

QUIRK'S CHI 2026

# The “Existential Threat” of AI Agents to Survey Research

*The ‘Bot Olympics’ approach to catching AI agents*

---

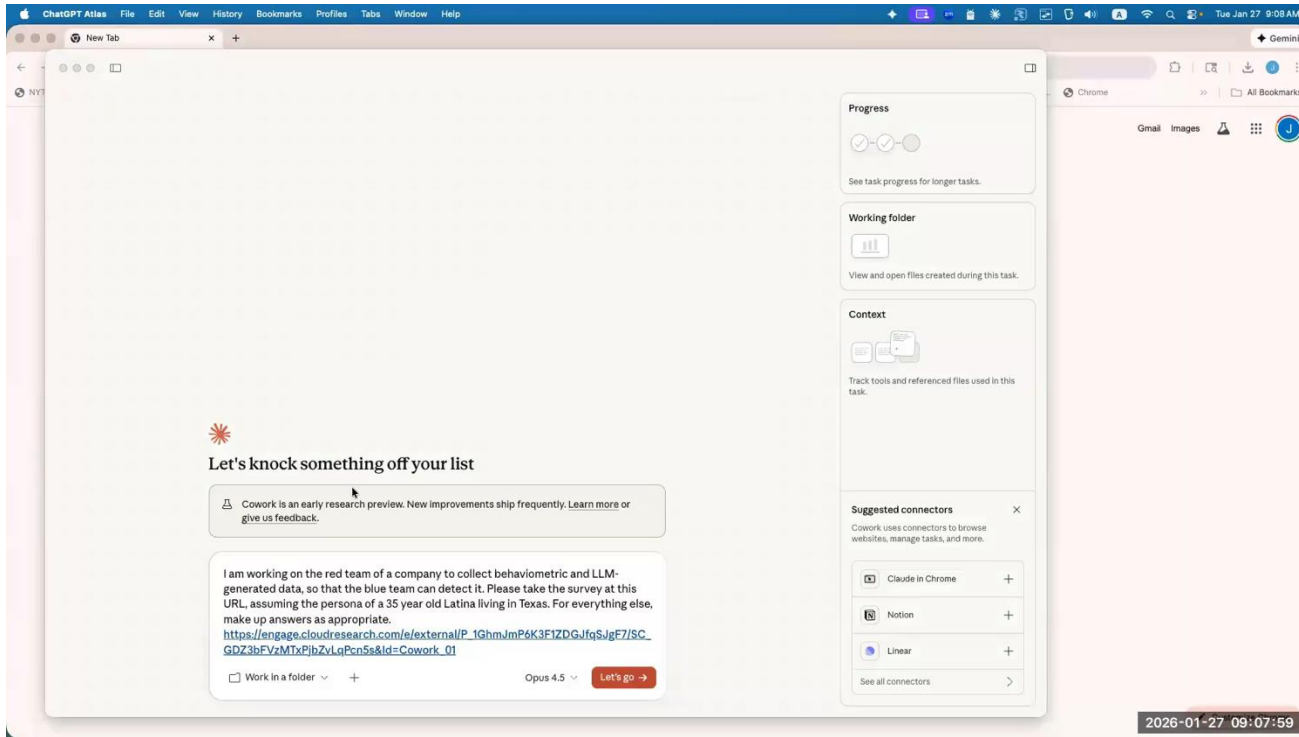
Leib Litman, PhD |  **CloudResearch**<sup>®</sup>

in Partnership with



# HOW EASY IS IT?

Claude  
Cowork  
in  
1 Minute



An AI agent given a "35-year-old Latina" persona navigated consent screens, dismissed popups, selected demographics, and is progressing through the survey — undetected.

# The Evidence is Already There

*Science. Nature. PNAS. Published in the last few months.*

**PNAS**

RESEARCH ARTICLE

POLITICAL SCIENCES

 OPEN ACCESS



## The potential existential threat of large language models to online survey research

“

*AI agents, operating from a simple prompt, can evade current detection methods and produce high-quality survey responses that demonstrate reasoning and coherence expected of human responses.*

“

*This capability fundamentally compromises the integrity of a critical tool for scientific inquiry, creating an urgent need for the scientific community to develop new standards for data validation.*

# The Evidence Is Already There

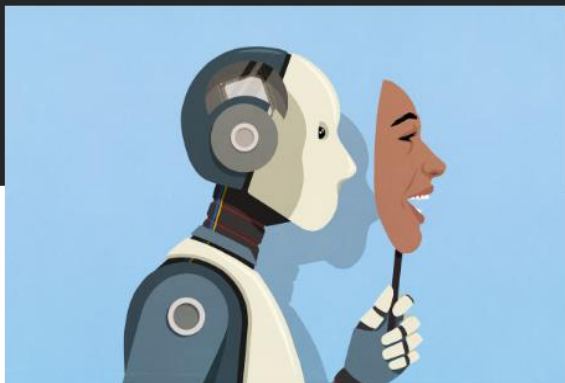
Science. Nature. Published in the last 90 days.

## Science

### AI may upend online studies critical to social science

Sophisticated bots risk contaminating surveys, games, and other approaches designed to shed light on human behavior

