

# Human Detection of Political Speech Deepfakes across Transcripts, Audio, and Video

Matthew Groh<sup>1,†,\*</sup>, Aruna Sankaranarayanan<sup>2,3,†</sup>, Nikhil Singh<sup>2</sup>, Dong Young Kim<sup>2</sup>, Andrew Lippman<sup>2</sup>, and Rosalind Picard<sup>2</sup>

<sup>1</sup>Kellogg School of Management, Northwestern University, Evanston, IL, USA

<sup>2</sup>Media Lab, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>3</sup>CSAIL, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>†</sup>These authors contributed equally.

\*To whom correspondence should be addressed. E-mail: matthew.groh@kellogg.northwestern.edu

## ABSTRACT

Recent advances in technology for hyper-realistic visual and audio effects provoke the concern that deepfake videos of political speeches will soon be indistinguishable from authentic video recordings. The conventional wisdom in communication theory predicts people will fall for fake news more often when the same version of a story is presented as a video versus text. We conduct 5 pre-registered randomized experiments with 2,215 participants to evaluate how accurately humans distinguish real political speeches from fabrications across base rates of misinformation, audio sources, question framings, and media modalities. We find base rates of misinformation minimally influence discernment and deepfakes with audio produced by the state-of-the-art text-to-speech algorithms are harder to discern than the same deepfakes with voice actor audio. Moreover across all experiments, we find audio and visual information enables more accurate discernment than text alone: human discernment relies more on how something is said, the audio-visual cues, than what is said, the speech content.

## Introduction

Recent advances in technology for algorithmically applying hyper-realistic manipulations to video are simultaneously enabling new forms of interpersonal communication and posing a threat to traditional standards of evidence and trust in media<sup>1–8</sup>. In the last few years, computer scientists have trained machine learning models to generate photorealistic images of people who do not exist<sup>9–11</sup>, inpaint people out of images<sup>12,13</sup>, clone voices based on a few samples<sup>14,15</sup>, modulate the lip movements of people in videos to make them appear to say something they have not said<sup>16,17</sup>, and create synthetic videos based on simple text prompts<sup>18</sup>. The synthetic videos’ false appearance of indexicality – the presence of a direct relationship between the photographed scene and reality<sup>19,20</sup> – has the potential to lead people to believe video-based messages that they otherwise would not have believed if the messages were communicated via text. This potential influence is particularly concerning because research demonstrates that videos, especially videos of an injustice, elicit more engagement and emotional reactions (e.g., anger, sympathy) than text descriptions displaying the same information<sup>21–23</sup> (although, see ref.<sup>24</sup>). Moreover, visual misinformation is common on social media<sup>25</sup> and the emotional and motivational influences of visual communication have been attributed to why misleading viral videos have provoked mob-violence<sup>26,27</sup>. While people are more likely to believe a real event occurred after watching a video of the event than reading a description of the event<sup>28</sup>, an open question remains: Does visual communication relative to text or audio increase the believability of *fabricated* events?

The realism heuristic<sup>27,29</sup> predicts “people are more likely to trust audiovisual modality [relative to text] because its content has a higher resemblance to the real world.” This prediction is relevant for many deepfake videos<sup>30</sup> and suggests fabricated video would be more believable than fabricated text conditional on the absence of obvious perceptual distortions. Yet there exists little direct empirical evidence for this heuristic applied to algorithmically manipulated video. In an experiment using 3 false news videos as stimuli, researchers found that stories presented as videos are perceived as more credible than stories presented as text or read aloud in audio form<sup>27</sup>. In contrast, in an experiment by Barari et al 2021 showing 6 political deepfake videos (videos manipulated by artificial intelligence to make someone say something they did not say) and 9 non-manipulated videos, researchers did not find differences between truth discernment rates in video, audio, and text<sup>31</sup>. Perhaps some of the participants did not take the videos’ “indexicality” as evidence of authenticity because participants were aware of how easily such videos could be manipulated. Alternatively, some participants may have noticed perceptual distortions or indicators of satire (e.g. facial expressions and comedic timing) in the videos used in Barari et al 2021, which would naturally lead one to believe the video has been manipulated. In experiments examining how people react to deepfake videos of politicians, researchers find

people are not more likely to report false memories after watching deepfake videos than reading the same false news as text<sup>32</sup>, people are more likely to feel uncertain than misled after viewing a deepfake of Barack Obama<sup>33</sup>, people consider a deepfake of a Dutch politician significantly less credible than the real video from which it was adapted<sup>34</sup> and the previously mentioned deepfake video of the Dutch politician is not more persuasive than the text alone<sup>35</sup>. In the experiment examining the deepfake of the Dutch politician, some respondents explained their credibility judgments by indicating audio-visual cues of how the message was communicated (e.g., unnatural mouth movements); others indicated inconsistency in the content of the message itself (e.g., contextually unrealistic speeches)<sup>34</sup>. One explanation for the current mixed evidence on the role of communication modalities in mediating people's ability to discern fabricated content is the large possibility space of how political speeches may appear in videos and the small number of stimuli in media effects research<sup>36</sup>.

Related research demonstrates how fake images can be persuasive and difficult to distinguish from real images. People rarely question the authenticity of images even when primed<sup>37</sup>. Images can increase the credibility of disinformation<sup>38</sup>. Images of synthetic faces produced by StyleGAN2<sup>10</sup> are indistinguishable to research participants from the original photos on which the StyleGAN2 algorithm was trained<sup>39</sup>. Moreover, research shows that non-probative and uninformative photos can lead people to believe false claims<sup>40</sup>, lead people to believe they know more than they actually know<sup>41</sup>, and promote "truthiness" by creating illusory truth effects<sup>42,43</sup>, which can lead people to believe falsehoods they previously knew to be falsehoods<sup>44,45</sup>. When it comes to ostensibly probative videos of political speeches, the question of whether people are more likely to believe an event occurred because they saw it as opposed to only read about it remains open.

In fact, today's algorithmically generated deepfakes are not yet consistently indistinguishable from real videos. On a sample of 166 videos from the largest publicly available dataset of deepfake videos to date<sup>46</sup>, people are significantly better than chance but far from perfect at discerning whether an unknown actor's face has been visually manipulated by a deepfake algorithm<sup>47</sup>. This finding is significant because it demonstrates that people can identify deepfake videos from real videos based solely on visual cues. However, some videos are more difficult than others to distinguish because of their blurry, dark, or grainy visual features. On a subset of 11 of the 166 videos, Kobis et al 2021 do not find that people can detect deepfakes better than chance<sup>48</sup>.

In another experiment with 25 deepfake videos and 4 real videos but only 94 participants, researchers found that the overall discernment accuracy is 51% and a media literacy training increases discernment accuracy by 24 percentage points for participants assigned to the training relative to the control group<sup>49</sup>.

People's capacity to identify multimedia manipulations raises questions: how do various kinds of fabricated media (e.g., synthesized audio and video of political speeches that never happened) alter the perceived credibility of misinformation, how do audience characteristics (e.g., reflective reasoning) moderate media effects, and how does the source and content of a message interact with the fabricated media and audience characteristics<sup>50</sup>? A growing field of misinformation science is beginning to address these questions. Research on news source quality demonstrates that people in the United States are generally accurate at identifying high and low-quality publishers<sup>51</sup> and the salience of source information does not appear to change how accurately people identify fabricated news stories<sup>52</sup>, manipulated images<sup>53</sup>, or false news headlines<sup>54,55</sup> although evidence on false news headlines is mixed<sup>56,57</sup>. Research on political false news content suggests an individual's tendency to rely on intuition instead of analytic thinking is a stronger factor than motivated reasoning in explaining why people fall for false news<sup>58</sup>, and similarly, people with more analytic cognitive styles worldwide are more accurate at discerning between authentic and fabricated political videos<sup>59</sup> and true and false headlines related to COVID-19<sup>60</sup>. In fact, people tend to be better at discerning truth from falsehood when evaluating news headlines that are concordant with their political partisanship relative to when evaluating news headlines that are discordant<sup>61</sup>. While the science of misinformation has generally focused on the messengers (the source credibility of publishers)<sup>62</sup> and the message of what is said (the media credibility of written articles and headlines)<sup>61</sup>, the relevance of audio-visual communication channels to the psychology of misinformation has received less attention<sup>63</sup> and is important for addressing the problem of misinformation<sup>64</sup>.

In this paper, we conduct 5 pre-registered experiments to evaluate how well people can distinguish between real and fabricated political speeches by well-known politicians and how communication modalities, contexts, audio sources, and base rates of fabricated content influence discernment. The stimuli include 32 videos from the Presidential Deepfake Dataset<sup>65</sup> and 12 additional videos<sup>31</sup> (see Figure 1 for a screenshot from each video). In total, we analyze data from 2,215 recruited participants in these 5 pre-registered experiments and an additional 41,313 non-recruited participants who participated in Experiment 1 but were not pre-registered.

We begin with Experiment 1, which addresses the question: How does media modality influence non-recruited participants' ability to discern real and fabricated political speeches? In Experiment 1, we present participants with 32 political speeches – half real and half fabricated – by Donald Trump and Joseph Biden that are randomized to be displayed via the 7 possible permutations of text, audio, and video: a transcript, an audio clip, a silent video, audio with subtitles, silent video with subtitles, video with audio, and video with audio and subtitles. By randomly assigning political speeches to these permutations of text, audio, and video modalities and asking participants to discern truth from falsehood (see Methods section for exact question wording and Figure 9 for a screenshot from the digital experiment), this experiment is designed to disentangle the degree to



**Figure 1. Screenshots of the Experiments’ Stimuli** Frames from the 10th second in the 32 videos from the Presidential Deepfakes Dataset (PDD) and the 12 other videos used in Barari et al 2021. The left column shows the real videos, and the right column shows deepfakes. The top 4 rows show frames from the PDD videos and the bottom 2 rows show frames from the other videos.

which participants attend to and consider the content of what is said and the audio-visual cues as to how it is said.

Experiment 2 builds upon Experiment 1 by enhancing and extending the stimuli and adapting the wording of the experiment. In experiment 2, we present participants with 20 videos randomly sampled from 60 videos of politicians, which include 12 videos used in Barari et al 2021<sup>31</sup> and videos of the same 32 political speeches from Experiment 1 where the 16 real videos are the same and the deepfakes are enhanced with the state-of-the-art algorithms in 2023<sup>66</sup> and include 16 deepfakes with voice actor audio and 16 deepfakes with audio produced by a text-to-speech algorithm fine-tuned on the presidents’ voices<sup>67</sup>. Experiment 2 offers an empirical investigation into human discernment of deepfakes videos with different sources of audio (audio from a voice actor or text-to-speech algorithm) and different contexts of videos (non-satirical presidential speeches in the Presidential Deepfake Dataset (PDD) videos and satirical speeches and explicit discussions of synthetic media in the other videos used in Barari et al 2021).

In Experiment 3, we examine how the base rate of fabrications influences participant accuracy by randomizing participants to a low or high base rate of fabricated political speeches. We present participants with 20 political speeches by Donald Trump and Joseph Biden that are randomized to appear as a transcript, a silent video, an audio clip, or video with audio. Experiment 3 provides a conceptual replication of experiment 1 with enhanced stimuli and an opportunity to evaluate the influence of base rates of misinformation.

In Experiment 4, we present participants with the 16 videos or audio clips of 16 real political speeches by Donald Trump and Joseph Biden and the same 16 real political speeches with audio produced by a voice actor; we ask participants if they can identify which stimuli are voiced by the US presidents and which are voiced by the voice actor. Experiment 4 offers an opportunity to evaluate how accurately participants can distinguish Donald Trump and Joseph Biden’s voice from a voice actor’s voice.

Finally, in Experiment 5, we present participants with 10 videos randomly sampled from the same 32 videos in Experiment 3, but we do not prime people with a direct question of authenticity. In contrast to the previous experiments, we ask participants “What comes to mind after watching the following video/listening to the following audio/reading the following quote?” This final experiment reveals how even when participants are not necessarily paying attention to authenticity they reveal suspicions of fabrications differently across media modalities.

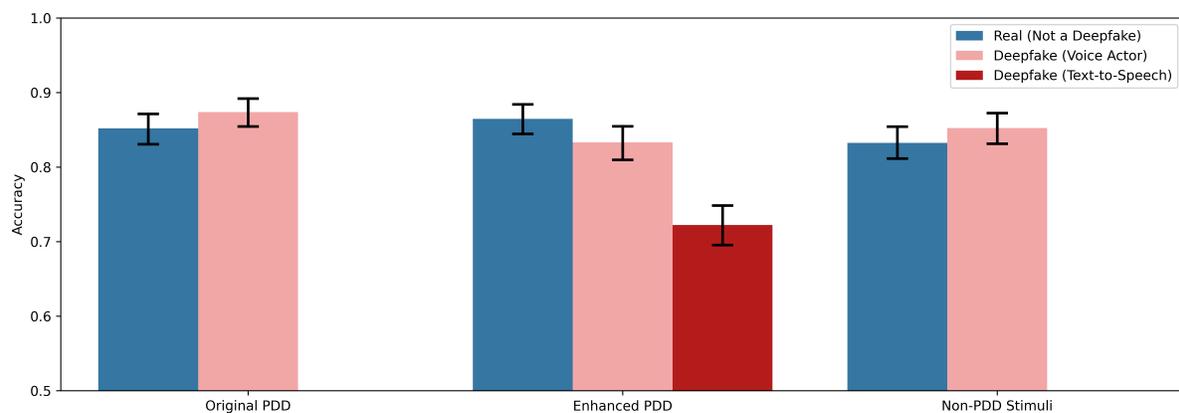
## Results

Experiments 1a, 2, 3, 4, and 5 involve participants recruited from Prolific and are pre-registered on aspredicted.org at the following URLs: [1a](#), [2](#), [3](#), [4](#), and [5](#). Experiment 1b is not pre-registered but includes 41,313 participants who discovered the experiment organically through search engines or the news. In the Methods section, we provide details on participants, the digital experiment interface, experimental stimuli, and randomization protocol. Throughout this paper, accuracy refers to human participants not machine learning models unless otherwise noted.

Experiment	Prereg	Feedback	Known Rate	Stimuli	Obs	Modalities	Rate of Fabrications
Experiment 1a	Yes	Yes	Yes	32 original PDD	32	Text, audio, and video	50%
Experiment 1b	No	Yes	Yes	32 original PDD	32	Text, audio, and video	50%
Experiment 2	Yes	No	No	48 enhanced PDD and 12 other videos	20	Video	60%
Experiment 3	Yes	No	No	32 enhanced PDD (TTS deepfakes)	20	Text, audio, and video	20% or 80%
Experiment 4	Yes	No	No	16 real PDD videos	16	Audio and video	50%
Experiment 5	Yes	No	No	32 enhanced PDD (TTS deepfakes)	10	Text, audio, and video	50%

**Table 1. Overview of the five experiments** “Prereg” indicates whether the experiment was pre-registered, “Feedback” indicates whether we give participants immediate feedback on whether a stimulus is fabricated or not, “Known Rate” indicates whether we informed participants the rate of fabrications, “Stimuli” refers to the stimuli used, “Obs” refers to the maximum number of observations provided by each participant, “Modalities” indicates the possible modalities in which the stimuli are presented, and “Rate of Fabrications” refers to the base rate of fabricated speeches, which was randomized in Experiment 3.

Figure 2 presents accuracy by participants in Experiment 1a and Experiment 2 on real and fabricated videos with audio from the Presidential Deepfakes Data (PDD) videos and 12 other videos previously examined in Barari et al 2021<sup>31</sup>. In Experiment 1a, participants correctly identified real PDD videos and deepfakes in 85% and 87% observations, respectively. In Experiment 2, participants correctly identified real PDD videos, the enhanced PDD voice actor deepfakes, enhanced PDD text-to-speech deepfakes, real other videos, and other deepfakes in 86%, 83%, 72%, 85%, and 83% of observations. As a baseline for comparison, random guessing on this task would lead to 50% accuracy. Participants are closer to random guessing than a perfect score on the enhanced PDD text-to-speech deepfakes but closer to a perfect score on the rest of the stimuli.

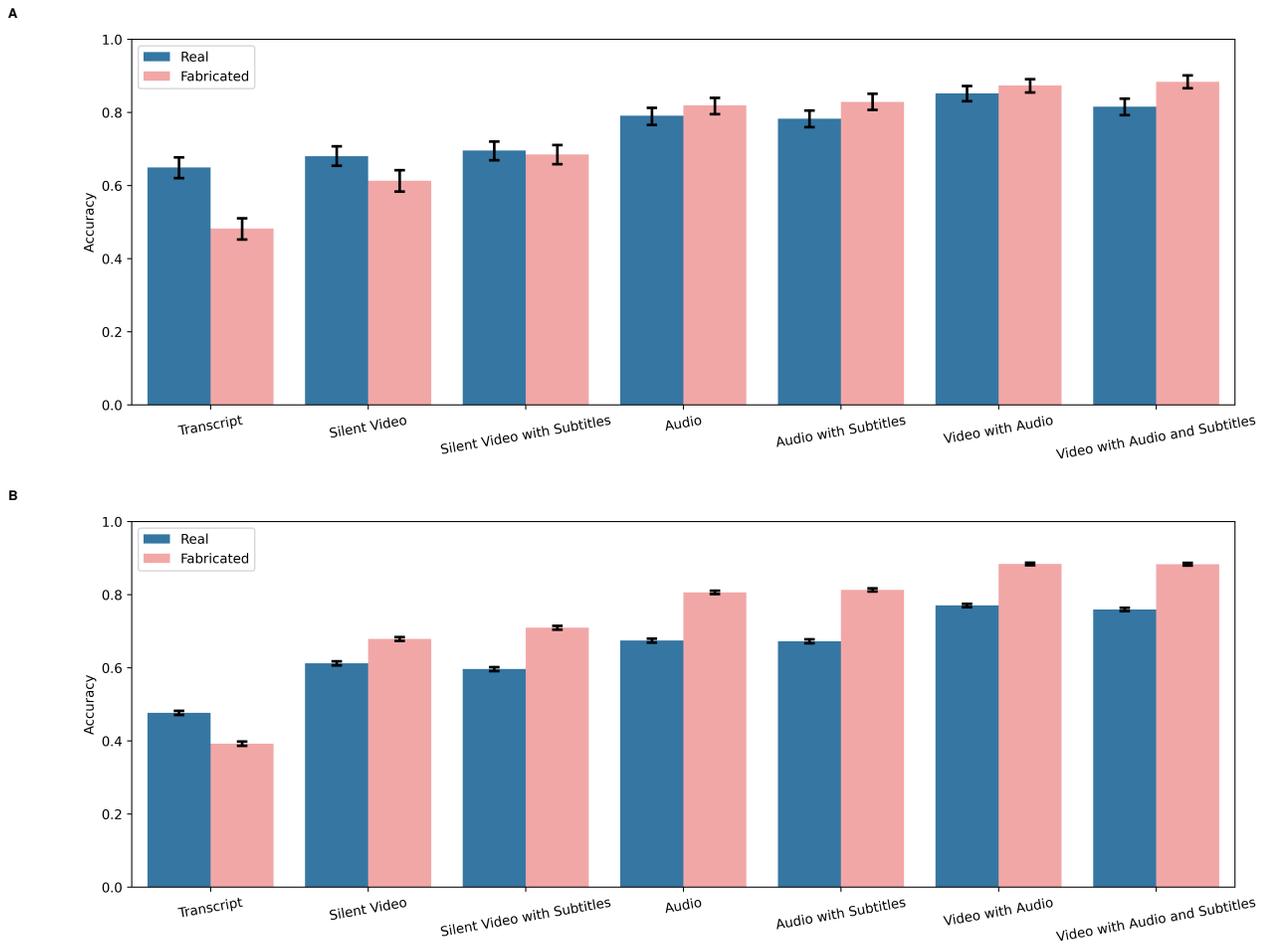


**Figure 2. Accuracy Distinguishing Real and Fabricated Speeches across Video Stimuli** Accuracy across the original Presidential Deepfakes Dataset (PDD) video stimuli in Experiment 1a, the enhanced PDD video stimuli in Experiment 2, and the non-PDD video stimuli in Experiment 2. The error bars represent 95% confidence intervals.

