

# Imagining a False Alibi Impairs Concealed Memory Detection With the Autobiographical Implicit Association Test

Phot Dhammapeera

University of Kent and Chulalongkorn University

Xiaoqing Hu

The University of Hong Kong and HKU-Shenzhen Institute of Research and Innovation, Shenzhen, China

Zara M. Bergström

University of Kent

Imagining counterfactual versions of past events can distort memory. In 3 experiments, we examined whether imagining a false alibi for a mock crime would make suspects appear less guilty in a concealed memory detection test, the autobiographical Implicit Association Test (aIAT), which aims to determine which of 2 autobiographical events are true. First, “guilty” participants completed a mock crime, whereas “innocent” participants completed an innocent act. Next, some of the guilty participants were asked to imagine a false alibi that corresponded to the innocent act. Finally, all groups completed the aIAT. Across experiments, we varied the type of aIAT used and also compared the effectiveness of the false alibi countermeasure when only imagined once, versus when it was repeatedly imagined over a week-long period. The aIAT accurately detected the mock crime as true for guilty participants without a false alibi, but was consistently less able to detect the mock crime as true for guilty participants who had imagined a false alibi. The findings suggest that if guilty suspects fabricate an alibi, this may create a memory for the alibi that appears to be true based on the aIAT, which is problematic for its real-life applications in concealed memory detection.

### Public Significance Statement

We found that rehearsing a false alibi can impair truth detection with a computerized test, the autobiographical implicit association test. This finding is important because it suggests the test is vulnerable to faking, and that real-life applications of this test are premature.


**Keywords:** memory, imagination, autobiographical Implicit Association Test, truth

**Supplemental materials:** <http://dx.doi.org/10.1037/xap0000250.supp>

Forensic memory detection aims to determine if a criminal suspect has concealed information stored in their memory that is indicative of guilt. Guilty suspects are expected to have unique

knowledge of the crime that would not be known by innocent suspects. Therefore, nonverbal markers of memory, such as memory-related brain activity (e.g., Allen, Iacono, & Danielson, 1992; Gamer, Klimecki, Bauermann, Stoeter, & Vossel, 2012; Rosenfeld, Angell, Johnson, & Qian, 1991; van Hooff, Brunia, & Allen, 1996), autonomic activity (Gamer, 2011; Lykken, 1959), or reaction times (RTs) and accuracy on indirect memory tests (Sartori, Agosta, Zogmaister, Ferrara, & Castiello, 2008; Verschuere & De Houwer, 2011), can be measured to detect if a suspect is concealing incriminating knowledge. Many of these methods can very accurately detect concealed information, at least in cooperative research participants with little motivation to hide their guilt (Granhag, Vrij, & Verschuere, 2015; Verschuere, Ben-Shakhar, & Meijer, 2011). However, one prominent concern is that real criminals may use countermeasure strategies to attempt to hide their guilt (e.g., Bergström, Anderson, Buda, Simons, & Richardson-Klavehn, 2013; Hu, Bergström, Bodenhausen, & Rosenfeld, 2015; Verschuere, Prati, & De Houwer, 2009; for a review, see Ben-Shakhar, 2011), threatening the validity of these tests in real-life

This article was published Online First September 26, 2019.

Phot Dhammapeera, School of Psychology, University of Kent, and Faculty of Psychology, Chulalongkorn University; Xiaoqing Hu, Department of Psychology, The State Key Laboratory of Brain and Cognitive Sciences, The University of Hong Kong, and HKU-Shenzhen Institute of Research and Innovation, Shenzhen, China;  Zara M. Bergström, School of Psychology, University of Kent.

We thank Jessica Amos, Ellie Anslow, Ashley Bailey, Chloe Brunskill, Rhiannon Chappell, Amber Gardner, Lucy Hendleman, Catalina Marin, Chloe Walker, and Eleanor Webster for help with data collection. This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Correspondence concerning this article should be addressed to Zara M. Bergström, School of Psychology, University of Kent, Canterbury CT2 7NP, United Kingdom. E-mail: [z.m.bergstrom@kent.ac.uk](mailto:z.m.bergstrom@kent.ac.uk)

settings. Considering the important societal, legal, and ethical implications of forensic memory detection, it is therefore critical to evaluate whether memory detection tests are susceptible to countermeasures. It is also important to assess which types of countermeasures are likely to be successful in order to ensure that memory detection tests are optimally designed to withstand evasion attempts.

The autobiographical implicit association test (aIAT, Sartori et al., 2008), is a computerized task that bears high promise in assessing the implicit truth value of autobiographical statements, which can therefore be used to detect concealed autobiographical memories. The aIAT measures RTs and accuracy in a simple sentence classification task as markers of whether an autobiographical event is true or false for an individual, and is thus considerably easier and cheaper to implement than physiology and brain activity-based techniques that necessitate specialist equipment and highly trained administrators. In a criminal context (e.g., Sartori et al., 2008), the aIAT involves presenting suspects with four different types of statements that suspects have to classify on two dimensions: logically true versus false, or crime-related versus innocent-related, by pressing two different buttons. Sentences for the first dimension are true or false for everyone taking the test (e.g., true: "I am in front of a computer" vs. false: "I am in a restaurant"), whereas the truth of sentences for the second dimension depend on whether the suspect has committed the crime or not (e.g., true if guilty/false if innocent: "I stole a ring" (a crime-related sentence) vs. false if guilty/true if innocent: "I bought a ring" (an innocence-related sentence)). In guilt congruent blocks, logically true and crime-related statements share one button, while logically false and innocent-related statements share another button. In guilt incongruent blocks, logically false and crime-related statements share one button, while logically true and innocent-related statements share another button. Guilty suspects are expected to respond faster and more accurately in guilt congruent than incongruent blocks due to crime-related sentences having implicit and automatic associations with the truth. Innocent suspects are expected to show the opposite pattern.

Many studies have shown very accurate memory detection using the aIAT (reviewed in Agosta & Sartori, 2013; however, see Suchotzki, Verschuere, Van Bockstaele, Ben-Shakhar, & Crombez, 2017 for evidence that the aIAT may be less effective than other RT-based memory detection paradigms). Moreover, the aIAT is not only able to detect which of two autobiographical events is more strongly associated with truth, but is also better at detecting true memories than false memories that the participant believes are true (Marini, Agosta, Mazzoni, Barba, & Sartori, 2012). Because of such promising results, the aIAT has already been applied in at least one real court case in Italy, where it was used by the defense team as part of a battery of tests to suggest that the defendant had memory impairments, which was accepted by the judge as evidence of diminished culpability and contributed to a reduced penalty for a convicted murderer (Sirgiovanni, Corbellini, & Caporale, 2016). In contrast, other research has shown that the aIAT may be susceptible to relatively simple countermeasures that guilty suspects can apply during the test, such as slowing down responses in the guilt congruent blocks (Verschuere et al., 2009) or speeding up responses in the guilt incongruent blocks (Hu, Rosenfeld, & Bodenhausen, 2012), especially when partici-

pants are allowed to practice in advance of the test. However, suspects who used such strategies may be caught out by selectively modifying their response times only during critical blocks but not during other, noncritical blocks (Agosta, Ghirardi, Zogmaister, Castiello, & Sartori, 2011). Thus, trying to beat the aIAT by directly altering response times may not be a particularly effective countermeasure, because such faking attempts may be detectable by unusual patterns of response times across different blocks (although see Hu et al., 2012).

An alternative strategy that guilty suspects could use for evading forensic memory detection is to intentionally modify their memories in advance of the test, in order to make these memories more consistent with innocence. A large body of evidence shows that memories for experienced events remain malleable after encoding and can be updated or inhibited at a later stage (e.g., Anderson & Hanslmayr, 2014; Dudai, 2012). Indeed, in several experiments we have found that by intentionally suppressing memories of committing a mock crime, guilty suspects were able to significantly reduce retrieval-related ERPs thus increasing the likelihood of appearing innocent on an EEG-based memory detection test (Bergström et al., 2013; Hu et al., 2015). Furthermore, suppression of mock crime memories weakened the associative strength between the crime and the truth so that guilty suspects also appeared more innocent on a later aIAT, even without engaging any intentional strategies during the aIAT itself (Hu et al., 2015). Thus, modifying memories in advance of a memory detection test may be an effective countermeasure strategy that is less detectable than on-line faking attempts during the test itself.

Whereas previous research showed that suspects can intentionally weaken incriminating memories to evade detection, another strategy by which guilty suspects could appear innocent is to intentionally store false information in memory that suggests innocence. It is well established that people can hold vivid memories for events that they have never experienced in real life (Loftus & Pickrell, 1995; Schacter, Guerin, & St. Jacques, 2011). Such memories can be created simply by imagining a novel event (Loftus, 2003) that becomes encoded as a memory representation with similar perceptual and conceptual features as a memory based on an experienced event, making true and false memories similar in terms of their neural and behavioral characteristics (Mitchell & Johnson, 2009). Consistent with this view, imagining performing simple actions (such as picking a specific card from a deck of playing cards) enhances implicit associations between the imagined event and the truth when contrasted with nonimagined events in an aIAT. Some research found this to be the case particularly when participants misremembered imagined actions as previously performed (Takarangi, Strange, Shortland, & James, 2013), whereas in other studies, aIAT truth detection of imagined actions was enhanced even when participants knew the imagined event did not occur in real life (Shidlovski, Schul, & Mayo, 2014; see also Mangiulli et al., 2018; Takarangi, Strange, & Houghton, 2015; Vargo, Petróczi, Shah, & Naughton, 2014). Furthermore, in a mock criminal context, asking people to deliberately memorize a hypothetical alternative version of a mock crime can weaken skin conductance responses associated with a true mock crime, and thereby impair memory detection with autonomic measures (Gronau, Elber, Satran, Breska, & Ben-Shakhar, 2015).

However, to our knowledge, no previous research has investigated whether guilty suspects can intentionally memorize false information indicative of innocence as a countermeasure strategy for evading guilt detection with the aIAT. In real life, guilty suspects may fabricate an untrue version of what they were doing at the time of the crime to use as a false alibi, and by doing so, they may encode this information into memory in a form that may share some characteristics with a true memory, which may potentially also distort or impair their memory for the true crime event (Otgaar & Baker, 2018). Recent research has shown that adopting a false alibi can impair identification of guilty suspects in deception detection paradigms (Foerster, Wirth, Herbort, Kunde, & Pfister, 2017; Suchotzki, Berlijn, Donath, & Gamer, 2018), but this issue has not been investigated with the aIAT. We addressed these issues in three experiments that used the aIAT to investigate whether imagining a false alibi impaired guilt detection by enhancing the implicit truth value of an alibi and/or decreasing the implicit truth value of a committed mock crime. We also investigated whether the alibi countermeasure was more effective when applied repeatedly over an extended time period compared with just in one brief session. To preempt the results, we found a consistent pattern across studies whereby the false alibi significantly impaired guilt detection with the aIAT, which seemed to be primarily driven by the alibi being detected as true rather than a substantial impairment of the original mock crime memory.

### Experiment 1

The first experiment was conducted in three stages. First, “guilty” participants carried out a mock crime which involved stealing a ring from a bag in a university staff office area, whereas “innocent” participants carried out an innocent act that involved going to the same office area but instead writing their e-mail address on a paper slip on a staff member’s door. Next, half of the guilty participants were instructed to imagine performing the innocent act with the explicit intention of using this as a false alibi in order to appear innocent. The other half of guilty participants and the innocent group performed an unrelated filler task. Finally, all three groups undertook an aIAT where the relative truth value of the mock crime and innocent/false alibi events were compared in all three groups.

We hypothesized that imagining a false alibi would create a memory for the imagined act, which may have some implicit associations with the truth even though participants knew their alibi was fake at an explicit level (Shidlovski et al., 2014). Imagining a fake alibi would thus lead to lower aIAT discrimination between the objectively true mock crime and the objectively false innocent act when this group was compared with the guilty group who did not imagine the alibi. If imagining an alibi as a countermeasure was completely successful at making guilty suspects appear innocent, aIAT performance for these guilty participants would be indistinguishable from the innocent group who actually conducted the innocent act in real life.

### Method

**Participants.** The design was based on our previous experiment which included 78 participants divided across three groups and found a large effect size (Cohen’s  $d = 0.78$ ) for reduced aIAT

memory detection in a suppression countermeasure group compared with a standard guilty group (Hu et al., 2015). That prior experiment was designed to have 0.8 power to detect a  $d = 0.8$  effect size, and we increased our sample size in the current study to further enhance statistical power, and therefore recruited 108 participants who were split into three groups, resulting in  $>0.9$  power to detect a  $d = 0.8$  effect size, or 0.8 power to detect a  $d = 0.7$  effect size (we decided a priori that we were primarily interested in detecting large effects of the alibi countermeasure on the aIAT, as only large countermeasure effects have substantial implications for practical applications involving guilt classification at the individual level). The participants were undergraduate students at the University of Kent who took part via a research participation scheme in return for course credits. Participants were randomly assigned to three experimental groups ( $N = 36$  in each): the guilty-alibi group (30 female, six male), the guilty-standard group (29 female, seven male), and the innocent group (28 female, eight male). Twenty additional participants were replaced due to technical problems or not following the instructions during the mock crime/innocent act (such as stealing the wrong object, or going to the wrong part of the building). Participants’ age ranged from 18–28 ( $M = 19.83$ ,  $SD = 1.62$ ). The groups did not significantly differ in terms of age,  $F(2, 104) = .80$ ,  $p = .451$ ,  $\eta_p^2 = .02$  gender ( $\chi^2(2) = .36$ ,  $p = .837$ ,  $\phi = .84$ ). All participants had English as their first language, had normal or corrected-to-normal vision, and had no diagnosis of dyslexia. The study was approved by the University of Kent Psychology Ethics committee.