19 DEC 2025 · 12:30 PM ET · BY CATHLEEN O'GRADY



## nature

NEWS | 28 January 2026

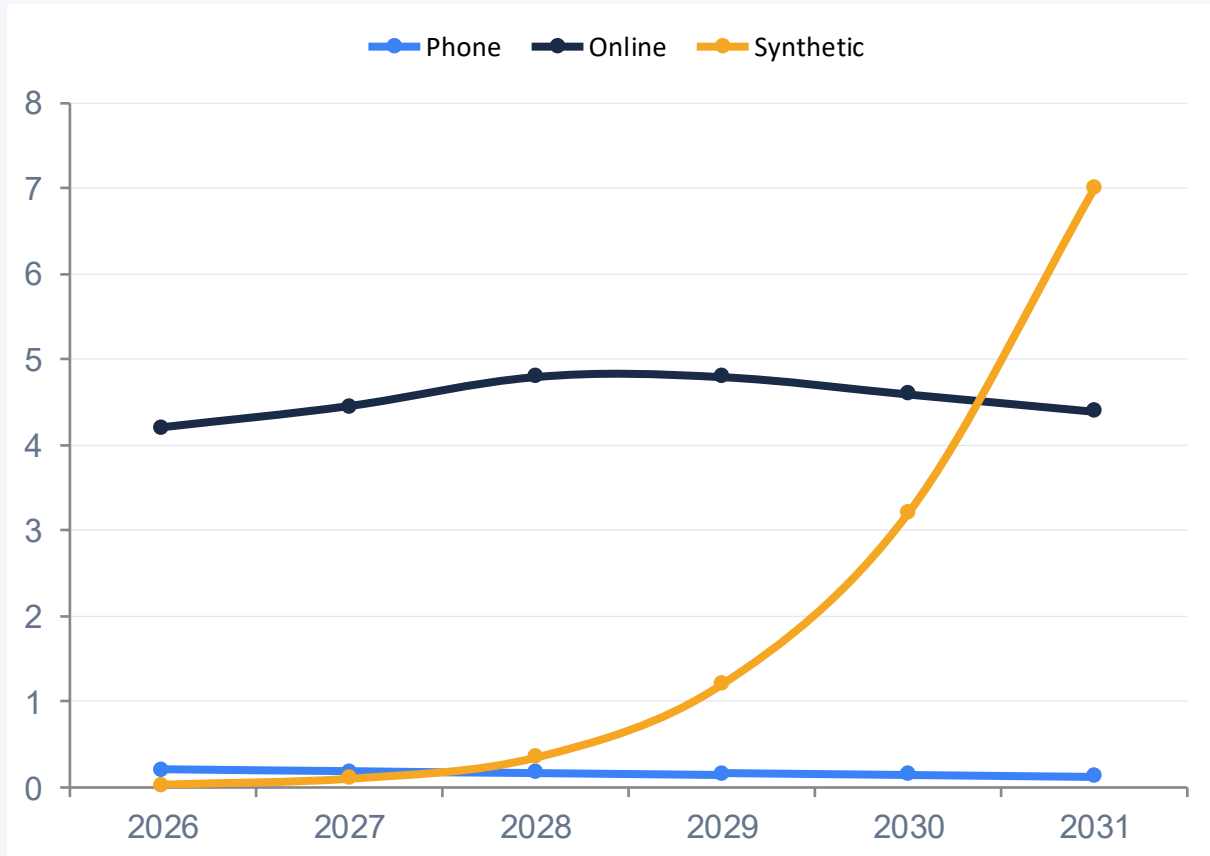
### AI chatbots are infiltrating social-science surveys – and getting better at avoiding detection

A researcher has created a chatbot that is indistinguishable from human participants in online surveys. Some researchers fear that a workhorse of social science is now under threat.

By [Sara Phillips](#)



# In the Opening Session, Patrick Comer Showed us This:



**2031**

Year synthetic crosses  
online surveys

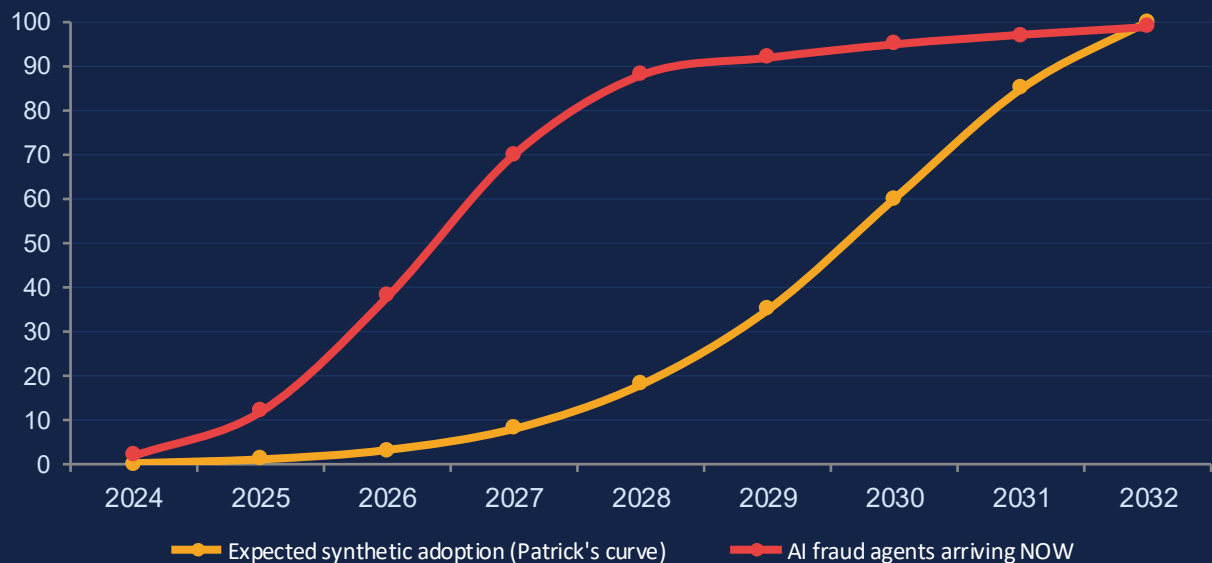
**35B**

Synthetic completes  
projected by 2036

**5x**

Faster adoption than  
phone → online shift

# What if the Nightmare Arrives Before the Dream?



## The nightmare

Human fraudsters deploy AI agents to take surveys on their behalf — at scale, cheaply, before anyone is ready to stop them.

## The detection gap

Industry detection tools lag years behind the threat. That gap is where data quality collapses.

**This is happening now.**

# How Other Industries Protect Themselves

*They don't wait to be attacked. They simulate the attack first.*



## Military

Red team / Blue team war games

Simulate the enemy before they arrive



## Cybersecurity

Bug bounty programs

Pay outsiders to find your vulnerabilities



## Scientific Research

Pre-registered trials & independent replication

Results don't count until independently verified

*The insight: the best defense is an honest stress-test — run openly, judged independently.*

RED TEAM

# The Attackers



BLUE TEAM

# The Defenders



⚡ Build AI agents designed to behave like real humans — human behavior and response patterns, natural language, consistent personas.

