

Prebunking Elections Rumors: Artificial Intelligence Assisted Interventions Increase Confidence in American Elections

Mitchell Linegar^{1*}, Betsy Sinclair², Sander van der Linden³
R. Michael Alvarez¹

¹Linde Center for Science, Society, and Policy, California Institute of Technology, Pasadena & 91125, USA.

²Department of Political Science, Washington University in St. Louis, St. Louis & 63130, USA.

³Department of Psychology, Cambridge University, Cambridge & CB2 3EB, UK.

*Corresponding author. Email: mlinegar@caltech.edu

Election Rumors and Conspiracies

- Widespread and widely believed
- Strong partisan trends
- Actively endorsed by political elites
- Undermine democracy, democratic participation
- Discourage peaceful transition of power

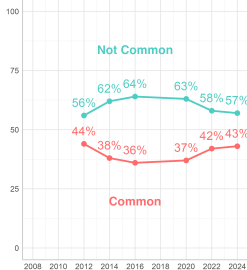
A Need for Tools to Combat Disinformation

- Enormous space of misinformation
- Debunking and prebunking can reduce belief in misinformation
- Individually pushing back against false allegations is costly and slow
- Enter AI

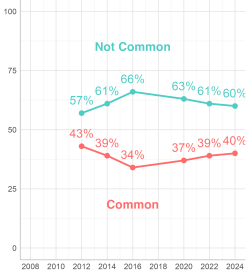
This paper

- Preregistered, two-wave experimental study of U.S. registered voters
- YouGov panel ($N = 4,293$)
- Goal: prebunk / inoculate against election disinformation
- Test five common and widespread election myths
- Use AI to automatically produce inoculation doses
- AI-generated prebunks reduce belief in election conspiracies, increase belief in election integrity

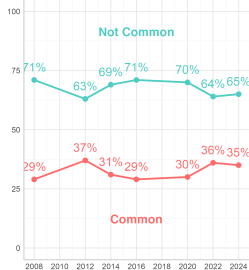
People voting an absentee ballot intended for another person



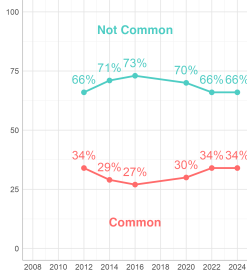
People voting who are not U.S. citizens



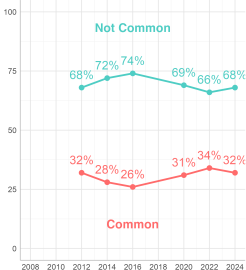
People pretending to be someone else when going to vote



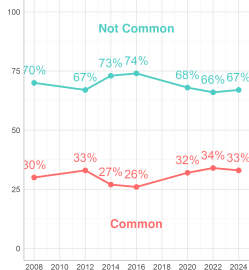
People voting more than once in an election



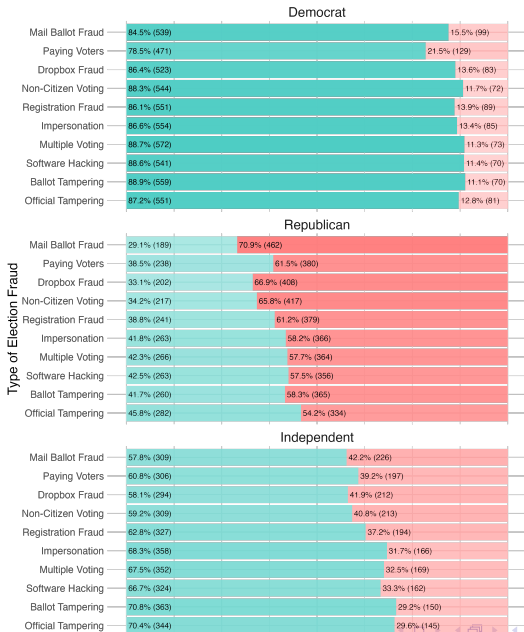
Officials fraudulently changing the reported vote count



People stealing or tampering with ballots that have been voted



Frequency of Beliefs About Various Types of Election Fraud by Party Identification (Weighted by Sampling Weights)



Hypotheses

- **H1:** Participants exposed to prebunking of a specific election-related rumor will report lower confidence in that rumor compared to the control group.
- **H2:** Participants exposed to prebunking of a specific election-related rumor will report higher confidence that their votes will be accurately counted in the next election compared to the control group.