**Materials, design, and procedure.** First, participants in the two guilty groups were required to go to a kitchen adjacent to staff offices in a university building, find a bag, and steal a box from inside the bag. They were explicitly asked to look and take note of what was inside the box (a ring), and then return with the box and its content to the experimental room. The word ring was not mentioned in the instructions so that the memory of the ring was gained solely from enacting the crime. Innocent participants were required to go to the same area in the building, but instead they were told to write their e-mail address on an appointment sign-up sheet on the door of a lecturer’s office. Thus, innocent participants were unaware of the mock crime.

Next, participants in the guilty-alibi group were provided with a fake alibi scenario, which was designed to help them appear innocent on the aIAT. Participants were told that they would soon take part in a test designed to detect their guilt, however they should aim to appear innocent by adopting the alibi. Participants were instructed that it was essential that they try to imagine the scenario as if it were true and that their memory for scenario details would later be tested. The alibi scenario was a short verbal description of the innocent act: “You were on your way to find your lecturer. On their door, there was a sheet of paper specifying that you could leave your e-mail address for the lecturer to get back to you. So you tore off a bit of paper and wrote your e-mail address and left it in the envelope provided and came back here. The envelope has since been destroyed so there is no evidence that your alibi is false.” Participants were told to close their eyes and vividly imagine the alibi for 2 min. Next, they were asked to describe the scenario in detail and answer a few questions about it. If they gave incorrect answers, the alibi story was repeated and the questions asked again until the correct answers were given. Participants in the guilty-standard and innocent groups were instead

required to carry out a filler task of solving Sudoku puzzles. They were given two puzzles as well as written instructions and told to do the best they could while they were timed for 5 min.

In the final stage, all participants took part in a seven-block computerized aIAT (Hu et al., 2015; Sartori et al., 2008). Participants were instructed that multiple sentences would appear on the screen and they would need to classify them as either logically true or false, or ring-related or e-mail-related by pressing left or right buttons on the keyboard. To avoid online attempts to modify the test result, they were not informed regarding how the test worked or how to alter their responses to appear innocent (cf. Agosta et al., 2011; Hu et al., 2012; Verschuere et al., 2009). The first block (20 trials) was a simple classification block that required participants to classify five true and five false sentences, with each sentence repeated twice in random order. Participants were instructed to press the left key “Z” for logically true sentences (e.g., “I am a research participant”) and the right key “M” for logically false sentences (e.g., “I am playing football”), based on what they were doing at that time. The labels “True” and “False” were displayed on the left and right sides of the screen respectively, to remind participants of the response-key mapping. The second block (20 trials) was a simple classification block that required participants to classify five sentences related to the guilty act (e.g., “I took a ring”) and five sentences related to the innocent act/alibi scenario (e.g., “I wrote my e-mail”). Participants were asked to press the left key “Z” for ring-related sentences and the right key “M” for e-mail-related sentences, and the labels “Ring” and “E-mail” were displayed on the left and right sides of the screen, respectively. Blocks 3 (20 trials) and 4 (40 trials) were critical double classification blocks which tested participants’ responses to guilt congruent sentence pairings, because logically true and autobiographically true sentences for the guilty groups were paired to the same response button. Participants were instructed to press “Z” if the sentence was logically true or ring-related and “M” if the sentence was logically false or e-mail-related, and the labels “True/Ring” and “False/E-mail” were displayed on the left and right sides of the screen, respectively. Block 5 (20 trials) was a practice reverse simple classification block, which reversed the key assignments for ring and e-mail-related sentences (“Z” for e-mail-related and “M” for ring-related sentences, with the left label changed to “E-mail” and the right label changed to “Ring”). Finally, Blocks 6 (20 trials) and 7 (40 trials) were also critical double classification blocks with the reversed keys, thus testing participants’ responses to guilt incongruent sentence pairings, because logically false and autobiographically true sentences for the guilty groups were paired to the same response button. Participants were instructed to press “Z” if the sentence was logically true or e-mail-related and “M” if the sentence was logically false or ring-related, and the labels “True/E-mail” and “False/Ring” were displayed on the left and right sides of the screen respectively. Faster RT and higher accuracy for guilt congruent blocks than guilt incongruent blocks indicate an association between the crime and the truth, whereas the reverse pattern indicate an association between the innocent act and the truth.

Half of the participants within each group conducted the blocks in the order described above, while Blocks 2–4 and 5–7 were swapped for the other half of participants in order to counterbalance the order of guilt congruent and guilt incongruent blocks. Thus, counterbalancing formats were balanced within groups and

matched across groups. For all blocks, sentences were presented on the screen in random order, and stayed on the screen until participants pressed a button. Participants were instructed to respond as quickly and accurately as possible, and if they pressed the incorrect button a red “X” appeared on the screen until the pressed the correct button.

**Data analysis.** The main measure of guilt in the aIAT is the D-score, which combines accuracy and RT into a single, standardized measure (Greenwald, Nosek, & Banaji, 2003; Sartori et al., 2008). We used the same formula to calculate D as in the most relevant previous studies (Hu et al., 2012, 2015). First, extreme RTs (<100 ms or >10,000 ms) were deleted. As in prior research, incorrect responses were given a 600ms penalty, and the mean RTs were calculated for the guilt congruent and guilt incongruent blocks separately, including the incorrect responses with the applied penalties. Finally, the mean RT difference between guilt congruent and guilt incongruent blocks was divided by the standard deviation of the RT distribution for correct trials only, from both blocks combined, in order to obtain the D-score. In the Experiment 1 version of the aIAT, a positive D-score indicated guilt because it suggests that participants associated sentences describing the mock crime with the truth, whereas a negative D-score indicated innocence because it suggests that participants associated sentences describing the innocent act with the truth.

Potential group differences in D-scores were analyzed with commonly used frequentist inferential tests from the GLM (ANOVA, *t* tests). Effect sizes were estimated using partial eta-squared for ANOVAs, and Cohen’s *d* for *t* tests. Cohen’s *d* for both paired and independent *t* tests was calculated as the difference between means divided by the pooled standard deviation rather than from the *t*-values to avoid inflating effect size estimates for paired *t* tests (Dunlap, Cortina, Vaslow, & Burke, 1996). As the key hypotheses relied on testing whether D-scores were above or below 0 within each group and whether there were pairwise group differences in D-scores, frequentist *t* tests for such differences were supplemented with Bayes factors ( $BF_{10}$ ) to evaluate the relative support for a difference ( $H_1$ ) versus no difference ( $H_0$ ). These were calculated with Bayesian *t* tests in JASP (JASP Team, 2019) using default priors (a Cauchy distribution with center = 0,  $r = .707$ ). The Bayes factors is a ratio that contrasts the likelihood that the data would occur under the alternative ( $H_1$ ) versus null ( $H_0$ ) hypotheses, with values over 1 indicating support for  $H_1$  and values below 1 indicating support for the  $H_0$ . Values close to 1 are only considered weakly/anecdotally supportive of one hypothesis over the other, whereas  $BF_{10} > 3$  are typically interpreted as substantial evidence in support of  $H_1$  over  $H_0$ , and  $BF_{10} < 0.33$  are interpreted as substantial evidence in support of  $H_0$  over  $H_1$  (see Wagenmakers, Wetzels, Borsboom, & van der Maas, 2011).

The aIAT was developed to diagnose guilt or innocence at the individual level, which is typically done by classifying individuals with positive D-scores as “guilty” and individuals with negative D-scores as “innocent” when contrasting a guilty versus innocent event in this way (Sartori et al., 2008). However, because such classification rates are dependent on choosing specific cut-offs and the optimal cut-off may vary across samples and experimental designs, we instead conducted a threshold-independent receiver operating characteristic (ROC) analysis to evaluate classification performance using areas under the curve (AUCs; following e.g., Hu et al., 2015, but see online supplemental materials for

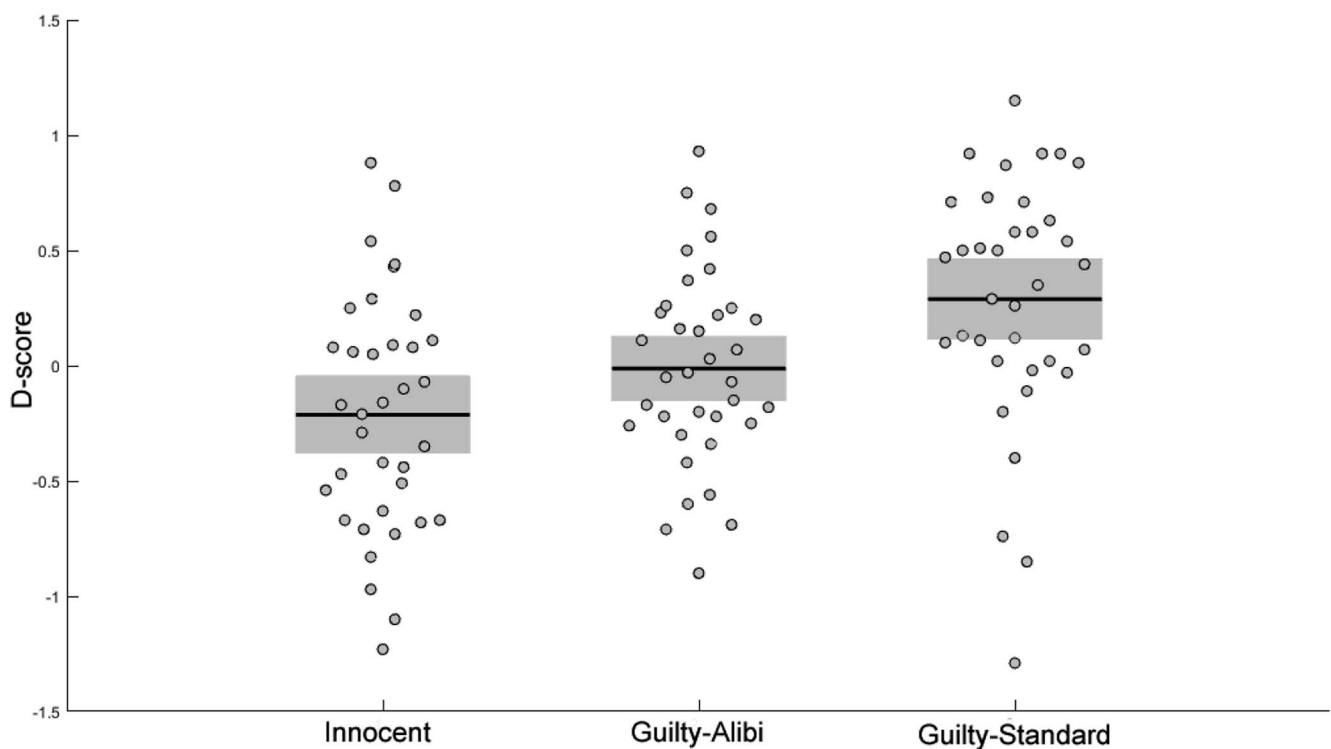
threshold-based classification). The AUCs reflect the accuracy with which a randomly chosen participant can be classified into the correct group (guilty or innocent), where .5 reflects chance classification and 1.0 reflects perfect classification.

In addition to analyzing the D-score, we also analyzed the raw RT and accuracy rates separately for the guilt-congruent versus incongruent blocks for each group. However, because these analyses only revealed patterns that were consistent with the main D-score findings, they are presented in the [online supplemental materials](#). Furthermore, in a final analysis, we also calculated a “faking index” (Agosta et al., 2011) that has been proposed as a method for detecting whether participants are showing unusual RT patterns that indicates countermeasure use. Therefore, we used the faking index to assess whether rehearsing a false alibi resulted in unusual RT patterns across aIAT blocks that could function as signals of guilt even when the main guilt measure (i.e., D-score) is disrupted by countermeasures. However, this analysis revealed that the faking index did not discriminate well between the groups, so these results are also presented in the [online supplemental materials](#). Individual level data for this project is available at <https://osf.io/wumdy/>.

