

The ABC’s of Who Benefits from Working with AI: Ability, Beliefs, and Calibration

Online Appendix

Andrew Caplin, David Deming, Shangwen Li, Daniel Martin,
Philip Marx, Ben Weidmann, and Kadachi Jiada Ye

A Robustness Analyses

A.1 Demographic Balance and Controls

Table 3 presents summary statistics by participant treatment group. This includes additional demographic information optionally reported by participants to the experimental platform. The table shows that there was balance across control and treated participant characteristics, with marginally significant differences in birthplace and language when we do not account for multiple hypothesis testing. An issue with optional demographic information – which for clarity is omitted from Table 3 – is that this data is incomplete. In particular, 41% of participants do not report education or employment status, though we find no evidence of imbalanced non-reporting across treatment assignment. Table 4 replicates our main Table 2 when we additionally control for the remaining, mostly non-missing demographic indicators (age, gender, ethnicity, birthplace, language) at baseline and interacted with treatment. This confirms that the results of Table 2 are essentially unchanged if we additionally control for participant demographics.

A.2 Calibration and Confidence

In this appendix we consider alternative specifications and implications of net confidence instead of calibration as a covariate. One possibility is to re-estimate our main specifications of Table 2, but distinguishing between under-confident and over-confident miscalibration. In such a specification, we fail to reject ($p = 0.34$) that the coefficients on

Table 3: Subject Summary Statistics

	Control	Treatment	P-value
Total Time (Minutes)	34.14 (13.42)	34.99 (14.13)	0.44
Bonus (\$)	4.27 (2.16)	4.29 (2.08)	0.90
IQ (0-14)	5.03 (2.80)	4.86 (2.79)	0.42
Age	40.74 (14.39)	40.45 (13.56)	0.79
Female	0.56 (0.50)	0.58 (0.49)	0.69
Nonwhite	0.39 (0.49)	0.35 (0.48)	0.35
Born Outside the U.S.	0.10 (0.30)	0.06 (0.24)	0.09
First Language Not English	0.06 (0.24)	0.03 (0.16)	0.02
Student	0.21 (0.41)	0.16 (0.37)	0.19
Employed Full-Time	0.55 (0.50)	0.48 (0.50)	0.13
Sample Size	239	493	

Note: The values in parentheses are the sample standard deviations. P-values measure the probability of finding a mean difference as large as the one observed, under the null hypothesis of no difference across groups.

under-confident and over-confident miscalibration are equal. However, such a regression is empirically under-powered for under-confidence given the small number of under-confident participants in the experiment, and so the overall results are not reported.

Instead, we explore the mechanism that under-confident individuals may commit too many Type I errors (following the AI when it is incorrect), while over-confident individuals may commit too many Type II errors (failing to follow the AI when it is correct). As also explained in the main text, in Table 5 we repeat our main specification with net confidence instead of calibration and disaggregating by whether the AI prediction was correct. To disaggregate by image attribute (AI correctness), this regression is at the participant-response level with standard errors clustered by participant. (Note that

Table 4: Regression: Skills and Working with AI, with Demographic Controls (Omitted)

Skills (Control Block)	Outcome:			
	Accuracy (Treated Block)			
	(1)	(2)	(3)	(4)
Treatment	5.01 (1.98)	6.06 (1.79)	5.62 (1.78)	6.24 (1.83)
Accuracy		3.90 (0.41)	3.80 (0.49)	3.44 (0.49)
Accuracy \times Treatment		-1.51 (0.53)	-2.08 (0.61)	-1.77 (0.61)
Calibration			0.21 (0.47)	0.09 (0.46)
Calibration \times Treatment			1.35 (0.56)	1.42 (0.57)
IQ				1.29 (0.45)
IQ \times Treatment				-1.01 (0.54)
Constant	69.09 (1.67)	66.64 (1.49)	66.69 (1.49)	65.90 (1.54)
Observations	702	702	702	702

Note: Observations are at the subject level. Robust standard errors are in parentheses. Skills are standardized within session.

previous regressions at the participant level with robust standard errors are equivalent to a pooled regression at the participant-response level with standard errors clustered by participant.) Consistent with the preceding Type I/II error explanation, we find that the less confident are relatively more disadvantaged for images where the AI is ex post incorrect, with a standardized treatment effect across net confidence of 1.93 (s.e.=1.34), whereas the more confident are more disadvantaged for images where the AI is ex post correct, with a standardized treatment effect across net confidence of -2.52 (s.e.=0.75). We note that the results where the AI is incorrect are also under-powered, given that the AI is usually correct. Also, a comparison of average baseline performance of 53% and 70% accuracy when the AI is respectively incorrect and correct reveals correlation in image difficulty across human participants and the AI.