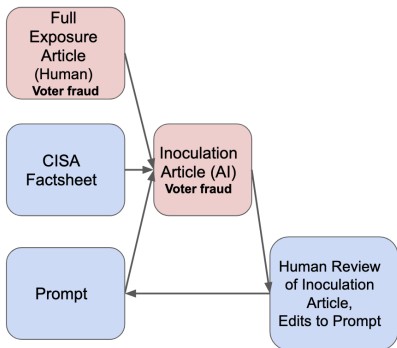
Experimental Design

- Participants answer pre-treatment questions
- Pick five salient election rumors, each participant assigned to one
- Choose Breitbart articles endorsing each rumor
- All participants read rumor-relevant article (“Full Exposure Article”)
- Prior to this, read either AI-written “Inoculation Article” or AI-written “Placebo Article”
- Participants answer post-treatment questions, and again one week later

AI-Written Inoculation Articles

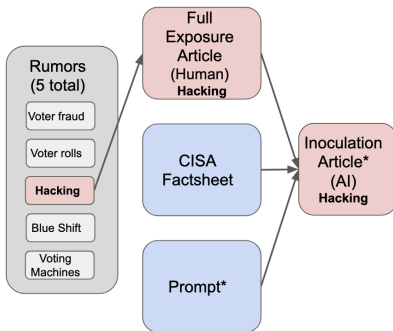
- Input: article endorsing false rumor, CISA fact sheet, prompt
- Human-in-the-loop process: iterate initial prompt until can produce satisfactory inoculation articles for a single rumor
- Use same prompt for all other rumors
- Randomly assign HITL and purely LLM generated articles

AI-Article Prompt Creation Process



Repeat until no edits to prompt
→ Prompt*

AI-Article Writing Process



Use same Prompt* for all rumors

Figure: Prompt and article writing procedure.

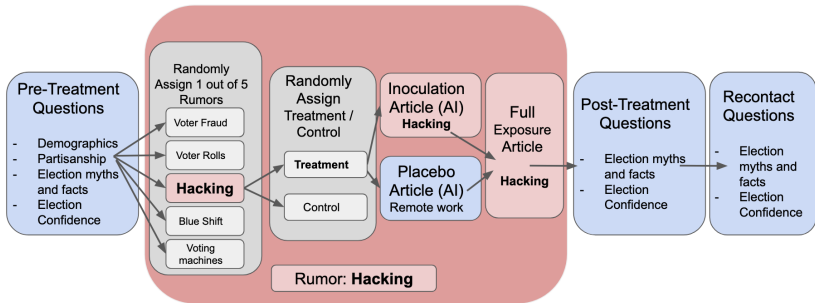
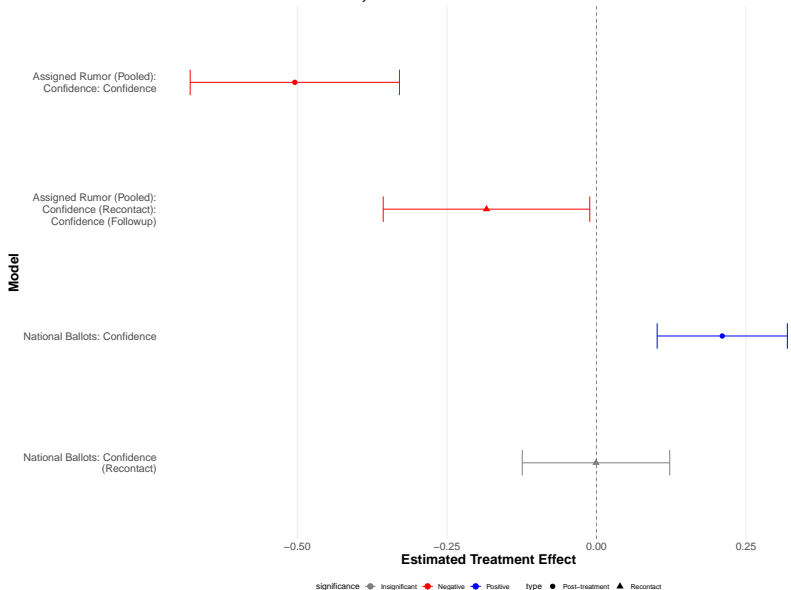


Figure: Experimental design. Blue: common to all participants, regardless of assigned rumor. Grey: where randomization occurs. Red: articles participants are assigned to.