## Results

Mean D-scores were in the expected direction, with the highest scores in the guilty standard group and the lowest scores in the

innocent group, and were significantly different between the three groups,  $F(2, 105) = 9.46, p < .001, \eta_p^2 = 0.15$  (see [Figure 1](#)). The innocent participants, who undertook the innocent act but did not have any knowledge of the mock crime, elicited D-scores below 0,  $t(35) = -2.48, p = .018, d = 0.41, BF_{10} = 2.55$ . Guilty-standard participants, who committed the mock crime but did not have any knowledge of the innocent act, elicited D-scores above 0,  $t(35) = 3.25, p = .003, d = 0.54, BF_{10} = 13.70$ . The guilty-alibi participants, who committed the mock crime and were also provided with an alibi scenario consistent with the innocent act, elicited D-scores nondistinguishable from 0,  $t(35) = 0.17, p = .87, d = 0.03, BF_{10} = 0.18$ . D-scores were higher in the guilty-standard group than the innocent group, strongly supported by both frequentist and Bayesian statistics,  $t(70) = 4.06, p < .001, d = 0.96, BF_{10} = 179.99$ . However, there was only a nonsignificant trend for higher D-scores in the guilty-alibi compared with the innocent group, and the Bayes factor was very close to 1 and thus inconclusive,  $t(70) = 1.80, p = .076, d = 0.43, BF_{10} = 0.97$ . Importantly, D-scores were significantly reduced in the guilty-alibi group compared with the guilty-standard group, and the Bayes factor indicated substantial evidence in favor of a difference ( $H_1$ ) compared with no difference ( $H_0$ ) between groups,  $t(70) = 2.66, p = .010, d = 0.62, BF_{10} = 4.55$ . These results indicate that, as expected, imagining a fake alibi consistent with innocence impaired memory detection with the aIAT.



*Figure 1.* D-scores for the three groups from the mock crime/innocent event aIAT in Experiment 1. Each dot indicates an individual score. The black lines shows the mean score and the gray boxes show the 95% confidence intervals of the mean. D-scores above 0 suggest guilt (that the mock crime-related sentences are associated with the truth) and D-scores below 0 suggest innocence (that the innocent-related sentences are associated with the truth). Scores are jittered along the x-direction for display purposes.

Because applied uses of the aIAT involves classifying individual suspects as guilty or innocent, we also conducted a ROC analysis to evaluate how accurately our participants could be classified based on their D-scores. This analysis showed that when comparing guilty-standard and innocent groups, D-score classification was significantly better than chance ( $AUC = .70$ ,  $SE = .06$ ,  $p = .004$ ), but comparing guilty-alibi and innocent groups, D-score classification was less accurate and not significantly different than chance ( $AUC = .62$ ,  $SE = .07$ ,  $p = .093$ ). Thus, individual classification rates also supported our prediction that imagining a false alibi would impair memory detection.

## Discussion

In Experiment 1, the aIAT showed relatively good discrimination between guilt and innocence in participants who did not employ countermeasures, consistent with previous findings (e.g., Agosta & Sartori, 2013; Sartori et al., 2008). However, the false alibi countermeasure reduced memory detection when compared with a standard guilty group who were not trying to evade the test, consistent with our predictions. Performance in the innocent group showed a stronger relative association between the innocent act and the truth than the mock crime and the truth, whereas performance in the guilty-standard group indicated the opposite relative association. Performance in the guilty-alibi group however was equivocal as to which scenario was truthful. This pattern indicates that imagining a fake alibi scenario likely created a memory for the imagined alibi act that had some implicit associations with the truth, even though participants knew their alibi was false at an explicit level (cf. Shidlovski et al., 2014; Takarangi et al., 2013, 2015). This account is consistent with more general findings that imagining an event can create a memory for that event that has similar perceptual and behavioral characteristics as memories based on true experiences (e.g., Loftus, 2003; Loftus & Pickrell, 1995; Mitchell & Johnson, 2009; Schacter et al., 2011). Presumably, because both the mock crime and the imagined alibi act had some associations with the truth, neither of the critical aIAT blocks were truly congruent or incongruent with their memories, leading to similar performance in both blocks.

The results are consistent with the explanation that imagining a false alibi increased the implicit truth value of that scenario, which thereby disrupted aIAT discrimination between the alibi and the mock crime. However, imagining a counterfactual version of an event may also interfere with the veridical memory of the event and decrease its implicit truth value (cf. Otgaar & Baker, 2018). Gronau et al. (2015) asked participants to learn a hypothetical crime scenario with various details that were different from a mock crime they had actually conducted. Results showed that learning a false version of the mock crime impaired explicit recall of true crime details, and furthermore, reduced skin-conductance markers of true crime memories. They argued that true crime memories may have become inhibited as a result of retrieval competition between true and false crime details, similarly to the retrieval-induced forgetting phenomenon (Anderson, Bjork, & Bjork, 2000, 1994; Anderson & Levy, 2007), or alternatively, that the memory for alibi information interfered with and blocked access to the memory for the true mock crime (see Anderson & Neely, 1996, for review). Because the aIAT in Experiment 1 measured the relative truth of the false alibi versus mock crime scenarios, we can

conclude that these scenarios had similar implicit truth values in the alibi countermeasure group. However, we cannot determine whether the lack of a difference was due to increased implicit truth value of the false alibi, or reduced implicit truth value of the mock crime, or a combination of both. This issue was addressed in the next experiment.

## Experiment 2

Experiment 2 used exactly the same false alibi manipulation, materials, and procedure as in Experiment 1, with the only change being that the final test involved a different aIAT design that contrasted the mock crime with an unexperienced event that was clearly different from the learned false alibi. Thus, this study investigated whether imagining a false alibi would still impair detection of the mock crime regardless of which other scenario it is compared to. If such a pattern was found, it would indicate that the implicit truth value of the original crime-related memory was weakened by rehearsing an alibi, because any reduction in mock crime detection in this aIAT could not be due to inflated implicit truth value of the imagined alibi event as this scenario was not used as a contrast in the test.

We hypothesized that if the alibi manipulation was successful at reducing the implicit truth value of the true mock crime memory, perhaps by reducing access to this memory through inhibition or an interference “blocking” mechanism (Anderson et al., 1994; Anderson & Levy, 2007; Gronau et al., 2015), then rehearsing an alibi should reduce detection of guilty suspects on the aIAT by lowering their D-scores when compared with guilty suspects who did not rehearse an alibi after committing the mock crime. As a consequence, the D-scores for guilty suspects who rehearsed an alibi should be more similar to the innocent group than to the guilty-standard group. Alternatively, if our previous finding was caused only by an increase in implicit truth value of the alibi scenario due to an imagination inflation-related process (e.g., Loftus & Pickrell, 1995; Shidlovski et al., 2014), then there should be no difference in aIAT performance between the guilty-alibi and guilty-standard groups as guilt detection rates in both groups should be equal, but both groups should have higher D-scores and be more likely to be detected as guilty than the innocent group.

## Method

**Participants.** The final sample consisted of 108 undergraduate students from the University of Kent who took part via a research participation scheme in return for course credits ( $M_{age} = 18.94$  years,  $SD = 1.98$ , age range = 18–36 years), maintaining the same statistical power as in Experiment 1. Twelve additional participants were excluded due to technical errors or failures to follow instructions. Participants were randomly assigned to three experimental groups ( $N = 36$  in each group): the guilty-alibi group (31 female, five male), the guilty-standard group (33 female, three male), and the innocent group (30 female, six male). The groups did not differ in age,  $F(2, 105) = 0.78$ ,  $p = .461$ ,  $\eta_p^2 = 0.02$ , nor gender,  $\chi^2(2) = 1.15$ ,  $p = .563$ ,  $\phi = 0.10$ . All participants had English as their first language, had normal or corrected-to-normal vision, and had no diagnosis of dyslexia. The study was approved by the University of Kent Psychology Ethics committee.

**Materials, design, and procedure.** The materials, design, and procedure was identical to Experiment 1 with one exception: The

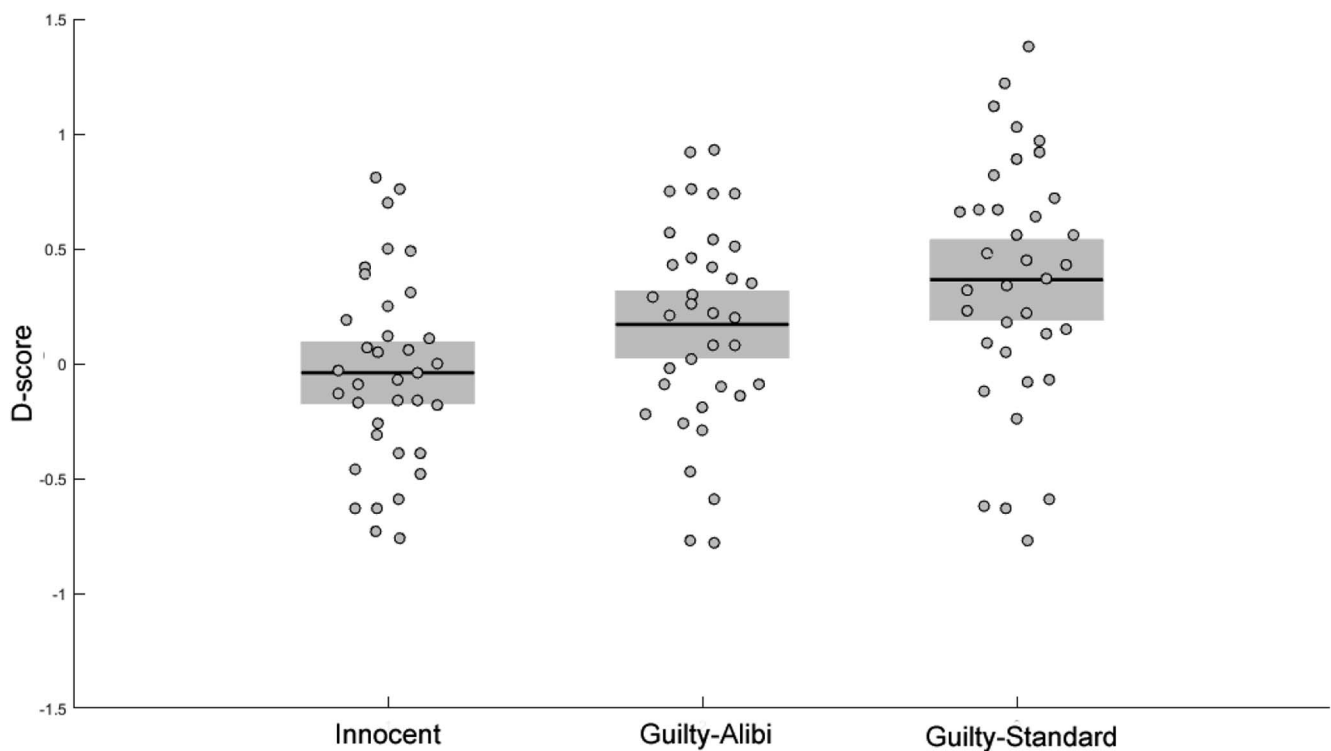
aIAT version was different. As in Experiment 1, the study was conducted in three stages. First, participants in the two guilty groups carried out a mock crime in which they required to go to an office block and steal a ring from a bag, while innocent participants carried out an innocent act, involving writing their e-mail address on a paper in the same area as the guilty participants. Next, half of the guilty participants were instructed to imagine performing the innocent act as a fake alibi with the explicit intention to use it as a strategy to appear innocent. The rest of participants performed a filler task. Finally, all three groups took an aIAT, which assessed which of two events had a stronger relative association with the truth. Importantly, instead of contrasting the mock crime and innocent act/false alibi directly, the aIAT in Experiment 2 contrasted the mock crime with a completely novel unexperienced event involving entering a lecturer's office and stealing a CD with exam questions on (henceforth referred to as the "exam" event, adapted from Sartori et al., 2008) that should not be associated with any truth value for any of the groups. All aspects of the aIAT task design and instructions were the same as in Experiment 1, with the only change being that sentences related to the alibi/innocent act were replaced with sentences related to the unexperienced event. As in Experiment 1, the order of the guilt congruent versus incongruent blocks was counterbalanced across participants, and an equal number of participants within each group received each order.

After the main experiment, all participants completed a questionnaire where they rated how they had experienced and con-

ducted the different tasks. They rated how nervous they had been while conducting the mock crime/innocent act (as applicable; on a 0–6 scale where 0 = *not nervous at all*; 6 = *extremely nervous*), and how often they were thinking about the mock crime/innocent act during the aIAT (0 = *not at all*; 6 = *all the time*). The two guilty groups also rated their motivation to beat the aIAT (0 = *not motivated at all*; 6 = *extremely motivated*), and answered open-ended questions on whether they used any strategy to intentionally distort the test. There were also two additional questions for guilty-alibi participants: how vividly they had been able to imagine the alibi (0 = *not vivid at all*; 6 = *extremely vivid*) and how often they were thinking about the alibi during the aIAT (0 = *not at all*; 6 = *all the time*).

