# SOCIAL REPUTATION AS ONE OF THE KEY DRIVER OF AI OVER-RELIANCE: AN EXPERIMENTAL TEST WITH CHATGPT-3.5

MATHIEU CHEVRIER

# Social Reputation as one of the Key Driver of AI Over-Reliance: An Experimental Test with ChatGPT-3.5

Mathieu Chevrier[*]

**GREDEG WP No. 2025-12**

**Abstract**

Understanding an agent's true competencies is crucial for a principal, particularly when delegating tasks. A principal may assign a task to an AI system, which is often perceived as highly competent, even in domains where its actual capabilities are limited. This experimental study demonstrates that participants mistakenly bet on ChatGPT-3.5's ability to solve mathematical tasks, even when explicitly informed that it only processes textual data. This overestimation leads participants to earn 67.2% less compared to those who rely on the competencies of another human. Overconfidence in ChatGPT-3.5 persists irrespective of task difficulty, time spent using ChatGPT-3.5, nor prior experience posing mathematical or counting questions to it mitigates this bias. I highlight that overconfidence in ChatGPT-3.5 is driven by the algorithm's social reputation. The more participants perceive ChatGPT-3.5 as socially trusted, the more they tend to rely on it.

**JEL Codes:** C92; D91.

**Keywords:** ChatGPT-3.5, Overconfidence, Competence, Social Reputation, Over-reliance, Laboratory experiment

# 1 Introduction

It is critical for a principal to correctly estimate what an agent is capable of doing - you wouldn't ask a trader to perform surgery or a musician to write a mathematical theorem. So why delegate to an AI agent a specific task for which it is not competent?

Research shows that when we feel close to someone, we can better understand their thoughts and feelings (Epley et al., 2004). In contrast, when someone seems socially distant, we tend to view them in broad, abstract terms (Trope and Liberman, 2010). Thus, I first assume that assessing the true competence of a human—whom we can understand and empathize with—is generally easier than evaluating that of a "black-box" algorithm. Furthermore, studies by Lee and See (2004) and Dietvorst et al. (2015) indicate that people often hold algorithms to extremely high standards, expecting them to deliver consistent performance. I secondly assume that principals may perceive algorithms competence as unaffected by task difficulty. However, algorithms optimized for efficiency, like ChatGPT-3.5, may struggle with complex tasks by prioritizing fluency over precision. This discrepancy can distort participants' perception of the algorithm's true capabilities.

To understand how individual factors shape confidence estimation, I third argue that personal experience and social reputation are key. Frequent use of an algorithm helps individuals recognize its limitations, as identifying its errors improves capability assessment. Conversely, the algorithm's social reputation, shaped by others' usage and praise, may foster overconfidence in its abilities.

In this experimental study, I used ChatGPT-3.5, the most well-known algorithm in January 2024, to examine participants' ability to evaluate its competence. Participants estimated whether ChatGPT-3.5 or another person could better count the number of 1s in 100-digit binary sequences, knowing ChatGPT-3.5 processes only textual data.

This study contributes to the literature on over-reliance on algorithms (Klingbeil et al., 2024; Chevrier et al., 2024). Participants consistently overestimate ChatGPT-3.5's abilities, perceiving it as unaffected by task difficulty. The social reputation of ChatGPT-3.5 is a key factor in explaining why individuals tend to over-rely on it.

## 2    Design

This study employs a between-subjects design with a single treatment variation. The primary objective is to examine how the principal (participants) estimate the true competence of an agent (thereafter called worker) in a counting task. In the baseline condition, the worker is another participant, who is randomly selected at the beginning of the experimental session.[1] In the treatment condition, the worker is ChatGPT-3.5.

