social Psychology*, *44*(5), 1368–1375.
883    Trivers, R. (2011). *The Folly of Fools: The Logic of Deceit and Self-Deception in Human*
884    *Life*. New York: Basic Books.

885     Van der Henst, J.-B., Carles, L., & Sperber, D. (2002). Truthfulness and relevance in telling
886     the time. *Mind & Language*, *17*(5), 457–466.
887     Yaniv, I. (1997). Weighting and trimming: Heuristics for aggregating judgments under
888     uncertainty. *Organizational Behavior and Human Decision Processes*, *69*(3), 237–249.
889     Yaniv, I. (2004). Receiving other people's advice: Influence and benefit. *Organizational*
890     *Behavior and Human Decision Processes*, *93*, 1–13.
891     Yaniv, I., & Foster, D. P. (1995). Graininess of judgment under uncertainty: An accuracy-
892     informativeness trade-off. *Journal of Experimental Psychology: General*, *124*(4), 424.
893