### Experiment 1a (501 participants, preregistered)

We designed Experiment 1a to address the following question: How does media modality influence participants’ ability to discern real and fabricated political speeches? In order to answer this question, we show participants 32 videos from the original Presidential Deepfakes Dataset (PDD), inform participants that half are real and half are fake, and ask participants to indicate their level of confidence that the stimulus is a fabricated political speech or not. After each response, we inform participants of whether the stimulus was real or fabricated.

We find that participants’ accuracy increases as they have access to additional communication modalities. In particular, accuracy by modality from lowest to highest starts with transcripts at 57% accuracy followed by silent videos without subtitles at 64% accuracy, silent videos with subtitles at 69% accuracy, audio (with and without subtitles) at 81% accuracy, video with audio and subtitles at 85% accuracy, and video with audio but no subtitles at 86% accuracy. Figure 3 shows accuracy across modalities for real and fabricated stimuli for Experiments 1a and 1b. In Table 2, we present the pre-registered regression analysis on “confidence score” which is a measure of accuracy weighted by participants’ confidence defined as the participant’s confidence (ranging from 50 to 100) if correct and 100 minus the participant’s confidence if incorrect. Based on ordinary least squares regressions with standard robust errors clustered at the participant level following Abadie et al (2017)<sup>68</sup>, we find results mirror accuracy by modality where transcripts have the lowest confidence score of 58% ( $p < 0.001$ ), silent videos are 7 percentage



**Figure 3. Accuracy Distinguishing Real and Fabricated Speeches across Media Modalities in Experiment 1** A. Accuracy across all permutations of text, audio, and video in Experiment 1a with 501 recruited participants. B. Accuracy across all permutations of text, audio, and video in Experiment 1b with 41,313 non-recruited participants. The error bars represent 95% confidence intervals.

points higher ( $p < 0.001$ ), silent videos with subtitles are 9 points higher ( $p < 0.001$ ), audio (with and without subtitles) is 19 points higher ( $p < 0.001$ ), and video with audio (with and without subtitles) is 25 points higher ( $p < 0.001$ ). In columns 2 and 3 of Figure 2, we present results for real and fabricated speeches by themselves, which shows that additional media modalities help participants identify fabrications as fabrications even more than additional media modalities help participants identify real speeches as not fabricated.

As a secondary analysis, we find the participants' accuracy increase by 2.1 percentage points ( $p < 0.001$ ) for each of the three Cognitive Reflection Test (CRT)<sup>69</sup> questions they answered correctly. In addition, we find the participants' accuracy increases by 0.12 percentage points ( $p < 0.001$ ) for every stimulus seen.

### Experiment 1b (41,313 participants, not preregistered)

Experiment 1b presents a robustness check of Experiment 1a and is identical to Experiment 1a except Experiment 1b is has two orders of magnitude more participants and is not pre-registered. The results of Experiment 1b directionally corroborate results from Experiment 1a. Specifically, accuracy by modality from lowest to highest starts with transcripts at 43% accuracy followed by silent videos (with and without subtitles) at 65% accuracy, audio (with and without subtitles) at 74% accuracy, video with audio and subtitles at 82% accuracy, and video with audio but no subtitles at 83% accuracy.

### Experiment 2 (302 participants, preregistered)

In order to evaluate the generalizability of the results from Experiment 1, we designed and curated an enhanced and extended set of stimuli. In addition, we adapted Experiment 2 such that participants are not informed about the base rate of deepfakes, participants are not informed of whether stimuli are real or fabricated until the end of the experiment when we debrief participants, and we slightly adjust the experimental interface (see Figure 10). These changes in Experiment 2 allow us to

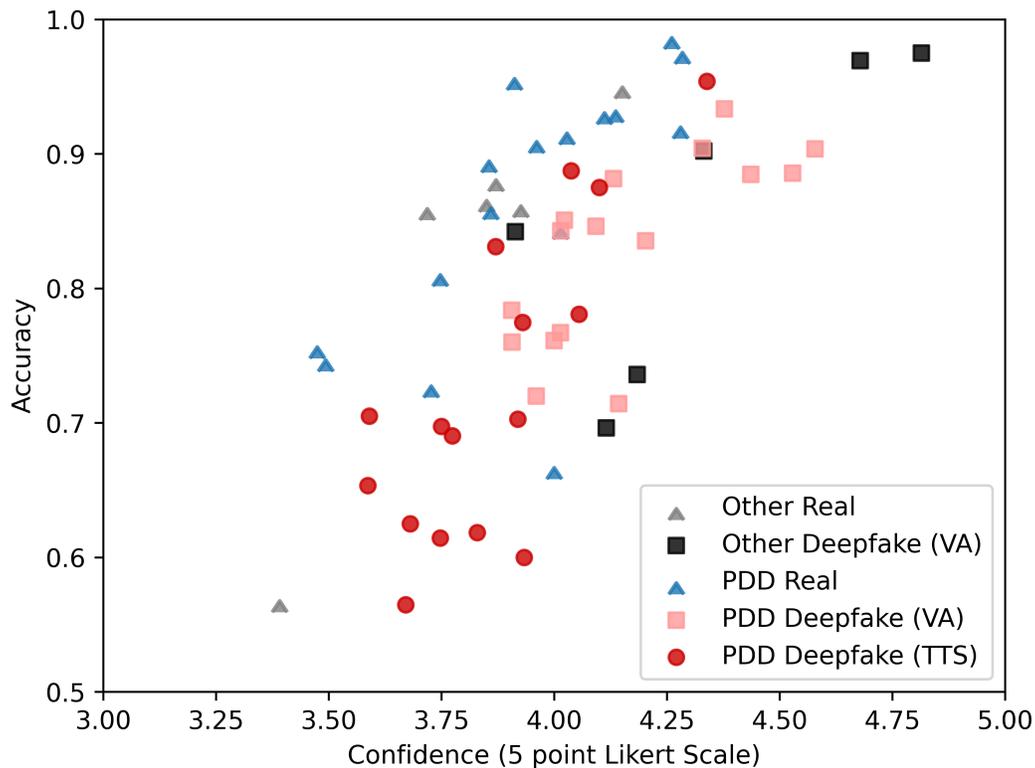
address the following question: How does manipulation methodologies and context influence participants’ ability to discern real and fabricated political speeches on an enhanced and extended set of stimuli?

In Experiment 2, we show participants 20 videos randomly sampled from the 60 enhanced PDD videos and videos used in Barari et al 2021 and ask participants whether they think the speech is fabricated, and how confident they are in their judgment. Each participant sees 4 real videos from the PDD, 4 voice actor deepfakes, 4 text-to-speech deepfakes, 4 real videos used in Barari et al 2021, and 4 voice actor deepfakes used in Barari et al 2021.

We find participants are similarly accurate at identifying real videos as real and voice actor deepfakes as fabricated, but participants’ accuracy is significantly lower on text-to-speech deepfakes than voice actor deepfakes. In Table 3, we present the pre-registered ordinary least square regressions on accuracy, which is a binary variable defined as 1 if participants accurately identify the stimulus and 0 otherwise. Specifically, we do not find a difference between accuracy on Barari et al 2021 deepfakes and enhanced PDD voice actor deepfakes ( $p = 0.685$ ), real PDD videos ( $p = 0.804$ ), or real videos used in Barari et al 2021 ( $p = 0.763$ ). However, we find that participants’ accuracy on text-to-speech deepfakes is 13 percentage points lower than their accuracy on the deepfakes used in Barari et al 2021 ( $p = 0.013$ ).

In a series of 16 pre-registered t-tests comparing PDD voice actor deepfakes with their text-to-speech counterparts, we find the accuracy on 5 out of 16 deepfakes is statistically significant and lower on text-to-speech videos than voice actor videos when controlling the false discovery rate using the Benjamini-Hochberg procedure<sup>70</sup>. Table 4 presents the accuracy rates across each of the 16 enhanced deepfakes from the two audio sources alongside p-values from the t-tests and the number of observations for each video.

Figure 4 presents the distribution of accuracy and confidence (as reported on a 5 point Likert scale) for each of the 60 videos in Experiment 2. In particular, this scatterplot demonstrates relatively high variance in accuracy across contexts and shows accuracy and confidence are positively correlated (the Pearson Correlation Coefficient between them is 0.66 ( $p < 0.001$ )). The real video with the lowest accuracy is a hot-mic of Obama speaking with Dimitri Medvedev<sup>71</sup>, and the deepfake with the lowest accuracy is Donald Trump (as voiced by a text-to-speech algorithm) speaking about his phenomenal respect for women. The deepfakes with the highest accuracy (97%) are the two deepfakes of Donald Trump used in Barari et al 2021.



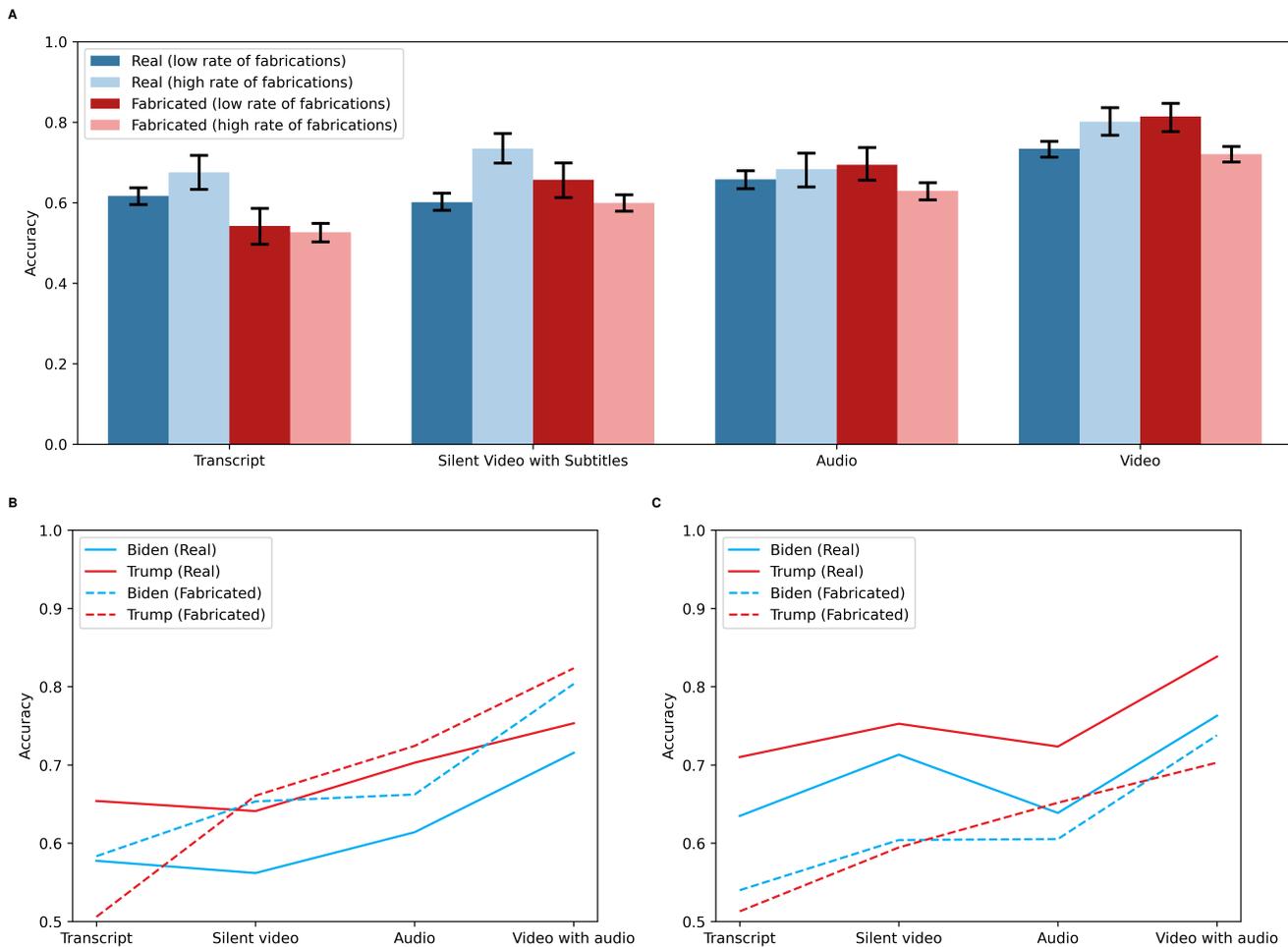
**Figure 4. Average Accuracy and Confidence for All Video Stimuli in Experiment 2** Scatter plot showing participants’ mean accuracy and confidence on each of the 60 videos in experiment 2. “PDD” indicates videos from the enhanced Presidential Deepfake Dataset and “Other” indicates videos are drawn from the same sample used in Barari et al 2021. “VA” indicates voice actor deepfakes and “TTS” indicates text-to-speech deepfakes.

We present the pre-registered secondary analysis in Table 5 where we examine correct confidence (a weighted measure of accuracy defined as  $1-(5-\text{confidence})/5$  if correct and  $-(\text{confidence})/5$  if incorrect), response time, and a binary variable for toggling the play/pause button. We find the results on correct confidence corroborate the main analysis on accuracy, and we do not find any effects of stimuli context on response time or toggling the play/pause button.

We do not find the participants' accuracy changes over the course of the number of videos watched ( $p = 0.828$ ).

### Experiment 3 (1,006 participants, preregistered)

In order to further evaluate the generalizability of the results in Experiment 1 and how base rates of misinformation may influence these results, we conduct Experiment 3 following the protocol in Experiment 2 where we show participants 20 political speeches randomly sampled from the 32 PDD political speeches and ask participants whether they think the speech is fabricated and how confident they are in their judgment. The deepfakes are all enhanced text-to-speech deepfakes from the PDD. We randomize participants to high and low base rate conditions where participants either see 16 or 4 fabricated speeches, respectively. Just like Experiment 2, we do not inform participants of the base rate of fabricated speeches and we do not inform participants of whether stimuli are real or fabricated until the end of the experiment when we debrief participants.



**Figure 5. Accuracy Distinguishing Real and Fabricated Speeches across Media Modalities and Base Rates in Experiment 3** A. Accuracy across all permutations of text, audio, and video and high and low-base rate conditions in Experiment 3. The error bars represent 95% confidence intervals. B. Low base rate condition in Experiment 3: Accuracy across on the four modalities. C. High base rate condition in Experiment 3: Accuracy across on the four modalities.

We find that the base rate of deepfakes has minimal influence on participants' overall accuracy and accuracy increases as participants have access to additional communication modalities. In Table 6, we present the pre-registered ordinary least square regressions on accuracy, which is a binary variable defined as 1 if participants accurately identify the stimulus and 0 otherwise. In columns 2 and 3 of Table 6, we find the high base rate of fakes leads to a 7.2 percentage point higher accuracy on real stimuli ( $p < 0.001$ ) and 5.8 percentage point lower accuracy on fabricated stimuli ( $p < 0.001$ ). When considering interactions in columns 4-6, we do not find silent videos with subtitles increase participants' overall accuracy beyond accuracy on transcripts ( $p = 0.532$ ), but we find silent videos with subtitles increase participants' accuracy on fabricated speeches by 11.4 percentage

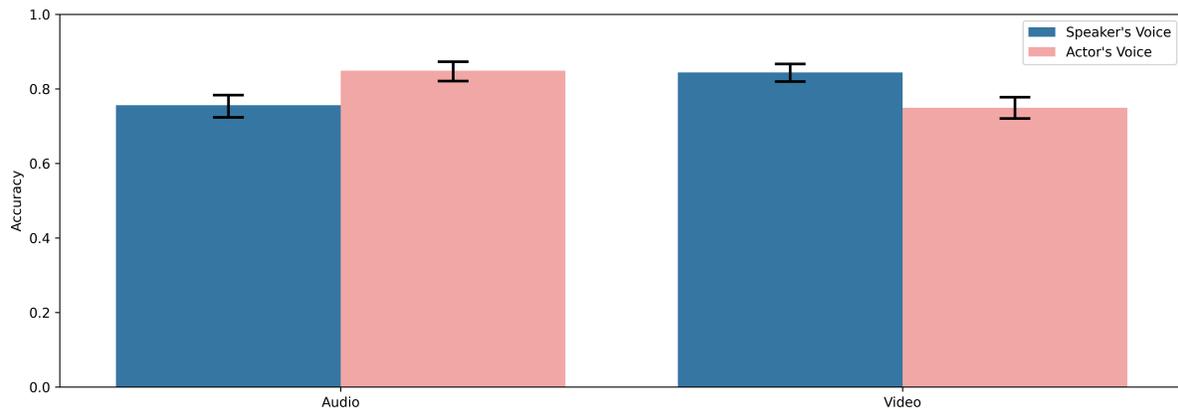
points ( $p = 0.006$ ). We find audio increases participants' overall accuracy by 6.3 percentage points ( $p = 0.003$ ) and specifically increases accuracy on real speeches by 4.1 percentage points ( $p = 0.072$ ) and fabricated speeches by 15.2 percentage points ( $p = 0.002$ ). We find the largest impact on accuracy from video with audio which increases overall accuracy by 14.8 percentage points ( $p < 0.001$ ) and specifically increases accuracy on real speeches by 11.7 percentage points ( $p < 0.001$ ) and fabricated speeches by 27.2 percentage points ( $p < 0.001$ ). We find minimal effects of interactions between modalities and base rate assignment. Figure 5 shows accuracy across base rates and modalities for real and fabricated stimuli for Experiment 3.