### 2.1    The Counting Task

The worker (participant or ChatGPT-3.5) counts the number of '1' digits in a 100-digit binary sequence. The task is repeated three times, each with a randomly assigned difficulty: Simple (10 ones), Intermediate (25 ones), and Complex (50 ones).[2]

### 2.2    Estimation Task

Before the counting task begins, participants receive a short description of the worker on their screen. The human worker is introduced as follows: *"Another participant, randomly chosen at the start of the experiment,has 20 seconds to count the sequences of numbers. For each correctly counted sequence, this participant will earn €1, up to a total of €3."*

In the ChatGPT-3.5 condition, participants read: *"I am ChatGPT-3.5, a language model developed by OpenAI. I am a computer program designed to understand and generate text in response to human queries. My ability to understand and generate text arises from advanced machine learning algorithms trained on large amounts of textual data. I do not possess consciousness or real understanding; I operate by analyzing statistical patterns in the data I was trained on (up to January 2022)."*

After viewing the description of the worker, participants are then asked to estimate the worker's true competence. To elicit this, I follow the steps shown in Figure 1. These steps are repeated for each sequence. In **Step 1** a sequence of 100 digits is randomly displayed on the participant's screen. In **Step 2** participants estimate whether the worker can accurately count the number of '1' digits in the presented sequence (Yes or No). Then, in **Step 3** participants state their confidence in the worker's accuracy on a scale from 50% to 100%. In **Step 4** social reputation is elicited by asking participants to estimate

---

[1]Each session consisted of 9 to 17 participants.

[2]Task complexity increases with Shannon entropy.

the percentage of other participants who also believe the worker can count the sequence accurately. In **Step 5** participants rate the complexity of the presented sequence, from 1 (very simple) to 7 (very complex).

## 2.3 Confidence and Social Reputation Elicitation

To incentivize the true reporting of both participants' confidence levels and social reputation of the worker's competence, I implement two distinct payment mechanisms. The Matching Probability (MP) mechanism (Massoni et al., 2014), is applied as follows (see Figure 1, **Step 3**). Let $p$ represent the participant's stated confidence (in percent) concerning the worker's ability to count the '1' digits. Two lotteries, $L_1$ and $L_2$, are then drawn uniformly from the integers 1 to 100. If $L_1 \leq p$, the participant's earnings hinge on whether their estimation in **Step 2** aligns with the worker's actual performance. Specifically, they receive €3 if they answered "Yes" ("No") and the worker was correct (incorrect), and €0 otherwise. If $L_1 > p$, the participant's payoff is determined by whether $L_1$ exceeds $L_2$. They earn €3 if $L_1 > L_2$ and €0 otherwise.

A variant of the Bayesian Truth Serum (Prelec, 2004) is used to elicit social reputation (see **Step 4**). Participants are asked to estimate the exact proportion of all other participants across every experimental session of the study who answered "Yes" in **Step 2** (i.e., who believed that the worker can count correctly). Participants who correctly guessed the exact percentage received €30 upon completion of all experimental sessions.
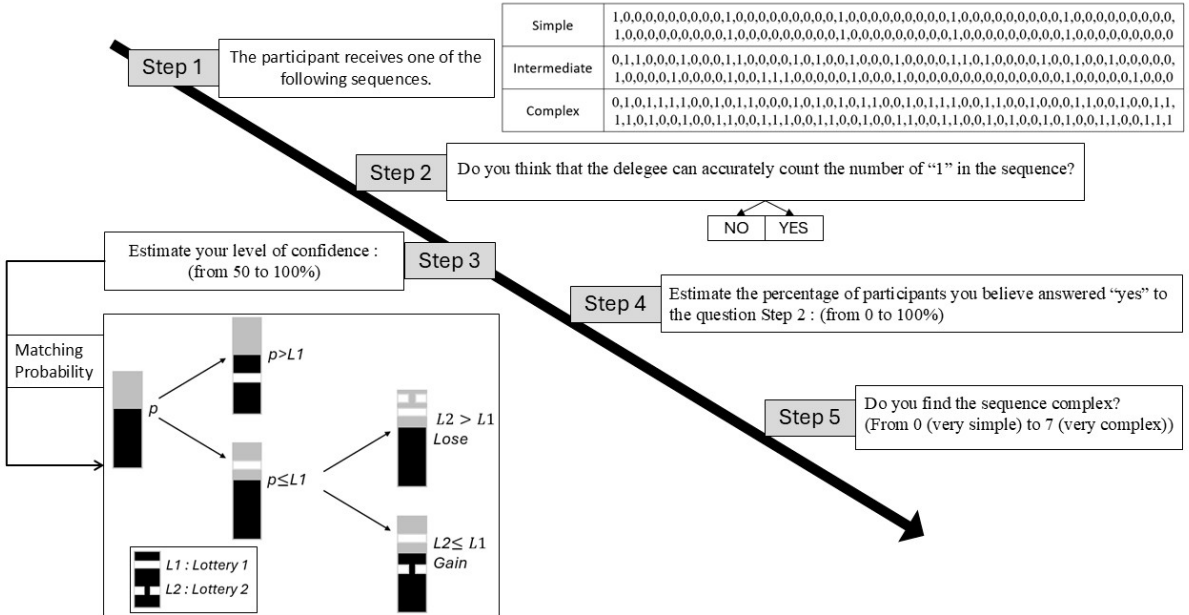
## 2.4 Experimental Protocol