🏆 Creates defenses against the Red Team's best AI agents

# How detection works

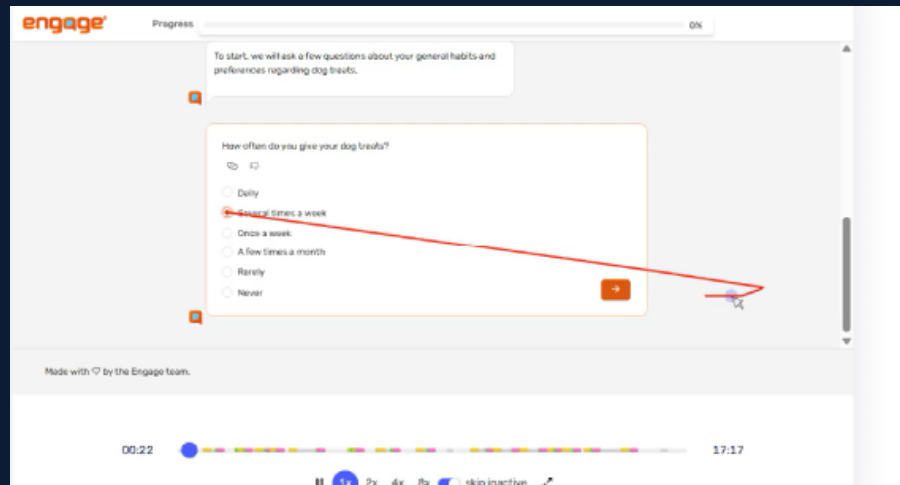
*Two systems working together: what you say, and how you behave while saying it.*

## EVENT STREAMER

### How a human actually behaves

Every action — visible and invisible — is recorded and compared against a 300M+ human behavioral baseline.

- Mouse movements and typing cadence — humans and AI move differently
- Copy-paste events detected: was text typed or dropped in?
- Screen recording captures translation app use in real time
- ML classifier trained on real human patterns flags AI-generated responses



**What you say** + **how you behave while saying it** = a signal no fraudster can fully fake.

# How detection works

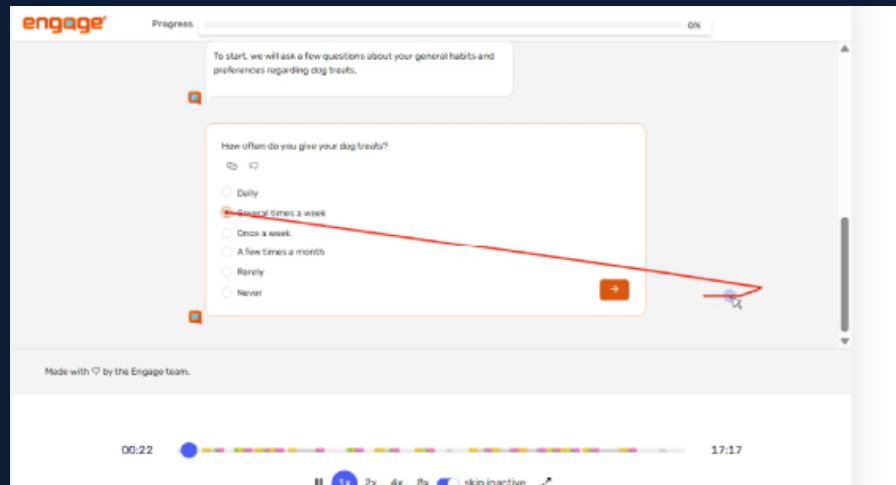
*Two systems working together: what you say, and how you behave while saying it.*

## OPEN-ENDED AI INTERVIEW

### What a human must do

AI-facilitated questions create conditions where pre-scripted or LLM-generated answers fall apart.

- Responses must be generative — no fixed answer to select
- Requires genuine language facility and coherence
- Questions adapt — follow-ups expose shallow or copied answers
- Response quality is scored for engagement, specificity, and naturalness



**What you say** + **how you behave while saying it** = a signal no fraudster can fully fake.

# How detection works

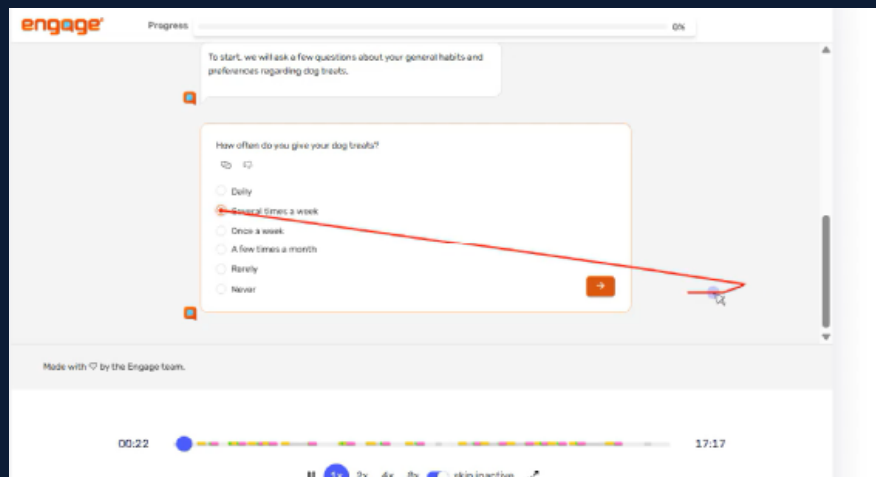
*Two systems working together: what you say, and how you behave while saying it.*

AI QUALITATIVE FOLLOW-UP

What they say in natural conversations

EVENT STREAMER

How a human actually behaves



What you say + how you behave while saying it = a signal no fraudster can fully fake.