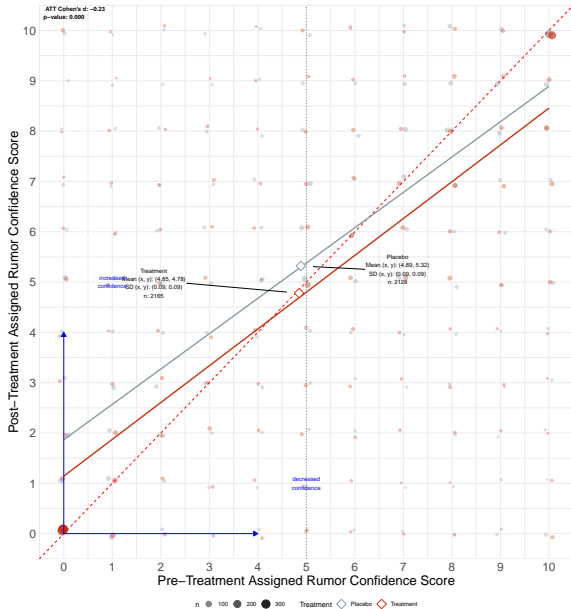
Treatment Effect Estimates

Confidence that: Election Rumors are True, Ballots are Accurately Counted Nationally

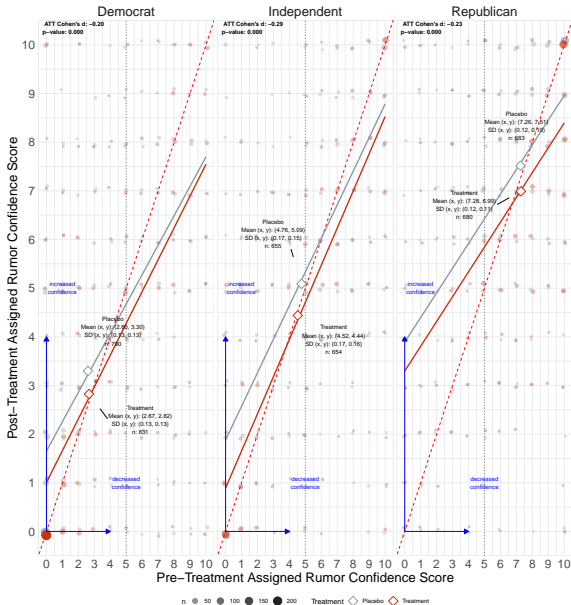
Measured Immediately Post-Treatment and After One Week



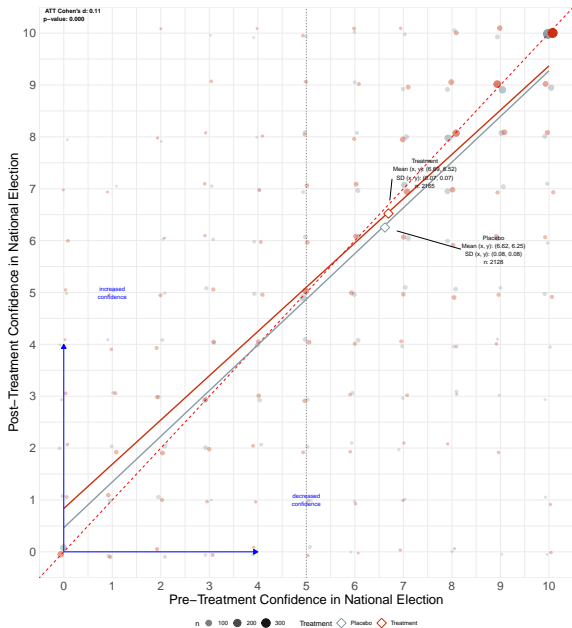
Confidence in Assigned Election Rumor: Pre-Treatment vs. Post-Treatment