I conducted experiments at the Laboratory of Experimental Economics of Nice in January 2024 [3], following an earlier short 30-minute experiment. At the end of the first experiment, participants were invited to join an optional, paid follow-up experiment.[4] I collected 100 observations for the ChatGPT-3.5 treatment and 79 for the Human treatment. In total, 40.2% of participants were male, with ages ranging from 18 to 25 years. Participants had already answered personal questions and completed the cognitive reflection test (Primi et al., 2016) in the earlier experiment. At the start of this experiment, participants were asked questions to assess their personal experience with ChatGPT-3.5

---

[3]The ethics committee of Côte d'Azur University has authorized the implementation of this experiment.

[4]8 participants declined due to time constraints.

Figure 1: Decision Making Timeline.

and perceived competence of the algorithm. After reading the instructions[5], participants completed a trial in which they received a hypothetical payoff to familiarize themselves with the MP mechanism before starting the main experiment.

At the end of the experiment, participants were invited one by one into a separate room to receive their payment. For each participant, one of the three sequences was randomly selected. According to the matching probability mechanism, participants either received their payoff through the lotteries (€3 or €0) or were informed about whether the worker had correctly counted the selected sequence. In the Human condition, participants learned immediately whether their predictions were correct. Participants did not observe an other participant performing the counting task, as that was conducted in parallel while they completed the experiment. In the ChatGPT-3.5 condition, the experimenter, in the participant's presence, asked ChatGPT-3.5 to count the randomly selected sequence. Both the experimenter and the participant could verify ChatGPT-3.5's result on the spot, and the participant's payment was determined by the correctness of their prediction about ChatGPT-3.5.[6]

---

[5]Complete instructions are available here: https://osf.io/rhsja/files/osfstorage/67d98f329cc590366a053aa2

[6]Across over 100 participants, the selected sequence was only counted correctly twice by ChatGPT-3.5. To maintain consistency and prevent learning across participants, all prompts given to ChatGPT-3.5 were cleared before each new participant arrived.

For the payment of the social reputation elicitation, at the end of all experimental sessions, one sequence from each treatment was randomly selected. One week later, participants whose estimates matched the true value received €30.

On average, the experiment lasts 10 minutes. All participants received a €5 show-up fee. Additionally, on average, participants earned €1.54 based on the Matching Probability rule, and two participants received a €30 bonus for social reputation elicitation.