Table 5: Regression: Split by Whether AI Correct

Skills (Control Block)	Outcome: Subject Response Correct (Treated Block)	
	(1)	(2)
	AI Incorrect	AI Correct
Treatment	-17.43 (1.07)	13.57 (0.61)
Accuracy	3.21 (0.94)	4.29 (0.56)
Accuracy \times Treatment	-1.63 (1.30)	-2.63 (0.82)
Net Confidence	-0.06 (0.91)	0.01 (0.50)
Net Confidence \times Treatment	1.93 (1.34)	-2.52 (0.75)
Constant	53.14 (0.78)	69.97 (0.44)
Observations	13,877	44,536

Note: Observations are at the subject-response level. Subject-clustered standard errors are in parentheses. Skills are standardized within session.

A.3 Measurement Error

To address the possibility of measurement error affecting our results, we repeat the preceding OLS analysis of Table 2 with an adaptation of the Obviously Related Instrumental Variables (ORIV) methodology of Gillen, Snowberg and Yariv (2019). Their basic idea is to instrument one imprecise measure with another imprecise measure of the same underlying trait. We adapt this to our experimental setting by randomly drawing a split of the 80 images used for skill estimation into an instrumented and instrumenting set, leading to instrumented and instrumenting estimates of each skill for each individual. To remove the noise introduced by a single split realization, we repeat this splitting procedure 10,000 times. We then repeat our estimation of regression 3 across the 10,000 splits while instrumenting our skill estimates. To correct the standard errors from artificially inflating sample size, we follow Gillen, Snowberg and Yariv (2019) and cluster standard errors at the participant level.

Table 6: Regression: ORIV

Skills (Control Block)	Outcome: Accuracy (Treated Block)			
	(1)	(2)	(3)	(4)
Treatment	6.93 (0.56)	5.45 (0.62)	5.34 (0.62)	5.11 (0.68)
Accuracy		11.24 (1.65)	11.21 (1.80)	12.17 (2.37)
Accuracy \times Treatment		-3.86 (1.98)	-4.82 (2.11)	-5.02 (2.73)
Calibration			0.05 (0.72)	0.13 (0.77)
Calibration \times Treatment			1.84 (0.84)	2.02 (0.89)
IQ				-1.05 (0.83)
IQ \times Treatment				0.01 (0.94)
Constant	65.46 (0.48)	66.58 (0.53)	66.58 (0.54)	66.74 (0.60)
Observations	7,320,000	7,320,000	7,320,000	7,320,000

Note: Observations are at the subject-simulation level. Subject-clustered standard errors are in parentheses. Skills are standardized within session and simulated split.

Table 6 repeats the analysis of Table 2 adopting the ORIV estimation approach for ability and calibration; we continue to include IQ as an exogenous control. Consistent with the simple univariate attenuation bias intuition, we find *larger* heterogeneous effects across ability and calibration when accounting for measurement error, which (though less precise) remain significant at conventional levels. For example, in Column (3) including only accuracy and calibration, we find standardized treatment effects across ability and calibration of -4.82 (s.e.=2.11) and 1.84 (s.e.=0.84) respectively, compared to -2.16 (s.e.=0.57) and 1.35 (s.e.=0.54) in the OLS results of Table 2. Interestingly, accounting for measurement error leads to negligible variation in treatment effects across IQ, suggesting that IQ may proxy for imprecisely measured accuracy in the OLS specification. In conclusion, however, accounting for measurement error only strengthens the magnitude of our main results regarding the relationship between calibration and the benefits

of working with AI.

A.4 Alternative Skill Measures

Table 7: Regression: Alternative Skill Measures

Skills (Control Block)	Outcome: Accuracy (Treated Block)			
	(1)	(2)	(3)	(4)
Treatment	8.22 (0.63)	7.34 (0.53)	7.28 (0.53)	7.40 (0.53)
Accuracy		4.54 (0.40)	4.64 (0.61)	4.36 (0.62)
Accuracy \times Treatment		-1.34 (0.52)	-2.95 (0.73)	-2.71 (0.73)
Calibration			-0.08 (0.57)	-0.26 (0.59)
Calibration \times Treatment			2.53 (0.72)	2.66 (0.74)
IQ				1.14 (0.50)
IQ \times Treatment				-0.90 (0.58)
Constant	70.96 (0.53)	71.62 (0.43)	71.61 (0.44)	71.49 (0.44)
Observations	732	732	729	729

Note: Observations are at the subject level. Robust standard errors are in parentheses. Skills are standardized within session.