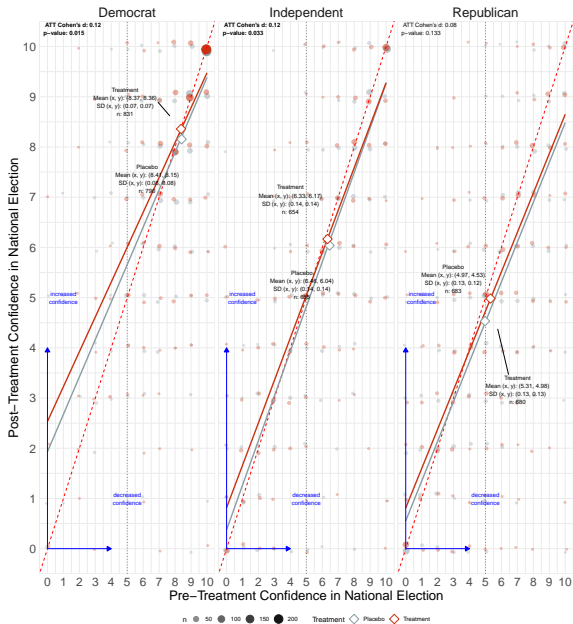
Confidence in Assigned Election Rumor: Pre-Treatment vs. Post-Treatment By Party Identification



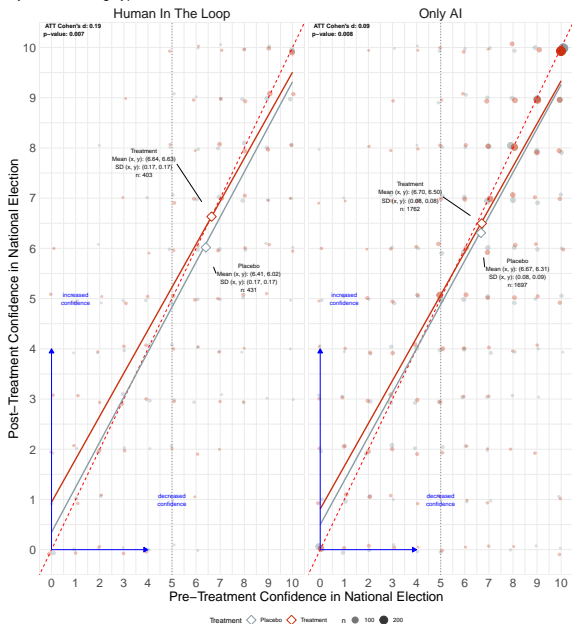
Confidence in National Election: Pre-Treatment vs. Post-Treatment



Confidence in National Election: Pre-Treatment vs. Post-Treatment By Party Identification

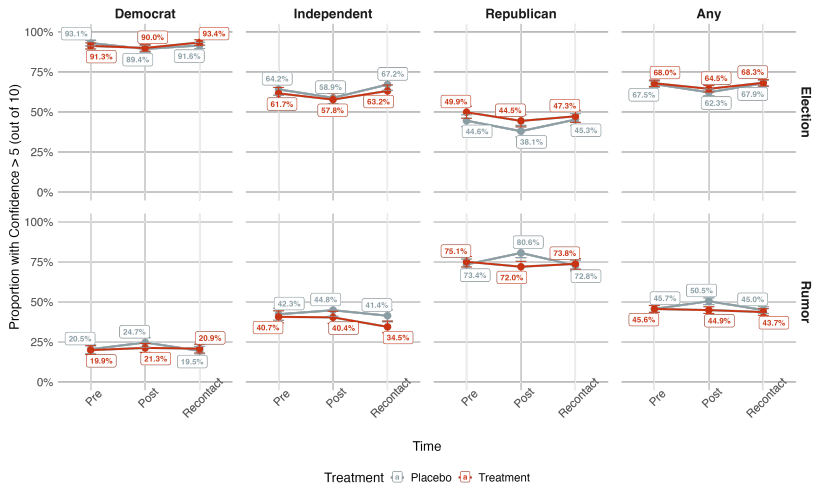


Confidence in National Election: Pre-Treatment vs. Post-Treatment By Article Writing Type



Confidence in Election Integrity and Election Rumors

By Party Identification Over Time, Weighted by Sampling Weights



Conclusion

- Presented experimental method for using AI to prebunk false election rumors and conspiracies
- Labor intensive, Human-in-the-Loop written articles perform similarly as those written purely by AI
- Prebunks are durably effective at reducing belief in specific rumors
- Prebunks are temporarily effective at increasing confidence in elections
- More work is needed to bolster long-term election confidence

Next Steps

- Are more “intensive” interventions more effective?
- Over the 2024 election: ran experiment comparing AI-powered chatbot conversations and AI-written inoculation articles
- Currently analyzing results
- Try the bot out at: <https://electionbot.chat>