In order to evaluate the robustness of media effects on individual speeches, we conduct 6 families of comparisons of 32 pre-registered t-tests comparing accuracy in one modality to accuracy in another modality. Tables 7, 8, 9, 10, 11, and 12 present the accuracy rates across each of the 32 political speeches alongside p-values from the t-tests between modalities and the number of observations for each political speech in each modality. We find the accuracy on 6, 9, and 21 out of 32 political speeches are statistically significant and lower on transcripts than silent videos, audio, and video with audio, respectively, when controlling the false discovery rate using the Benjamini-Hochberg procedure<sup>70</sup>. Likewise, we find the accuracy on 4 and 19 out of 32 political speeches are statistically significant and lower on silent video than audio and video with audio. Finally, the accuracy on 12 out of 32 political speeches is statistically significant and lower on audio than video with audio.

We present the pre-registered secondary analysis in Table 13 where we examine correct confidence, response time, and a binary variable for toggling the play/pause button. We find the results on correct confidence corroborate the main analysis on accuracy. We do not find any effects of the base rate condition on response time or toggling the play/pause button ( $p = 0.381$  and  $p = 0.700$ , respectively). We find silent video takes participants an additional 6.7 seconds when viewing silent video relative to audio ( $p < 0.001$ ), and we find silent video and video with audio leads participants to toggle the play and pause 15.1 and 3.0 percentage points ( $p < 0.001$  and  $p < 0.001$ ) more than audio, respectively.

We do not find that participants' accuracy changes over the course of the number of stimuli seen in either condition ( $p = 0.231$ ) or high base rate condition by itself ( $p = 0.881$ ).

#### Experiment 4 (206 participants, preregistered)



**Figure 6. Accuracy Distinguishing Speakers' and Actors' Voices across Video Stimuli in Experiment 4** Accuracy across all audio and video in Experiment 4. The error bars represent 95% confidence intervals.

In order to further identify the role of manipulated audio in participants' ability to distinguish between real and fabricated content across the previous experiments, we designed Experiment 4 to address the following question: How does media modality influence recruited participants ability to distinguish between a well-known speaker's real voice and an actor's voice? In Experiment 4, we show participants 16 real PDD political speeches and ask participants whether they think the voice is the speaker's or a voice actor, and how confident they are in their judgment. By focusing on only real videos, we examine the role of perfectly realistic visual information to influence accuracy. Just like Experiment 2 and 3, we do not inform participants of the base rate of voice actor audio and we do not inform participants of whether stimuli are the actual speakers' or voice actor's voice until the end of the experiment when we debrief participants.

We find that that participants are more accurate at identifying voice actors' audio than real speakers' audio, more accurate on real speakers' video than real speakers' video, but less accurate on voice actor video with audio than voice actor audio by itself. In Figure 6 and Table 14, we present the pre-registered ordinary least square regression analysis on accuracy. Specifically, participants are 75.6% accurate on audio by the actual speakers and we find voice actor audio increases accuracy by 9.2 percentage points ( $p = 0.006$ ), video with audio increases accuracy by 8.8 percentage points ( $p = 0.004$ ), but the combination of video with audio by voice actors lowers accuracy by 18.7 ( $p < 0.001$ ) percentage points such that it is the same level of

accuracy as participants obtain on audio (without video) by the actual speakers. Despite identical base rates of voice actor audio in the video and audio conditions and similar overall accuracy rates of 80% in both video and audio conditions, participants were biased to identifying audio stimuli as an actor's voice in 55% of observations and identifying video stimuli as an actor's voice in 45% of observations. This suggests the voice actor's audio consistently matched with the real video footage and audio with video increased participants' beliefs that all audio (both real and fake) was authentic relative to when participants were listening to audio only.

We present the pre-registered secondary analysis in Table 13 where we examine correct confidence, response time, and a binary variable for toggling the play/pause button. We find the results on correct confidence corroborate the main analysis on accuracy. We do not find effects on toggling the play/pause button, but we find that voice actor audio leads participants to take an additional 1.5 seconds ( $p < 0.001$ ) beyond what they take for audio by the real speakers.

For each additional stimulus seen, we find the participants' accuracy increases by 0.3 percentage points ( $p = 0.011$ ); participants slightly improve based on seeing and hearing real and fabricated stimuli.

### Experiment 5 (200 participants, preregistered)

In contrast to Experiments 1 through 4 where participants are explicitly asked to consider the veracity of the stimulus, Experiment 5 is designed not to alert participants to the dependent variable of interest. Instead in Experiment 5, participants are asked, "What comes to mind after watching the following video/listening to the following audio/reading the following quote?" Instead of the custom website at detectfakes.media.mit.edu used in the previous experiments, Experiment 5 is hosted on Qualtrics and we changed the initial instructions as follows: "This is an MIT research project. You will be shown quotes, audio, and video files that one might expect to see on social or digital media. You will be requested to share your thoughts and opinions after reading, listening, or viewing each of these media files." Just like Experiments 2 through 4, we do not inform participants of the base rate of voice actor audio and we do not inform participants of whether stimuli are real or fabricated until the end of the experiment when we debrief participants.

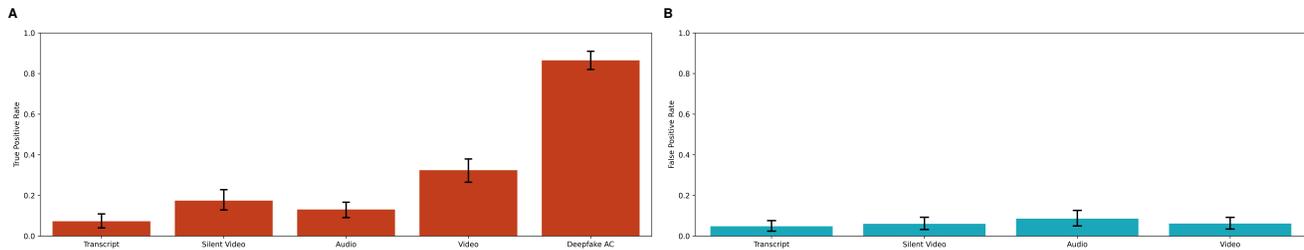
In Experiment 5, each participant sees a forced choice attention check, the deepfake attention check, and 10 randomly sampled stimuli from Experiment 3 (32 real and fabricated speeches from the PDD where the 16 deepfakes are the text-to-speech deepfakes). In order to evaluate the effect of recent exposure to a deepfake, we randomly assign participants to see the deepfake attention check video at the start or end of the experiment. Additional details about the design of Experiment 5 and how free text responses are annotated into a binary variable for suspicion of a fabrication are shared in the Methods section.

Similar to Experiment 1, we find that participants' accurate suspicion of a fabrication increases as participants have access to additional communication modalities. Further, we find that the false positive rate (inaccurate suspicion of a fabrication) is not associated with communication modalities. Table 18 presents pre-registered regression results for Experiment 5 and column 2 of Table 18 shows that relative to text transcripts, silent video increases the true positive rate by 10.3 percentage points ( $p = 0.001$ ), audio increases the true positive rate by 6.3 percentage points ( $p = 0.022$ ), video increases the true positive rate by 25 percentage points ( $p < 0.001$ ), and priming (assignment to the attention check deepfake displayed as the first stimulus as opposed to the last) increases the true positive rate by 7.4 percentage points ( $p = 0.057$ ). In contrast, the same regression on the false positive rate does not have significant values on modality conditions (the p-values for silent video, audio, and video are 0.540, 0.194, and 0.621, respectively) but increase in the true positive rate after priming is statistically significant and increases the false positive rate by 6.3 percentage points ( $p = 0.001$ ). In Figure 7, we present the true positive and false positive rates for suspicion of fabrications in the experiment. In 32.4% of observations of deepfake videos with audio, participants' responses were judged to be suspicions of fabrications whereas the rate of suspicions was only 6 percentage points in real videos with audio. Audio, by itself, is much less likely to reveal suspicions: 13.0% of observations of text-to-speech audio are suspected whereas 8.6% of observations of real audio are suspected. Finally, transcripts are the least likely modality to reveal suspicions; 7.3 percent of observations of fabricated transcripts show suspicions compared to 4.8 percent of observations of real transcripts.

We note that the vast majority of participants' responses (86.5%) to the obvious, attention check deepfake indicate participants suspected the video included a fabrication. Additionally, we note we did not pre-register heterogeneous analyses based on age but we find strong evidence that accurate suspicion of a deepfake is correlated with age. In particular, the true positive rate across participants' ages ranges from 31% for 32 people in their 20s, 43% for 54 people in their 30s, 31% for 48 people in their 40s, 25% for 24 people in their 50s, and 20% for people in their 60s and beyond.

## Discussion

This paper provides evidence, via multiple pre-registered randomized experiments with 2,215 participants that visual and auditory communication modalities increase people's ability to distinguish authentic political speeches from fabricated political speeches. In the context of **authentic** speeches, we provide corroborating evidence for the conventional wisdom around the "seeing is believing" narrative (the realism heuristic that suggests people will tend to trust video over text<sup>29</sup> and results from Wittenberg et al 2021 showing people "are more likely to believe an event occurred when it is presented in video versus textual



**Figure 7. True Positive and False Positive Rate in Suspicion of Fabrications in Experiment 5** A. True positive rate (accurate suspicion of a fabrication) across modalities and the attention check deepfake where a man’s face and glasses transform in front of the viewer. B. False positive rate (inaccurate suspicion of a fabrication) across modalities. The error bars represent 95% confidence intervals.

form<sup>28</sup>); people are significantly more accurate at identifying authentic speeches as authentic when the speeches include audio and visual modalities as opposed to only text (although note Wittenberg et al 2021 finds minimal effects of video on persuasiveness). However with respect to **fabricated** content, the results from Experiments 1-3 and 5 add considerable nuance to the seeing is believing narrative: people are significantly more accurate at identifying fabricated speeches as fabricated when the speeches include audio and visual modalities as opposed to only text. In other words, we find participants are significantly more accurate at distinguishing between authentic and fabricated political videos than transcripts. Moreover, Experiment 5 demonstrates that this continues to be the case even when people are not directly asked about the authenticity of a speech or primed to consider accuracy.

These results are based on an experiment with a stimuli set that is much larger than most stimuli sets for the psychology of media effects research<sup>36</sup> and deepfake detection<sup>31,33,34</sup>, but it is important to add a caveat that we focused on a single context, political speeches, and a combination of algorithms, the deepfake lip-syncing wav2lip algorithm and the DeepFaceLab library, which are very effective at manipulating a person who is facing forward and already speaking into a convincing deepfake video. While we present evidence that adds considerable nuance to the media effects literature on communication modalities, future work may consider additional nuances by exploring heterogeneity based on other kinds of deepfake manipulations like face swapping and head puppetry<sup>72</sup>, contexts that require more sophistication to produce a convincing deepfake (e.g. where a person is moving, turning their head, and interacting with other people), who is being manipulated<sup>73</sup>, and contexts immediately relevant to current events (for example in March 2023, fake arrest images of Donald Trump were released on social media leading up to his indictments by the Manhattan District Attorney’s office and the United States Department of Justice<sup>74</sup>).

The results from these experiments cannot simply be explained by the deepfake manipulations being too obvious or unrealistic. Figure 2 illustrates accuracy across the political speeches displayed as video with audio and reveals that while participants are relatively highly accurately at identifying voice actor deepfakes, participants only identify text-to-speech deepfakes in 72% of observations, which is closer to random guessing than a perfect score. Moreover, human discernment is influenced by a number of factors and we find 10 text-to-speech PDD deepfakes, 2 voice actor PDD deepfakes, and 2 deepfake videos used Barari et al 2021 are accurately identified in less than 75% of observations. The participants’ low accuracy offers evidence that visual artifacts and inconsistencies created by the lip syncing deepfake manipulations are not readily apparent to most people, and as such, these videos represent a reasonable stimuli set for examining how well people can distinguish real from deepfake videos and how communication modalities, audio sources, and base rates of misinformation influence discernment.

People distinguish authentic from fabricated videos based on perceptual cues from video and audio and considerations about the content (e.g., the degree to which what is said matches participants’ expectations of what the speaker would say, which is known as the expectancy violation heuristic<sup>75</sup>). With the message content alone, participants are only slightly better than random guessing at 57% accuracy in Experiment 1a and 58% in Experiment 3. With perceptual information from video and the message content via subtitles, participants are slightly more accurate at 66% accuracy in Experiment 1a and 62% accuracy in Experiment 3. With information from audio only, participants are more accurate at 80.5% in Experiment 1a and 65% in Experiment 3. Finally, with information from both video and audio, participants are even more accurate at 82% accuracy in Experiment 1a and 74% accuracy in Experiment 3. Our finding that participants are more accurate at distinguishing between real and fabricated on voice actor audio than silent video with subtitles aligns with the social psychology literature demonstrating people tend to rely on auditory information more than visual information for both discerning sincerity<sup>76</sup> and ascribing authorship of a script to a human (as opposed to a computer)<sup>77</sup>. Another factor that could be supporting improved detection of deepfakes in the audiovisual regime is a beneficial effect of multisensory integration. Complementary audio-visual information has been shown to improve accuracy on perceptual decision-making tasks compared with visual information only, by amplifying post-sensory decision evidence<sup>78</sup>. However, the low accuracy of participants in distinguishing the speaker’s audio from text-to-speech audio trained on the speaker suggests social cues oriented towards speech in digital interfaces will

need to adapt to fabricated audio that is nearly indistinguishable from real audio. Overall, the experiment's results show that as participants have access to more information via audio and video, they are better able to distinguish whether a political speech has been fabricated.

Political deepfakes are most dangerous when people are least expecting information to be manipulated, and these experiments examine the influence of the base rate of misinformation on participant discernment. In Experiment 1, 50% of content is fake, and we explicitly inform participants of this base rate. In Experiment 3, the base rate of misinformation is randomized to be 20% or 80% of stimuli and we do not inform participants of this base rate. We find the high base rate of fakes compared to the low base rate leads participants to a 7.2 percentage point higher accuracy on real stimuli and 5.8 percentage point lower accuracy on fabricated stimuli, which are both statistically significant. In other words, participants responded that both real and fabricated stimuli are fake less often in the condition with the high base rate of misinformation than the low base rate. One explanation for this difference may be that people generally do not expect a very high base rate of fakes, which may lead people to respond that a stimulus is fake less often than they would in a more balanced setting. While false news is relatively rare in today's media ecosystem and is approximated to make up less than a fraction of a percent<sup>79,80</sup> of news, the capacity for generative AI to create misinformation is expanding<sup>81</sup> and future base rates of misinformation may be higher.

The political danger of fabricated videos may not be the average algorithmically produced deepfake but rather a single, highly polished, and extremely convincing video. For example, hyper-realistic deepfakes like the Tom Cruise deepfakes on TikTok (see <https://www.tiktok.com/@deeptomcruise>) are produced by visual effects artists using artificial intelligence algorithms but also using traditional video editing software and highly trained look-alike actors. While these hyper-realistic deepfakes may still contain manipulation artifacts (e.g., unattached earlobes that do not match Tom Cruise's attached earlobes<sup>82</sup>), future work on the psychology of multimedia misinformation may consider hyper-realistic videos produced by visual effects studios in addition to algorithmically manipulated videos. Experiment 4 offers insights on the future influence of hyper realistic videos by demonstrating that perfectly realistic video (the real videos paired with voice actor audio) leads people to believe audio is more likely to be authentic than when listening to audio by itself.

These experiments are useful to study how people discern multimedia information when attending to questions of accuracy, but they are less useful in understanding how people will share misinformation they consume on social media. People are generally highly accurate in discerning the veracity of news headlines yet share false news headlines because their attention is not focused on accuracy<sup>83</sup>. In fact, Epstein et al 2023 show that simply considering whether to share news on social media decreases people's accuracy at truth discernment<sup>84</sup>. Similarly, our findings in experiment 5 reveal that priming participants with a video showcasing visual effects manipulations lead people to express slightly more suspicions in free responses than participants who have not yet seen such a showcase video. These results support recent research showing that educational material on common misinformation techniques can improve people's ability to discern trustworthy from untrustworthy videos<sup>85</sup>.

On social media, video-based misinformation will often be designed to incorporate characteristics (e.g., fear, disgust, surprise, novelty) that divert people's focus from accuracy and make content go viral<sup>86-89</sup>. Given that multimedia misinformation may be both easier to discern and more frequently shared on social media than text-based media, more research needs to be done to understand how people allocate attention while browsing the Internet<sup>90</sup>. Our findings in Experiment 5 (which presents an environment much more similar to social media than the previous experiments) suggest that many people pay attention to the question of authenticity and remark on suspicions even without being asked about a stimulus' authenticity.

It is important to keep in mind that discernment – how accurately people discern misinformation – is different than belief – how much people report they believe misinformation. It is possible (though quite peculiar) that someone could be highly accurate at discerning truth from falsehood while also tending to believe the fabricated content and not believe the true content. For example, research on false news headlines and articles finds that people are better at discerning news concordant with their political leanings than discordant news while also believing concordant news more often than discordant news<sup>61</sup>.

The findings that videos of political speeches are easier to distinguish as authentic or fabricated than text transcripts highlights the need to re-introduce and explain the oft-forgotten second half of the “seeing is believing” adage. In 1732, the old English adage appears as: “Seeing is believing but feeling is the truth.”<sup>91</sup> Here, “feeling” does not refer to emotion but rather direct experience. Since the advent of photography, people across society have generally understood that what we see in a photograph is not always the truth and further assessment is often necessary<sup>92-94</sup>.

In this paper, we examined a bounded question – how well can ordinary people discern (and how often do they suspect) whether or not a short soundbite of a political speech by a well-known politician in text, audio, or video has been fabricated – and we find that more information via communication modalities – text transcripts vs. silent, subtitled video vs. video with audio – enables people to more accurately discern fabricated and real political speeches. These results are particularly relevant for the design of content moderation systems for flagging misinformation on social media. In particular, we suggest content moderation flags include explanations that address which component part of a video appears to be fabricated. These explanations could allow people to appropriately allocate attention to the content<sup>95</sup> or perceptual cues (e.g., low-level pixel

features, high-level semantic features, and biometric-based features<sup>96</sup>) when trying to assess the content’s authenticity.