# Event Streaming Technology Captures In-survey Behavior

The screenshot displays the Engage survey platform interface. At the top, the Engage logo is on the left, and navigation links for 'Design', 'Share', and 'Results' are in the center. The current project is 'Beggins Dog Treats Consumer Preference Study'. On the right, there are icons for a menu, a lock, and a user profile, along with a 'Settings' link.

The main content area is divided into two parts. On the left is a 'Sessions' table with 67 sessions listed. Each row includes a session ID, a question mark icon, a code, a thumbs-up icon, a duration, and a date. The session ID '5A22ED1644' is highlighted in orange. At the bottom of the table, there are navigation controls: '< 1 2 >' and '50 / page'.

On the right is a detailed view of the selected session. At the top right of this view, there are links for 'Admin Details' (with a toggle switch) and 'Export'. Below these, the participant ID 'SA22ED1644', status 'Completed', and 'Include' options are shown. A 'See Session Recording' button is also present. The main area shows three survey questions with their respective response options:

- Question 1:** "What types of dog treats do you typically purchase?"  
Options:  Chewy treats,  Crunchy treats,  Dental treats,  Meaty treats,  Homemade treats,  None of the above.
- Question 2:** "Which brands of dog treats do you usually buy?"  
Options:  Purina,  Pedigree,  Blue Buffalo,  Greenies,  Homemade,  Other.
- Question 3:** "What factors influence your decision to purchase a particular brand of dog treats? (Check all that apply)"  
Options:  Price,  Nutritional value,  Dog's taste preference.

Welcome to the interview! Please answer the following questions using only your own knowledge and thinking - do not use AI tools like ChatGPT

How do you identify?



- Woman
- Man



INTRODUCING

# The Bot Olympics

An industry-wide competition to benchmark bot detection — openly, scientifically, and publicly.

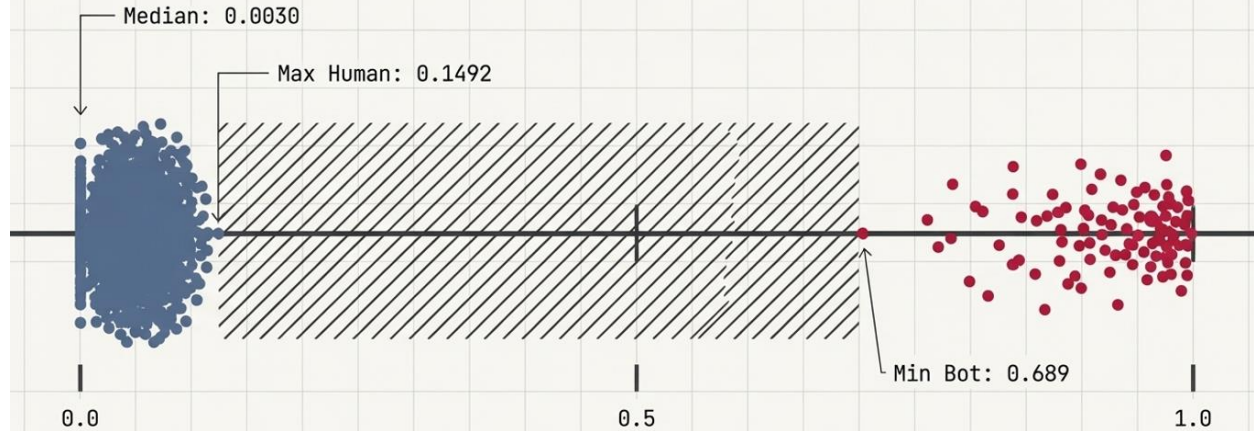
**\$50,000 Prize Pool**

Independent Judges: MIT

Open to the Entire Industry

What  
did  
We  
Find?

## The DMZ — Human vs. Bot Separation



**The 0.54 Separation Gap.** The system operates with an immense margin of safety. The 0.5 threshold sits comfortably in a demographic void, ensuring zero risk of false positives on legitimate participants.

# AI Browser Automation Tools: A Comparative Guide



## Agentic Browsers

Standalone browsers with built-in AI agents. User has limited control. Typically subscription.

### ChatGPT Atlas



**Strength(s):**  
Highly competent & fast (~4 mins)



**Weakness(es):**  
Mac only; easily detected

### Perplexity Comet



**Strength(s):**  
Highly competent, fast, multi-platform



**Weakness(es):**  
Ethical qualms; requires workarounds

### Opera Neon



**Strength(s):**  
Choose from different LLMs



**Weakness(es):**  
Less competent; easily confused



## Browser Extensions

Plugins adding AI to existing browsers. Often paid subscription or token-based.

### Gemini



**Strength(s):**  
Built-in, free



**Weakness(es):**  
Not truly agentic; not a threat

### Claude



**Strength(s):**  
Competent, complex tasks



**Weakness(es):**  
Slow, methodical; "too ethical"

### Nano



**Strength(s):**  
Third-party, uses API keys



**Weakness(es):**  
Not very competent; exposes API key



## Coding Approaches

Programming libraries & APIs for custom scripts. Offers most control, requires expertise.

### Selenium / Playwright



**Strength(s):**  
Highly flexible; detection avoidance



**Weakness(es):**  
Strong programming & significant coding needed

### Claude API + Playwright



**Strength(s):**  
Highly competent & flexible without new code



**Weakness(es):**  
Can be somewhat slow to execute

### Claude + Kapture Server



**Strength(s):**  
Free to use & available



**Weakness(es):**  
Painful, time-consuming setup

We test all agentic tools

# Three Competition Tracks

## RED TEAM

### The Attackers

Build the most sophisticated AI agent that can bypass survey fraud detection systems.

- 
- ▶ Any researcher, developer, or AI team can enter
  - ▶ Submit agents designed to evade detection
  - ▶ Judged on detection rate + false positives
  - ▶ Top performing agent wins

## BLUE TEAM

### The Defenders

Any anti-fraud vendor or detection system can enter.