## 3    Results

### 3.1    Dependent Variables

The main metric in the data analysis is the "calibration" benchmark (Lichtenstein et al., 1977). It is defined as $\frac{1}{3n}\sum_{i=1}^{n}\sum_{j=1}^{3}\big|\text{Confidence Level}_{i,j} - \text{Performance}_{i,j}\big|$, where 'i' indexes the participants and 'j' indexes the trials. The 'Confidence Level' is the participant's percentage estimate of the worker's ability to complete the task (capture in Figure 1, **Step 2**).[7] The 'Performance' is the worker's success in the counting task for the given trial. Overconfidence arises when average confidence surpasses observed performance, and underconfidence occurs otherwise. Smaller gaps indicate better calibration. Depending on their estimation, participants earn either €3 or €0, with higher expected payoffs when stated confidence closely matches the worker's actual performance. Calibration is also computed for each sequence to assess how task difficulty affects participants' responses (see Figure 2). Social reputation is measured in percentage terms by the question depicted in Figure 1 (**Step 4**). Finally, I run a panel regression with random effects on "calibration per sequence" and a logit panel regression on participants' "success", where success equals 1 if a participant's estimation is correct (see Table 1).

### 3.2    Data Analysis

On average, participants earned 67.2% less than those evaluating another participant (ChatGPT-3.5 = €0.81, Human = €2.47, Rank-Sum Test, p < 0.0001). Participants' show near-perfect calibration in assessing another human's competence (+2.4%), but significantly higher overconfidence when evaluating ChatGPT-3.5 (+81.3%) (Rank-Sum Test, p < 0.0001). Regressions (1) and (2) in Table 1 further confirm this overconfidence
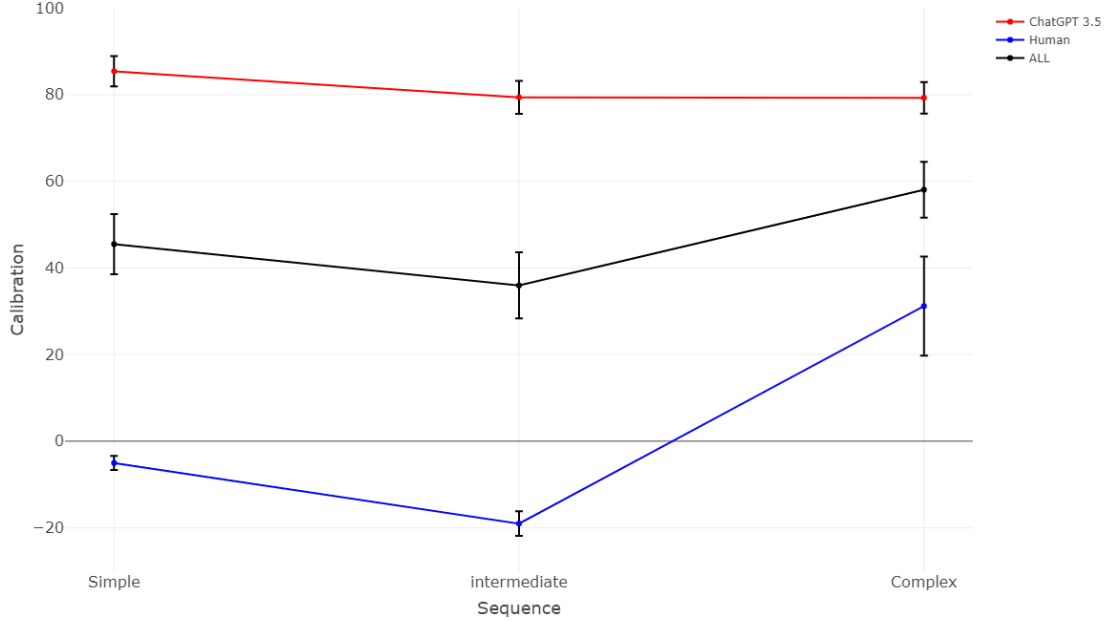
---

[7] If a participant states that the worker cannot complete the task, the confidence level is adjusted by reversing the score: |100% − Confidence Level|.