## Results

The mean standardized D-score indices of guilt (Greenwald et al., 2003; Hu et al., 2015) were significantly different between the groups,  $F(2, 105) = 6.73, p = .002, \eta_p^2 = 0.11$  (see Figure 2). Innocent participants, who had no knowledge of neither the mock crime nor the novel "exam" event, obtained a D-score that was not significantly different from 0 as expected, and the Bayes factor showed relatively stronger evidence for no difference than a difference,  $t(35) = -0.57, p = .569, d = 0.10, BF_{10} = 0.21$ . Guilty-standard participants, who committed the mock crime and did not have any knowledge of the exam event, elicited D-scores significantly above 0, strongly supported by a very large Bayes



*Figure 2.* D-scores for the three groups from the mock crime/unexperienced event aIAT in Experiment 2. The black lines shows the mean scores and the gray boxes show the 95% confidence intervals of the mean. D-scores above 0 suggest guilt (that the ring-related sentences are associated with the truth). D-scores close to 0 suggest that the events were equally associated with the truth, but because the test did not include a truly "innocent" event, innocence cannot be detected in this aIAT version.

factor,  $t(35) = 4.10$ ,  $p < .001$ ,  $d = 0.68$ ,  $BF_{10} = 115.39$ . The guilty-alibi participants, who also committed the mock crime and did not have any knowledge about the exam event, also elicited D-scores significantly above 0, but with only anecdotal support for a difference from the Bayes factor,  $t(35) = 2.28$ ,  $p = .029$ ,  $d = 0.38$ ,  $BF_{10} = 1.75$ . D-scores were significantly lower in the innocent group than guilty-standard,  $t(70) = 3.59$ ,  $p < .001$ ,  $d = 0.85$ ,  $BF_{10} = 46.66$ ; and guilty-alibi groups,  $t(70) = 2.06$ ,  $p = .043$ ,  $d = 0.49$ ,  $BF_{10} = 1.48$ . There was also a nonsignificant trend toward lower D-scores in the guilty-alibi than guilty-standard group, but for this test the Bayes factor was weakly more supportive of no group difference than a difference,  $t(70) = 1.67$ ,  $p = .099$ ,  $d = 0.39$ ,  $BF_{10} = 0.80$ . These results indicate that imagining a false alibi does not abolish the implicit truth value of the true crime memory since the mean D-score was significantly above 0 in the guilty-alibi group, and there was now only a weak, nonsignificant tendency, and no Bayesian support for reduced aIAT memory detection in this group compared to the standard guilty condition.

A threshold-independent ROC analysis to evaluate classification performance showed that when comparing guilty-standard and innocent groups, D-score classification was significantly better than chance ( $AUC = .73$ ,  $SE = .06$ ,  $p = .001$ ). Comparing guilty-alibi and innocent groups, D-score classification was lower, but also better than chance ( $AUC = .64$ ,  $SE = .07$ ,  $p = .043$ ). The D-score classification results thus indicated that rehearsing an alibi did not fully impair the original memory of the mock crime because these participants could still be detected as guilty, yet there was a subtle numerical reduction in guilt classification for guilty-alibi participants.

Ten participants (four innocent, three guilty-standard, and three guilty-alibi) were excluded from the questionnaire analysis due to missing responses. The results revealed no differences between guilty-standard ( $M = 2.76$ ,  $SD = 1.60$ ) and guilty-alibi ( $M = 2.60$ ,  $SD = 1.46$ ) groups in nervousness during the mock crime,  $t(64) = 0.40$ ,  $p = .689$ ,  $d = 0.10$  and the extent to which they thought about the mock crime during the aIAT ( $M = 3.21$ ,  $SD = 1.53$ ;  $M = 3.52$ ,  $SD = 1.17$ , respectively;  $t(64) = 0.90$ ,  $p = .372$ ,  $d = 0.23$ ). However, there was a significant difference between guilty groups in their motivation to beat the test: The guilty-alibi ( $M = 4.15$ ,  $SD = 1.18$ ) group was more motivated to appear innocent than the guilty-standard group ( $M = 3.45$ ,  $SD = 1.35$ ;  $t(62) = 2.24$ ,  $p = .029$ ,  $d = 0.56$ ). The innocent group reported being significantly less nervous while conducting the innocent task than the guilty groups were while conducting the mock crime (innocent  $M = 1.78$ ,  $SD = 1.60$ ; innocent vs. guilty-alibi:  $t(63) = 2.17$ ,  $p = .033$ ,  $d = 0.55$ ; innocent vs. guilty-standard:  $t(63) = 2.46$ ,  $p = .017$ ,  $d = 0.62$ ). They also thought less about the innocent act during the aIAT than the two guilty groups thought about the mock crime during the aIAT (innocent  $M = 1.00$ ,  $SD = 1.50$ ; innocent vs. guilty-alibi:  $t(63) = 7.53$ ,  $p < .001$ ,  $d = 1.90$ ; innocent vs. guilty-standard:  $t(63) = 5.87$ ,  $p < .001$ ,  $d = 1.48$ ), as would be expected because there were no sentences related to the innocent act in this aIAT version. Exploratory correlation analyses were also conducted to investigate whether any of the self-report measures correlated with performance in the aIAT, but there were no significant correlations.

## Discussion

Experiment 2 assessed whether imagining a false alibi reduces the implicit truth value of the true crime memory, in line with previous findings that have shown that learning counterfactual details after a mock crime can impair true memories of the crime (Gronau et al., 2015). In Experiment 1, the results showed that the aIAT was unable to determine whether an experienced mock crime or an imagined false alibi was true. However, the aIAT design did not permit us to test whether this lack of discrimination was caused by increased truth value of the imagined alibi or decreased truth value of the mock crime, or a combination of both. In Experiment 2, we therefore contrasted the mock crime with a novel event that had been neither experienced nor imagined in an aIAT, in order to assess the implicit truth value of the mock crime memory independent of the alibi memory. In this study, the mock crime was still detected despite participants previously imagining a false alibi, suggesting that the alibi had not impaired the true memory of the crime to a substantial extent.

As expected in Experiment 2, the mean D-score of innocent participants was close to 0, suggesting that neither event was strongly associated with the truth in this group. Both guilty groups scored above 0, indicating that they associated the mock crime with the truth more than the unexperienced event. Therefore, it appears that the low discrimination between the experienced mock crime and imagined alibi in Experiment 1 was mainly driven by the alibi manipulation increasing the implicit truth value of the imagined scenario, rather than a reduction of implicit truth value of the mock crime memory. This finding contrasts with other research that has suggested that rehearsing a false alibi can cause it to become a default response such that when a cue triggers a memory about a crime, that memory is automatically inhibited to facilitate a false alibi response (Foerster et al., 2017), and that thinking counterfactually can impair memories for the event that actually occurred (Petrocelli & Crysel, 2009; see also Otgaar & Baker, 2018).

One possible reason why the true mock crime memory was unimpaired in Experiment 2 might be that the alibi manipulation was only implemented through one brief rehearsal and imagination phase. Thus, the effect of the alibi manipulation may not have been as strong as in real-life situations where suspects may prepare and imagine an alibi repeatedly and over a long-time period before the interrogation. If participants were able to rehearse/imagine the alibi in this way, it may be more likely to impair the true memory of the mock crime, either by increased retroactive interference or by inhibition of the crime memory representation itself (e.g., Gronau et al., 2015). Previous research has suggested that when multiple memories are associated to the same cue, repeatedly retrieving one memory in the face of competitive activation of another memory can cause the nonselected memory to become inhibited (Anderson et al., 1994). Likewise, repeatedly pushing an unwanted memory out of mind by thinking of a substitute thought may interfere with (Bergström, de Fockert, & Richardson-Klavehn, 2009) retrieval of the original memory, or even inhibit it (Benoit & Anderson, 2012). The literature on motivated forgetting suggests that such impairments of unwanted memories are gradual and increase with repetition (e.g., Anderson & Green, 2001), predicting that a true crime memory might only become impaired if a false alibi is repeatedly retrieved. Likewise, the retroactive interference theory suggests that repeatedly rehearsing



one memory associated to a cue may strengthen that association, which can block access to other associated memories without those memories being inhibited (see Anderson & Neely, 1996). Thus, multiple theoretical accounts suggest that repeated and temporally extended imagination of an alibi should be more likely to impair access to the original crime memory, as addressed in the next experiment.

### Experiment 3

Experiment 3 was designed to replicate and extend on findings from the previous studies, with particular focus on whether repeated rehearsal of a false alibi over an extended time period might be more effective at impairing the true memories compared to a single brief alibi intervention just before the aIAT. In the previous two experiments all experimental phases were conducted in the same session; participants first conducted a mock crime, then immediately learned and imagined the false alibi, which was followed by the aIAT. We therefore added a time delay of 1 week between the mock crime and test, which made the design more realistic and enabled us to investigate the effect of repeated and distributed false alibi rehearsal on aIAT memory detection.

The experimental design was similar to the previous studies, except that it was conducted in two sessions 1 week apart, and included an additional experimental group. Furthermore, in the second session, all participants completed three versions of the aIAT that contrasted the mock crime versus the innocent/alibi event (same aIAT as in Experiment 1), the mock crime versus an unexperienced event (same aIAT as in Experiment 2), and the alibi versus the unexperienced event (a new aIAT version to assess the implicit truth value of the innocent act/alibi independently of the mock crime). Similarly to previous experiments, participants first conducted either an innocent act or a mock crime, depending on which group they were assigned to. All participants then came back for the aIAT session a week later. In one countermeasure group (“guilty-alibi”), participants conducted a mock crime during the first session, then left and returned a week later at which point they learned and imagined the false alibi immediately before the aIATs. In the other countermeasure group (“guilty-alibi with home training”), participants learned and imagined the false alibi during the first session immediately after conducting the mock crime, and were also required to repeat this imagination task at home once a day for a week before returning to complete the aIATs. These two countermeasure groups were compared against innocent and guilty-standard groups, as in the previous two studies.

We expected that participants who carried out an innocent act should be detected as innocent and participants who committed a mock crime without learning an alibi should be detected as a guilty across the relevant aIAT versions. However, participants who learned the false alibi would be less likely to be detected as guilty than the standard guilty group. If imagining a false alibi leads to gradual strengthening of the false alibi information in memory and/or gradual impairment of the true memory with repetition, then extended rehearsal of a false alibi for a week before the test should be particularly effective at reducing detection of guilty suspects.

### Method

**Participants.** The final sample consisted of 144 undergraduate students from the University of Kent who took part via a

research participation scheme in return for course credits ( $M_{age} = 19.13$  year,  $SD = 1.57$ , age range = 18–34 years). Twenty-eight additional participants were excluded due to technical errors, failures to follow instructions, or failure to attend both sessions. Participants were randomly assigned to one of the four groups ( $N = 36$  in each group): innocent (30 female, six male), guilty-standard (30 female, six male), guilty-alibi (27 female, six male), and guilty-alibi with home training (HT; 31 female, five male). Thus, this experiment maintained the same statistical power as the previous two experiments for pairwise comparisons between groups. The groups did not differ in terms of age,  $F(3, 140) = 0.74, p = .531, \eta_p^2 = .02$ , nor gender,  $\chi^2(3) = 1.69, p = .639, \phi = 0.11$ . All participants had English as their first language, had normal or corrected-to-normal vision, and had no diagnosis of dyslexia. The study was approved by the University of Kent Psychology Ethics committee.

**Materials, design, and procedure.** First, participants in all three guilty groups committed a mock crime involving going to a staff office area and stealing a ring, whereas participants in the innocent group completed an innocent task involving writing their e-mail address on a sign-up sheet in the same area (both these tasks were kept identical to Experiments 1 and 2). Next, all participants were dismissed and asked to come back to the laboratory after a week, except the guilty-alibi with HT group. The latter group were given instructions to perform an extra task after completing the mock crime. They first learned and imagined a false alibi, which described the innocent act, using the same materials and procedure as in Experiments 1 and 2. Next, they were given a home training task, which required them to access an Internet link in order to rehearse the false alibi once every day in the intervening 6 days until the test day. When they accessed the link, they were asked to read a description of the alibi (using the same text as used on the first day) and imagine themselves completing the described actions as vividly and accurately as possible. After that, they were asked to write down a detailed description of the scenario they had imagined and rate how vivid their imagination of the alibi had been. Participants were only included in the final sample if they had completed this task as instructed.

After a week, all participants came back to the lab to complete the rest of the study. Participants in innocent and guilty-standard group were asked to complete a filler task (solving Sudoku puzzles), while the two alibi groups rehearsed the alibi (describing the innocent act). For the guilty-alibi group, this was the first time they learned that they needed to use an alibi to appear innocent and found out the details of the alibi/innocent act, whereas for the guilty-alibi with HT group it was another chance to rehearse the alibi they had learned and repeatedly rehearsed during the preceding week. Finally, all participants completed three versions of the aIAT: (a) contrasting the mock crime versus the innocent/alibi event (same aIAT as in Experiment 1); (b) contrasting the mock crime versus the unexperienced event involving stealing an exam (same aIAT as in Experiment 2); and (c) contrasting the innocent/alibi versus the unexperienced event (a novel aIAT version used to assess whether the innocent event would be detected as true after rehearsing a false alibi). The aIAT task design, sentences, and instructions were identical to those used in the previous studies, with the only changes being the added new Version 3, and that all participants undertook all three versions. The order of aIAT congruent/incongruent blocks and order of versions was fully coun-

terbalanced across participants to prevent order effect confounds, and counterbalancing formats were equally distributed within each of the four groups.

After the experiment, participants were asked to complete a questionnaire, which was similar to the one used in Experiment 2 with a few additional questions about details of the innocent act or mock crime. For the innocent group, participants were required to give answers relating to details of the innocent act and give ratings on a scale from 0 to 6 regarding their behavior and experience during the initial act and the aIAT (e.g., in how much detail they could remember the act, their motivation to beat the aIAT, and the extent to which they thought about the act during the aIAT). The guilty groups were asked to provide answers regarding details of the mock crime and provide various ratings on a 0–6 scale regarding their nervousness during the mock crime, their motivation to beat the aIAT, the extent to which they thought about the mock crime during the aIAT, and whether they had intentionally used any strategy to distort the test, including the extent to which they thought about the alibi scenario during the aIAT and how vividly they had imagined an alibi (for the guilty-alibi groups only).