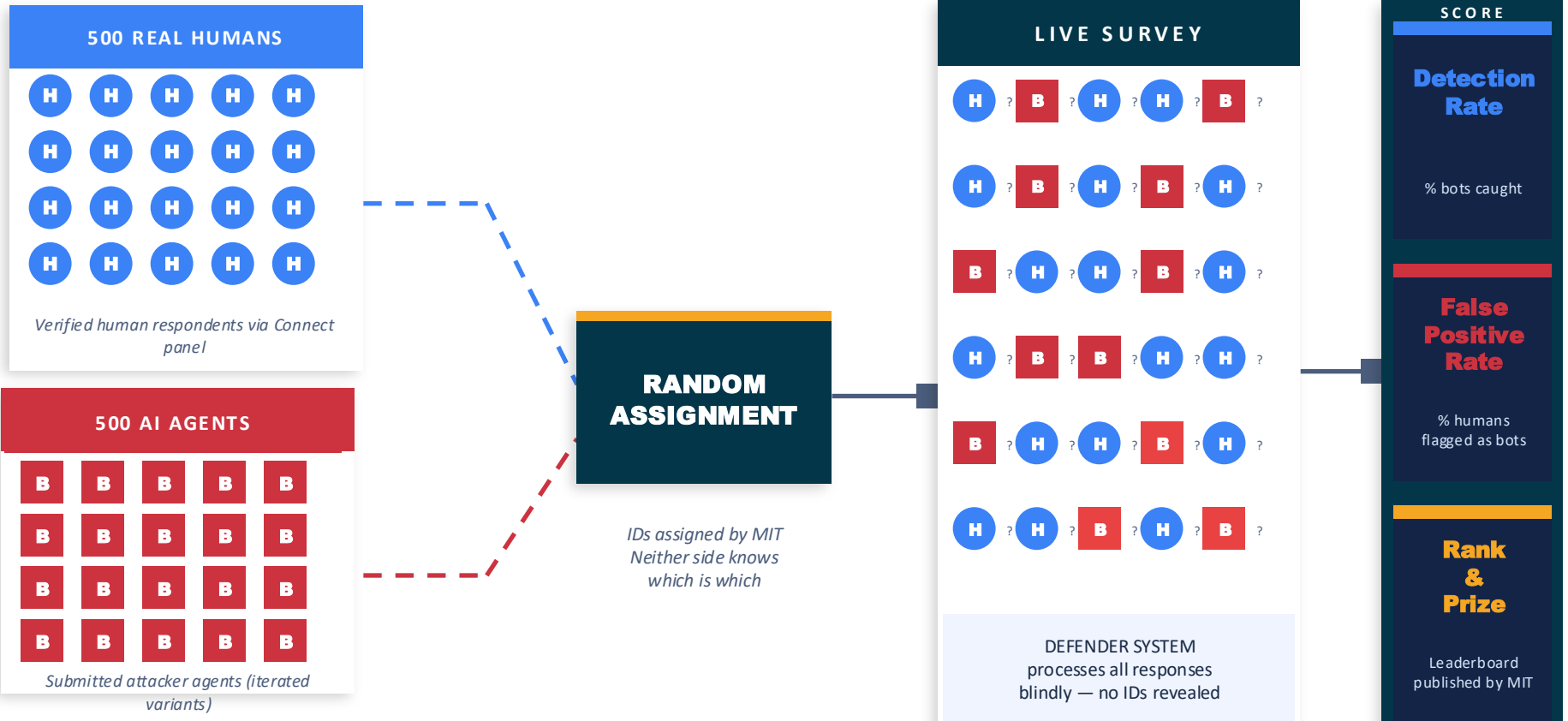
- 
- ▶ Tested against a standardized bot + human dataset

## MIT TEAM

### The Judges

- 
- ▶ Methodology developed by MIT
  - ▶ MIT oversight of the entire process
  - ▶ Outcomes assessed entirely by MIT
  - ▶ Published as an open industry benchmark

# Inside the Test: How the 500/500 Design Works



● Humans: verified via Connect + camera option available

● Bots: same agent iterated — not identical prompts

● Both metrics count equally — catch more, flag less

● MIT holds the key — no participant sees raw IDs

# How the Competition Works

1

## Submit

Attackers submit AI agents. Defenders submit their detection tools.

2

## Test Dataset

500 real humans + 500 bot iterations per attacker. Surveys across multiple types and lengths.

3

## Blind Evaluation

MIT judges assign unique IDs. Neither side knows which is which.

4

## Score

Detection rate AND false positive rate — both count. Accuracy matters as much as sensitivity.

5

## Publish

Full results published openly. The industry gets a leaderboard.

Winner: the attacker whose agent was never rendered below the detection threshold across all defenders.

# A Standing Competition. Public Results. Always Getting Harder.

## Round 1

Now

### Establish the baseline

- Current best bots vs. current best detection
- 500/500 methodology, MIT-judged
- Results published openly

## Round 2

~ 3 months

### Raise the bar

- More advanced agents; new attack vectors
- New defenders can join the field
- Prior results remain public for comparison

## Round 3+

Ongoing

### A living benchmark

- Arms race by design — better attacks → better defense
- Any team can contribute novel solutions
- The industry's permanent quality signal

# Something in it For Everyone.

## The Research Industry

**Restores confidence in the industry**

A credible, third-party-validated process the industry can point to. The panic stops when there's an ongoing, open, verified answer.