As a final robustness check, we repeat the main analysis of Table 2 using alternative measures of ability and calibration. As an alternative measure of ability, we consider the area under the receiver operator characteristic curve (AUC) (Bamber, 1975; Hanley and McNeil, 1982), a common summary measure of diagnostic ability that extends the Blackwell order (Chan, Gentzkow and Yu, 2022). Intuitively, AUC is interpretable as the probability that a random positive instance has a higher report than a random negative instance. AUC is a useful complement to accuracy because it measures how well predictions are ranked, and not just the performance of a binary classifier induced by the symmetric classification threshold 0.5. On the other hand, changes in accuracy are more

readily interpretable.

For our alternative measure of calibration, we use the coefficient on reported probabilities in a logistic regression of true states S_j on reported probabilities B_{ij} across the images used to estimate skills for each individual i . This corresponds to a one-parameter version of the standard model of [Grether \(1980\)](#) that allows for arbitrary degrees of participant over- and under-confidence. To facilitate interpretation, we take the negative absolute value of the natural log of the aforementioned coefficient as our measure of calibration. The log transformation maintains our earlier interpretation of negative values as under-confident, zero values as calibrated, and positive values as over-confident; further taking the negative absolute value is analogous to our main measure of calibration [\(2\)](#) and allows to interpret higher values with better calibration. Because the logistic function is not defined for reported beliefs of certainty 0, 1, such responses would drop out of estimation even though they suggest over-confidence; therefore we replace these responses with the nearest interior responses, respectively 0.01 and 0.99. Another complication is that our alternative calibration measure is undefined for participants with a negative coefficient in the logistic regression, corresponding to a negative association between reported probabilities and true states and thus indicative of severe inattention. We drop 3 such participants from evaluation.

Table [7](#) repeats the analysis of Table [2](#) adopting the alternative measures for ability and calibration; we continue to include IQ as an exogenous control. While the magnitude of coefficients is not directly comparable to Table [2](#) given the different outcome (AUC instead of accuracy), we can still compare the relative magnitude of standardized coefficients on the interactions of ability, calibration, and IQ with treatment. Doing so, we find a significant and larger effect of calibration. For example, in the full specification (Column 4 of Table [7](#)), we estimate a standardized treatment effect across calibration of 2.66 (s.e.=0.74), which is nearly equal in magnitude to the analogous ability effect of -2.71 (s.e.=0.73). Interestingly, as in the preceding ORIV specification, statistical significance of the IQ effect (coef=-0.90, s.e.=0.58) is not robust to the alternative specification, although the coefficient is similar in magnitude to the analogous estimate in Table 2

(coef=-1.01, s.e.=0.51). In conclusion, the alternative skill specification supports our main results regarding the importance of calibration.

B Experimental Details

B.1 Overview of Caffe Model

The guesses of our AI assistant are based on the “Convolutional Architecture for Fast Feature Embedding,” or Caffe, which is a machine learning model specifically developed for deep learning applications, particularly in image recognition tasks. It allows models to learn statistical patterns directly from large sets of labeled images.

The structure of a Caffe model consists of several key components. Images are initially provided as input data. These images then pass through convolutional layers, which systematically identify and extract visual features such as edges, textures, shapes, and other essential details from the input data. After feature extraction, pooling layers reduce the dimensionality of these features by simplifying and summarizing them, thus helping the model to efficiently handle vast amounts of information. Finally, fully-connected layers aggregate and interpret these condensed features to perform accurate predictions or classifications.

During the training phase, the model adjusts and estimates millions of parameters, known as weights, to minimize prediction errors through iterative learning. This extensive parameter adjustment process enables the model to improve its accuracy continuously. Once trained, the Caffe model can effectively classify new images or recognize visual patterns by leveraging its learned parameters, resulting in high precision and performance in image-based tasks.

B.2 Image Selection

From the images in the desired age ranges, we chose images taken between 2010 and 2014 to ensure they had a high enough resolution and so that participants would be familiar with the outfits of the people in the images. From the remaining images, we removed images that were unclear or contained multiple faces. To further refine the set of images, we conducted two pilots with 100 participants each. For the first pilot, we randomly split 200 images into two subsets, and approximately 50 participants were piloted for each subset. Based on this pilot, we removed images if they were too easy, too difficult, or

too certain. We defined the image as too easy if 95% of participants answered correctly; too difficult if fewer than 30% of participants answered correctly; and too certain if more than 30% of participants reported less than 1 % or above 99%. These exclusion criteria reduced the number of images to 124, so we conducted a second pilot. In the second pilot, we added 196 new images to create a set of 320 images, and increased the restrictiveness of our exclusion criteria. An image was deemed too easy if the average accuracy was above 90% while an image was deemed too difficult if the average accuracy was below 20%. Based on the second round, we selected 227 out of 320 images. Finally, we selected 160 of the 227 suitable images so that AI confidence scores would be well calibrated.