(coefficient: 24.69, p = 0.013, IC [5.112, 44.286]).

**Result 1:** *Participants show significant overconfidence in ChatGPT-3.5's capabilities, resulting in earnings that are 67.2% lower than those of participants interacting with another human.*

Figure 2: calibration per sequence across worker with confidence intervals (p = 0.95).



Calibration is computed for each sequence. A calibration equal to 0 means that the participant's confidence level is well calibrated. Below (higher than) 0 means that the participant is underconfident (overconfident).

The complexity control variable in **Step 5** (Figure 1) confirms that participants accurately perceived differences in sequence difficulty (simple: 2.0, intermediate: 3.8, complex: 4.8; Rank-Sum Tests, p < 0.0001). Figure 2 shows that calibration level per sequence vary with task difficulty. Regardless of difficulty, participants consistently display high overconfidence in ChatGPT-3.5's abilities compared to other participants. Notably, overconfidence is greater for simple tasks (85.4%) than for intermediate (79.3%) and complex (79.2%) tasks (Rank-Sum Tests, p < 0.01). In contrast, participants assessing another human exhibit better calibration across all levels, with more accurate estimations for simpler tasks.

**Result 2:** *Task difficulty minimally influences participants' calibration regarding ChatGPT-3.5.*

## 3.3 Which factors drive calibration ?

Using participant individual factors related to ChatGPT-3.5 usage[8] and social reputation[9], I examine calibration quality. Among 100 participants, 90% are students aged 18–25, with 71% having previously used ChatGPT-3.5 and 64% using it at least one hour per week. Of these, 39% have used it for calculations and 29% for counting tasks. Table 1 shows no significant relationship between hours, computation, and count, suggesting prior interactions with ChatGPT-3.5 do not help participants assess its competence. Perceived competence of ChatGPT-3.5 also fails to predict true AI competence, while perceived ChatGPT-3.5 computation ability, the control variable, is significant (coefficient: 2.623, $p < 0.029$, CI [0.262, 4.984]). Thus, participants' perception of ChatGPT-3.5 computation is key to predicting calibration quality. However, social reputation as well significantly influences participants' acceptance of ChatGPT-3.5.

Previous studies have shown that a positive reputation increases individuals' confidence in other participants (Bohnet and Huck, 2004). In addition, online marketplace feedback signals social reputation and bolsters users' confidence in these platforms (Bolton et al., 2004). In this study, the higher ChatGPT-3.5's perceived social reputation, the more participants' confidence increases, making them more prone to overconfidence. In regressions (1) and (2), reputation increases overestimation of ChatGPT-3.5's capabilities compared to a human (coefficient: 0.693, $p < 0.0001$, CI [0.454, 0.931]). In regressions (3) and (4), reputation explains overconfidence (coefficient: 0.303, $p < 0.0001$, CI [0.186, 0.419]). Regressions (5) and (6) further reveal that stronger social reputation effects correlate with higher error rates (coefficient: -0.194, $p < 0.0001$, CI [-0.278, -0.110]).

**Result 3:** *Participants' overconfidence toward ChatGPT-3.5 is driving by the social reputation of the AI.*

## 4 Discussion

Over-reliance on AI is a critical societal issue. For example, 42% of IT firms use AI systems that are often inaccurate for recruitment (Lytton, 2024). Experts suggest that over 44.7% of autopilot crashes could be avoided if drivers followed the system's rec-

---

[8]"If you have already used ChatGPT-3.5, could you estimate the number of hours you use it per week? (If not, enter 0)", "Have you ever asked ChatGPT-3.5 to perform calculations for you?", "Have you ever asked ChatGPT-3.5 to count something for you?".

[9]"Across the tasks performed by ChatGPT-3.5, do you perceive it as competent overall?", "Do you perceive ChatGPT-3.5 as competent in performing computations?"