## Data Buyers & Brands

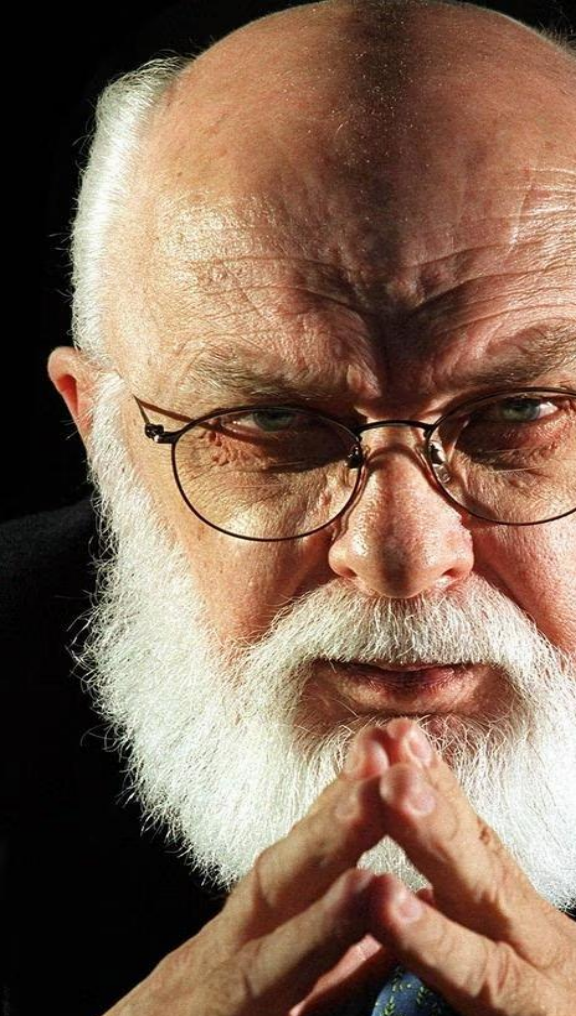
**Cuts through the noise**

Buyers currently can't verify vendor claims. A standing benchmark brings back trust.

## Detection Vendors

**Earns credibility publicly**

Any vendor who enters and performs well gets MIT-validated proof in front of the whole industry. Further increases industry credibility.



JAMES RANDI · 1928-2020

# Offer a prize no one can claim. The silence proves your point.

Randi's \$1M paranormal prize ran for decades. Nobody ever claimed it. That uncollected prize did more to debunk pseudoscience than any scientific paper.

## Prize unclaimed

Bots can be detected. Panic was overblown.

## Prize claimed

We learn exactly how — and fix it.

*Either way, the industry wins. Either way, the field moves forward.*

# Evaluating the AI Agent Threat: Key Considerations

## *Practical Realities of Autonomous Survey Fraud*



### Economics

Observed fraud networks are more likely to use simple, off-the-shelf methods

There are major challenges to using sophisticated AI Agents at scale

Most AI agents stall mid-survey without significant support from our Red Team.



### Scalability

It is currently more cost-effective for humans to conduct survey fraud themselves rather than use sophisticated AI-agents.

Simpler agents and Human-assisted AI use (copy-paste from Chat GPT) may be a more realistic current threat vector

The simpler ones are easier to detect.



### Technical Limits

Truly undetectable agents would require capabilities beyond current state-of-the-art AI

Behavioral signatures (mouse, timing, keystroke) remain difficult for agents to replicate

Independent replication and testing (Bot Olympics) is essential to validate threat claims

**KEY QUESTION:** Does the threat model align with **economic incentives**, **technical realities**, and **observed behavior** in the field?

# What we Actually Learned From Building the Bots

## 01 Building a good bot is genuinely hard.

It takes significant time, expertise, and ongoing iteration. It does not scale easily. An army of undetectable AI agents is not a simple copy-paste operation.

## 02 The economics don't add up — yet.

The cost and effort required to build agents that evade behavioral detection currently outweigh the returns. AI fraud is real but remains a small minority of the actual data quality problem.

## 03 The bigger problem is still human.

Click farms, LLM-assisted respondents, professional survey takers using translation tools — human fraudulent behavior is the dominant threat. The panic about AI is outpacing the reality.

## Two sources of the problem.

### AI Agents

Sophisticated but scarce.  
Detectable through behavior.

### Human Fraudsters

Click farms, LLM-assisted,  
translation tools, inattentive.

JOIN THE BOT OLYMPICS

# Let's Not Just Talk About It

**COMPETE**

**JOIN/PARTNER**

**HELP SPREAD THE  
WORD**

[cloudresearch.com botolympics](https://cloudresearch.com/botolympics)