B.3 Image Splits

In each session, we fixed whether an image would appear in a control block or treated block. To increase robustness, we changed the way that images were selected to be in treated or control blocks between two sessions. For the first session, we used the following method to split images. First, all images were ordered lexicographically based on the true label and then the AI prediction of being over 21. Following this order, we alternately assigned each image to treated and control blocks. This produced a set of images containing 40 images of each label in both the treated and control blocks. For the second session, we used the same order based on true label and AI confidence score to form adjacent pairs of images. For each pair, we randomly assigned one image to the control block and the other image to the treated block. We then selected a split that achieved two objectives: (1) having a non-decreasing calibration curve for AI scores in the treated group and (2) having the accuracy and confidence levels of humans in the pilot sessions be well-balanced across treated and control images.

B.4 Screenshots

Consent

This is an economics research study. The purpose is to learn about how individuals make decisions. If you choose to participate, you will make a series of computer-based choices which should take **30-35 minutes**.

In addition to a fixed payment, you can earn a sizable bonus (based on your choices) within 7 days of submission. The study is entirely anonymous. Any data produced cannot be linked back to you. Data may be used to publish in an academic journal and may be shared with other researchers in the future.

The study is entirely voluntary. You may choose to end your participation at any time by closing your browser. If you end early, you will (1) forfeit your fixed payment; (2) forfeit your bonus; and (3) your data will be automatically deleted.

Economics studies have a strict standard of no deception. The study will operate exactly as described in the instructions.

If you have questions or concerns, you may contact the study's Protocol Director, Daniel Martin, at danielmartin@ucsb.edu. If you are concerned about your rights as a participant, you may contact the University of California, Santa Barbara (UCSB) Human Subjects Committee (HSC) at hsc@research.ucsb.edu or call at (805) 893-3807 or (805) 893-4290.

Authorization

- Yes, I would like to participate in this study, I am a resident of the United States, and I am 18 or older
- No, I cannot participate

To make sure you have read the consent page, the next button will be activated after 10 seconds.

Next

Figure 4: Consent form.

Before we begin the study, please answer the following two questions.

Please enter the correct answer to this math problem.

$100/2=$

- 20
- 30
- 50

If you are reading these instructions carefully, DO NOT select the correct answer. Instead, please select the LAST option.

$2+5+3=$

- 10
- 12
- 14

Next

Figure 5: Attention checks.

Welcome

This is a study about classifying images.

The whole session will last for **30 to 35 minutes**. The participation payment you will receive for completing the session is **\$5**, and there are two possible additional bonus payments of **\$5** and **\$1**. Your performance during the session will determine the likelihood you receive the bonus payments.

You will see some instructions over the next few pages, please read them carefully.

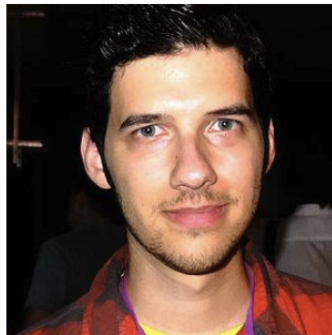
Next

Figure 6: Welcome page.

Instructions

In this session, you will be shown **images** of people one at a time.

For each image, you will be asked the **probability** (from 0% to 100%) that you think the person was over 21 years old when the image was taken. For example, if you were presented with the following image, you would be asked the probability you think the person was over 21 years old.



After seeing 160 images, you will be asked to complete a very brief survey and a short, unrelated additional task.

To make sure you have read the instructions, the next button will be activated after 10 seconds.

Next

Figure 7: Instructions for the experiment.

Instructions

Payment

The **participation payment is \$5**, which you will get if you **finish** the session.

The **main possible bonus payment is an additional \$5**. At the end of the session, one of the images will be selected at random to determine your bonus payment. The payment rule for the randomly selected image is designed so that you maximize your overall likelihood of getting the bonus payment when you **truthfully report your belief that the person in the image was over 21**.

Another possible bonus payment is \$1, which will be determined by your performance on the unrelated additional task at the end of the session.

To make sure you have read the instructions, the next button will be activated after 10 seconds.

Next

Figure 8: Details about payment.

Instructions

Helpful tips

1. All images were taken between 2010 and 2014.
2. Half (50%) of the 160 images are of people who were under 21 years old (that is, 20 years and younger) and the other half are people who were over 21 years old (that is, 22 years and older). None of the images are of people 21 years old.
3. For each image, you will have **60 seconds** to enter the probability you think the person in the image was over 21 years old.
4. If you **do not submit** your answer within the time limit, you will have a 0% chance of winning the bonus payment if the image is randomly selected for payment.
5. You will not be able to go back to previously seen images.