## Results

**Mock crime/innocent event aIAT.** The mock crime/innocent version of the aIAT directly contrasted the mock crime (ring) with the innocent/alibi (e-mail) event, and was identical to the aIAT

used in Experiment 1. In this test, positive D-scores (Greenwald et al., 2003; Hu et al., 2015) are indicative of guilt because they suggest participants associate the mock crime with the truth whereas negative D-scores are indicative of innocence because they suggest participants associate the innocent event with the truth. The mean D-scores were significantly different between the groups,  $F(3, 140) = 6.78, p < .001, \eta_p^2 = 0.13$  (see Figure 3). The mean D-score of the innocent group was not significantly different from 0, with the Bayes factor indicating (weak) relative support of no difference over a difference,  $t(35) = -1.30, p = .20, d = 0.22, BF_{10} = 0.39$ , inconsistent with the predictions and suggesting that the innocent event was on average not detected as true in this group. The guilty-standard group, however, did obtain a D-score that was significantly above 0 with strong supporting evidence from the Bayes factor,  $t(35) = 3.75, p < .001, d = 0.63, BF_{10} = 47.32$ , indicating successful guilt detection in this group. The guilty-alibi group who committed a mock crime and learned a false alibi just prior to the test, however, had a mean score significantly below 0,  $t(35) = -2.06, p = .047, d = 0.34, BF_{10} = 1.18$ , thus appearing more innocent than guilty, although the Bayes factor was only weakly supportive of a difference from 0 in this group. In contrast, the guilty-alibi with HT group, who committed a mock crime and then repeatedly rehearsed a false alibi for a week before the test, did not have a mean D-score that differed from 0,  $t(35) = 1.01, p = .320, d = 0.17, BF_{10} = 0.29$ . Independent *t* tests revealed that the mean D-score of the innocent group was signif-

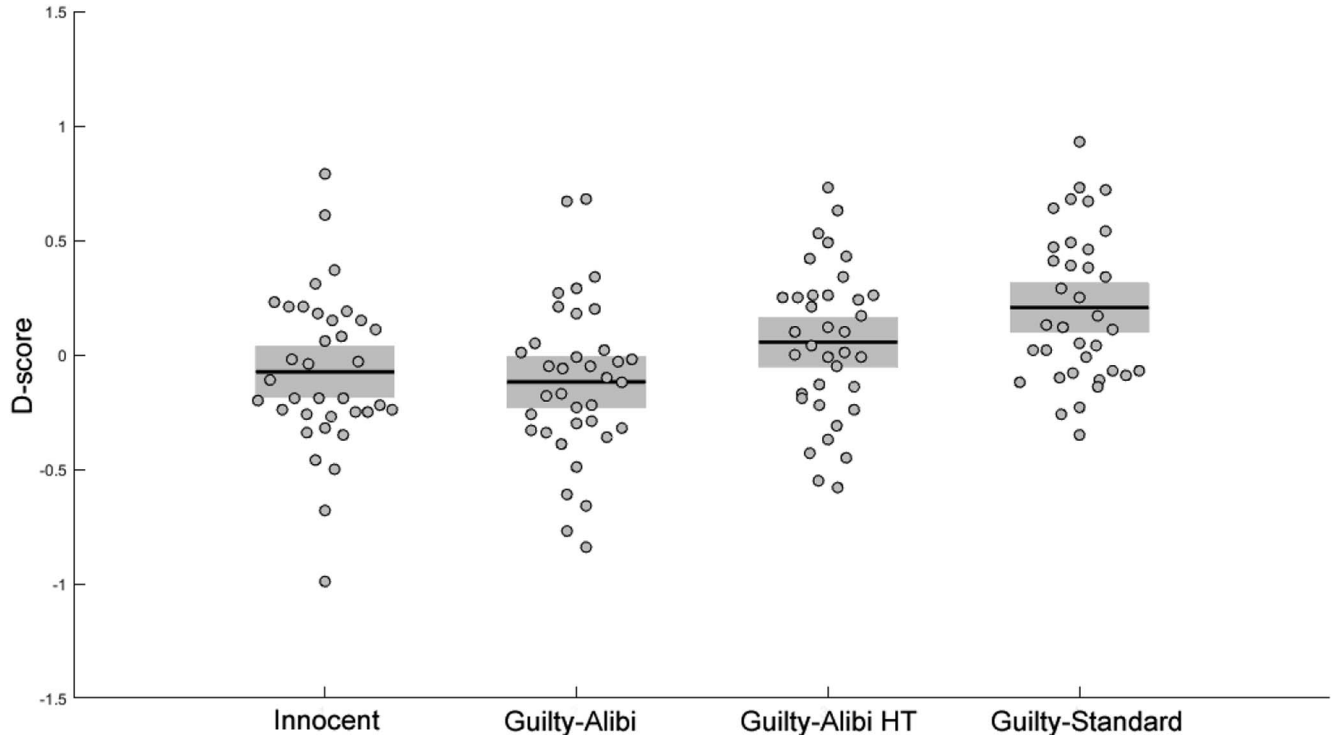


Figure 3. D-scores for the four groups from the mock crime/innocent event aIAT in Experiment 3. The black lines shows the mean score and the gray boxes show the 95% confidence intervals of the mean. D-scores above 0 suggest guilt (that the ring-related sentences are associated with the truth) and D-scores below 0 suggest innocence (that the e-mail-related sentences are associated with the truth). HT = home training.

icantly lower than in the guilty-standard group,  $t(70) = 3.54, p < .001, d = 0.83, BF_{10} = 40.34$ , while there were no differences between the innocent and either of the alibi groups (innocent vs. guilty-alibi:  $t(70) = 0.54, p = .59, d = 0.13, BF_{10} = 0.28$ ; innocent vs. guilty-alibi with HT:  $t(70) = 1.64, p = .10, d = 0.39, BF_{10} = 0.76$ ). However, the mean D-score of the guilty-standard group was significantly higher than in the guilty-alibi group, with strong support for a difference from the Bayes factor,  $t(70) = 4.08, p < .001, d = 0.96, BF_{10} = 194.05$ , but only trend level higher than in the guilty-alibi with HT group with only anecdotal Bayesian support for a difference,  $t(70) = 1.95, p = .056, d = 0.46, BF_{10} = 1.21$ . Surprisingly, the mean D-score of the guilty-alibi with HT group was significantly higher than the guilty-alibi group with anecdotal Bayesian support for a difference between the two alibi groups,  $t(70) = 2.19, p = .03, d = 0.52, BF_{10} = 1.84$ , suggesting that extended training with the alibi actually made it a *less* effective strategy for appearing innocent on this aIAT version.

A threshold-independent ROC analysis to evaluate classification performance showed that when comparing innocent and guilty-standard groups, D-score classification was significantly better than chance ( $AUC = .72, SE = .060, p = .001$ ). However, D-score classification was not accurate when comparing innocent and guilty-alibi groups ( $AUC = .54, SE = .069, p = .581$ ), nor when comparing innocent and guilty-alibi with HT groups, although the latter was at trend-level ( $AUC = .62, SE = .067, p = .073$ ).

So in sum, the mock crime/innocent aIAT largely replicated the findings from Experiment 1: Guilty participants who did not use

countermeasures could be detected as guilty, whereas imagining a false alibi led to lower detection rates. However, this countermeasure was most effective when applied only once immediately before the aIAT, contrary to our predictions that extended and repeated alibi rehearsal would enhance the effectiveness of this strategy. Also somewhat surprising was that detection of innocent participants was relatively poor compared with Experiment 1.

**Mock crime/unexperienced event aIAT.** The mock crime/unexperienced event version of the aIAT contrasted the mock crime (ring) with an event that none of the groups had experience nor knowledge of (exam), and was identical to the aIAT version used in Experiment 2. In this test, positive D-scores are indicative of guilt because they suggest that participants associate the mock crime with the truth, whereas D-scores around 0 suggest that participants associate both events equally strongly with the truth (i.e., they associate either both, or neither event with the truth). Because none of the two events is indicative of innocence there is no result that would be diagnostic of innocence in this aIAT version, and no groups were predicted to show negative D-scores. In this test, there was only a trend toward differences between the groups in mean D-scores,  $F(3, 140) = 2.50, p = .062, \eta_p^2 = 0.05$  (see Figure 4), suggesting that this aIAT version did not discriminate between the groups as well as the mock crime/innocent event aIAT (as would be expected since there should be less variability between groups when the test is designed to only produce scores either around 0 or above, and no negative scores). The mean D-scores of guilty-standard,  $t(35) = 3.99, p < .001, d = 0.67$ ,

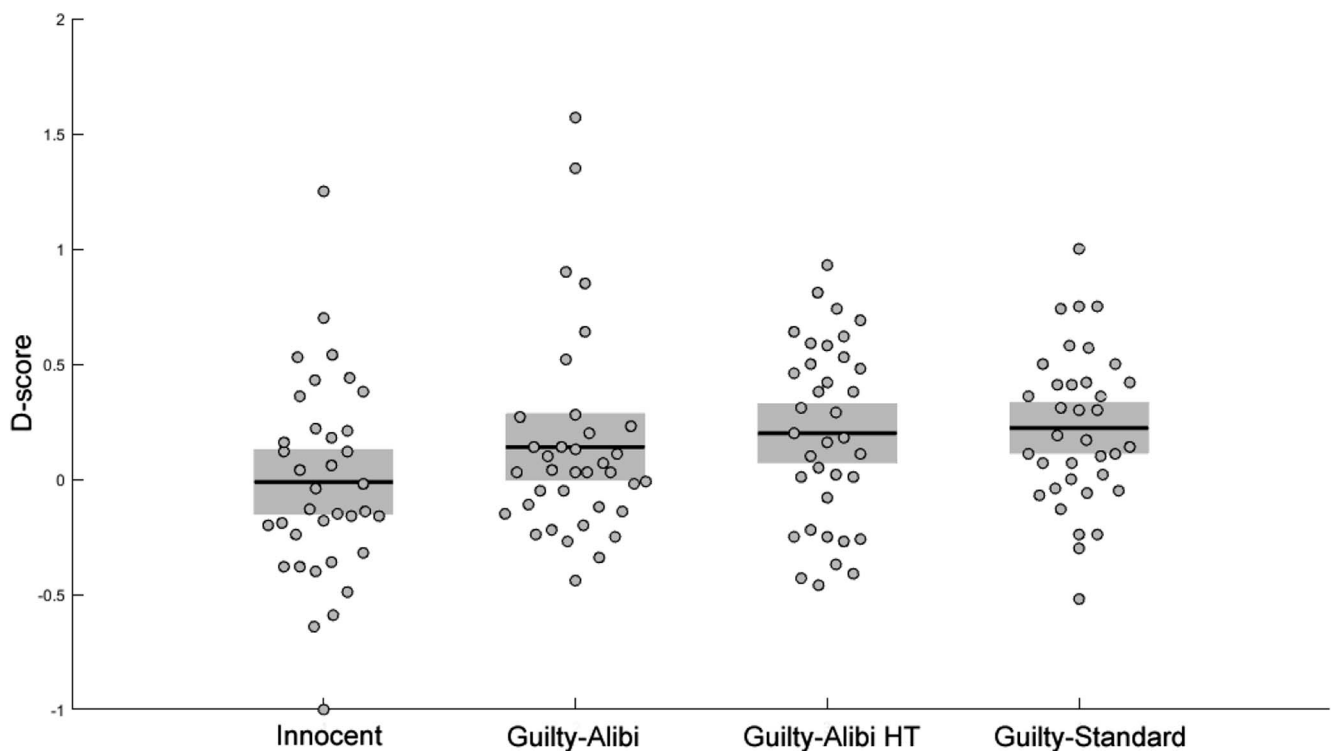


Figure 4. D-scores for the four groups from the mock crime/unexperienced event aIAT in Experiment 3. The black lines shows the mean score and the gray boxes show the 95% confidence intervals of the mean. D-scores above 0 suggest guilt (that the ring-related sentences are associated with the truth). D-scores close to 0 suggest that the events were equally associated with the truth, but because the test did not include a truly “innocent” event, innocence cannot be classified in this aIAT version. HT = home training.

$BF_{10} = 87.79$ , and guilty-alibi with HT group,  $t(35) = 3.07$ ,  $p = .004$ ,  $d = 0.51$ ,  $BF_{10} = 8.97$ , were significantly above 0, supported by large Bayes factors. However, the mean D-scores for innocent,  $t(35) = -0.17$ ,  $p = .868$ ,  $d = 0.03$ ,  $BF_{10} = 0.18$ , and guilty-alibi groups,  $t(35) = 1.91$ ,  $p = .064$ ,  $d = 0.32$ ,  $BF_{10} = 0.91$ , were not significantly different from 0, with the Bayesian evidence more in favor of no difference than a difference. Independent  $t$  tests revealed that the mean D-score of the innocent group was significantly lower than in the guilty-standard group,  $t(70) = 2.59$ ,  $p = .01$ ,  $d = 0.61$ ,  $BF_{10} = 4.04$ , and the guilty-alibi with HT groups,  $t(70) = 2.19$ ,  $p = .03$ ,  $d = 0.52$ ,  $BF_{10} = 1.84$ . However, no significant differences between the groups emerged from the other pairwise comparisons (innocent vs. guilty-alibi:  $t(70) = 1.49$ ,  $p = .14$ ,  $d = 0.35$ ,  $BF_{10} = 0.63$ ; guilty-standard vs. guilty-alibi:  $t(70) = 0.89$ ,  $p = .38$ ,  $d = 0.21$ ,  $BF_{10} = 0.34$ ; guilty-standard vs. guilty-alibi with HT:  $t(70) = 0.27$ ,  $p = .79$ ,  $d = 0.06$ ,  $BF_{10} = 0.25$ ; guilty-alibi vs. guilty-alibi with HT:  $t(70) = 0.61$ ,  $p = .55$ ,  $d = 0.14$ ,  $BF_{10} = 0.29$ ).