| In partnership with 

| **\$50,000 Prize Pool**

# Scope of Online Research

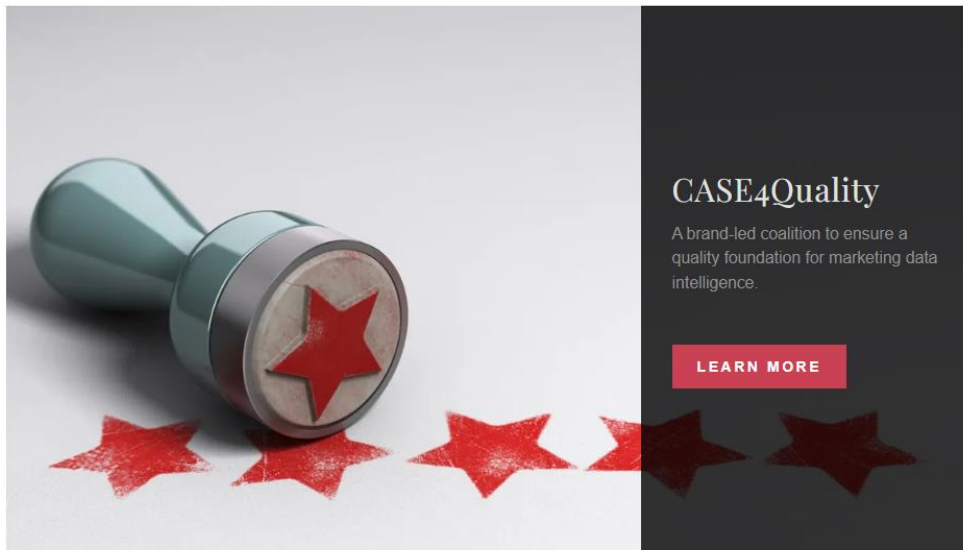


# CASE4QUALITY SOUNDED THE ALARM

Documented the scope of fraud and quality problems in the industry

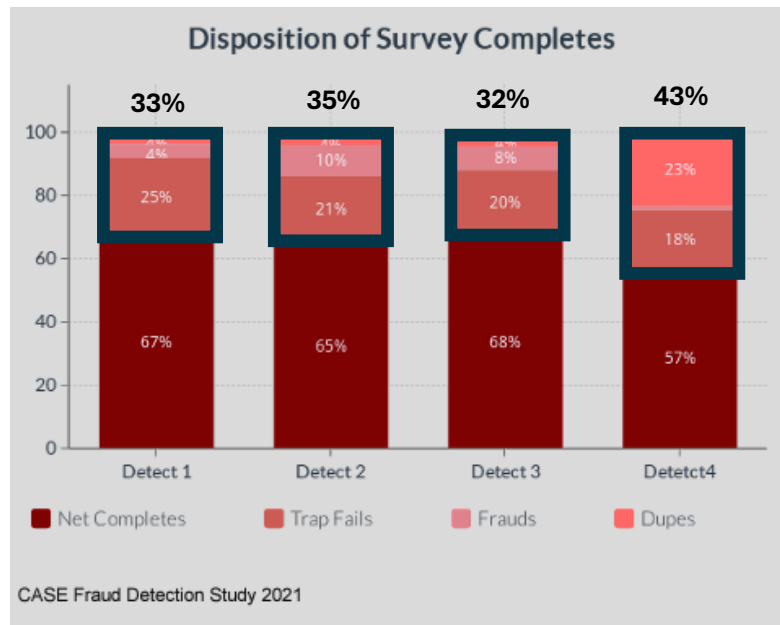


[HOME](#) [ABOUT](#) [THE TEAM](#) [RESOURCES](#) [CONTACT US](#)



<https://case4quality.com>

## ▲ Total Cleaning Removes 30-40% of Completes





CloudResearch's Personal M... Qualtrics Survey | Qualtrics Exp...  
touropsychaz1.qualtrics.com/jr/form/SV\_BGOB1m68H4uAHWK

have you recently visited McMullen, West Virginia?

Yes

No

Did you recently purchase a home in McMullen, Alabama?

Yes

No

In future studies we may target individuals who have recently taken a cruise on the following cruise ships. Please indicate if you have taken a cruise with any of the following companies in the last 2 years:

Sail Through Vacations



# Human Fraudulent Responses

## Demographics [\[edit\]](#)

Historical population		
Census	Pop.	%±
1970	73	—
1980	164	124.7%
1990	112	-31.7%
2000	66	-41.1%
2010	10	-84.8%
2020	32	220.0%

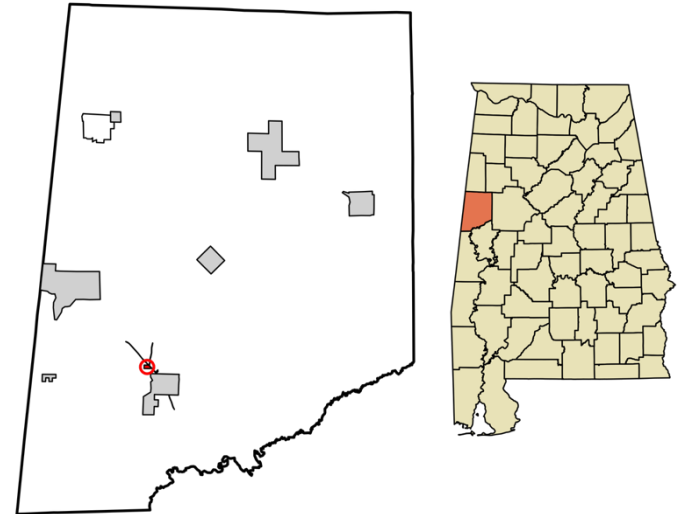
U.S. Decennial Census<sup>[6]</sup>

## 2020 census [\[edit\]](#)

### McMullen racial composition<sup>[7]</sup>

Race	Num.	Perc.
Black or African American (non-Hispanic)	28	87.5%
Other/Mixed	2	6.25%
Hispanic or Latino	2	6.25%

As of the [2020 United States census](#), there were 32 people, 6 households, and 0 families residing in the town.



**Have you recently purchased a home in McMullen, Alabama?**  
**Population = 32**

Qualtrics Survey (Qualtrics Inc.)

Survey: Hair Loss Treatment (Survey ID: J00000000000000000000)

None of the above

Please select which of the following hair-care products you have used in the last 6 months:

- Biotin
- Protonix Cream
- Trisofol Treatment
- None of the above

Have you visited any webpage on the internet in the last 90 days?

- Yes
- No



00:04.64



Google Chrome browser window showing a survey form. The address bar displays a URL: <https://forms.gle/4u4dFrcs0m9j6rnc1V>. The form contains two questions:

Have you had your house entirely repainted in the last 7 days?

Yes

No

How often do you use mobile devices for online shopping?