Before beginning, you have **4 practice** images, which will not impact your payment.

To make sure you have read the instructions, the next button will be activated after 10 seconds.

Next

Figure 9: Helpful tips and reminders.

What is the probability that the person in this image was **over** 21 years old?



Time left to complete this page: 52

Click anywhere on the bar to move around your choice.

Under 21


Over 21

Submit

Practice 1 out of 4

Figure 10: Practice round example (before clicking the slider bar).

What is the probability that the person in this image was **over** 21 years old?



Time left to complete this page: 35

Click anywhere on the bar to move around your choice.

Under 21 Over 21

Probability over 21 in image is
41%

Submit

Practice 1 out of 4

Figure 11: Practice round example (after clicking the slider bar).

Begin Experiment

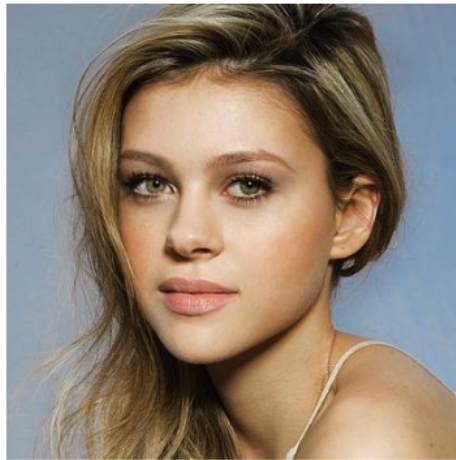
The practice images are now complete. You will now proceed with the experiment.

There are a total of **160 images**, and you can pause every **20 images**. You will not be able to go back to previous images during the session.

Next

Figure 12: Before incentivized rounds.

What is the probability that the person in this image was **over** 21 years old?



Time left to complete this page: 52

Click anywhere on the bar to move around your choice.

Under 21

Over 21

Submit

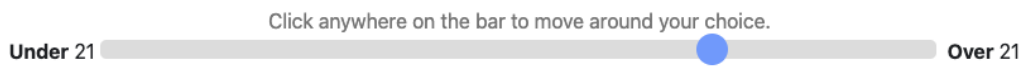
Round 1 out of 160

Figure 13: Incentivized round example (before clicking the slider bar).

What is the probability that the person in this image was **over** 21 years old?



Time left to complete this page: 44



Probability over 21 in image is

74%

Submit

Round 1 out of 160

Figure 14: Incentivized round example (after clicking the slider bar).

Artificial Intelligence (AI) Assistant

For the next 8 images, **you will see an AI Assistant guess.**

The AI Assistant provides a guess for the probability of the person in the image was over 21.

The AI Assistant's guesses are more accurate than the average human, but worse than the most skilled humans.

The AI Assistant is calibrated. For example, when the AI Assistant guesses 60%, the person in the image was over 21 about 60% of the time (and under 21 the other 40% of the time). [Click here to learn more about calibration.](#)

Calibration

To make sure you have read the instructions, the next button will be activated after 10 seconds.

Next

Figure 15: Introduction to AI Assistant.

Calibration



For each image, you will receive a probability from the AI Assistant that was trained using similar images. For example, on a particular image you might be told: "AI assistant's guess is **72%**."

The probability that an image was over 21 for a given probability is given by the figure below. For example, for a score between 50 and 60, approximately 53% of the people in the images were above 21.