Table 1: Panel regression with random effect for "calibration" and logit panel regression with random effect for "success" across all treatments.

| Dependent Variable | calibration | | | | success | |
|---|---|---|---|---|---|---|
| Data | ALL | ALL | ChatGPT-3.5 | ChatGPT-3.5 | ChatGPT-3.5 | ChatGPT-3.5 |
| | (1) | (2) | (3) | (4) | (5) | (6) |
| Difficulty | 6.387**** | 6.387**** | -1.235** | -1.232** | 0.112 | 0.066 |
| | (1.651) | (1.668) | (0.533) | (0.536) | (0.486) | (0.494) |
| ChatGPT-3.5 | 24.700** | 24.712** | | | | |
| | (9.993) | (10.060) | | | | |
| Reputation | -0.186* | -0.187* | 0.294**** | 0.294**** | -0.178**** | -0.185**** |
| | (0.101) | (0.099) | (0.060) | (0.060) | (0.045) | (0.047) |
| ChatGPT-3.5 × Reputation | 0.693**** | 0.696**** | | | | |
| | (0.122) | (0.122) | | | | |
| Hours | | | -0.225 | -0.450 | 0.103 | 0.075 |
| | | | (0.525) | (0.687) | (0.216) | (0.226) |
| Computation | | | -3.596 | -1.036 | -0.520 | -1.439 |
| | | | (4.371) | (3.128) | (1.210) | (1.305) |
| Count | | | 3.245 | 2.691 | 0.775 | 1.738 |
| | | | (6.494) | (4.510) | (1.617) | (1.375) |
| Competent | | | -0.574 | -0.170 | -0.250 | -0.407 |
| | | | (1.673) | (1.471) | (0.526) | (0.535) |
| Competent in Computation | | | 2.623** | 2.192** | -1.197** | -1.219*** |
| | | | (1.204) | (1.031) | (0.468) | (0.453) |
| Constant | 3.684 | -4.680 | 53.489**** | 46.413**** | 14.137*** | 16.243*** |
| | (9.837) | (14.364) | (6.942) | (12.624) | (4.374) | (5.740) |
| Control | NO | YES | NO | YES | NO | YES |
| Observations | 537 | 537 | 300 | 300 | 300 | 300 |
| Number of obs | 179 | 179 | 100 | 100 | 100 | 100 |
| Wald Chi2 | 1538 | 1927 | 46.31 | 52.99 | 20.81 | 22.02 |
| R2 (Overall) | 0.6907 | 0.6987 | 0.2239 | 0.2540 | | |
| p-value | 0 | 0 | 7.87e-09 | 4.13e-07 | 0.0980 | 0 |

*Note:* Robust standard errors at the individual level are reported in parentheses. "Calibration" is the confidence level minus the actual success rate for each decision. "Success" is a binary variable equal to 1 if the estimation is correct, otherwise 0. "Difficulty" equals 1 (2) [3] if the sequence is simple (intermediate) [complex]. "ChatGPT-3.5" is equal to 1 in ChatGPT-3.5 treatment. "Reputation" ranges from 0 to 100. "Hours" refers to the number of hours participants use ChatGPT-3.5 per week. "Computation" and "Count" are binary variables equal to 1 if the participant has previously asked ChatGPT-3.5 to compute or count something, respectively. Perceived competence is assessed using "competent" and "competent in counting" on a Likert scale ranging from 1 (not competent) to 7 (very competent). Individual control variables include the order of the randomized digit sequence. The CRT, which ranges from 0 to 6. "Man" is a binary variable equal to 0 for women and 1 for men. "Years of study" range from 0 to 8. The "area of study", "Risk" (Dohmen et al., 2011) and "control"(Wolff et al., 2022) are considered. None of the individual control variables are significant. **** $p<0.001$, *** $p<0.01$, ** $p<0.05$, * $p<0.1$

ommendations, such as not sleeping and keeping a hand on the wheel (Duncan, 2024). In our laboratory experiment, regardless of task difficulty, participants who exhibited overconfidence in ChatGPT-3.5 earned 67.20% less than those who delegated the task to a human worker. This overconfidence is driven by the ChatGPT-3.5's social reputation rather than by individual knowledge of the algorithm.[10] Our findings indicate that social reputation plays a key role in AI adoption and can lead to dangerous over-reliance.