Threshold independent ROC analyses showed that D-score classification performance was above chance when comparing the innocent and guilty-standard groups ( $AUC = .68$ ,  $SE = .064$ ,  $p = .009$ ) and when comparing the innocent and guilty-alibi with HT groups ( $AUC = .62$ ,  $SE = .065$ ,  $p = .037$ ). However, classification performance was not accurate when comparing the innocent and guilty-alibi groups ( $AUC = .59$ ,  $SE = .068$ ,  $p = .207$ ).

To summarize, results of the mock crime/unexperienced event aIAT in Experiment 3 replicated the findings from Experiment 2 that guilty participants who did not use countermeasures could be detected as guilty when compared with an innocent group. Consistent with results from the mock crime/innocent event version in Experiment 3, the mock crime/unexperienced event aIAT also indicated that whereas the guilty-alibi with HT group could be detected as guilty, the guilty-alibi group without HT appeared less guilty (they were not significantly different from the innocent group in any analysis). This pattern again suggests that the false alibi countermeasure was most effective when applied only once immediately before the aIAT, contrary to our predictions. However, the effects of the alibi manipulation were weaker on this version of the aIAT compared with the mock crime/innocent event aIAT, because the guilty-alibi group did not show a significant reduction in D-score compared with the guilty-standard group. Thus, consistent with Experiments 1 and 2, the alibi manipulation was most effective when the mock crime and alibi were directly contrasted, and was less effective when the mock crime was contrasted with an unexperienced event.

**Innocent/unexperienced event aIAT.** The innocent/unexperienced event version of the aIAT contrasted the innocent/alibi event (involving writing an e-mail) with an event that none of the groups had experience nor knowledge of (stealing an exam) in order to assess whether the innocent/alibi event would be detected as true for any of the groups. That is, would learning and rehearsing a false alibi lead that scenario to be detected as true, or would it only be detected as true for the innocent group who had actually conducted the act? In this test, positive D-scores are indicative of innocence because they suggest that participants associate the e-mail event with the truth, whereas D-scores around 0 suggest that participants associate both events equally strongly with the truth (i.e., they associate either both, or neither event with the truth). Because neither of the two events is indicative of guilt, there is no

result that would be diagnostic of guilt in this aIAT version, and no groups were predicted to show negative D-scores. In this test, the mean D-score of the guilty-standard group was not different from 0,  $t(35) = 0.09$ ,  $p = .928$ ,  $d = 0.02$ ,  $BF_{10} = 0.18$ , as expected, because this group had no knowledge of either event. In contrast, the guilty-alibi,  $t(35) = 2.28$ ,  $p = .029$ ,  $d = 0.38$ ,  $BF_{10} = 1.73$ , and guilty-alibi with HT groups,  $t(35) = 2.23$ ,  $p = .033$ ,  $d = 0.37$ ,  $BF_{10} = 1.58$ , did score significantly above 0, suggesting that the alibi was detected as if true on average in these groups (although with only weak support from the Bayes factor). Surprisingly however, the innocent group's mean D-score was not significantly above 0 and the Bayes factor indicated relative support for no difference from 0,  $t(35) = 0.40$ ,  $p = .687$ ,  $d = 0.07$ ,  $BF_{10} = 0.19$ , showing a failure of the test to detect the innocent event even though it was actually true for that group. There was also no overall significant difference between the groups in mean D-scores,  $F(3, 140) = 1.95$ ,  $p = .124$ ,  $\eta_p^2 = 0.04$  (see Figure 5), suggesting that this aIAT version did not discriminate between the groups well. Comparing differences in mean D-score between groups using independent  $t$  tests, there were nonsignificant trends toward more positive D-scores in the two alibi groups than in the guilty-standard group (guilty-alibi vs. guilty-standard:  $t(70) = 1.81$ ,  $p = .08$ ,  $d = 0.43$ ,  $BF_{10} = 0.98$ ; guilty-alibi with HT vs. guilty-standard:  $t(70) = 1.79$ ,  $p = .08$ ,  $d = 0.42$ ,  $BF_{10} = 0.95$ ) but none of the other differences approached significance and the Bayesian analysis indicated relatively more support for no difference than a difference for all comparisons (innocent vs. guilty-standard:  $t(70) = 0.37$ ,  $p = .72$ ,  $d = 0.09$ ,  $BF_{10} = 0.26$ ; innocent vs. guilty-alibi:  $t(70) = 1.36$ ,  $p = .18$ ,  $d = 0.32$ ,  $BF_{10} = 0.54$ ; innocent vs. guilty-alibi with HT:  $t(70) = 1.35$ ,  $p = .18$ ,  $d = 0.32$ ,  $BF_{10} = 0.53$ ; guilty-alibi vs. guilty-alibi with HT:  $t(70) = 0.01$ ,  $p = .99$ ,  $d < 0.01$ ,  $BF_{10} = 0.24$ ).

Threshold-independent ROC analyses revealed that D-score classification based on the innocent/unexperienced event aIAT was inaccurate. Comparing the innocent group with the guilty-standard group, classification performance was at chance ( $AUC = .52$ ,  $SE = .069$ ,  $p = .787$ ), and it was only slightly better but still not significant when comparing innocent participants with guilty-slibi ( $AUC = .59$ ,  $SE = .067$ ,  $p = .177$ ) and guilty-alibi with HT ( $AUC = .59$ ,  $SE = .068$ ,  $p = .169$ ).

Thus, in this aIAT version, we found very poor detection of the participants who had actually performed the innocent act, whereas imagining a false alibi seemed to have slightly increased detection of this false scenario as true in the two alibi groups. However, because the groups were not significantly different from each other in mean D-scores or classification rates, this slight increase in the alibi groups was not reliable.

**Postexperiment questionnaire.** Results from the final questionnaire are shown in Table 1. The innocent group rated their memory of the innocent act as less vivid than the three guilty groups rated their memory for the mock crime act (innocent vs. guilty-standard:  $t(70) = 3.46$ ,  $p = .001$ ,  $d = 0.83$ ; innocent vs. guilty-alibi:  $t(70) = 3.39$ ,  $p = .001$ ,  $d = 0.81$ ; innocent vs. guilty-alibi with HT:  $t(70) = 4.45$ ,  $p < .001$ ,  $d = 1.06$ ) and they also reported that they remembered fewer details of the act (innocent vs. guilty-standard:  $t(70) = 4.42$ ,  $p < .001$ ,  $d = 1.06$ ; innocent vs. guilty-alibi:  $t(70) = 5.20$ ,  $p < .001$ ,  $d = 1.24$ ; innocent vs. guilty-alibi with HT:  $t(70) = 4.93$ ,  $p < .001$ ,  $d = 1.18$ ). The innocent group also reported having been less nervous during the innocent act than the three guilty groups

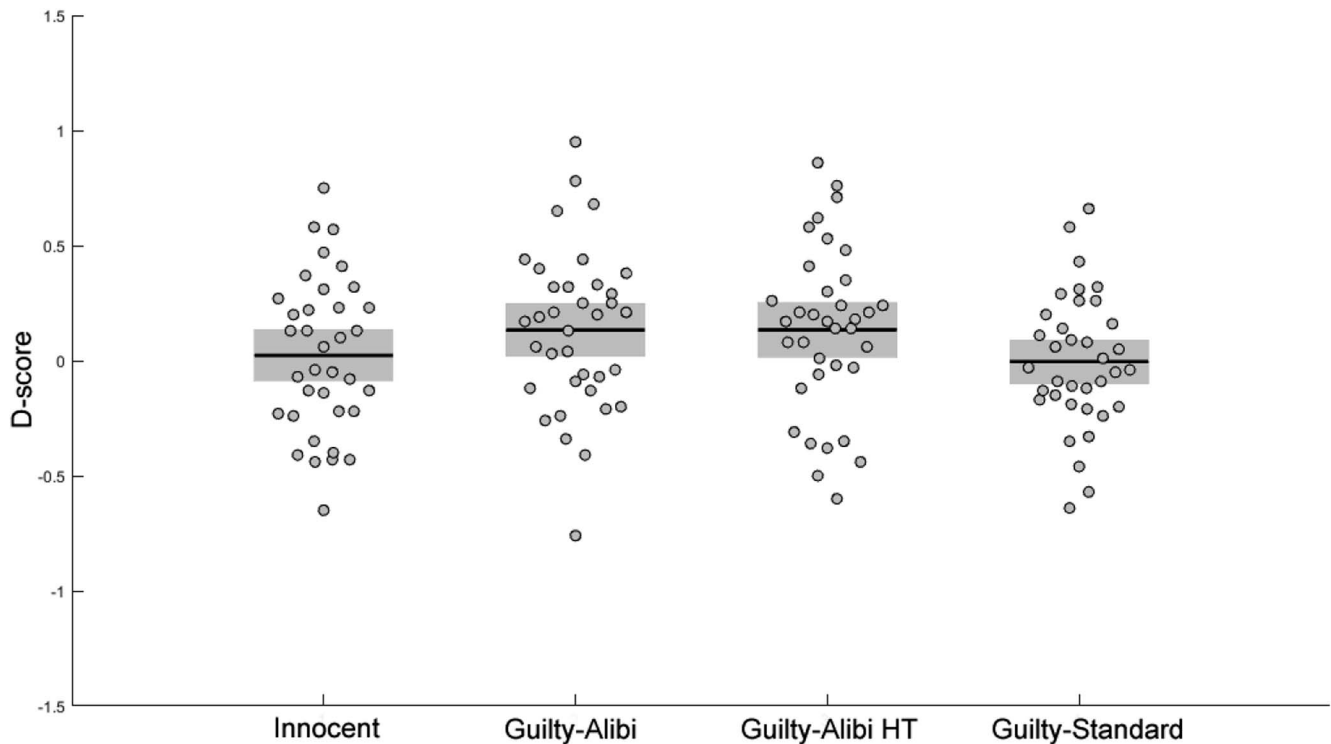


Figure 5. D-scores for the four groups from the innocent/unexperienced event aIAT in Experiment 3. The black lines shows the mean score and the gray boxes show the 95% confidence intervals of the mean. D-scores above 0 suggest innocence (that the e-mail-related sentences are associated with the truth). D-scores close to 0 suggest that the events were equally associated with the truth, but because the test did not include a truly “guilty” event, guilt cannot be classified in this aIAT version. HT = home training.

were when they committed the mock crime (innocent vs. guilty-standard:  $t(70) = 2.80, p = .007, d = 0.67$ ; innocent vs. guilty-alibi:  $t(70) = 2.13, p = .037, d = 0.51$ ; innocent vs. guilty-alibi with HT:  $t(70) = 3.83, p < .001, d = 0.92$ ), and reported thinking about the innocent act less during the aIATs than the three guilty groups thought about the mock crime during the aIATs (innocent vs. guilty-standard:  $t(70) = 3.85, p < .001, d = 0.92$ ; innocent vs. guilty-alibi:  $t(70) = 2.13, p = .037, d = 0.51$ ; innocent vs. guilty-alibi with HT:  $t(70) = 3.95, p < .001, d = 0.94$ ). There were no significant differences between the three guilty groups on any of those questions (all  $ps > 0.14$ ).

The alibi groups and the innocent group were all more motivated to appear innocent on the aIATs than the guilty-standard group (guilty-standard vs. innocent:  $t(70) = 2.04, p = .045, d = 0.49$ ; guilty-standard vs. guilty-alibi:  $t(70) = 3.09, p = .003, d = 0.74$ ; guilty-standard vs. guilty-alibi with HT:  $t(70) = 2.83, p = .006, d = 0.68$ ), but did not differ between each other in levels of motivation (all  $ps > 0.39$ ). With regards to the alibi-specific questions, there were no differences between the alibi groups in terms of how much they were thinking of the alibi during the aIATs,  $t(70) = 0.75, p = .46, d = 0.18$ , but the guilty-alibi with HT group reported being able to

Table 1

Mean and Standard Deviations of Self-Reported Ratings on the Final Questionnaire for the Four Groups

Questionnaire item	Innocent	Guilty-alibi	Guilty-alibi with HT	Guilty-standard
Remember detail of the act	3.39 (1.25)	4.64 (.72)	4.64 (.87)	4.53 (.91)
Vividness of the act memory	3.50 (1.76)	4.36 (.83)	4.69 (.98)	4.44 (1.03)
Nervousness during the act	1.67 (1.29)	2.33 (1.37)	3.05 (1.76)	2.69 (1.79)
Thinking about the act during aIAT	1.58 (1.56)	2.50 (1.68)	3.11 (1.71)	3.08 (1.75)
Motivation to beat the aIAT	3.86 (1.50)	4.14 (1.22)	4.14 (1.50)	3.11 (1.58)
Imagine detail of the alibi	—	3.94 (1.33)	4.57 (.70)	—
Vividness of the alibi imagination	—	3.89 (1.47)	4.57 (.88)	—
Thinking about the alibi during aIAT	—	2.83 (1.68)	3.14 (1.78)	—

Note. The scale had 7 points (0–6), and lower scores always indicate less of the item being measured (e.g. less vividness/nervousness/motivation, etc.) and higher scores always indicate more of the item being measured (e.g. more vividness/nervousness/motivation, etc.). The “act” refers to the act conducted in the first session (i.e. either mock crime or innocent act, depending on group). aIAT = autobiographical implicit association test; HT = home training.

imagine the alibi scenario in more details,  $t(70) = 2.48$ ,  $p = .016$ ,  $d = 0.59$ , and more vividly than the guilty-alibi group,  $t(70) = 2.36$ ,  $p = .021$ ,  $d = 0.56$ . Exploratory correlation analyses were also conducted to investigate whether any of the self-report measures correlated with performance in the aIAT, but there were no significant correlations.