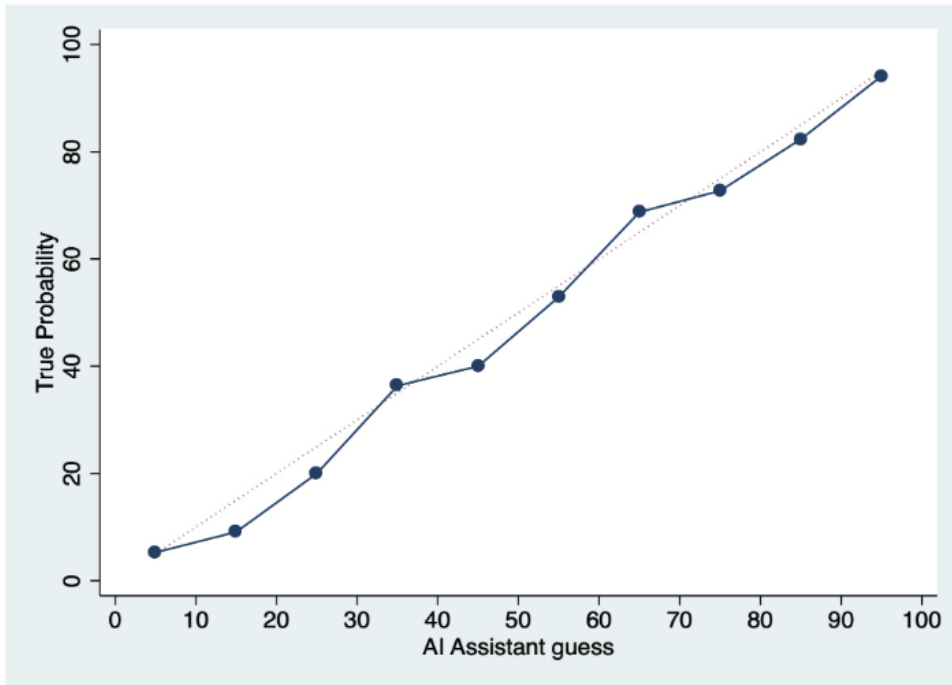


Figure 16: Explanation of calibration.

You have completed 20 out of 160 images

For the next set of 20 images, you will see an **AI Assistant guess**.

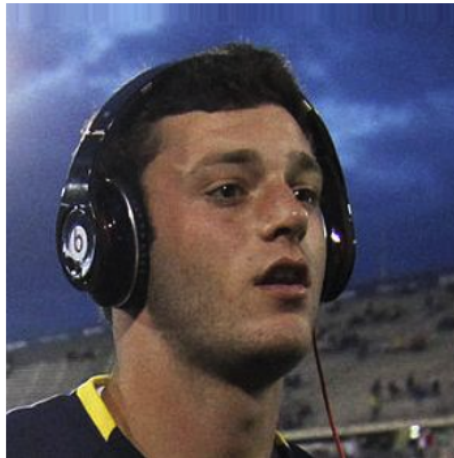
Recall that the **AI Assistant's guesses are more accurate than the average human, but worse than the most skilled human.**

Click next to continue.

Next

Figure 17: Page shown to treatment participants before treated blocks begin.

What is the probability that the person in this image was **over** 21 years old?



Time left to complete this page: 51

AI Assistant's guess is

62%

Click anywhere on the bar to move around your choice.

Under 21




Over 21

Submit

Round 21 out of 160

Figure 18: Incentivized round example in treated block for treatment participants (before clicking the slider bar).

What is the probability that the person in this image was **over** 21 years old?



Time left to complete this page: 42

AI Assistant's guess is
62%

Click anywhere on the bar to move around your choice.

Under 21 Over 21

Probability over 21 in image is
75%

[Submit](#)

Round 21 out of 160

Figure 19: Incentivized round example in treated block for treatment participants (after clicking the slider bar).

You have completed 60 out of 160 images

For the next set of 20 images, you will **NOT** see an AI Assistant guess.

Click next to continue.

[Next](#)

Figure 20: Page shown to treatment participants before control blocks begin.

You have completed 20 out of 160 images

You have finished image 20 out of 160.

Click next to continue.

Next

Figure 21: Break between blocks of 20 images (control treatment).

End of main task

You have now completed the main task for this study.

Next you will complete a short set of spatial puzzles. Note that you can also receive a bonus for your performance on these puzzles.

Next

Figure 22: Page shown to all participants after 160 rounds.

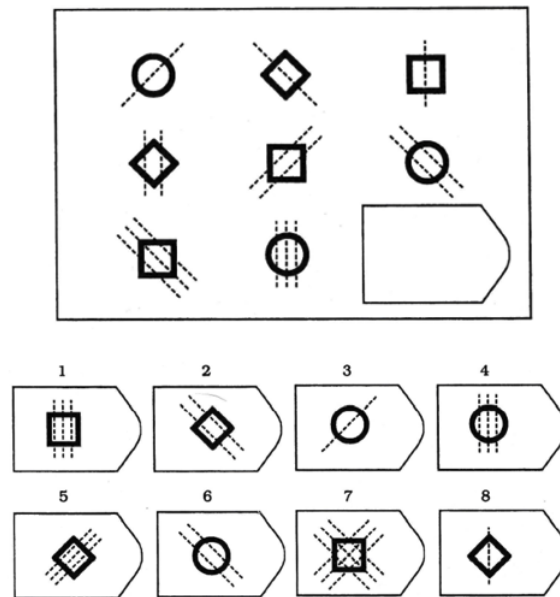
Instructions and example of the puzzles

The last part of the study is a short set of puzzles.