Finally, these findings offer insights into political communication and communication theory more generally; there is more to how humans form beliefs than the “seeing is believing” narrative would suggest because people can pay attention and seek out inconsistencies to both what is said and how something is said.

## Methods

### Consent and Ethics

This research complied with all relevant ethical regulations and the Massachusetts Institute of Technology’s Committee on the Use of Humans as Experimental Subjects approved this study as Exempt Category 3 – Benign Behavioral Intervention. This study’s exemption identification numbers are E-3105, E-3354, E-4735, and E-5493. For experiments 1 through 4, all participants were presented with an informed consent statement: “Detect Fakes is an MIT research project. All guesses will be collected for research purposes. All data for research are collected anonymously. For questions, please contact [detectfakes@mit.edu](mailto:detectfakes@mit.edu). If you are under 18 years old, you need consent from your parents to use Detect Fakes.”

For participants in experiment 5, the informed consent and instructions statement was presented as follows: “This is an MIT research project. You will be shown quotes, audio, and video files that one might expect to see on social or digital media. You will be requested to share your thoughts and opinions after reading, listening, or viewing each of these media files. All response data is collected anonymously for research purposes... Participation is voluntary, and you may only participate if you are 18 years of age or older. For questions, please contact [arunas@mit.edu](mailto:arunas@mit.edu).”

Before beginning any of the experiments, all participants from Prolific were also provided a research statement, “The findings of this study are being used to shape science. It is very important that you honestly follow the instructions requested of you on this task, which should take a total of 15 minutes. Check the box below based on your promise:” with two options, “I promise to do the tasks with honesty and integrity, trying to do them uninterrupted with focus for the next 15 minutes.” or “I cannot promise this at this time.” Participants who responded that they could not do this at this time were re-directed to the end of the experiment.

In Experiment 1, we immediately debriefed participants on which political speeches are real and which are fabricated after each video seen. In Experiments 2 through 5, we debrief participants on which political speeches are real and which are fabricated at the end of the experiment. In order to limit the potential for these deepfakes to be taken out of their research context, we created a public website showing the deepfakes signed using the C2PA protocol to indicate these videos are partially AI generated. If deepfakes were taken out of context, people can reference these signed deepfakes to identify them as fabrications designed for research. Figure 8 presents a screenshot showing how we overlaid deepfake videos’ metadata for ease of reference.

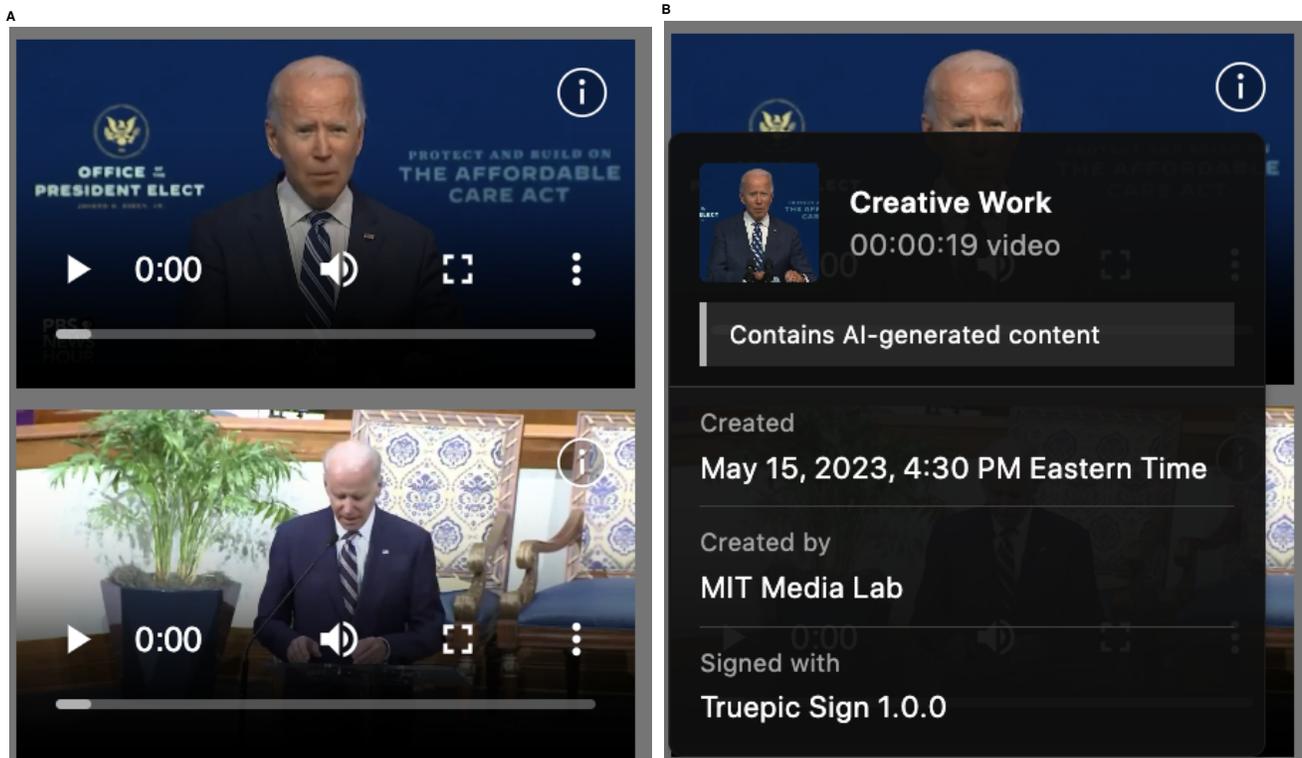
For participants in Experiment 1a recruited from Prolific, we compensated participants at a rate of \$9.78 an hour and provided bonus payments of \$5 to the top 1% of participants in terms of accuracy. In Experiment 1b, we did not compensate participants financially because they arrived at the website via organic links on the Internet. In Experiments 2 through 5, all participants are recruited from Prolific and compensated at a rate of \$12.00 an hour.

### Digital Experiment Interface

In experiments 1 through 4, we hosted multimedia stimuli – transcripts, audio, and video of authentic and fabricated political speeches – on a custom designed website called Detect Fakes, which was hosted at <https://detectfakes.media.mit.edu/>. In these experiments, we asked participants to identify stimuli as fabricated and non-fabricated. In experiment 5, we used qualtrics and asked participants “What comes to mind after watching the following video/listening to the following audio/reading the following quote?”

#### **Experiments 1a and 1b**

First, we collected informed consent, presented participants with instructions and an attention check, and then we showed participants a short political speech and ask “Did [Joseph Biden/Donald Trump] say that?” followed by “Please [read/listen/watch] this [transcript/audio clip/video] from [Joseph Biden/Donald Trump] and share how confident you are that it is fabricated. Remember half the media snippets we show are real and half are fabricated.” Figure 9 in the Supplementary Information section presents a screenshot of the user interface, which shows participants were instructed to move a slider to report their confidence from 50% to 100% that a stimulus is fabricated (or 50% to 100% that a stimulus is not fabricated). After each response, we presented feedback to participants to inform them whether the stimulus was actually fabricated. Then, we presented participants with another stimulus (selected at random) and the process repeated until participants viewed all 32 stimuli or decided to leave the experiment.



**Figure 8. Screenshot of Debrief with Metadata** A. Screenshot of videos with information icon in the top right corner. B. Screenshot of metadata content displayed after clicking the information icon.

### **Experiments 2, 3 and 4**

Similar to Experiment 1, we first collected informed consent and presented participants with instructions and an attention check. Next, we showed participants a short political speech and asked “Fabricated or Authentic?” followed by “Please [read/listen/watch] this [transcript/audio clip/video] from [speaker] and share whether you think the speech is fabricated, and how confident you are in this judgment.” In Experiment 4, we edited the follow up to “Please watch and listen to this [audio clip/video] from [speaker] and share whether you think the voice is [speaker’s] or a voice actor, and how confident you are in this judgment.” Figure 10 in the Supplementary Information section presents a screenshot of the user interface, which shows participants chose between two options – Real or Fake – and reported their confidence along a Likert scale from 1=Not confident at all to 2=Slightly confident to 3=Somewhat confident to 4=Fairly confident to 5=Completely confident. The experiment interface prevented participants from selecting the real or fake radio button until they have watched as least 15 seconds of the audio or video, prevented participants from selecting their confidence rating until they have selected a radio button, and prevented participants from submitting their response until they have selected a radio button and confidence rating.

### **Experiment 5**

In Experiment 5, we began by collecting informed consent and presenting participants with instructions and a directed choice attention check. Instead of directly asking participants to identify whether the video is real or fake, Experiment 5 asks, “What comes to mind after watching the following video/listening to the following audio/reading the following quote?” As shown in Figure 11, participants are presented with a text response box. In order to proceed to the next stimulus, the interface requires participants to write at least 50 characters and spend at least 30 seconds on each stimulus. Experiment 5 is hosted on Qualtrics (as opposed to our custom website described in the Methods section) to prevent participants from knowing the dependent variable of interest in this experiment.

In order to create the key dependent variables for this experiment, three authors of this article independently annotated 2,200 participants’ responses for whether participants expressed suspicion that the stimulus was fabricated or not. The dependent variable represents majority agreement between the three annotators. Cohen’s Kappa between annotators on PDD stimuli ranged from 0.56 to 0.85; the maximum Cohen’s Kappa between annotators is 0.72 on text, 0.79 on silent video, 0.87 on audio, and 0.90 on video with audio. Two annotators identified suspicion in 11% of observations and the third annotator identified suspicion in 22% of observations. Across the 2,000 observations, 9% are identified by all three as revealing suspicion, 3% are identified by two of three, 11% are identified by one, and 77% are identified by all three as not revealing suspicion. Examples

of revealing suspicion by two authors but not a third include: “I’m not sure this is a real statement made by Trump” and “I’m not sure if Biden actually said this or not, but it sounds convincing enough.” Likewise, examples that don’t count as revealing suspicion because only one author annotated it as suspicion include: “If what he says were true then that would be a very beneficial thing.” and “I’m not completely sure what Biden is trying to say here, I would want more information.”

## **Experiment Stimuli**

All real and deepfake stimuli are listed in Table 19 in the Appendix.

### ***Original Presidential Deepfakes Dataset***

In Experiments 1a and 1b, the stimuli were drawn from the Presidential Deepfake Dataset (PDD)<sup>65</sup>. The PDD consists of 32 videos showing two United States presidents – Donald Trump and Joseph Biden – making political speeches. Half the videos are original videos that had not been altered or manipulated. The other half had been fabricated to make the politicians appear to say something that they have not said. The fabricated videos were produced by writing a fabricated script, recording professional voice-actors reading the script, and applying a deepfake lip-syncing algorithm<sup>16</sup> to real videos of Joseph Biden and Donald Trump to make it appear as if the politicians actually gave such a fabricated speech. The mean duration of the videos is 21 seconds and all videos are recorded at 30 frames per second. The PDD is balanced across three dimensions: (1) videos that had and had not been fabricated, (2) videos of Joseph Biden and Donald Trump, and (3) videos of the two politicians making concordant and discordant speeches with what the general public believes were the politicians’ political views.

In order to validate the concordance and discordance of speeches, we conducted an independent survey where 84 participants who passed an attention check rated each of the 32 transcripts for how well the political speeches match either politicians’ political views. Participants were instructed “For each statement, we want you to rank how closely the statement matches your understanding of President Joseph Biden or President Donald Trump’s political views” and asked to provide a judgment on a 5-point Likert scale from “Strongly Disagree” (-2) to “Strongly Agree” (2) that “This statement matches President [Joseph Biden’s/Donald Trump’s] political viewpoint: [statement].” Participants’ responses confirm that speeches designed to be concordant and discordant with the two politicians views were indeed concordant and discordant with the average participants’ perception of the politicians’ views. The Z-values of participants responses to concordant and discordant speeches are -0.25 and 0.21, respectively, and this difference is statistically significant with  $p < 0.001$  based on a t-test.

In Experiment 1, we transformed each of the original videos from the PDD into 7 different forms of media: a transcript, an audio clip, a silent video, audio with subtitles, silent video with subtitles, video with audio, and video with audio and subtitles. As a result, there were 7 modality conditions, 32 unique speeches, and 224 unique stimuli. In the digital experiment, the transcript appears as HTML text and the six other forms of media content appear in a video player. The audio clip shows a black screen in the video player and the audio clip with subtitles shows a black screen with subtitles at the bottom.

### ***Enhanced Presidential Deepfakes Dataset***

The enhanced PDD included the same real videos from the PDD and two sets of 16 enhanced deepfakes with voice actor audio and audio from a text-to-speech algorithm. The 16 enhanced deepfakes included source videos more amenable to deepfake lip-syncing manipulations than the initial source videos for the PDD deepfakes, visual touch ups with the DeepFaceLab library, audio from both a voice actor and a text-to-speech algorithm trained on the speakers’ voices, and additional audio engineering to add in background noise and subtle acoustic elements to make the audio appear real. We provide additional details on the deepfake generation process below in the section titled “Enhanced Deepfake Generation Process.”

### ***Other Stimuli***

We included 6 additional deepfakes and 6 additional real videos in Experiment 2, which are used in Barari et al 2021. The deepfakes included two deepfakes of Bernie Sanders and Hilary Clinton from the Agarwal et al 2019 dataset<sup>5</sup>, which involve face swapping of the politicians’ faces onto videos of Saturday Night Live skits by Larry David and Kate McKinnon. The four other deepfakes can be found on Youtube and showed Boris Johnson endorsing his opponent, Trump announcing we have eradicated AIDS, Trump announcing deepfakes would make it easy to make him say ridiculous things, and Obama announcing we are in an era in which our enemies can make it look like anyone is saying anything.

## **Enhanced Deepfake Generation Process**

### ***Source Video Curation***

We selected 16 videos from YouTube of Donald Trump and Joseph Biden giving speeches where the politician is facing forward and no other people are in the frame. As we curated videos, we selected videos based on the following: the distance of the speaker from the camera, which controls for the size of the politician’s face, the speaker’s facial expressions and body language such that they match the sentiment in the script, the stationarity of the face position such that there are minimal clues from yaw and pitch rotation or zooming or panning, and situational artifacts like microphones, news ribbons, and other contextual clues that match the content of the politician’s speech. For example, ‘b-04.mp4’ features President Biden speaking about the

meaninglessness of war by referencing the death of his son. We included several contextual clues in the video frame such as the logos of the Army National Guard and Air National Guard, the name of President Biden’s son on the pedestal, and a sombre expression on President Biden’s face with tears streaming down his face. Further, we selected speeches where the cadence of the audio qualitatively matches the cadence of the fabricated audio. We found that the Wav2Lip algorithm, used to synthesize the lip sync, generated more authentic videos when the visual frames in which the politicians speak and those in which they pause approximately align with the speaking and non-speaking portions of the audio frames, respectively. We also ensured that we used different source videos to generate each deepfake.

We calculated the accutance of a randomly selected frame of each video in our dataset to measure the perceived ‘sharpness’ of the video. This helped us compare the relative sharpness or blurriness of the fabricated videos to the real ones, and ensure that the observed accuracy is not a consequence of the visual quality of the videos. The accutance measures the contrast between adjacent pixels, particularly along the edges of objects or color transitions in the frame. The accutance is a measure that is more closely based on human visual perception of the sharpness. The real Biden videos have an average accutance of 26.29 and the text-to-speech and voice actor deepfakes of Biden have comparable accutance values of 24.71 and 24.69, respectively. Similarly, the real Trump videos have an average accutance of 34.071 and the text-to-speech and voice actor deepfakes of Trump have comparable accutance values of 36.129 and 35.81, respectively. Interestingly, the accutance of the real Biden videos is lower than that of the fabricated videos, while the accutance of the real Trump videos is higher than that of the fabricated videos. We also show, quantitatively, that none of the videos are blurry, given the range of the accutance values observed.

The width of faces in all the videos are comparable. The face width in the Biden videos varies between 176px - 335px for real videos and between 190px - 200px for the fabricated videos. The width of faces in the Trump videos varies between 136px - 250px for real videos and between 132px - 241px for the fabricated videos.

### ***Speech Content Fabrication***

We selected statements made by political pundits and fact-checking sources that resembled President Donald Trump and Joseph Biden’s communications style, but misrepresented their political talking points. As such, we fabricated statements that are concordant or discordant with the featured politician’s views on a topic. In order to create convincing fabricated speeches, we employed tactics commonly found in fabricated news such as misrepresenting situational or temporal contexts (i.e. presenting a view held by the politician in the past as that being held by them in the present), presenting statements made by someone aligned or misaligned with the politician’s viewpoints as statements made by them, and manipulating the situational context in the video to present a statement that was not made by the politician but is relevant to the context. We restricted the topics referenced by the fabricated statements to Obama, gun control, birth control, minority groups, social security, and school reform to ensure topical parity.

### ***Voice Actor Audio***

The audio for the PDD deepfake stimuli was created by voice actors (J-L Cauvin and Austin Nasso for the original deepfakes in Experiment 1 and Colin Cassidy for the enhanced deepfakes in Experiment 2) to read the fabricated speech transcripts in their most realistic impression of Donald Trump and Joseph Biden. The enhanced voice actor videos were further improved with realistic acoustic reverberation effects, background noise, and microphone filtering as described in the section titled “Audio Engineering.”

### ***Text-to-Speech Audio***

We generated 16 fabricated audios using Eleven Labs<sup>67</sup> text-to-speech product, which took text and source audio as inputs and generated a synthetic audio with speech from the text input and stylistic aspects such as pitch, volume, tone, and inflection that matched with the source audio.

We selected source audio from speeches by President Biden and Trump with minimal background noise, no additional voices, and emotional tones similar to the ones we were interested in synthesizing. We used audio snippets between 1-5 minutes in length. For example, in one of our audios of President Biden talking about the pointlessness of war (b-04-t.mp4), we introduced sombre tones using an audio of President Biden speaking at John McCain’s memorial service. Similarly, the audio supporting the right of women to abortion contained passionate and angry tones created using President Biden’s passionate speech on the “extremist threat to democracy”<sup>97</sup>.

In order to introduce pauses, emphasis, features such as yelling, sombre tones, mirth, sarcasm, scorn, and other emotions, we manipulated the input text with various punctuation symbols. We did this by adding periods, commas, and semi-colons to introduce pauses and using capitalized text and exclamations to control volume. Asterisks, and quoted words, placed emphasis on certain portions of the speech and question marks allowed us to add query-tones. Eleven Labs also allowed user control through two additional input attributes – ‘Stability’ and ‘Clarity + Similarity Enhancement’. Reducing ‘Stability’ improved the expressiveness in the voice but also led to frequently changing inflection points in the audio. Reducing the ‘Clarity + Similarity Enhancement’ attribute minimized the effects of background noise and other undesirable features in the source audio. Using

lower values of the ‘Stability’ attribute allowed us to create anger, and loud phrases in the synthesized audio. Lower values of the ‘Stability’ attribute produced more realistic synthesis of President Donald Trump’s voice, but higher values contributed to realistic synthesis of President Biden’s voice.