Therefore, in sum, the questionnaire data from Experiment 3 suggested that the innocent group had poorer memory of the innocent act than the guilty groups' memory of the mock crime, whereas repeated and extended rehearsal of the alibi scenario in the guilty-alibi with HT group led to improved ability to imagine the alibi scenario when compared with the guilty-alibi group. Furthermore, the innocent and alibi groups were more motivated to appear innocent on the aIATs than the guilty-standard group.

## Discussion

The aim of this study was to further investigate the effect of rehearsing alibi as a countermeasure on the aIAT (Agosta & Sartori, 2013; Sartori et al., 2008). Previous research suggested that rehearsing a counterfactual scenario to what actually happened during a mock crime can impair access to the true memory (Gronau et al., 2015). In Experiment 3, we investigated whether learning and imagining a false alibi prior to the aIAT would impair the original memory for a mock crime and/or increase the implicit truth value of the alibi itself, and whether these effects would be particularly enhanced when the alibi was repeatedly rehearsed and imagined over an extended time period, in line with theoretical accounts of retrieval interference and inhibition (see, e.g., Anderson & Green, 2001; Anderson & Hanslmayr, 2014; Anderson & Neely, 1996). Such extended and repeated practice of an alibi might be expected to occur in real life, because a guilty criminal might adopt a false alibi and then practice it extensively prior to an investigation several days, weeks, or months later.

The results indicated that in the aIAT that tested the relative strength of the mock crime versus innocent act/alibi, the mock crime was possible to detect after a week delay in guilty-standard participants. However, this aIAT could not distinguish which of the two events were true for innocent participants, nor for the guilty-alibi with HT groups. Interestingly, in the guilty-alibi group that did not receive home training, the test result was more indicative of innocence than guilt. In the aIAT that tested the relative strength of the mock crime versus an unexperienced event, results suggested that the mock crime was possible to detect in guilty-standard and guilty-alibi with HT groups, while it was undetectable in innocent and guilty-alibi groups. In the aIAT that tested the relative strength of the innocent/alibi act versus an unexperienced event, none of the groups showed strong evidence of innocence and this aIAT showed poor discrimination between all groups.

Our findings thus indicate that the strongest effect of the alibi countermeasure was in the guilty-alibi participants who learned and imagined a fabricated alibi one week *after* the mock crime and just prior to the test, without repeated rehearsal. In this group, the results suggested that they associated the imagined false alibi event more with the truth relative to the objectively true mock crime event. Moreover, the aIAT that contrasted the mock crime with an unexperienced event was not able to distinguish which of the two events was true for these guilty participants, suggesting that access to the mock crime memories may have been impaired in this

group. Thus, the effect of the alibi countermeasure in this group was even stronger than the findings in Experiments 1 and 2, where the alibi group did not show significant associations between the alibi and truth (Experiment 1) and they also showed evidence of associating the mock crime with truth when contrasted with the unexperienced event (Experiment 2). These differences across studies may be due to differences in the relative strength of the memory representations for the alibi information versus the mock crime. Mental simulation of the alibi event just before the aIAT may have caused this imagined event memory to be more vivid or salient than the true memory of the mock crime, which may have been weaker in this experiment than in the previous two studies due to the longer time delay between the event and the test. Because of the relatively weak memory for the mock crime, the alibi countermeasure may have been more effective at obscuring detection of that memory than in the previous two studies (cf. Gronau et al., 2015, for related findings with psychophysiological memory detection).

Surprisingly, a different result pattern was observed in the guilty participants who received repeated alibi training for a week before the aIATs. We predicted that extended rehearsal of an imagined alibi would be particularly effective at inducing blocking by retroactive interference or competitive inhibition of the true memory (e.g., Anderson & Hanslmayr, 2014; Anderson & Neely, 1996), and that this group would therefore be more likely to appear innocent compared with a group who only imagined the alibi once just before the test. However, we found the opposite result—although the extended alibi training did reduce memory detection on the aIAT version that directly contrasted the mock crime with the alibi, the magnitude of this reduction was smaller than in the alibi group without extended training. Furthermore, in the aIAT version that contrasted the mock crime with an unexperienced event, the mock crime was still detected as true in the extended training group. These results suggest that extensive and repeated rehearsal of the false alibi did not impair the original mock crime memory, rather, it may have actually strengthened that memory. The home training task may have had an ironic effect of reminding participants of the mock crime and leading the memory for the crime to become strengthened as a result, consistent with prior findings that repeated reminders can enhance automatic influences of memories, which can produce ironic effects when such enhancement affects behavior in unwanted ways (Jacoby, 1999). Future research should assess whether alibi-induced ironic strengthening of the true crime memory can be avoided by explicitly training participants to suppress thoughts of the mock crime while completing the alibi imagination task, which might be an effective strategy for reducing mock crime memory strength while simultaneously strengthening memory for the alibi (cf. Anderson & Green, 2001; Bergström et al., 2013; Hu et al., 2015).

Another surprising finding in Experiment 3 was that none of the aIAT versions detected the innocent act as true for participants in the innocent group despite them actually having conducted the act in real life. In contrast, the mock crime could be detected in the guilty participants who did not use countermeasures. This difference may be related to the 1-week delay that we introduced between the initial act and the aIATs, which may have weakened innocent participants' memory of the innocent act more than it weakened guilty participants' memory of the mock crime. In line with this suggestion, the innocent participants rated their memories

of their act as less vivid and detailed than the guilty participants' ratings of the mock crime memories, and also reported that they had been less nervous while conducting the act than the guilty participants were when conducting the mock crime. This pattern of results suggest that the mock crime memories were associated with higher emotional arousal, which is known to enhance the subjective vividness of memories and their durability over time (Kensinger, 2009). This finding is interesting as it converges with other evidence that memories of recent, familiar events are more detectable in the aIAT than memories of distant, less familiar events (Takarangi et al., 2015) in pointing toward a role of subjective memory quality in aIAT accuracy—the test may only be able to detect memories that are subjectively detailed and vivid, and any factors that reduce memory quality may also reduce the test's effectiveness. It also suggests general limitations with laboratory studies that investigate memory detection with mock crimes, because memories of mock crimes may differ substantially from real criminal memories in terms of emotional arousal. Future research should investigate whether countermeasures can be used against aIAT memory detection of real autobiographical memories that are emotionally arousing.

Overall, the results of Experiment 3 support our hypothesis that rehearsing a false alibi before an aIAT may distort the test results, but they also show that the effectiveness of this strategy depends on how the alibi countermeasure is used, and also on how the aIAT is designed.

### General Discussion

The aIAT has been promoted as an accurate tool for determining which of two autobiographical events are true, with promising applications in forensic memory detection (Agosta & Sartori, 2013; Sartori et al., 2008). However, a growing body of research has revealed potential countermeasures that guilty suspects can adopt to make themselves appear innocent, such as intentionally altering their responses during the test itself (Agosta et al., 2011; Hu et al., 2012; Verschuere et al., 2009), or suppressing their incriminating memories in advance of the test (Hu et al., 2015). We tested whether a novel countermeasure that has recently been applied in physiological memory detection (Gronau et al., 2015) and deception detection paradigms (Foerster et al., 2017; Suchotzki et al., 2018) would also be effective at reducing detection using the aIAT. Specifically, we assessed whether instructing guilty suspects to intentionally store false information in memory would enable those suspects to appear innocent on the test. In line with our predictions, imagining a false alibi impaired memory detection with the aIAT so that the test could no longer distinguish between the objectively true mock crime memory and the objectively false alibi, and this finding was replicated with a large effect size in two experiments. Consistent with previous research (e.g., Agosta & Sartori, 2013; Sartori et al., 2008), our results showed relatively good discrimination between guilt and innocence in participants who did not employ countermeasures. However, the false alibi countermeasure significantly reduced memory detection when compared with a standard guilty group who were not trying to evade the test.

Across experiments, the strongest and most consistent effect of the alibi manipulation occurred on the aIAT version that directly contrasted the mock crime with the alibi to assess their relative

truth value, whereas there were only weaker, less consistent effects on the aIAT that contrasted the mock crime with an unexperienced novel event to detect the truth value of the mock crime itself. This pattern indicates that the effectiveness of the alibi strategy was primarily driven by increased detection of the alibi as true, rather than decreased detection of the mock crime as true. Imagining a false alibi may have created a memory for the alibi scenario that had some implicit associations with the truth, even though participants knew their alibi was false at an explicit level. This account converges with more general findings that imagining an event can create a memory for that event that has similar perceptual and behavioral characteristics as memories based on true experiences (e.g., Loftus, 2003; Loftus & Pickrell, 1995; Mitchell & Johnson, 2009; Schacter et al., 2011), and previous findings that imagining simple actions can increase detection of those actions as true in the aIAT, either by inducing misremembering that imagined actions were actually performed (Takarangi et al., 2013) or sometimes even despite participants knowing the imagined action did not actually happen (Shidlovski et al., 2014).

Our findings thus converge with other research that have found dissociations between explicit and implicit measures of truth (Shidlovski et al., 2014). It has been suggested that these dissociations occur because people can make contrary implicit and explicit evaluations of truth, which may help them deceive both themselves and others (Shidlovski et al., 2014). However, an alternative and more parsimonious explanation is that the aIAT does not actually measure implicit associations between events and the truth, but instead is simply sensitive to the relative salience of different events. In line with this view, asking participants to rehearse and imagine the alibi may have increased the relative salience of this event compared with the mock crime or unexperienced event (cf. Rothermund & Wentura, 2004). Regardless of which account is correct, this uncertainty regarding what the aIAT measures is in our view a fundamental problem for using the aIAT in real criminal cases (see Sirgiovanni et al., 2016)—if researchers do not know what the test is measuring, how can using the test be justified when a false result may have direct real-life consequences? Clearly, practical applications of the aIAT are premature until further research has clarified what the test actually measures, and in what situations it will produce accurate results.

Although our key finding that the false alibi countermeasure reduced the aIAT's ability to discriminate between a true mock crime and a false alibi was strong and robust, our sample sizes and designs were not optimized to detect more subtle changes in guilt detection between groups. For example, there were nonsignificant trends toward differences between groups in several other comparisons (e.g., alibi vs. innocent groups in Experiment 1) that could have been informative if we had increased the statistical power of the design. Likewise, these other group comparisons sometimes produced inconclusive Bayes factors that were not clearly supportive of the alternative nor the null hypothesis, which indicates that the sample sizes were too small to discriminate between these competing hypotheses using Bayesian analyses (see Lakens, Mclatchie, Isager, Roedel, & Dienes, 2018 for discussion). This limitation should be addressed by employing larger sample sizes in future research to better understand variations in truth detection of autobiographical events with the aIAT.

To conclude, we show that imagining a false alibi impaired memory detection with the aIAT because it was unable to distin-

guish between a true mock crime and a false alibi. This finding raises serious concerns for potential real-life applications of this test as a forensic tool with lying, uncooperative suspects. In real life, guilty suspects may spontaneously fabricate false alibis, and investigators may want to use the aIAT to compare the truth value of a suspect's alibi with the crime they are accused of. Our results suggest that such real-life applications may be unsuccessful due to suspects inadvertently modifying their memories by fabricating a false alibi. Furthermore, memories of unethical behavior such as crimes may be particularly susceptible to modification because forgetting immoral acts allow people to maintain a positive self-concept (Kouchaki & Gino, 2016; although see Stanley, Yang, & De Brigard, 2018). Thus, guilty suspects may have several strong motivations to change their memories for self-serving reasons, which in turn may enable them to appear innocent on forensic memory detection tests.