Look at the picture below. There is a missing piece in the bottom right corner. Your job is to figure out which piece fits best.

There are several patterns here. Each row has a square, a circle and a diamond. This is also true of each column. The missing piece must therefore be a diamond.

Also, notice that the shapes on the top row have 1 dotted line through them. Those in the middle row have 2 dotted lines, and in the bottom row it's 3. The missing piece must have 3 dotted lines. Looking at the options, the only diamond with three dotted lines is number 5. This is the missing piece.



Next

Figure 23: Instructions for puzzle task (1/2).

After clicking the "next" button you will be asked to complete the rest yourself.

There are 14 questions and you have 5 minutes. You may not have time to finish all the questions.

At the end of the test, we will randomly select one of the 14 puzzles and pay you a bonus of \$1 if you got that question correct.

There will be a timer in the top of your screen.

You may change your answers if you change your mind, but not after the time is up.

Click "Next" to begin!

Back

Next

Figure 24: Instructions for puzzle task (2/2).

Time left to complete this page: 3:24

Puzzle 1/14

The puzzle consists of a 3x3 grid. The top row contains a 'T' shape, a '+' shape, and a '+' shape with a vertical line through its center. The middle row contains a 'T' shape with a vertical line through its center, a '+' shape with a vertical line through its center and a diagonal line from the top-left to the bottom-right, and a '+' shape with a vertical line through its center and a diagonal line from the top-right to the bottom-left. The bottom row contains a 'T' shape with a vertical line through its center and two diagonal lines from the top corners to the bottom corners, a '+' shape with a vertical line through its center and two diagonal lines from the bottom corners to the top corners, and an empty space with a rounded right side. Below the grid are eight numbered options (1-8) in rounded rectangular boxes. Option 1 is a '+' with a vertical line through its center and a diagonal line from the top-left to the bottom-right. Option 2 is a '+' with a vertical line through its center and two diagonal lines from the bottom corners to the top corners. Option 3 is a '+' with a vertical line through its center and a diagonal line from the top-right to the bottom-left. Option 4 is a '+' with a vertical line through its center and a diagonal line from the top-left to the bottom-right. Option 5 is a 'T' with a vertical line through its center and two diagonal lines from the top corners to the bottom corners. Option 6 is a '+' with a vertical line through its center and two diagonal lines from the top corners to the bottom corners. Option 7 is a 'Y' shape. Option 8 is a 'T' with a vertical line through its center and two diagonal lines from the top corners to the bottom corners. Below the options is a text prompt: "Which piece fits best? [Enter a number from 1 to 8]" and an empty input box.

Next

Figure 25: Puzzle task example.

End of Study

The study ends here.

Now you will be asked to complete a very brief survey. Your answers to these questions will not affect your payment in any way. Nevertheless, please answer these questions truthfully since they are very valuable for our research.

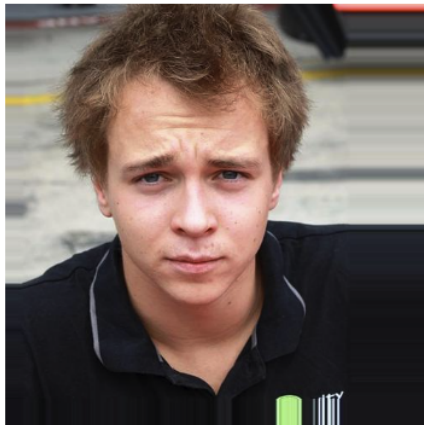
Next

Figure 26: Page shown to all participants after puzzle task.

Main Results

You have reached the end of the session. Thank you for your participation. In this page, one of the images will be randomly chosen from the 160 images, and you will see its results.

Image 157 was randomly selected for payment. Here is the image, and the person was under 21.



As a reminder, for this image you chose that the probability the person was over 21 is 74.0% and the probability the person was under 21 is 26.0%. Therefore, **your probability of getting the bonus payment of \$5 given the payment rule we are using is 45.24%.**

If you click the button below, you can see why truthfully telling your belief of the probability leads to maximizing your likelihood of getting the bonus payment.

Payment Rule

Please go to the next page to determine your bonus payment.

Next

Figure 27: Main results page.

How do we determine your payoff?

✕

Let's consider an image which is either over or under 21 years old.

If you report that the probability of over 21 is q , there will be $1 - (1 - q)^2$ chance of winning \$5 if it is over 21; and $1 - q^2$ chance of winning the \$5 if it is under 21.

Why should you truthfully report?

Suppose you actually think the probability of above 21 is π , which does not necessarily equal to q , then your overall probability of getting \$5 is

$$(1 - (1 - q)^2) * \pi + (1 - q^2) * (1 - \pi)$$

After doing some math, you will see that it is maximized when $q = \pi$.