### **Lip-Syncing and Face Swapping**

The enhanced deepfakes were created from our set of curated source videos using a 2 step process. First, the source video was morphed to sync the featured politician’s lips to the fabricated audio. This synthesized video was then visually polished using a face-swap deepface algorithm.

**Lip syncing** The audio stream of a curated video was first removed to produce a silent video. The silent video was then cut up into a snippet that is the same length as the fabricated audio featuring the visual frames to be used in the final deepfake. We used two Wav2Lip models pretrained on the LRS-2 dataset<sup>98</sup> – the first model prioritized lip sync accuracy over the visual quality of the synthesized lip sync, the second model prioritized visual quality over lip sync accuracy. The video-audio snippet pair was fed into each of the Wav2Lip models to produce two synthetic videos where the lip sync in the video was matched to the fabricated audio.

**Deepfake generation** We used the DeepFaceLab library<sup>66</sup> to create our set of enhanced deepfakes. The library used the face-swap deepfake technique to superimpose facial features from a source video onto a destination video. In our case, each video synthesized by the Wav2Lip model was considered the destination video. The source video was created for each President by concatenating a set of videos featuring them in a diversity of lighting conditions and facial poses. Each concatenated video was about 90 minutes long. We created facesets using these concatenated videos such that for each second of video, we obtained 20 faces using the S3FD face detector<sup>99</sup>. Faces in the faceset were then rotated about the yaw to squarely face the camera to standardize and smooth facial landmarks. A face segmentation algorithm<sup>100</sup> was further applied to ablate hair and other facial props from the faces in the faceset. Deepfake models were created for each assembled faceset using the training process outlined by the DeepFaceLab project. We trained the simplest model for 600K iterations, to use multiple 64x64 patches of the faces in the faceset to produce the synthesized output. The model was used to visually enhance the face in each of the two Wav2Lip videos, and the more realistic output was selected out of the two.

### **Audio Engineering**

In order to make the voice actor and text-to-speech audio congruent with the videos, we engineered the audio clips to mimic aspects of the real videos such as acoustic reverberation, background noise, and microphone filtering. Without these efforts, the audio production methods (by a voice actor recording in a controlled environment with a high-quality microphone or by a speech synthesis model which does not explicitly account for the visually-depicted environment) could readily provide auditory cues as to the clips’ origins.

We observed from real video clips of the speakers that typical background noise in these political videos consisted of a few key components. One was *room tone*, or the baseline noise in different environments. This is commonly added in contexts like film and journalism, to provide realistic environmental context and continuity. Other components were more specific to the settings of political speaking: background voices (e.g. *walla* or audience reactions), camera shutters, and sometimes other unidentified noises.

We collected a variety of such sounds from Freesound.org<sup>101</sup>. To each video clip, we added a suitable room tone as a minimal environmental factor. In some cases, we judged that additional noises were not likely, and simply filtered the room tone to resemble that from realistic videos. In many cases, we then proceeded to introduce other sound effects including camera shutters and background voices, tailoring each arrangement to the specific video. We then replicated these background layers for both voice actor and text-to-speech versions of a given video clip.

Acoustic reverberation refers to the prolongation of sounds in physical environments caused by an accumulation of reflections. Our expectations about acoustic reverberation are linked to visual perception<sup>102, 103</sup>; for instance, imagine the same speech in a large auditorium vs. in a small office. To simulate the acoustic effect of the depicted visual environments, we introduced artificial reverberation with reference to the real videos. For example, a large rally would likely contain reflections from far-away surfaces and thus have a longer reverberation than an announcement from an office.

Beginning with a suite of reverberation simulations which match the range of different depicted environments, we selected 1-2 options for each video as most appropriate, and then blended the applied reverberation with the dry auditory mix to yield a realistic result. We applied the same reverberation to voice-actor and text-to-speech audio, with one exception: for the clip b-04-t, we added a little extra reverberation to the text-to-speech voice signal (as the reverberation applied to the voice-actor audio did not translate realistically without this). One factor worth noting here is that the text-to-speech clips, despite being dry in the sense that they did not contain audible *late* reverberation, did already contain some of the acoustic characteristics of the original environments (such as early reflections). This might be due to the audio that they were modeled on, as it is challenging to get studio-quality anechoic recordings of political figures.

To mimic the effects of different microphones, speaker-microphone spatial relationships, and broadcast audio pipelines, we applied equalization to the voice actor audio clips and the background layers. For each speaker audio, we first used an EQ matching process to mimic the frequency response of a real audio clip. Since this process is often error-prone due to the difficulty of estimating the right transfer function, we then manually tweaked the effect of the match EQ and added various additional equalization. As an example, if a speaker is far from a microphone, or significantly off-axis in terms of orientation, then high frequencies will be attenuated and *proximity effects* might also result in less low-frequency content than a very close microphone, resulting in a mid-frequency focus. We heuristically equalized all voice-actor audio clips with this in mind. We did not apply any equalization to the text-to-speech audio clips; as mentioned previously, these already seemed to implicitly model such aspects of real clips.

There were other aspects of real audio clips that were difficult to synthetically recreate. One example of this is plosives, wherein consonants like "p" and "b" result in bursts of air that sound as low-frequency pops on many microphones. Recreating these convincingly is challenging, and so we did not include any such effects. We also applied other very subtle effects in some cases, to move the voice-actor audio closer to the real audio, including very slight formant shifting (for closer timbre matching), and very slight mid- and high-frequency excitation to recreate artifacts of the original speeches, in addition to the equalization.

For the voice-actor recordings that pertained to the *real* videos, we needed to time-align the voice-actor audio with the original audio in order for the auditory and visual components to be synchronized. This overall timing encoded minor factors like the onset of the speech, but then also the speed and cadence which can vary arbitrarily and result in significant misalignment. This process mimics aspects of so-called Automated Dialog Replacement (ADR) wherein dubbed audio in film settings needs to be re-aligned with the original recordings. Though there are tools that partially automate this, it is traditionally done by hand by moving and stretching spoken words and syllables until they align well. We followed this process, carefully aligning the voice-actor audio with the originals. Though this process introduces artifacts due to the audio stretching, we tuned the stretching algorithm and its parameters to minimize them to the extent practically possible.

## Preregistration

Experiments 1a, 2, 3, 4, and 5 involved participants recruited from Prolific and were pre-registered on aspredicted.org at the following URLs: Experiment 1a, Experiment 2, Experiment 3, Experiment 4, and Experiment 5. With one exception, there were no deviations from the pre-registration. The one exception is we do not conduct analyses towards the third main hypothesis in Experiment 1a about motivated reasoning because this hypothesis was unfeasible to address with this experimental setup.

Experiment 1b was not pre-registered and included 41,313 participants who discovered the experiment organically through search engines or news media.

## Participants

In Experiments 1a, 2, 3, 4, and 5, we recruited participants via the Prolific platform<sup>104</sup>. In each of these experiments, participants responded to a baseline survey, which consisted of questions on political preferences, experience with deepfakes, trust in media and politics. Prolific provided basic demographic data including participants' self-reported sex, ethnicity, and age. In experiment 1a but not Experiments 2-4, we included three questions from the Cognitive Reflection Test (CRT)<sup>69</sup>. In Experiment 1b, we collected data from participants who visited the experiment organically but did not collect demographic, political identity, or pre-experiment questions for these non-recruited participants.

In Experiments 1a, 2, 3, and 4, 95%, 83%, 84%, and 89% of participants provided responses to the complete set of stimuli, which results in 99.9%, 98.7%, 97.5%, and 97.6% of the expected data in each experiment. The missing data appears to be missing due to intermittent or slow network connection issues where participants could proceed without their data getting submitted to the server. In Experiment 1b where participants were not recruited and were not asked to complete a pre-specified number of responses, 15% of participants provided responses to the complete set of stimuli.

We exclude participants from participating in multiple experiments. However, 9 participants who participated in Experiment 2 also participated in Experiment 4. The results in Experiment 4 are robust to both including and excluding these 9 participants.

In these experiments, we do not find consistent differences based on sex, and we do not report sex-based analyses because we did not pre-register sex-based analyses nor do we have theoretical grounds for suspecting differences across sex.

### **Experiment 1a – March 20 to April 8, 2021**

In Experiment 1a, we recruited 501 participants from the United States who successfully passed the attention check and provided 16,011 observations. The demographic distribution of participants along sex and age is: 50% male, 49% female, and 1% unknown; 60% 18 to 35, 37% 36-64, and 3% over 64. We do not have data on participants' race or ethnicity. The sample of 509 recruited participants is balanced across political identities: 255 recruited participants self-report as Democrats, and the other 246 recruited participants self-report as Republicans. In response to a pre-experiment question on participants' experience with deepfakes, fewer than 1% of participants responded that they have created their own deepfakes, 73% of participants have seen a few to several examples of deepfakes, and 27% of participants have yet to see their first deepfake or don't know their

experience with deepfakes. In terms of trust and confidence in media, 43% of participants report a fair amount or a great deal of trust and confidence in media and 57% of participants report not very much or none at all. On the topic of following news on government and public affairs, 81% of participants report following the news most or some of the time and 19% of participants report following only now and then or hardly at all. In this experiment, 44 participants fail the attention check and 8 participants withdrew. We do not find statistically significant differences in the failure rate on the attention check across political identities<sup>105</sup>.

#### ***Experiment 1b – March 19 2021 to June 30 2022***

In Experiment 1b, 41,313 participants visited the experiment, passed the attention check, and provided 416,901 observations. According to data from Google Analytics, 76% of these participants participated from outside the United States. 5,106 individuals participated in the experiment during the pre-registration window from March 4, 2021 to June 1, 2021. These participants found the website organically and completed 44,461 trials. Between June 1, 2021 and July 1, 2022, an additional 67,576 individuals (70% of whom visited from outside the United States) completed 566,343 trials. We include participants in Experiment 1b who participated outside the pre-registered window because we had an unexpectedly very large sample, which is due to around a thousand of participants visiting the website each week and ten thousand participants visiting the website after it was posted to a website called Hacker News in March 2022. In total, 31,369 of these non-recruited participants failed the attention check.

#### ***Experiment 2 – June 9, 2023***

In Experiment 2, we recruited 302 participants from the United States who successfully passed the attention check and provided 5,964 observations. The demographic distribution of participants along sex, age, and ethnicity is: 54% male, 44% female, and 2% unknown; 46% 18 to 35, 46% 36-64, 5% over 64, and 2% unknown; and 70% White, 11% Black, 7% Asian, 5% mixed, 4% other, and 2% unknown. With respect to political beliefs, 63% of participants self-report their political preference as democratic, 16% as equally democratic and republican, and 22% as republican, and similarly, 61% of participants report voting for Joseph Biden in 2020, 20% of participants report voting for Donald Trump in 2020, and the rest decline to answer or report voting for another candidate. In response to a pre-experiment question on participants' experience with deepfakes, fewer than 1% of participants responded that they have created their own deepfakes, 85% of participants have seen a few to several examples of deepfakes, and 13% of participants have yet to see their first deepfake or don't know their experience with deepfakes. In terms of trust and confidence in media, 49% of participants report a fair amount or a great deal of trust and confidence in media and 51% of participants report not very much or none at all. On the topic of following news on government and public affairs, 69% of participants report following the news most or some of the time and 31% of participants report following only now and then or hardly at all. In this experiment, 30 participants fail the attention check.

#### ***Experiment 3 – June 15, 2023 to June 20, 2023***

In Experiment 3, we recruited 1006 participants from the United States who successfully passed the attention check and provided 19,812 observations. The demographic distribution of participants along sex, age, and ethnicity is: 48% male, 51% female, and 1% unknown; 33% 18 to 35, 54% 36-64, 13% over 64, and 1% unknown; and 77% White, 13% Black, 6% Asian, 2% mixed, and 2% other. With respect to political beliefs, 62% of participants self-report their political preference as democratic, 12% as equally democratic and republican, and 26% as republican, and similarly, 59% of participants report voting for Joseph Biden in 2020, 22% of participants report voting for Donald Trump in 2020, and the rest decline to answer or report voting for another candidate. In response to a pre-experiment question on participants' experience with deepfakes, fewer than 1% of participants responded that they have created their own deepfakes, 88% of participants have seen a few to several examples of deepfakes, and 13% of participants have yet to see their first deepfake or don't know their experience with deepfakes. In terms of trust and confidence in media, 45% of participants report a fair amount or a great deal of trust and confidence in media and 55% of participants report not very much or none at all. On the topic of following news on government and public affairs, 75% of participants report following the news most or some of the time and 25% of participants report following only now and then or hardly at all. In this experiment, 59 participants fail the attention check.

#### ***Experiment 4 – June 14, 2023***

In Experiment 4, we recruited 206 participants from the United States who successfully passed the attention check and provided 3,215 observations. The demographic distribution of participants along sex, age, and ethnicity is: 49% male, 48% female, and 3% unknown; 38% 18 to 35, 52% 36-64, 5% over 64, and 5% unknown; and 74% White, 7% Black, 8% Asian, 4% mixed, 2% other, and 3% unknown. With respect to political beliefs, 65% of participants self-report their political preference as democratic, 11% as equally democratic and republican, and 24% as republican, and similarly, 58% of participants report voting for Joseph Biden in 2020, 21% of participants report voting for Donald Trump in 2020, and the rest decline to answer or report voting for another candidate. In response to a pre-experiment question on participants' experience with deepfakes, 1% of participants responded that they have created their own deepfakes, 87% of participants have seen a few to several examples of deepfakes,

and 12% of participants have yet to see their first deepfake or don't know their experience with deepfakes. In terms of trust and confidence in media, 39% of participants report a fair amount or a great deal of trust and confidence in media and 61% of participants report not very much or none at all. On the topic of following news on government and public affairs, 76% of participants report following the news most or some of the time and 24% of participants report following only now and then or hardly at all. In this experiment, 14 participants fail the attention check and 1 participant withdrew.

### **Experiment 5 – December 22, 2023**

In Experiment 5, we recruited 200 participants from the United States who successfully passed the directed choice attention check and provided 2,200 observations. The demographic distribution of participants along sex, age, and ethnicity is: 47.5% male, 50.5% female, and 2% unknown; 46% 18 to 35, 46% 36-64, 7% over 64, and 1% unknown; and 55% White, 18% Black, 11% Asian, 6% mixed, 6.5% other, and 3.5% unknown. In this experiment, 3 participants were excluded (2 participants fail the directed choice attention check and 1 participant responded with the same unrelated response to all questions) and 21 participants withdrew before completing the experiment.

### **Randomization**

In all experiments, we randomized the order of the political speeches and each participant encounters each political speech only once. In Experiment 1, participants engaged with up to 32 unique political speeches. We randomized the display of the political speech as one of the seven modality conditions. In Experiment 2, participants engaged with a random sample of 20 unique videos from a pool of 60 videos, which consisted of 4 videos from each of the following 5 stimuli categories: 16 real videos from the PDD dataset, 16 deepfakes with audio from voice actors from the PDD dataset, 16 deepfakes with audio from a text-to-speech algorithm from the PDD dataset, 6 real videos used in Barari et al 2021, and 6 deepfakes used in Barari et al 2021<sup>31</sup>). In Experiment 3, participants engaged with a random sample of 20 unique political speeches from a pool of 32 unique political speeches, which consisted of 16 real videos and 16 deepfakes with audio from a text-to-speech algorithm from the PDD dataset. We randomized the display of the political speeches as one of the four modality conditions, and we also randomize the base-rate of deepfakes seen. In Experiment 4, participants engaged with 32 unique political speeches, which consisted of 16 real videos with the original audio and 16 real videos with voice actor audio. Finally, in Experiment 5, participants saw 10 stimuli randomly sampled from the same set as Experiment 3.

### **Data Availability**

All stimuli and data collected are available on Research Box at [https://researchbox.org/1723&PEER\\_REVIEW\\_passcode=EGVULE](https://researchbox.org/1723&PEER_REVIEW_passcode=EGVULE). The original sources for the videos are described in Table 19 in the Appendix.

### **Data Availability**

All code produced to analyze the data collected are available on Research Box at [https://researchbox.org/1723&PEER\\_REVIEW\\_passcode=EGVULE](https://researchbox.org/1723&PEER_REVIEW_passcode=EGVULE).

### **Acknowledgements**

The authors would like to acknowledge funding and support for signing videos from Truepic, funding from MIT Media Lab member companies, thank Colin Cassidy, J-L Cauvin, Austin Nasso for providing voice impressions, thank the following users who contributed sounds from Freesound.org including aaronstar, aleclubin, cmilan, funwithsound, jgarc, johnsonbrandediting, klankbeeld, macohibs, mzui, noisecollector, peridactyloptrix, speedygonzo, zabuhailo, thank David Rand, Gordon Pennycook, Rahul Bhui, Yunhao (Jerry) Zhang, Ziv Epstein, and members of the Affective Computing lab at the MIT Media Lab and the Human Cooperation lab at MIT Sloan School of Management for helpful feedback on early versions of this manuscript, Anna Murphy, Shreya Kalyan, Theo Chen, and Alicia Guo for excellent research assistance, Craig Ferguson for helpful feedback on technical development, and Stephan Lewandowsky, Simon Clark, and anonymous reviewers for insightful reviews and important questions to consider in experiments 2 through 5.

### **Author Contributions**

M.G. conceived the experiments, A.S. and D.K. curated and created the deepfakes, N.S. performed audio engineering, A.S., M.G., and N.S. conducted the experiments, A.S. and M.G. analyzed the results, A.S., M.G., and N.S. wrote the manuscript, and A.S., A.L., M.G., N.S. and R.P. reviewed and edited the manuscript.

## Competing Interests

The authors declare funding for participant recruitment from Truepic and declare no competing interests.