## References

- Agosta, S., Ghirardi, V., Zogmaister, C., Castiello, U., & Sartori, G. (2011). Detecting fakers of the autobiographical IAT. *Applied Cognitive Psychology, 25*, 299–306. <http://dx.doi.org/10.1002/acp.1691>
- Agosta, S., & Sartori, G. (2013). The autobiographical IAT: A review. *Frontiers in Psychology, 4*, 519. <http://dx.doi.org/10.3389/fpsyg.2013.00519>
- Allen, J. J., Iacono, W. G., & Danielson, K. D. (1992). The identification of concealed memories using the event-related potential and implicit behavioral measures: A methodology for prediction in the face of individual differences. *Psychophysiology, 29*, 504–522. <http://dx.doi.org/10.1111/j.1469-8986.1992.tb02024.x>
- Anderson, M. C., Bjork, E. L., & Bjork, R. A. (2000). Retrieval-induced forgetting: Evidence for a recall-specific mechanism. *Psychonomic Bulletin & Review, 7*, 522–530. <http://dx.doi.org/10.3758/BF03214366>
- Anderson, M. C., Bjork, R. A., & Bjork, E. L. (1994). Remembering can cause forgetting: Retrieval dynamics in long-term memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 20*, 1063–1087. <http://dx.doi.org/10.1037/0278-7393.20.5.1063>
- Anderson, M. C., & Green, C. (2001). Suppressing unwanted memories by executive control. *Nature, 410*, 366–369. <http://dx.doi.org/10.1038/35066572>
- Anderson, M. C., & Hanslmayr, S. (2014). Neural mechanisms of motivated forgetting. *Trends in Cognitive Sciences, 18*, 279–292. <http://dx.doi.org/10.1016/j.tics.2014.03.002>
- Anderson, M. C., & Levy, B. J. (2007). Theoretical issues in inhibition: Insights from research on human memory. In D. Gorfein & C. MacLeod (Eds.), *Inhibition in cognition* (pp. 81–102). Washington, DC: American Psychological Association. <http://dx.doi.org/10.1037/11587-005>
- Anderson, M. C., & Neely, J. H. (1996). Interference and inhibition in memory retrieval. In E. L. Bjork & R. A. Bjork (Eds.), *Memory: Handbook of perception and cognition* (2nd ed., pp. 237–313). San Diego, CA: Academic Press. <http://dx.doi.org/10.1016/B978-012102570-0/50010-0>
- Ben-Shakhar, G. (2011). Countermeasures. In B. Verschuere, G. Ben-Shakhar, & E. Meijer (Eds.), *Memory detection: Theory and application of the concealed information test* (pp. 200–214). Cambridge, UK: Cambridge University Press. <http://dx.doi.org/10.1017/CBO9780511975196.012>
- Benoit, R. G., & Anderson, M. C. (2012). Opposing mechanisms support the voluntary forgetting of unwanted memories. *Neuron, 76*, 450–460. <http://dx.doi.org/10.1016/j.neuron.2012.07.025>
- Bergström, Z. M., Anderson, M. C., Buda, M., Simons, J. S., & Richardson-Klavehn, A. (2013). Intentional retrieval suppression can conceal guilty knowledge in ERP memory detection tests. *Biological Psychology, 94*, 1–11. <http://dx.doi.org/10.1016/j.biopsycho.2013.04.012>
- Bergström, Z. M., de Fockert, J. W., & Richardson-Klavehn, A. (2009). ERP and behavioural evidence for direct suppression of unwanted memories. *NeuroImage, 48*, 726–737. <http://dx.doi.org/10.1016/j.neuroimage.2009.06.051>
- Dudai, Y. (2012). The restless engram: Consolidations never end. *Annual Review of Neuroscience, 35*, 227–247. <http://dx.doi.org/10.1146/annurev-neuro-062111-150500>
- Dunlap, W. P., Cortina, J. M., Vaslow, J. B., & Burke, M. J. (1996). Meta-analysis of experiments with matched groups or repeated measures designs. *Psychological Methods, 1*, 170–177. <http://dx.doi.org/10.1037/1082-989X.1.2.170>
- Foerster, A., Wirth, R., Herbold, O., Kunde, W., & Pfister, R. (2017). Lying upside-down: Alibis reverse cognitive burdens of dishonesty. *Journal of Experimental Psychology: Applied, 23*, 301–319. <http://dx.doi.org/10.1037/xap0000129>
- Gamer, M. (2011). Detecting concealed information using autonomic measures. In B. Verschuere, G. Ben-Shakhar, & E. Meijer (Eds.), *Memory detection: Theory and application of the concealed information test* (pp. 27–45). Cambridge, UK: Cambridge University Press. <http://dx.doi.org/10.1017/CBO9780511975196.003>
- Gamer, M., Klimecki, O., Bauermann, T., Stoeter, P., & Vossel, G. (2012). fMRI-activation patterns in the detection of concealed information rely on memory-related effects. *Social Cognitive and Affective Neuroscience, 7*, 506–515. <http://dx.doi.org/10.1093/scan/nsp005>
- Granhag, P. A., Vrij, A., & Verschuere, B. (2015). *Detecting deception: Current challenges and cognitive approaches*. West Sussex, UK: Wiley.
- Greenwald, A. G., Nosek, B. A., & Banaji, M. R. (2003). Understanding and using the implicit association test: I. An improved scoring algorithm. *Journal of Personality and Social Psychology, 85*, 197–216. <http://dx.doi.org/10.1037/0022-3514.85.2.197>
- Gronau, N., Elber, L., Satran, S., Breska, A., & Ben-Shakhar, G. (2015). Retroactive memory interference: A potential countermeasure technique against psychophysiological knowledge detection methods. *Biological Psychology, 106*, 68–78. <http://dx.doi.org/10.1016/j.biopsycho.2015.02.002>
- Hu, X., Bergström, Z. M., Bodenhausen, G. V., & Rosenfeld, J. P. (2015). Suppressing unwanted autobiographical memories reduces their automatic influences evidence from electrophysiology and an implicit autobiographical memory test. *Psychological Science, 26*, 1098–1106. <http://dx.doi.org/10.1177/0956797615575734>
- Hu, X., Rosenfeld, J. P., & Bodenhausen, G. V. (2012). Combating automatic autobiographical associations: The effect of instruction and training in strategically concealing information in the autobiographical implicit association test. *Psychological Science, 23*, 1079–1085. <http://dx.doi.org/10.1177/0956797612443834>
- Jacoby, L. L. (1999). Ironic effects of repetition: Measuring age-related differences in memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 25*, 3–22. <http://dx.doi.org/10.1037/0278-7393.25.1.3>
- JASP Team. (2019). JASP (Version 0.10.2) [Computer software]. Retrieved from <https://jasp-stats.org/>
- Kensinger, E. A. (2009). Remembering the details: Effects of emotion. *Emotion Review, 1*, 99–113. <http://dx.doi.org/10.1177/1754073908100432>
- Kouchaki, M., & Gino, F. (2016). Memories of unethical actions become obfuscated over time. *Proceedings of the National Academy of Sciences of the United States of America, 113*, 6166–6171. <http://dx.doi.org/10.1073/pnas.1523586113>
- Lakens, D., Mclatchie, N., Isager, P. M., Roedel, E. V., & Dienes, Z. (2018). Improving Inferences about null effects with Bayes factors and equivalence tests. *The Journals of Gerontology: Series B*. Advance online publication. <http://dx.doi.org/10.1093/geronb/gby065>
- Loftus, E. F. (2003). Make-believe memories. *American Psychologist, 58*, 867–873. <http://dx.doi.org/10.1037/0003-066X.58.11.867>



- Loftus, E. F., & Pickrell, J. E. (1995). The formation of false memories. *Psychiatric Annals*, *25*, 720–725. <http://dx.doi.org/10.3928/0048-5713-19951201-07>
- Lykken, D. T. (1959). The GSR in the detection of guilt. *Journal of Applied Psychology*, *43*, 385–388. <http://dx.doi.org/10.1037/h0046060>
- Mangiulli, I., Lanciano, T., Jelicic, M., van Oorsouw, K., Battista, F., Curci, A., & Curci, A. (2018). Can implicit measures detect source information in crime-related amnesia? *Memory*, *26*, 1019–1029. <http://dx.doi.org/10.1080/09658211.2018.1441421>
- Marini, M., Agosta, S., Mazzoni, G., Barba, G. D., & Sartori, G. (2012). True and false DRM memories: Differences detected with an implicit task. *Frontiers in Psychology*, *3*, 310. <http://dx.doi.org/10.3389/fpsyg.2012.00310>
- Mitchell, K. J., & Johnson, M. K. (2009). Source monitoring 15 years later: What have we learned from fMRI about the neural mechanisms of source memory? *Psychological Bulletin*, *135*, 638–677. <http://dx.doi.org/10.1037/a0015849>
- Otgaar, H., & Baker, A. (2018). When lying changes memory for the truth. *Memory*, *26*, 2–14. <http://dx.doi.org/10.1080/09658211.2017.1340286>
- Petrocelli, J. V., & Crysel, L. C. (2009). Counterfactual thinking and confidence in blackjack: A test of the counterfactual inflation hypothesis. *Journal of Experimental Social Psychology*, *45*, 1312–1315. <http://dx.doi.org/10.1016/j.jesp.2009.08.004>
- Rosenfeld, J. P., Angell, A., Johnson, M., & Qian, J. H. (1991). An ERP-based, control-question lie detector analog: Algorithms for discriminating effects within individuals' average waveforms. *Psychophysiology*, *28*, 319–335. <http://dx.doi.org/10.1111/j.1469-8986.1991.tb02202.x>
- Rothermund, K., & Wentura, D. (2004). Underlying processes in the implicit association test: Dissociating salience from associations. *Journal of Experimental Psychology: General*, *133*, 139–165. <http://dx.doi.org/10.1037/0096-3445.133.2.139>
- Sartori, G., Agosta, S., Zogmaister, C., Ferrara, S. D., & Castiello, U. (2008). How to accurately detect autobiographical events. *Psychological Science*, *19*, 772–780. <http://dx.doi.org/10.1111/j.1467-9280.2008.02156.x>
- Schacter, D. L., Guerin, S. A., & St Jacques, P. L. (2011). Memory distortion: An adaptive perspective. *Trends in Cognitive Sciences*, *15*, 467–474. <http://dx.doi.org/10.1016/j.tics.2011.08.004>
- Shidlovski, D., Schul, Y., & Mayo, R. (2014). If I imagine it, then it happened: The implicit truth value of imaginary representations. *Cognition*, *133*, 517–529. <http://dx.doi.org/10.1016/j.cognition.2014.08.005>
- Sirgiiovanni, E., Corbellini, C., & Caporale, C. (2016). A recap on Italian neurolaw: Epistemological and ethical issues. *Mind & Society*. Advance online publication. <http://dx.doi.org/10.1007/s11299-016-0188-1>
- Stanley, M. L., Yang, B. W., & De Brigard, F. (2018). No evidence for unethical amnesia for imagined actions: A failed replication and extension. *Memory & Cognition*, *46*, 787–795. <http://dx.doi.org/10.3758/s13421-018-0803-y>
- Suchotzki, K., Berlijn, A., Donath, M., & Gamer, M. (2018). Testing the applied potential of the Sheffield Lie Test. *Acta Psychologica*, *191*, 281–288. <http://dx.doi.org/10.1016/j.actpsy.2018.10.011>
- Suchotzki, K., Verschuere, B., Van Bockstaele, B., Ben-Shakhar, G., & Crombez, G. (2017). Lying takes time: A meta-analysis on reaction time measures of deception. *Psychological Bulletin*, *143*, 428–453. <http://dx.doi.org/10.1037/bul0000087>
- Takarangi, M. K., Strange, D., & Houghton, E. (2015). Event familiarity influences memory detection using the aIAT. *Memory*, *23*, 453–461. <http://dx.doi.org/10.1080/09658211.2014.902467>
- Takarangi, M. K., Strange, D., Shortland, A. E., & James, H. E. (2013). Source confusion influences the effectiveness of the autobiographical IAT. *Psychonomic Bulletin & Review*, *20*, 1232–1238. <http://dx.doi.org/10.3758/s13423-013-0430-3>
- van Hooff, J. C., Brunia, C. H., & Allen, J. J. (1996). Event-related potentials as indirect measures of recognition memory. *International Journal of Psychophysiology*, *21*, 15–31. [http://dx.doi.org/10.1016/0167-8760\(95\)00043-7](http://dx.doi.org/10.1016/0167-8760(95)00043-7)
- Vargo, E. J., Petróczi, A., Shah, I., & Naughton, D. P. (2014). It is not just memory: Propositional thinking influences performance on the autobiographical implicit association test. *Drug and Alcohol Dependence*, *145*, 150–155. <http://dx.doi.org/10.1016/j.drugalcdep.2014.10.008>
- Verschuere, B., Ben-Shakhar, G., & Meijer, E. (Eds.). (2011). *Memory detection: Theory and application of the concealed information test*. Cambridge, UK: Cambridge University Press. <http://dx.doi.org/10.1017/CBO9780511975196>
- Verschuere, B., & De Houwer, J. (2011). Detecting concealed information in less than a second: Response latency-based measures. In B. Verschuere, G. Ben-Shakhar, & E. Meijer (Eds.), *Memory detection: Theory and application of the concealed information test* (pp. 46–62). Cambridge, UK: Cambridge University Press. <http://dx.doi.org/10.1017/CBO9780511975196.004>
- Verschuere, B., Prati, V., & De Houwer, J. (2009). Cheating the lie detector: Faking in the autobiographical Implicit Association Test. *Psychological Science*, *20*, 410–413. <http://dx.doi.org/10.1111/j.1467-9280.2009.02308.x>
- Wagenmakers, E.-J., Wetzels, R., Borsboom, D., & van der Maas, H. L. (2011). Why psychologists must change the way they analyze their data: The case of psi: Comment on Bem (2011). *Journal of Personality and Social Psychology*, *100*, 426–432. <http://dx.doi.org/10.1037/a0022790>

Received February 27, 2019

Revision received August 22, 2019

Accepted August 22, 2019 ■