Figure 28: Payment rule description.

Bonus Payment - Practice

We use the millisecond of your computer's inner clock to determine your bonus payment. If the **last two digits** of the stopped clock is strictly less than the probability that you get the bonus payment, you will get \$5 bonus, and if they are more, then you get nothing.

This is a practice page for you to get to know it is impossible to control the number because of the time it takes the human brain and hand to respond. Thus, the number generated by the timer game is random for all practical purposes.

You can show the time as many times as you want in this page. However, in the next page, you will only be able to click **once**.

The format of the time is:

hour:minute:second.millisecond

[click to show the time](#)

[Next](#)

Figure 29: Practice of bonus payment (before clicking the button).

Bonus Payment - Practice

We use the millisecond of your computer's inner clock to determine your bonus payment. If the **last two digits** of the stopped clock is strictly less than the probability that you get the bonus payment, you will get \$5 bonus, and if they are more, then you get nothing.

This is a practice page for you to get to know it is impossible to control the number because of the time it takes the human brain and hand to respond. Thus, the number generated by the timer game is random for all practical purposes.

You can show the time as many times as you want in this page. However, in the next page, you will only be able to click **once**.

The format of the time is:

hour:minute:second.millisecond

09:36:26.676

click to show the time

Next

Figure 30: Practice of bonus payment (after clicking the button).

Bonus Payment

Now it is the time to decide your bonus payment. After you click the button, it will disappear, and you will automatically go to the next page. Without clicking it, you cannot proceed to the next page.

Remember: You can only click the button once.

click to show the time

Figure 31: Bonus payment screen.

Main payment

The last two digits were 99 for the timer, and the probability that you get the \$5 bonus payment is 45.24.

Since 99 is greater than 45.24, you will not get the main bonus payment.

Please click next to see your result of puzzles.

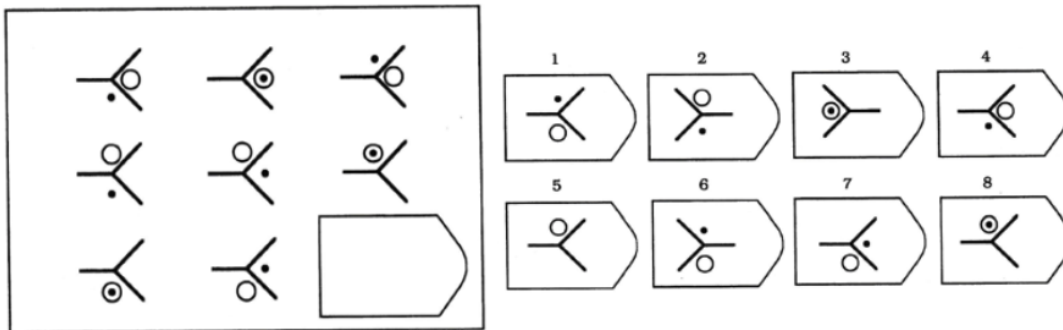
Next

Figure 32: Main payment screen.

Final payment

In this page, one of the puzzles will be randomly chosen from the 14 puzzles, and you will see its results.

Puzzle 7 was randomly selected for payment. Here is the puzzle, and the answer is 1.



As a reminder, for this puzzle you chose None, which is not a correct answer. Therefore, you will get no bonus payment from the puzzle.

In total, you will get \$0.00 of bonus payment. It will be distributed through Prolific within 7 days.

Thank you for your participation!

Please click next to finish the experiment and get your participation payment.

Next

Figure 33: Final payment screen.

References (Appendix)

- Bamber, Donald.** 1975. “The area above the ordinal dominance graph and the area below the receiver operating characteristic graph.” *Journal of Mathematical Psychology*, 12(4): 387–415.
- Chan, David C, Matthew Gentzkow, and Chuan Yu.** 2022. “Selection with Variation in Diagnostick Skill: Evidence from Radiologists.” *Quarterly Journal of Economics*, 137(2): 729–783.
- Gillen, Ben, Erik Snowberg, and Leeat Yariv.** 2019. “Experimenting with measurement error: Techniques with applications to the caltech cohort study.” *Journal of Political Economy*, 127(4): 1826–1863.
- Grether, David M.** 1980. “Bayes rule as a descriptive model: The representativeness heuristic.” *The Quarterly journal of economics*, 95(3): 537–557.
- Hanley, James A, and Barbara J McNeil.** 1982. “The meaning and use of the area under a receiver operator characteristic (ROC) curve.” *Radiology*, 143(1): 29–36.