## References

1. Hancock, J. T. & Bailenson, J. N. The social impact of deepfakes. *Cyberpsychology, Behav. Soc. Netw.* **24**, 149–152 (2021).
2. Chesney, B. & Citron, D. Deep fakes: A looming challenge for privacy, democracy, and national security. *Calif. L. Rev.* **107**, 1753 (2019).
3. Paris, B. & Donovan, J. Deepfakes and cheap fakes. *United States Am. Data & Soc.* (2019).
4. Leibowicz, C., McGregor, S. & Ovadya, A. The Deepfake Detection Dilemma: A Multistakeholder Exploration of Adversarial Dynamics in Synthetic Media. *arXiv:2102.06109 [cs]* (2021). Accessed 2021-07-07.
5. Agarwal, S. *et al.* Protecting world leaders against deep fakes. In *CVPR workshops*, vol. 1 (2019).
6. Pataranutaporn, P. *et al.* Ai-generated characters for supporting personalized learning and well-being. *Nat. Mach. Intell.* **3**, 1013–1022 (2021).
7. Guess, A. M. & Lyons, B. A. Misinformation, disinformation, and online propaganda. *Soc. media democracy: The state field, prospects for reform* 10–33 (2020).
8. Boháček, M. & Farid, H. Protecting world leaders against deep fakes using facial, gestural, and vocal mannerisms. *Proc. Natl. Acad. Sci.* **119**, e2216035119 (2022).
9. Karras, T., Laine, S. & Aila, T. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4401–4410 (2019).
10. Karras, T. *et al.* Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8110–8119 (2020).
11. Nichol, A. *et al.* Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).
12. Groh, M., Epstein, Z., Obradovich, N., Cebrian, M. & Rahwan, I. Human detection of machine-manipulated media. *Commun. ACM* **64**, 40–47 (2021).
13. Suvorov, R. *et al.* Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2149–2159 (2022).
14. Arik, S. O., Chen, J., Peng, K., Ping, W. & Zhou, Y. Neural voice cloning with a few samples. *arXiv preprint arXiv:1802.06006* (2018). Accessed 2021-07-07.
15. Luong, H.-T. & Yamagishi, J. Nautilus: a versatile voice cloning system. *IEEE/ACM Transactions on Audio, Speech, Lang. Process.* **28**, 2967–2981 (2020).
16. Prajwal, K. R., Mukhopadhyay, R., Namboodiri, V. P. & Jawahar, C. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International Conference on Multimedia*, MM '20, 484–492 (Association for Computing Machinery, New York, NY, USA, 2020).
17. Lahiri, A., Kwatra, V., Frueh, C., Lewis, J. & Bregler, C. Lipsync3d: Data-efficient learning of personalized 3d talking faces from video using pose and lighting normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2755–2764 (2021).
18. Hong, W., Ding, M., Zheng, W., Liu, X. & Tang, J. Cogvideo: Large-scale pretraining for text-to-video generation via transformers, DOI: [10.48550/ARXIV.2205.15868](https://doi.org/10.48550/ARXIV.2205.15868) (2022).
19. Peirce, C. S. *Peirce on signs: Writings on semiotic* (UNC Press Books, 1991).
20. Messaris, P. & Abraham, L. The role of images in framing news stories. In *Framing public life*, 231–242 (Routledge, 2001).
21. Glasford, D. E. Seeing is believing: communication modality, anger, and support for action on behalf of out-groups. *J. Appl. Soc. Psychol.* **43**, 2223–2230 (2013).
22. Yadav, A. *et al.* If a picture is worth a thousand words is video worth a million? differences in affective and cognitive processing of video and text cases. *J. Comput. High. Educ.* **23**, 15–37 (2011).

23. Appiah, O. Rich media, poor media: The impact of audio/video vs. text/picture testimonial ads on browsers' evaluations of commercial web sites and online products. *J. Curr. Issues & Res. Advert.* **28**, 73–86 (2006).
24. Powell, T. E., Boomgaarden, H. G., De Swert, K. & de Vreese, C. H. Video killed the news article? comparing multimodal framing effects in news videos and articles. *J. broadcasting & electronic media* **62**, 578–596 (2018).
25. Garimella, K. & Eckles, D. Images and misinformation in political groups: Evidence from whatsapp in india. *Harv. Kennedy Sch. Misinformation Rev.* (2020).
26. Goel, V., Raj, S. & Ravichandran, P. How whatsapp leads mobs to murder in india. *The New York Times* **18** (2018).
27. Sundar, S. S., Molina, M. D. & Cho, E. Seeing is believing: Is video modality more powerful in spreading fake news via online messaging apps? *J. Comput. Commun.* **26**, 301–319 (2021).
28. Wittenberg, C., Tappin, B. M., Berinsky, A. J. & Rand, D. G. The (minimal) persuasive advantage of political video over text. *Proc. Natl. Acad. Sci.* **118** (2021).
29. Sundar, S. S. The main model: A heuristic approach to understanding technology effects on credibility. *Digit. Media, Youth, Credibil.* (2008).
30. Hancock, J. T., Naaman, M. & Levy, K. Ai-mediated communication: definition, research agenda, and ethical considerations. *J. Comput. Commun.* **25**, 89–100 (2020).
31. Barari, S., Lucas, C. & Munger, K. Political Deepfake Videos Misinform the Public, But No More than Other Fake Media. preprint, Open Science Framework (2021). DOI: [10.31219/osf.io/cdfh3](https://doi.org/10.31219/osf.io/cdfh3). Accessed 2021-01-14.
32. Murphy, G. & Flynn, E. Deepfake false memories. *Memory* **30**, 480–492 (2022).
33. Vaccari, C. & Chadwick, A. Deepfakes and disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news. *Soc. Media+ Soc.* **6**, 2056305120903408 (2020).
34. Dobber, T., Metoui, N., Trilling, D., Helberger, N. & de Vreese, C. Do (microtargeted) deepfakes have real effects on political attitudes? *The Int. J. Press.* **26**, 69–91 (2021).
35. Hameleers, M., van der Meer, T. G. & Dobber, T. You won't believe what they just said! the effects of political deepfakes embedded as vox populi on social media. *Soc. Media+ Soc.* **8**, 20563051221116346 (2022).
36. Reeves, B., Yeykelis, L. & Cummings, J. J. The use of media in media psychology. *Media Psychol.* **19**, 49–71 (2016).
37. Kasra, M., Shen, C. & O'Brien, J. F. Seeing is believing: How people fail to identify fake images on the web. In *Extended abstracts of the 2018 CHI conference on human factors in computing systems*, 1–6 (2018).
38. Hameleers, M., Powell, T. E., Van Der Meer, T. G. & Bos, L. A picture paints a thousand lies? the effects and mechanisms of multimodal disinformation and rebuttals disseminated via social media. *Polit. Commun.* **37**, 281–301 (2020).
39. Nightingale, S. J. & Farid, H. Ai-synthesized faces are indistinguishable from real faces and more trustworthy. *Proc. Natl. Acad. Sci.* **119** (2022).
40. Cardwell, B. A., Henkel, L. A., Garry, M., Newman, E. J. & Foster, J. L. Nonprobative photos rapidly lead people to believe claims about their own (and other people's) pasts. *Mem. & cognition* **44**, 883–896 (2016).
41. Cardwell, B. A., Lindsay, D. S., Förster, K. & Garry, M. Uninformative photos can increase people's perceived knowledge of complicated processes. *J. Appl. Res. Mem. Cogn.* **6**, 244–252 (2017).
42. Newman, E. J., Jalbert, M. C., Schwarz, N. & Ly, D. P. Truthiness, the illusory truth effect, and the role of need for cognition. *Conscious. Cogn.* **78**, 102866 (2020).
43. Newman, E. J., Garry, M., Bernstein, D. M., Kantner, J. & Lindsay, D. S. Nonprobative photographs (or words) inflate truthiness. *Psychon. Bull. & Rev.* **19**, 969–974 (2012).
44. Fazio, L. K., Brashier, N. M., Payne, B. K. & Marsh, E. J. Knowledge does not protect against illusory truth. *J. Exp. Psychol. Gen.* **144**, 993 (2015).
45. Ecker, U. K. *et al.* The psychological drivers of misinformation belief and its resistance to correction. *Nat. Rev. Psychol.* **1**, 13–29 (2022).
46. Dolhansky, B. *et al.* The DeepFake Detection Challenge (DFDC) Dataset. *arXiv:2006.07397 [cs]* (2020). Accessed 2021-05-15.
47. Groh, M., Epstein, Z., Firestone, C. & Picard, R. Deepfake detection by human crowds, machines, and machine-informed crowds. *Proc. Natl. Acad. Sci.* **119** (2022).

48. Köbis, N., Doležalová, B. & Soraperra, I. Fooled twice—people cannot detect deepfakes but think they can. *iScience* **24** (2021).
49. Tahir, R. *et al.* Seeing is believing: Exploring perceptual differences in deepfake videos. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 1–16 (2021).
50. Lee, E.-J. & Shin, S. Y. Mediated misinformation: Questions answered, more questions to ask. *Am. Behav. Sci.* **65**, 259–276 (2021).
51. Pennycook, G. & Rand, D. G. Fighting misinformation on social media using crowdsourced judgments of news source quality. *Proc. Natl. Acad. Sci.* **116**, 2521–2526 (2019).
52. Austin, E. W. & Dong, Q. Source v. content effects on judgments of news believability. *Journalism quarterly* **71**, 973–983 (1994).
53. Shen, C. *et al.* Fake images: The effects of source, intermediary, and digital media literacy on contextual assessment of image credibility online. *New media & society* **21**, 438–463 (2019).
54. Dias, N., Pennycook, G. & Rand, D. G. Emphasizing publishers does not effectively reduce susceptibility to misinformation on social media. *Harv. Kennedy Sch. Misinformation Rev.* **1** (2020).
55. Jakesch, M., Koren, M., Evtushenko, A. & Naaman, M. The role of source, headline and expressive responding in political news evaluation. *Headl. Expressive Responding Polit. News Eval. (December 5, 2018)* (2018).
56. Nadarevic, L., Reber, R., Helmecke, A. J. & Köse, D. Perceived truth of statements and simulated social media postings: an experimental investigation of source credibility, repeated exposure, and presentation format. *Cogn. Res. Princ. Implic.* **5**, 1–16 (2020).
57. Kim, A., Moravec, P. L. & Dennis, A. R. Combating fake news on social media with source ratings: The effects of user and expert reputation ratings. *J. Manag. Inf. Syst.* **36**, 931–968 (2019).
58. Pennycook, G. & Rand, D. G. Lazy, not biased: Susceptibility to partisan fake news is better explained by lack of reasoning than by motivated reasoning. *Cognition* **188**, 39–50 (2019).
59. Appel, M. & Prieszel, F. The detection of political deepfakes. *J. Comput. Commun.* **27**, zmac008 (2022).
60. Arechar, A. A. *et al.* Understanding and reducing online misinformation across 16 countries on six continents. *PsyArXiv* (2022). Accessed 2022-02-22.
61. Pennycook, G. & Rand, D. G. The psychology of fake news. *Trends cognitive sciences* (2021).
62. Lazer, D. M. J. *et al.* The science of fake news. *Science* **359**, 1094–1096 (2018).
63. Dan, V. *et al.* Visual mis- and disinformation, social media, and democracy. *Journalism & Mass Commun. Q.* **98**, 641–664 (2021).
64. Calo, R., Coward, C., Spiro, E. S., Starbird, K. & West, J. D. How do you solve a problem like misinformation? *Sci. advances* **7**, eabn0481 (2021).
65. Sankaranarayanan, A., Groh, M., Picard, R. & Lippman, A. The presidential deepfakes dataset. In *Proceedings of the AIOFAI Workshop at the International Joint Conference on Artificial Intelligence* (2021).
66. Perov, I. *et al.* Deepfacelab: Integrated, flexible and extensible face-swapping framework. *arXiv preprint arXiv:2005.05535* (2020).
67. ElevenLabs. Elevenlabs.
68. Abadie, A., Athey, S., Imbens, G. & Wooldridge, J. When Should You Adjust Standard Errors for Clustering? *arXiv:1710.02926 [econ, math, stat]* (2017). Accessed 2021-12-15.
69. Frederick, S. Cognitive reflection and decision making. *J. Econ. perspectives* **19**, 25–42 (2005).
70. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. Royal statistical society: series B (Methodological)* **57**, 289–300 (1995).
71. Goodman, J. D. Microphone catches a candid obama (2012).
72. Lyu, S. DeepFake Detection: Current Challenges and Next Steps. *arXiv:2003.09234 [cs]* (2020). ArXiv: 2003.09234.
73. Bryan, C. J., Tipton, E. & Yeager, D. S. Behavioural science is unlikely to change the world without a heterogeneity revolution. *Nat. Hum. Behav.* **5**, 980–989, DOI: [10.1038/s41562-021-01143-3](https://doi.org/10.1038/s41562-021-01143-3) (2021).
74. Vincent, J. Ai image generator midjourney stops free trials but says influx of new users to blame (2023).

75. Metzger, M. J., Flanagin, A. J. & Medders, R. B. Social and heuristic approaches to credibility evaluation online. *J. communication* **60**, 413–439 (2010).
76. Barasch, A., Schroeder, J., Zev Berman, J. & Small, D. Cues to sincerity: How people assess and convey sincerity in language. *ACR North Am. Adv.* (2018).
77. Schroeder, J. & Epley, N. Mistaking minds and machines: How speech affects dehumanization and anthropomorphism. *J. Exp. Psychol. Gen.* **145**, 1427 (2016).
78. Franzen, L., Delis, I., Sousa, G. D., Kayser, C. & Piliastides, M. G. Auditory information enhances post-sensory visual evidence during rapid multisensory decision-making. *Nat. Commun.* **11** (2020).
79. Allen, J., Howland, B., Mobius, M., Rothschild, D. & Watts, D. J. Evaluating the fake news problem at the scale of the information ecosystem. *Sci. Adv.* **6**, eaay3539 (2020).
80. Watts, D. J., Rothschild, D. M. & Mobius, M. Measuring the news and its impact on democracy. *Proc. Natl. Acad. Sci.* **118** (2021).
81. Epstein, Z. *et al.* Art and the science of generative ai. *Science* **380**, 1110–1111 (2023).
82. Agarwal, S. & Farid, H. Detecting deep-fake videos from aural and oral dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 981–989 (2021).
83. Pennycook, G. *et al.* Shifting attention to accuracy can reduce misinformation online. *Nature* **592**, 590–595 (2021).
84. Epstein, Z., Sirlin, N., Arechar, A., Pennycook, G. & Rand, D. The social media context interferes with truth discernment. *Sci. Adv.* **9**, eabo6169 (2023).
85. Roozenbeek, J., van der Linden, S., Goldberg, B., Rathje, S. & Lewandowsky, S. Psychological inoculation improves resilience against misinformation on social media. *Sci. Adv.* **8**, eabo6254 (2022).
86. Berger, J. & Milkman, K. L. What makes online content viral? *J. marketing research* **49**, 192–205 (2012).
87. Vosoughi, S., Roy, D. & Aral, S. The spread of true and false news online. *Science* **359**, 1146–1151, DOI: [10.1126/science.aap9559](https://doi.org/10.1126/science.aap9559) (2018).
88. Brady, W. J., Wills, J. A., Jost, J. T., Tucker, J. A. & Van Bavel, J. J. Emotion shapes the diffusion of moralized content in social networks. *Proc. Natl. Acad. Sci.* **114**, 7313–7318 (2017).
89. Brady, W. J., Crockett, M. J. & Van Bavel, J. J. The mad model of moral contagion: The role of motivation, attention, and design in the spread of moralized content online. *Perspectives on Psychol. Sci.* **15**, 978–1010 (2020).
90. Lazer, D. Studying human attention on the internet. *Proc. Natl. Acad. Sci.* **117**, 21–22 (2020).
91. Fuller, T. *Gnomologia: adagies and proverbs; wise sentences and witty sayings, ancient and modern, foreign and British*, vol. 1 (B. Barker, 1732).
92. Messaris, P. *Visual persuasion: The role of images in advertising* (Sage, 1997).
93. Farid, H. Digital doctoring: how to tell the real from the fake. *Significance* **3**, 162–166 (2006).
94. King, D. *The commissar vanishes: The falsification of photographs and art in Stalin's Russia* (Metropolitan Books New York, 1997).
95. Lai, V. & Tan, C. On Human Predictions with Explanations and Predictions of Machine Learning Models: A Case Study on Deception Detection. *Proc. Conf. on Fairness, Accountability, Transpar.* 29–38, DOI: [10.1145/3287560.3287590](https://doi.org/10.1145/3287560.3287590) (2019). ArXiv: 1811.07901.
96. Agarwal, S. *et al.* Watch those words: Video falsification detection using word-conditioned facial motion. *arXiv preprint arXiv:2112.10936* (2021). Accessed 2022-02-14.
97. Gambino, L. Biden warns US democracy imperiled by Trump and Maga extremists. <https://www.theguardian.com/us-news/2022/sep/01/biden-speech-philadelphia-extremist-republicans-threaten-democracy> (2022). [Online; accessed 29-June-2023].
98. Son Chung, J., Senior, A., Vinyals, O. & Zisserman, A. Lip reading sentences in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 6447–6456 (2017).
99. Zhang, S. *et al.* S3fd: Single shot scale-invariant face detector. In *Proceedings of the IEEE international conference on computer vision*, 192–201 (2017).
100. Iglovikov, V. & Shvets, A. Terausnet: U-net with vgg11 encoder pre-trained on imagenet for image segmentation. *arXiv preprint arXiv:1801.05746* (2018).

101. Font, F., Roma, G. & Serra, X. Freesound technical demo. *Proc. 21st ACM international conference on Multimed.* (2013).
102. McCreery, A. & Calamia, P. Cross-modal perception of room acoustics. *The J. Acoust. Soc. Am.* **120**, 3150–3150 (2006).
103. Valente, D. L. & Braasch, J. Subjective scaling of spatial room acoustic parameters influenced by visual environmental cues. *The J. Acoust. Soc. Am.* **128**, 1952–1964 (2010).
104. Palan, S. & Schitter, C. Prolific.ac—A subject pool for online experiments. *J. Behav. Exp. Finance* **17**, 22–27, DOI: [10.1016/j.jbef.2017.12.004](https://doi.org/10.1016/j.jbef.2017.12.004) (2018).
105. Berinsky, A. J., Margolis, M. F. & Sances, M. W. Separating the shirkers from the workers? making sure respondents pay attention on self-administered surveys. *Am. J. Polit. Sci.* **58**, 739–753 (2014).

	<i>Dependent variable: Confidence Score</i>		
	All	Real	Fabricated
Constant	57.66*** (0.83)	63.48*** (1.08)	51.40*** (1.25)
Silent Video	6.56*** (1.20)	3.15* (1.53)	10.45*** (1.83)
Silent Video with Subtitles	8.73*** (1.12)	2.49 (1.52)	15.41*** (1.69)
Audio	19.44*** (1.17)	11.35*** (1.48)	27.99*** (1.68)
Audio with Subtitles	19.45*** (1.05)	10.50*** (1.36)	28.72*** (1.58)
Video with Audio	25.23*** (1.10)	17.63*** (1.42)	33.30*** (1.61)
Video with Audio and Subtitles	24.77*** (1.08)	14.95*** (1.49)	34.78*** (1.53)
Number of Individuals	501	501	501
Observations	16,011	8,004	8,007
$R^2$	0.07	0.04	0.12

*Note:*

\* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

**Table 2. Main Analysis for Experiment 1** Pre-registered main analysis for experiment 1a displaying ordinary least squares regressions with robust standard errors clustered on participant. Confidence score is the dependent variable, which is accuracy weighted by participants' confidence defined as the participant's confidence (ranging from 50 to 100) if correct and 100 minus the participant's confidence if incorrect. In each column, the independent variables indicate the modality with transcripts as the holdout. The first column shows all videos, the second column shows only real videos, and the third column shows only fabricated videos.