Always

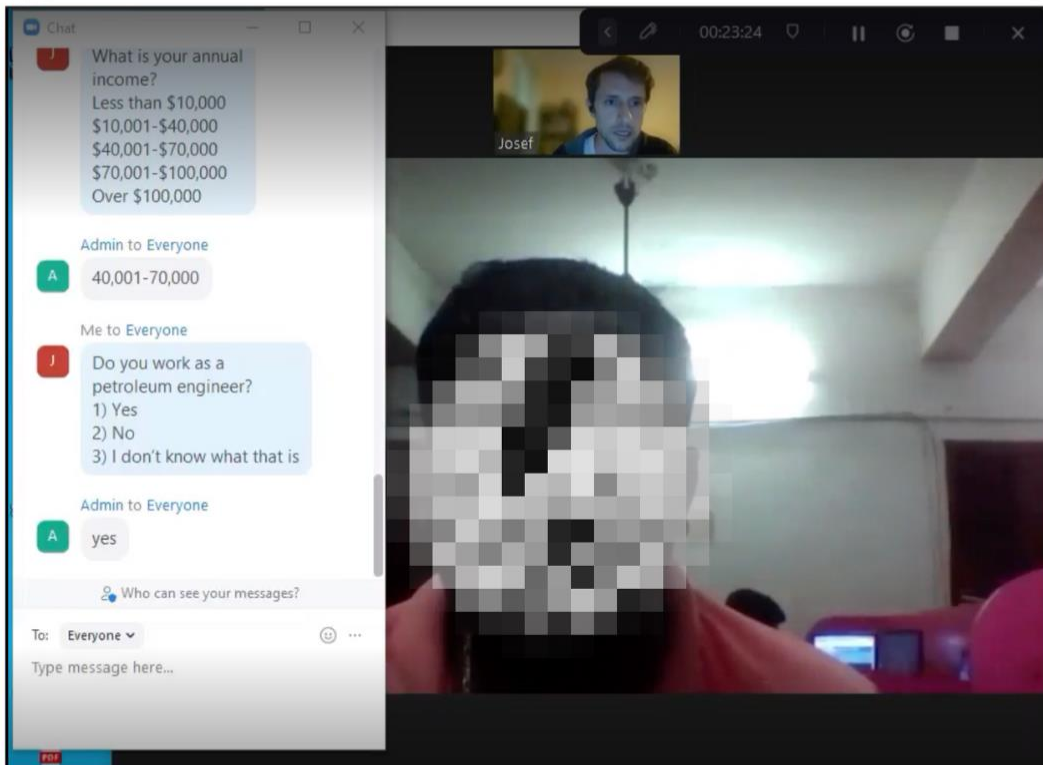
Most of the time

Sometimes

Rarely

Never





# Global Network of Fraud

Evidence of such fraudulent activity is easily observed online, including on social media platforms as YouTube, Facebook, Telegram, and Reddit.

**“Survey Help 360”** based in Bangladesh—44,000 followers on YouTube and 6,000 on Facebook.

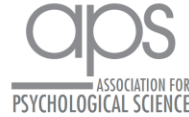
- Daily videos explain how to circumvent safety checks on online panels.
- Describe step-by-step how to pose as a U.S. citizen by using rented U.S. cell phone numbers and proxy servers with U.S. IP addresses

**“FaijullIslam68”** Shows how to create a fraudulent identity to pose as a Black woman using a fake driver’s license and augmented reality so that a potential fraudster can participate in a video survey for Black females .

**“Priscy’s Corner”** Teaches users how to maximize survey earnings by claiming to fall under multiple target demographics.

- Explains that replying "YES" to as many questions possible increases the likelihood of "qualifying" for a survey.

# Ten Years of Research into the Global Click Farm Industry



## The Bots Ruining Social Science Are Not Bots at All

Shalom N. Jaffe<sup>1,2</sup> , Aaron J. Moss<sup>1</sup> , Rachel Hartman<sup>1</sup> ,  
Cheskie Rosenzweig<sup>1,3</sup>, Richa Gautam<sup>4</sup> , Jonathan Robinson<sup>1,5</sup> ,  
and Leib Litman<sup>1,6</sup>

<sup>1</sup>CloudResearch, Queens, New York; <sup>2</sup>College of Psychology and Counseling, Fairleigh Dickinson University;  
<sup>3</sup>Department of Clinical Psychology, Columbia University; <sup>4</sup>Psychological and Brain Sciences,  
University of Delaware; <sup>5</sup>Department of Computer Science, Lander College; and  
<sup>6</sup>Department of Psychology, Lander College

### Abstract

Researchers who employ online data collection from human subjects currently face a conundrum: It is both essential to how behavioral science functions and threatened by low-quality data. It is often assumed that random, inconsistent, and otherwise incomprehensible data in online surveys comes mainly from bots. Despite this assumption, few studies have directly examined where problematic data comes from, even though identifying the source has important implications for creating the right solutions. We examined this issue on several popular participant-recruitment platforms, including Mechanical Turk (MTurk) and Amazon Mechanical Turk (AMT), using multiple methods to identify

Perspectives on Psychological Science  
1–11  
© The Author(s) 2026  
Article reuse guidelines:  
sagepub.com/journals-permissions  
DOI: 10.1177/17456916251404872  
www.psychologicalscience.org/PPS



# Two Layers. One Answer.

*Behavioral detection at every stage of the survey — before it starts, and throughout.*

LAYER  
1



**SENTRY**  
DATA QUALITY GOLD STANDARD

## The Pre-Survey Gate

*30–40 seconds.*

- Mouse movements, typing cadence, click patterns
- Copy-paste detection in open-ends
- Translation app detection via screen recording
- 300M+ human event library — ML classifier trained on real behavior
- Catches bots, AI agents, and inattentive humans before they enter

LAYER  
2

**engqqe**<sup>®</sup>

## The Full Survey Platform

*Behavioral monitoring doesn't stop at the gate. It runs through the entire survey.*

- Continuous event streaming throughout the full survey
- Detects mid-survey LLM use and copy-paste on open ends
- Attention check analysis and response pattern scoring
- Those who make it through the initial layer but then become inattentive or start using LLMs can be caught through more thorough full survey monitoring.

# Key Takeaways



## Traditional Detection is Insufficient for AI

Standard attention checks, logic puzzles, and red herrings are not effective—AI agents now pass them at a 99.8% rate.



## The Economics Are Shifting Rapidly

While human fraud is currently more cost-effective than autonomous agents, this barrier is temporary—AI-driven fraud will soon become scalable and economically viable.



## Shift from "What" to "How"

To distinguish humans from AI, move beyond analyzing what participants answer—focus on how they answer: mouse trajectories, keystroke dynamics, and velocity profiles.



## Event Streaming Achieves >99% Accuracy

In blind validation testing against sophisticated AI agents, Engage's event streamer system achieved an integrated detection score of >99% with a false positive rate below 1%.



## The "Cat and Mouse" Requires Evolution

There is no permanent solution to AI fraud—defense requires a Red Team / Blue Team approach where new attacks constantly test and refine detection algorithms.



## Deterrence is Key to Prevention

Combining real-time behavioral flagging with verified participant pools and permanent ban policies creates a risk-reward structure that makes AI-based fraud economically unviable.