**Declarations** Conflicts of interest: The author declares that he has no conflict of interest.

---

[10]Fu and Hanaki (2024) also show that reliance on ChatGPT is not explained by individual knowledge.

# Bibliography

Bohnet, I. and Huck, S. (2004). Repetition and reputation: Implications for trust and trustworthiness when institutions change. American economic review, 94(2):362–366.

Bolton, G., Katok, E., and Ockenfels, A. (2004). How effective are online reputation mechanisms? an experimental study. Management Science, 50(11):1587–1602.

Chevrier, M., Corgnet, B., Guerci, E., and Rosaz, J. (2024). Algorithm credulity: Human and algorithmic advice in prediction experiments. Available at SSRN 4828701.

Dietvorst, B. J., Simmons, J. P., and Massey, C. (2015). Algorithm aversion: people erroneously avoid algorithms after seeing them err. Journal of experimental psychology: General, 144(1):114.

Dohmen, T., Falk, A., Huffman, D., Sunde, U., Schupp, J., and Wagner, G. G. (2011). Individual risk attitudes: Measurement, determinants, and behavioral consequences. Journal of the european economic association, 9(3):522–550.

Duncan, I. (2024). How u.s. safety regulators have struggled to get a grip on tesla's autopilot. The Washington Post.

Epley, N., Keysar, B., Van Boven, L., and Gilovich, T. (2004). Perspective taking as egocentric anchoring and adjustment. Journal of personality and social psychology, 87(3):327.

Fu, Y. and Hanaki, N. (2024). Do people rely on chatgpt more than their peers to detect fake news? Technical report, ISER Discussion Paper.

Klingbeil, A., Grützner, C., and Schreck, P. (2024). Trust and reliance on ai—an experimental study on the extent and costs of overreliance on ai. Computers in Human Behavior, 160:108352.

Lee, J. D. and See, K. A. (2004). Trust in automation: Designing for appropriate reliance. Human factors, 46(1):50–80.

Lichtenstein, S., Fischhoff, B., and Phillips, L. D. (1977). Calibration of probabilities: The state of the art. In Decision Making and Change in Human Affairs: Proceedings of the Fifth Research Conference on Subjective Probability, Utility, and Decision Making, Darmstadt, 1–4 September, 1975, pages 275–324. Springer.

Lytton, C. (2024). Ai hiring tools may be filtering out the best job applicants. BBC.

Massoni, S., Gajdos, T., and Vergnaud, J.-C. (2014). Confidence measurement in the light of signal detection theory. Frontiers in psychology, 5:1455.

Prelec, D. (2004). A bayesian truth serum for subjective data. science, 306(5695):462–466.

Primi, C., Morsanyi, K., Chiesi, F., Donati, M. A., and Hamilton, J. (2016). The development and testing of a new version of the cognitive reflection test applying item response theory (irt). Journal of Behavioral Decision Making, 29(5):453–469.

Trope, Y. and Liberman, N. (2010). Construal-level theory of psychological distance. Psychological review, 117(2):440.

Wolff, W., Bieleke, M., Englert, C., Bertrams, A., Schüler, J., and Martarelli, C. S. (2022). A single item measure of self-control–validation and location in a nomological network of self-control, boredom, and if-then planning. Social Psychological Bulletin, 17:1–22.

# DOCUMENTS DE TRAVAIL GREDEG PARUS EN 2025
## *GREDEG Working Papers Released in 2025*