<i>Dependent variable: Accuracy</i>			
	(1)	(2)	(3)
Constant	0.852*** (0.044)	0.832*** (0.050)	0.833*** (0.020)
Voice Actor Videos (PDD)	-0.019 (0.047)		
Text-to-Speech Videos (PDD)	-0.130* (0.052)		-0.111*** (0.033)
Real Videos (PDD)	0.012 (0.050)	0.032 (0.055)	
Real Videos (Barari et al)	-0.020 (0.066)		
Number of Participants	302	302	302
Observations	5,964	2,383	2,390
R <sup>2</sup>	0.018	0.002	0.018

*Note:* \*p<0.05; \*\*p<0.01; \*\*\*p<0.001

**Table 3. Main Analysis for Experiment 2** Pre-registered main analysis for experiment 2 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. Accuracy is the dependent variable, which is a binary variable defined as 1 if participants accurately identify the stimuli as real and fake and 0 otherwise. In each column, the independent variables indicate the condition from which videos are drawn. In the first columns, the holdout condition is the fake videos used in Barari et al 2021. In the second column, the holdout condition is real videos used in Barari et al 2021. In the third column, the holdout condition is fake voice actor videos from the PDD.

Filename	Accuracy (TTS)	Accuracy (Voice Actor)	P-value	Significance (B-H)	Obs (TTS)	Obs (Voice Actor)
t-02	0.61	0.85	0.002	True	83	65
t-03	0.70	0.90	0.002	True	76	73
t-04	0.62	0.84	0.002	True	76	79
t-05	0.78	0.93	0.004	True	73	90
t-06	0.56	0.77	0.007	True	85	73
b-01	0.69	0.84	0.028	False	84	70
b-04	0.70	0.85	0.032	False	61	87
b-07	0.89	0.76	0.037	False	80	75
t-00	0.77	0.89	0.080	False	71	70
b-06	0.62	0.76	0.084	False	72	67
b-05	0.60	0.71	0.184	False	60	63
t-01	0.95	0.90	0.251	False	65	83
b-02	0.70	0.78	0.262	False	74	74
t-07	0.83	0.88	0.344	False	77	78
b-03	0.65	0.72	0.382	False	75	75
b-00	0.88	0.88	0.901	False	80	76

**Table 4. Text-to-Speech and Voice Actor Comparisons in Experiment 2** Pre-registered analysis comparing text-to-speech deepfakes videos to voice actor deepfake videos with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations for each text-to-speech deepfake and each voice actor deepfake.

	Correct Confidence	Response Time	Plays/Pauses
	(1)	(2)	(3)
Constant	0.650*** (0.086)	8.238*** (0.843)	0.125*** (0.017)
Voice Actor Videos (PDD)	-0.047 (0.090)	-0.768 (0.854)	-0.004 (0.017)
Text-to-Speech Videos (PDD)	-0.273** (0.099)	0.970 (0.870)	0.023 (0.017)
Real Videos (PDD)	-0.037 (0.096)	-0.113 (0.873)	-0.012 (0.017)
Real Videos (Barari et al)	-0.106 (0.120)	1.915 (0.981)	0.042 (0.034)
Number of Participants	302	302	302
Observations	5,964	5,964	5,964
$R^2$	0.024	0.009	0.003

Note:

\* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

**Table 5. Secondary Analysis for Experiment 2** Pre-registered secondary analysis for experiment 2 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. In column 1, “correct confidence” is the dependent variable, which is defined as  $1 - (5 - \text{confidence})/5$  if participants accurately identify the stimuli as real and fake and  $-(\text{confidence})/5$  otherwise. In column 2, windsorized marginal response time is defined as the response time minus the duration of the stimulus windsorized at the 5% and 95% values. In column 3, the dependent variable is a binary variable for playing or pausing the video more than once. In each column, the independent variables indicate the condition from which videos are drawn and the holdout condition is the fake videos used in Barari et al 2021.

	<i>Dependent variable: Accuracy</i>					
	All	Real	Fabricated	All	Real	Fabricated
	(1)	(2)	(3)	(4)	(5)	(6)
Constant	0.587*** (0.033)	0.614*** (0.041)	0.577*** (0.036)	0.602*** (0.033)	0.617*** (0.040)	0.542*** (0.049)
Audio	0.074*** (0.018)	0.034 (0.024)	0.112*** (0.025)	0.063** (0.021)	0.041 (0.023)	0.152** (0.048)
Silent Video with Subtitles	0.042* (0.018)	0.000 (0.013)	0.081** (0.031)	0.010 (0.017)	-0.016 (0.016)	0.114** (0.042)
Video with Audio	0.165*** (0.021)	0.119*** (0.025)	0.209*** (0.030)	0.148*** (0.024)	0.117*** (0.026)	0.272*** (0.049)
High Base Rate	-0.018 (0.030)	0.072*** (0.016)	-0.058*** (0.016)	-0.046 (0.033)	0.058* (0.023)	-0.016 (0.029)
Audio with Audio * High Base Rate				0.021 (0.023)	-0.033 (0.035)	-0.049 (0.043)
Silent Video * High Base Rate				0.062* (0.027)	0.075* (0.034)	-0.042 (0.035)
Video with Audio * High Base Rate				0.032 (0.028)	0.009 (0.028)	-0.078 (0.045)
Number of Participants	1008	1008	1008	1008	1008	1008
Observations	19,812	9,707	10,105	19,812	9,707	10,105
R <sup>2</sup>	0.017	0.014	0.026	0.017	0.016	0.027

Note:

\*p<0.05; \*\*p<0.01; \*\*\*p<0.001

**Table 6. Main Analysis for Experiment 3** Pre-registered main analysis for experiment 3 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. Accuracy is the dependent variable, which is a binary variable defined as 1 if participants accurately identify the stimuli as real and fake and 0 otherwise. In each column, the independent variables indicate the stimuli's modality, assignment to the high-base rate of fakes condition, and an interaction between modality and high base rate. The held out conditions are transcripts and assignment to low base rate of fakes.

Filename	Accuracy (Transcript)	Accuracy (Silent Video)	P-value	Significance (B-H)	Obs (Transcript)	Obs (Silent Video)
t-01	0.55	0.80	0.000	True	187	155
b-01	0.41	0.59	0.001	True	143	150
t-00	0.65	0.82	0.001	True	159	153
b-06	0.41	0.59	0.001	True	142	147
t-06	0.32	0.50	0.002	True	141	149
b-00	0.80	0.66	0.003	True	144	148
t-02	0.36	0.51	0.009	True	163	152
b-05	0.58	0.71	0.014	False	157	163
b-05-p	0.85	0.73	0.016	False	144	160
t-05	0.41	0.54	0.018	False	148	156
b-04	0.45	0.59	0.018	False	150	154
b-02-p	0.32	0.43	0.061	False	143	153
b-07	0.53	0.63	0.068	False	182	157
b-00-p	0.55	0.63	0.160	False	153	136
t-07-p	0.75	0.69	0.189	False	156	171
t-07	0.74	0.68	0.281	False	130	154
t-04-p	0.46	0.52	0.313	False	162	157
b-03	0.69	0.64	0.362	False	148	163
t-03-p	0.70	0.66	0.513	False	161	162
b-04-p	0.74	0.72	0.621	False	171	131
b-02	0.46	0.49	0.623	False	149	146
t-00-p	0.54	0.56	0.674	False	142	163
t-02-p	0.70	0.72	0.705	False	147	165
b-06-p	0.78	0.79	0.737	False	151	153
t-05-p	0.64	0.66	0.753	False	165	139
b-03-p	0.55	0.57	0.768	False	159	149
b-01-p	0.54	0.55	0.788	False	165	160
t-06-p	0.68	0.69	0.790	False	174	156
b-07-p	0.35	0.33	0.791	False	150	149
t-01-p	0.84	0.85	0.795	False	155	162
t-04	0.61	0.60	0.874	False	168	135
t-03	0.41	0.42	0.954	False	163	171

**Table 7. Transcripts and Silent Videos Comparisons in Experiment 3** Pre-registered analysis for experiment 3 comparing transcripts to silent videos with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations.

Filename	Accuracy (Transcript)	Accuracy (Audio Only)	P-value	Significance (B-H)	Obs (Transcript)	Obs (Audio Only)
t-01	0.55	0.74	0.000	True	187	151
t-06-p	0.68	0.85	0.000	True	143	129
b-02-p	0.32	0.56	0.000	True	159	151
t-05	0.41	0.66	0.000	True	142	165
t-03	0.41	0.72	0.000	True	141	164
t-02	0.36	0.60	0.000	True	144	143
b-01	0.41	0.59	0.001	True	163	165
t-05-p	0.64	0.79	0.004	True	157	151
b-07	0.53	0.68	0.006	True	144	149
b-04-p	0.74	0.60	0.007	True	148	159
t-01-p	0.84	0.92	0.022	False	150	150
t-07	0.74	0.84	0.029	False	143	164
b-04	0.45	0.57	0.042	False	182	144
t-02-p	0.70	0.59	0.049	False	153	154
b-06-p	0.78	0.86	0.072	False	156	139
t-06	0.32	0.41	0.092	False	130	147
b-00	0.80	0.86	0.155	False	162	139
b-07-p	0.35	0.43	0.162	False	148	143
t-04	0.61	0.67	0.225	False	161	170
t-00-p	0.54	0.61	0.243	False	171	161
b-06	0.41	0.47	0.338	False	149	174
b-05-p	0.85	0.82	0.591	False	142	145
t-00	0.65	0.68	0.609	False	147	161
b-03-p	0.55	0.58	0.613	False	151	163
b-02	0.46	0.49	0.641	False	165	161
t-03-p	0.70	0.67	0.697	False	159	162
b-05	0.58	0.56	0.728	False	165	161
b-03	0.69	0.71	0.755	False	174	135
t-07-p	0.75	0.76	0.864	False	150	160
b-00-p	0.55	0.54	0.872	False	155	148
b-01-p	0.54	0.53	0.890	False	168	165
t-04-p	0.46	0.46	0.945	False	163	156

**Table 8. Transcripts and Audios Comparisons in Experiment 3** Pre-registered analysis for experiment 3 comparing transcripts to audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations.

Filename	Accuracy (Transcript)	Accuracy (Video with Audio)	P-value	Significance (B-H)	Obs (Transcript)	Obs (Video with Audio)
t-05-p	0.64	0.87	0.000	True	187	151
t-00-p	0.54	0.80	0.000	True	143	168
b-07	0.53	0.83	0.000	True	159	158
t-03	0.41	0.65	0.000	True	142	154
t-05	0.41	0.81	0.000	True	141	166
t-01	0.55	0.84	0.000	True	144	167
b-04	0.45	0.75	0.000	True	163	156
t-02	0.36	0.68	0.000	True	157	146
t-06-p	0.68	0.86	0.000	True	144	177
b-02-p	0.32	0.62	0.000	True	148	144
b-02	0.46	0.71	0.000	True	150	169
b-01	0.41	0.78	0.000	True	143	132
b-00-p	0.55	0.80	0.000	True	182	149
b-05	0.58	0.76	0.001	True	153	159
b-01-p	0.54	0.72	0.001	True	156	160
t-07	0.74	0.89	0.001	True	130	146
t-00	0.65	0.81	0.001	True	162	166
t-06	0.32	0.49	0.001	True	148	148
b-07-p	0.35	0.52	0.003	True	161	134
b-06	0.41	0.56	0.009	True	171	145
b-00	0.80	0.90	0.012	True	149	175
b-04-p	0.74	0.84	0.042	False	142	161
t-01-p	0.84	0.91	0.051	False	147	176
t-04-p	0.46	0.57	0.057	False	151	151
b-06-p	0.78	0.86	0.077	False	165	168
t-07-p	0.75	0.83	0.101	False	159	146
t-04	0.61	0.69	0.107	False	165	157
t-03-p	0.70	0.62	0.183	False	174	140
b-03-p	0.55	0.60	0.394	False	150	174
t-02-p	0.70	0.73	0.666	False	155	160
b-05-p	0.85	0.86	0.819	False	168	152
b-03	0.69	0.70	0.916	False	163	137

**Table 9. Transcripts and Videos Comparisons in Experiment 3** Pre-registered analysis for experiment 3 comparing transcripts to video with audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations.

Filename	Accuracy (Silent Video)	Accuracy (Video with Audio)	P-value	Significance (B-H)	Obs (Silent Video)	Obs (Video with Audio)
b-00	0.66	0.90	0.000	True	155	151
t-06-p	0.69	0.86	0.000	True	150	168
b-01	0.59	0.78	0.000	True	153	158
t-05-p	0.66	0.87	0.000	True	147	154
b-02	0.49	0.71	0.000	True	149	166
t-05	0.54	0.81	0.000	True	148	167
t-03	0.42	0.65	0.000	True	152	156
t-00-p	0.56	0.80	0.000	True	163	146
t-07	0.68	0.89	0.000	True	160	177
b-07	0.63	0.83	0.000	True	156	144
b-07-p	0.33	0.52	0.001	True	154	169
b-00-p	0.63	0.80	0.001	True	153	132
b-02-p	0.43	0.62	0.001	True	157	149
b-04	0.59	0.75	0.001	True	136	159
b-01-p	0.55	0.72	0.002	True	171	160
t-02	0.51	0.68	0.002	True	154	146
t-07-p	0.69	0.83	0.004	True	157	166
b-05-p	0.73	0.86	0.010	True	163	148
b-04-p	0.72	0.84	0.011	True	162	134
t-04	0.60	0.69	0.089	False	131	145
t-01-p	0.85	0.91	0.108	False	146	175
b-06-p	0.79	0.86	0.167	False	163	161
t-01	0.80	0.84	0.299	False	165	176
b-03	0.64	0.70	0.315	False	153	151
t-04-p	0.52	0.57	0.374	False	139	168
b-05	0.71	0.76	0.381	False	149	146
t-03-p	0.66	0.62	0.496	False	160	157
b-06	0.59	0.56	0.534	False	156	140
b-03-p	0.57	0.60	0.568	False	149	174
t-02-p	0.72	0.73	0.955	False	162	160
t-00	0.82	0.81	0.963	False	135	152
t-06	0.50	0.49	0.966	False	171	137

**Table 10. Silent Videos and Video Comparisons in Experiment 3** Pre-registered analysis for experiment 3 comparing silent video to video with audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations.

Filename	Accuracy (Silent Video)	Accuracy (Audio Only)	P-value	Significance (B-H)	Obs (Silent Video)	Obs (Audio Only)
b-00	0.66	0.86	0.000	True	155	151
t-03	0.42	0.72	0.000	True	150	129
t-06-p	0.69	0.85	0.001	True	153	151
t-07	0.68	0.84	0.002	True	147	165
b-05	0.71	0.56	0.005	True	149	164
t-00	0.82	0.68	0.006	True	148	143
t-05-p	0.66	0.79	0.012	False	152	165
t-02-p	0.72	0.59	0.015	False	163	151
t-05	0.54	0.66	0.020	False	160	149
b-02-p	0.43	0.56	0.022	False	156	159
b-04-p	0.72	0.60	0.024	False	154	150
b-06	0.59	0.47	0.027	False	153	164
b-05-p	0.73	0.82	0.051	False	157	144
t-01-p	0.85	0.92	0.054	False	136	154
b-07-p	0.33	0.43	0.082	False	171	139
t-02	0.51	0.60	0.108	False	154	147
b-00-p	0.63	0.54	0.125	False	157	139
t-06	0.50	0.41	0.139	False	163	143
t-07-p	0.69	0.76	0.142	False	162	170
b-06-p	0.79	0.86	0.157	False	131	161
t-04	0.60	0.67	0.187	False	146	174
t-01	0.80	0.74	0.191	False	163	145
b-03	0.64	0.71	0.222	False	165	161
b-07	0.63	0.68	0.290	False	153	163
t-04-p	0.52	0.46	0.344	False	139	161
t-00-p	0.56	0.61	0.438	False	149	162
b-01-p	0.55	0.53	0.676	False	160	161
b-04	0.59	0.57	0.763	False	156	135
t-03-p	0.66	0.67	0.782	False	149	160
b-03-p	0.57	0.58	0.827	False	162	148
b-01	0.59	0.59	0.924	False	135	165
b-02	0.49	0.49	0.970	False	171	156

**Table 11. Silent Videos and Audios Comparisons in Experiment 3** Pre-registered analysis for experiment 3 comparing silent video to audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations.

Filename	Accuracy (Audio Only)	Accuracy (Video with Audio)	P-value	Significance (B-H)	Obs (Audio Only)	Obs (Video with Audio)
b-04-p	0.60	0.84	0.000	True	151	151
b-00-p	0.54	0.80	0.000	True	129	168
b-01	0.59	0.78	0.000	True	151	158
b-01-p	0.53	0.72	0.000	True	165	154
b-02	0.49	0.71	0.000	True	164	166
t-00-p	0.61	0.80	0.000	True	143	167
b-05	0.56	0.76	0.000	True	165	156
b-04	0.57	0.75	0.001	True	151	146
b-07	0.68	0.83	0.003	True	149	177
t-05	0.66	0.81	0.003	True	159	144
t-00	0.68	0.81	0.006	True	150	169
t-02-p	0.59	0.73	0.013	True	164	132
t-01	0.74	0.84	0.023	False	144	149
t-04-p	0.46	0.57	0.065	False	154	159
t-05-p	0.79	0.87	0.080	False	139	160
b-07-p	0.43	0.52	0.116	False	147	146
b-06	0.47	0.56	0.117	False	139	166
t-06	0.41	0.49	0.135	False	143	148
t-02	0.60	0.68	0.137	False	170	134
t-07-p	0.76	0.83	0.143	False	161	145
t-07	0.84	0.89	0.184	False	174	175
t-03	0.72	0.65	0.186	False	145	161
b-00	0.86	0.90	0.288	False	161	176
b-02-p	0.56	0.62	0.308	False	163	151
t-03-p	0.67	0.62	0.333	False	161	168
b-05-p	0.82	0.86	0.447	False	162	146
t-04	0.67	0.69	0.703	False	161	157
b-03-p	0.58	0.60	0.727	False	135	140
t-01-p	0.92	0.91	0.781	False	160	174
b-03	0.71	0.70	0.839	False	148	160
t-06-p	0.85	0.86	0.904	False	165	152
b-06-p	0.86	0.86	0.964	False	156	137

**Table 12. Audios and Videos Comparisons in Experiment 3** Pre-registered analysis for experiment 3 comparing audio to video with audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations.

	Correct Confidence (1)	Response Time (2)	Plays/Pauses (3)
Constant	0.152** (0.050)	10.813*** (0.764)	0.098*** (0.008)
Audio	0.115*** (0.031)		
Silent Video	0.019 (0.027)	6.744*** (1.839)	0.151*** (0.012)
Video with Audio	0.254*** (0.036)	1.909 (1.122)	0.030*** (0.009)
High Base Rate	-0.061 (0.049)	-0.817 (0.933)	-0.004 (0.011)
Audio with Audio * High Base Rate	0.026 (0.035)		
Silent Video * High Base Rate	0.093* (0.042)	-1.259 (2.057)	0.009 (0.015)
Video with Audio * High Base Rate	0.058 (0.045)	0.818 (1.354)	0.019 (0.014)
Number of Participants	1008	1008	1008
Observations	19,812	14,838	14,838
R <sup>2</sup>	0.023	0.003	0.032

Note:

\*p<0.05; \*\*p<0.01; \*\*\*p<0.001

**Table 13. Secondary Analysis in Experiment 3** Pre-registered secondary analysis for experiment 3 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. In column 1, “correct confidence” is the dependent variable, which is defined as  $1 - (5 - \text{confidence})/5$  if participants accurately identify the stimuli as real and fake and  $-(\text{confidence})/5$  otherwise. In column 2, windsorized marginal response time is defined as the response time minus the duration of the stimulus windsorized at the 5% and 95% values. In column 3, the dependent variable is a binary variable for playing or pausing the video more than once. In each column, the independent variables indicate the stimuli’s modality, assignment to the high-base rate of fakes condition, and an interaction between modality and high base rate. The held out conditions are transcripts in columns 1 and audio in columns 2 and 3 and assignment to low base rate of fakes in all columns.

<i>Dependent variable: Accuracy</i>	
	Correct (1)
Constant	0.756*** (0.031)
Voice Actor	0.092** (0.034)
Video with Audio	0.088** (0.030)
Voice Actor * Video with Audio	-0.187*** (0.044)
Number of Participants	206
Observations	3,215
$R^2$	0.014

*Note:* \*p<0.05; \*\*p<0.01; \*\*\*p<0.001

**Table 14. Main Analysis in Experiment 4** Pre-registered main analysis for experiment 4 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. Accuracy is the dependent variable, which is a binary variable defined as 1 if participants accurately identify the stimuli as real and fake and 0 otherwise. In each column, the independent variables indicate whether the stimulus was from the voice actor, whether the stimulus was shown as an audio or video with audio, and an interaction between these two variables.

Filename	Accuracy (Voice Actor)	Accuracy (Real)	P-value	Significance (B-H)	Obs (Voice Actor)	Obs (Real)
b-05-p	0.51	0.91	0.000	True	45	62
b-06-p	0.66	0.95	0.000	True	54	56
t-03-p	0.93	0.55	0.000	True	54	53
t-01-p	0.71	0.96	0.001	True	54	55
t-06-p	0.72	0.94	0.001	True	65	47
b-00-p	0.60	0.85	0.002	True	57	54
t-07-p	0.69	0.92	0.002	True	47	60
b-04-p	0.78	0.98	0.003	True	57	53
t-05-p	0.72	0.89	0.026	True	57	52
b-02-p	0.72	0.87	0.063	False	62	47
t-04-p	0.82	0.69	0.105	False	60	49
b-07-p	0.89	0.79	0.141	False	67	42
b-03-p	0.80	0.67	0.147	False	56	51
t-02-p	0.68	0.80	0.189	False	47	63
t-00-p	0.86	0.85	0.844	False	43	66
b-01-p	0.80	0.80	0.925	False	48	61

**Table 15. Voice Actor and Original Speaker in Videos Comparisons in Experiment 4** Pre-registered analysis comparing real videos with voice actor audio to real videos with the original speaker's audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations for each text-to-speech deepfake and each voice actor deepfake.

Filename	Accuracy (Voice Actor)	Accuracy (Real)	P-value	Significance (B-H)	Obs (Voice Actor)	Obs (Real)
b-02-p	0.93	0.51	0.000	True	49	44
t-00-p	0.93	0.64	0.001	True	46	47
b-07-p	0.87	0.70	0.038	False	54	37
t-04-p	0.88	0.70	0.040	False	51	41
b-00-p	0.82	0.64	0.051	False	47	44
t-02-p	0.78	0.59	0.056	False	43	48
b-03-p	0.82	0.68	0.119	False	48	45
t-07-p	0.85	0.93	0.221	False	47	46
t-05-p	0.86	0.79	0.339	False	46	45
t-03-p	0.76	0.72	0.639	False	47	43
b-06-p	0.83	0.87	0.657	False	49	44
t-06-p	0.90	0.92	0.728	False	34	57
b-05-p	0.81	0.83	0.811	False	42	47
b-01-p	0.83	0.81	0.829	False	37	56
b-04-p	0.85	0.86	0.866	False	41	52
t-01-p	0.83	0.84	0.926	False	52	42

**Table 16. Voice Actor and Original Speaker in Audios Comparisons in Experiment 4** Pre-registered analysis comparing voice actor audio to the original speaker’s audio with p-values from t-tests and statistical significance based on controlling the false discovery rate using the Benjamini-Hochberg procedure. The last two columns indicate the number of observations for each text-to-speech deepfake and each voice actor deepfake.

	Correct Confidence (1)	Response Time (2)	Plays/Pauses (3)
Constant	0.425*** (0.053)	7.847*** (0.509)	0.068*** (0.015)
Voice Actor	0.186** (0.059)	-1.539*** (0.408)	-0.001 (0.015)
Video with Audio	0.154** (0.051)	-0.658 (0.695)	0.007 (0.020)
Voice Actor * Video with Audio	-0.325*** (0.075)	0.998 (0.600)	0.020 (0.022)
Number of Participants	206	206	206
Observations	3,215	3,215	3,215
R <sup>2</sup>	0.016	0.005	0.002

Note:

\*p<0.05; \*\*p<0.01; \*\*\*p<0.001

**Table 17. Secondary Analysis in Experiment 4** Pre-registered secondary analysis for experiment 4 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. In column 1, “correct confidence” is the dependent variable, which is defined as 1 - (5-confidence)/5 if participants accurately identify the stimuli as real and fake and -(confidence)/5 otherwise. In column 2, windsorized marginal response time is defined as the response time minus the duration of the stimulus windsorized at the 5% and 95% values. In column 3, the dependent variable is a binary variable for playing or pausing the video more than once. In each column, the independent variables indicate whether the stimulus was from the voice actor, whether the stimulus was shown as an audio or video with audio, and an interaction between these two variables.

	Accurate Suspicion of Fake			Inaccurate Suspicion of Fake		
	(1)	(2)	(3)	(4)	(5)	(6)
Constant	0.073*** (0.017)	0.035 (0.026)	0.056* (0.026)	0.048** (0.017)	0.017 (0.017)	0.031 (0.018)
Silent Video	0.102** (0.032)	0.103** (0.032)	0.040 (0.054)	0.012 (0.022)	0.014 (0.022)	-0.024 (0.019)
Audio	0.058* (0.028)	0.063* (0.027)	0.072 (0.038)	0.038 (0.029)	0.037 (0.029)	0.032 (0.029)
Video	0.252*** (0.034)	0.251*** (0.034)	0.213*** (0.037)	0.013 (0.025)	0.011 (0.025)	0.001 (0.022)
Primed		0.074 (0.039)	0.032 (0.038)		0.063*** (0.019)	0.035 (0.030)
Primed * Silent Video			0.127* (0.064)			0.079 (0.041)
Primed * Audio			-0.027 (0.056)			0.010 (0.047)
Primed * Video			0.076 (0.071)			0.022 (0.043)
Number of Participants	200	200	200	200	200	200
Observations	1,018	1,018	1,018	982	982	982
$R^2$	0.060	0.069	0.076	0.003	0.020	0.024

Note:

\* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

**Table 18. Main Analysis for Experiment 5** Pre-registered main analysis for experiment 5 displaying ordinary least squares regressions with robust standard errors clustered on participants and stimuli. Accurate suspicion of a fake (true positive) is the dependent variable in columns 1 through 3 and inaccurate suspicion of a fake (false positive) is the dependent variable in columns 4 through 6. The independent variables indicate the stimuli's modality (Silent Video, Audio, and Video), assignment to seeing the deepfake attention check as the first stimulus (Primed), and interactions between these independent variables. The held out condition in this regression is the stimuli appearing as text transcripts.

Source	Exp.	Politician	Deepfake	Audio	Resolution	Speech ID	Video Source
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	640x360	t-00-p	<a href="#">59:06 to 59:27</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-01-p	<a href="#">00:10 to 00:30</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-02-p	<a href="#">00:12 to 00:33</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-03-p	<a href="#">03:19 to 03:38</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-04-p	<a href="#">00:00 to 00:21</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-05-p	<a href="#">00:17 to 00:37</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-06-p	<a href="#">00:16 to 00:38</a>
PDD (Youtube)	1-5	Donald Trump	0	Original/Actor	1280x720	t-07-p	<a href="#">00:18 to 00:40</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	1024x576	b-00-p	<a href="#">00:09 to 00:34</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	1280x720	b-01-p	<a href="#">05:43 to 05:55</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	640x360	b-02-p	<a href="#">01:00 to 01:15</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	1280x720	b-03-p	<a href="#">00:02 to 00:21</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	1280x720	b-04-p	<a href="#">28:55 to 29:15</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	1280x720	b-05-p	<a href="#">06:21 to 06:48</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	720x720	b-06-p	<a href="#">00:00 to 00:20</a>
PDD (Youtube)	1-5	Joseph Biden	0	Original/Actor	640x360	b-07-p	<a href="#">00:08 to 00:29</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1920x1080	b-00	<a href="#">02:21 to 02:45</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1280x720	b-01	<a href="#">00:38 to 01:02</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1920x1080	b-02	<a href="#">00:15 to 00:40</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1920x1080	b-03	<a href="#">00:02 to 00:22</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1280x720	b-04	<a href="#">01:13 to 01:42</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1280x720	b-05	<a href="#">10:40 to 11:02</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1920x1080	b-06	<a href="#">00:24 to 00:51</a>
PDD (Youtube)	1-3, 5	Joseph Biden	1	TTS/Actor	1920x1080	b-07	<a href="#">13:13 to 13:33</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-00	<a href="#">00:02 to 00:23</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-01	<a href="#">00:00 to 00:21</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-02	<a href="#">00:20 to 00:38</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-03	<a href="#">24:53 to 25:14</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1280x720	t-04	<a href="#">00:00 to 00:22</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-05	<a href="#">00:10 to 00:35</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-06	<a href="#">21:34 to 00:21</a>
PDD (Youtube)	1-3, 5	Donald Trump	1	TTS/Actor	1920x1080	t-07	<a href="#">00:00 to 00:25</a>
Barari (Agarwal)	2	Bernie Sanders	1	Voice Actor	1920x1080	o-bernie_1	<a href="#">00:00 to 00:17</a>
Barari (Youtube)	2	Boris Johnson	1	Voice Actor	1280x720	o-boris_1	<a href="#">00:00 to 00:20</a>
Barari (Agarwal)	2	Hilary Clinton	1	Voice Actor	1920x1080	o-hilary_1	<a href="#">00:00 to 00:18</a>
Barari (Youtube)	2	Donald Trump	1	Voice Actor	720x408	o-trump_1	<a href="#">00:19 to 00:38</a>
Barari (Youtube)	2	Donald Trump	1	Voice Actor	1280x720	o-trump_2	<a href="#">00:02 to 00:22</a>
Barari (Youtube)	2	Barack Obama	1	Voice Actor	1280x720	o-obama_1	<a href="#">00:00 to 00:20</a>
Barari (Youtube)	2	Donald Trump	0	Original	1280x720	o-trump_91	<a href="#">00:00 to 00:20</a>
Barari (Youtube)	2	Donald Trump	0	Original	1280x720	o-trump_92	<a href="#">00:23 to 00:43</a>
Barari (Youtube)	2	Elizabeth Warren	0	Original	1280x720	o-warren_9	<a href="#">00:00 to 00:20</a>
Barari (Youtube)	2	Joseph Biden	0	Original	1280x720	o-biden_91	<a href="#">00:31 to 00:51</a>
Barari (Youtube)	2	Barack Obama	0	Original	640x360	o-obama_91	<a href="#">00:02 to 00:18</a>
Barari (Youtube)	2	Barack Obama	0	Original	1280x720	o-obama_92	<a href="#">00:14 to 00:30</a>

**Table 19. Video Details and Sources** This table presents details on 44 real and deepfake video stimuli. The Presidential Deepfakes Dataset (PDD) includes 16 real videos and 16 deepfakes (with both voice actor and text-to-speech (TTS) audio). The other 12 videos, which were evaluated in Barari et al 2021<sup>31</sup> include 6 real videos and 6 deepfakes with voice actor audio (4 of these deepfakes are publicly available online and 2 of these deepfake are available upon request from the authors of Agarwal et al 2019<sup>5</sup>.) The voice actor deepfakes appear only in Experiment 1 and 2, the TTS deepfakes appear only in Experiment 2 and 3, and the real videos with actor audio appear only in Experiment 4. The “Video Source” column presents links to the source videos for each of the deepfakes.

[Instructions](#)

## Did Joseph Biden say that?

Please watch and listen to this video from Joseph Biden and share how confident you are that it is fabricated. Remember half the media snippets we show are real and half are fabricated.



I have seen/heard/read this before.

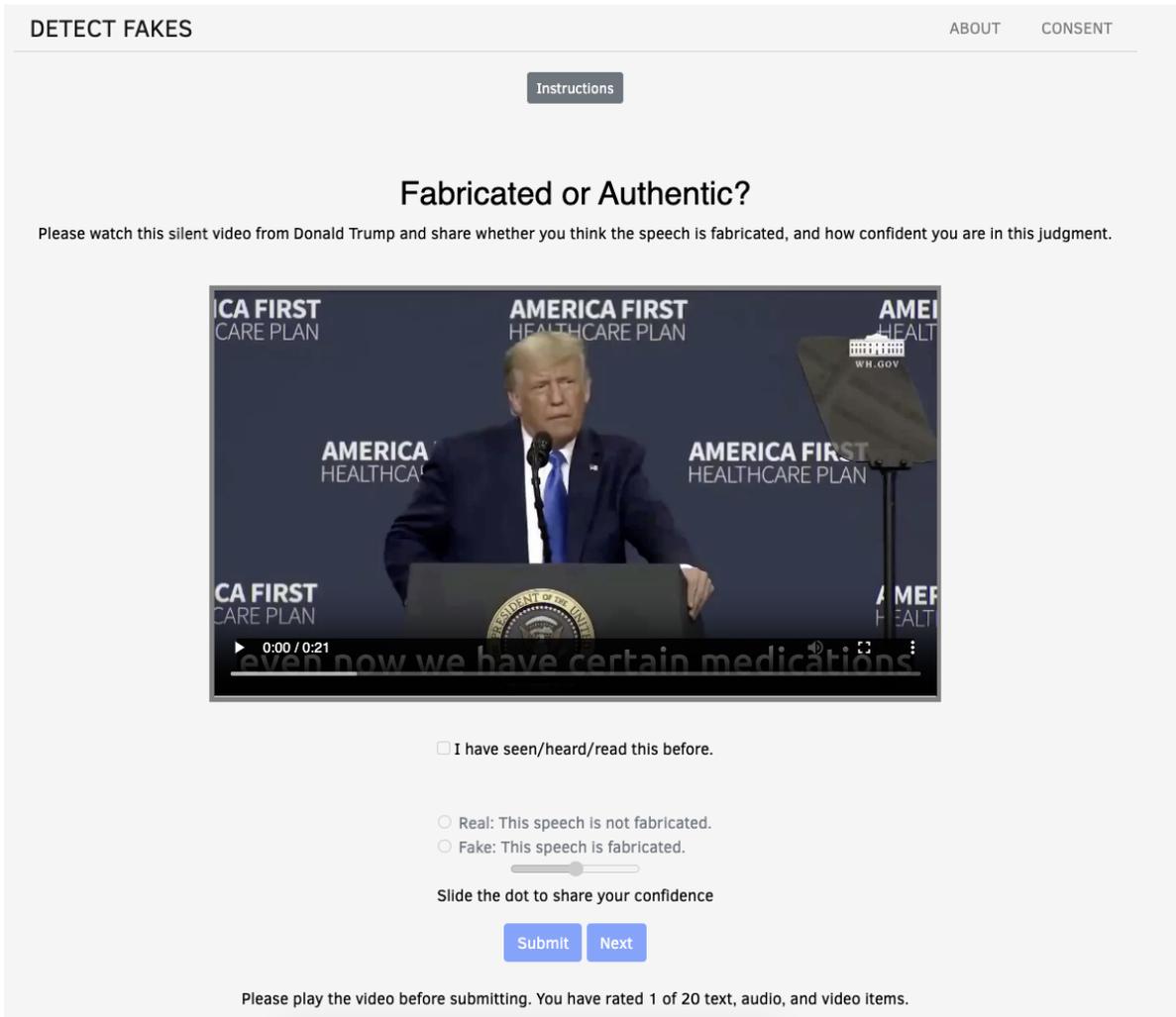


Just as likely real or fabricated

[Submit](#)[Next](#)

Please play the video before submitting. You have rated 1 of 32 text, audio, and video examples.

**Figure 9.** Screenshot of user interface in Experiment 1



**Figure 10. Screenshot of user interface in Experiment 2 and 3** The user interface for Experiment 4 is nearly identical to the user interface in Experiment 2 and 3 with the following two changes: First, the instructions at the top say, “Please watch and listen to this video from Donald Trump and share whether you think the voice is Donald Trump’s or a voice actor, and how confident you are in this judgment.” Second, the multiple choice options are “Real: The audio is from the speaker” and “Fake: The audio is from a voice actor.”

What comes to mind after watching the following video?



Please share your observations below:



Powered by Qualtrics [↗](#)

**Figure 11. Screenshot of user interface in Experiment 5** The user interface for Experiment 5 is built with Qualtrics and does not reveal that the dependent variable of interest in this experiment is the suspicion of a fabrication.