

RED TEAM REPORT

OPERATION 1 OF 3

**ANYONE WITH £100 AND A LAPTOP
CAN CREATE 10 HOURS OF
DEEPPFAKED 'FAKE NEWS'**

**UK PUBLIC MUST EXPECT 'FAKE
NEWS' AVALANCHE FOR GENERAL
ELECTION**

**DEEPPFAKE TECH IMPROVING AND COST
DROPPING EACH MONTH: PUBLIC
AWARENESS IS THE FIRST STEP
TOWARDS SAFEGUARDING DEMOCRACY**

Red Team Report

Part I: AI-Generated Media

Fenimore Harper Communications has conducted a 'Red Team' exercise on the generation of misinformation using AI to understand what techniques may be used by people. We have focused on the widely available, cheap-to-use tools that are likely to have the biggest influence on the public at large.

Executive Summary

- **Plummeting costs** of AI 'deep-fakes' have **greatly increased the risk of AI-powered misinformation** playing a role in the UK General Election.
- A 'Red Team' operation found that over **1,200 'deep-fake' misinformation short form clips** (over **10 hours of content**) can now be created **within an hour for just £100**.
- **Operation outputs: 9 deep-faked videos** of Prime Minister Rishi Sunak, Opposition Leader Keir Starmer and mainstream news presenters produced for **less than £1 in under 30 minutes**.
- Most widely available platforms allow users to create **high-fidelity deep-fakes of news presenters** within minutes.
- Widely available large language models (such as ChatGPT) can be made to produce **harmful misinformation with less than 200 characters of 'prompt hacking'**.
- Open source models allow for **voice clones of political candidates to be created within minutes**, despite being blocked by major platforms.

Introduction

In January 2024, Fenimore Harper uncovered [over 100 deep-fakes of Prime Minister Rishi Sunak](#) running as paid adverts on Meta's advertising platform.

These adverts were a financially-motivated scam. By leveraging the authority of the Prime Minister (and others, such as Elon Musk and newsreader Sophie Raworth), scammers were able to direct unsuspecting users to a fake investment platform which would attempt to steal their money.

Since then, we have seen hundreds more deep-faked adverts that have continued to run. Many of these also use the strategy of building short fake-news reports, stitching together deep-faked video, dubbed authentic video and manipulated news tickers.

In addition, politically-motivated misinformation has caused serious harm since the publication of our first report.

- A [voice-clone 'robocall' imitating US President Joe Biden](#) targeted New Hampshire voters and encouraged them to 'stay at home'.
- A [deep-faked video of a France24 newsreader](#) suggested that the Ukrainian government had planned an assassination attempt on French President Emmanuel Macron.
- [Fake clips of Slovakia's liberal party leader Michal Šimečka](#) appeared to show him discussing vote rigging and raising the price of beer.

These examples were quickly debunked by mainstream news media outlets, but had already been shared widely after their initial creation.

Back in 2023, [we warned of an 'explosion' of AI-powered misinformation](#). While a handful of false claims can be debunked by fact-checkers and news organisations on a daily basis, the challenge becomes much harder when the number of deep-faked clips on platforms reaches thousands per day.

Our Approach

What is a 'Red Team' operation?

A 'Red Team' operation refers to a security technique where a group pretends to be an enemy, in order to identify potential weaknesses and vulnerabilities. The term originates from the military, in which a 'Red Team' would play the role of an adversary in training exercises.

Red Team operations are useful in the context of misinformation and AI, as the techniques and technologies involved are rapidly evolving. [The White House backed an AI Red-Teaming exercise last year](#).

By attempting to create compelling misinformation ourselves, we can gain a better understanding of the challenges involved, and the potential impact of these techniques.

A 'Red Team' operation is **not** an endorsement of the techniques described, nor is it an instruction manual for creating harmful information. We have deliberately obscured the platforms and exact techniques used to create information. This report should be seen as a starting point for further discussion and action, rather than a definitive statement on the issue.

With that said, our goals for this report are to:

1. Show what scale of voter manipulation is currently possible with widely available and cheap tools.
2. Understand what tactics may be used to manipulate voter opinion in the UK.
3. Spread awareness of synthetic media in order to boost the public's resilience to manipulation.

PART 1: TACTICS

Misinformation campaigns, particularly those aimed at influencing political outcomes, often employ a range of tactics designed to manipulate public opinion and undermine trust in targeted parties or candidates.

FIG 1: COMMON MISINFORMATION TACTICS

<i>Emotional Manipulation</i>	<i>Using content that evokes strong emotions such as fear, anger, or excitement to bias the audience against a party. Playing on the human tendency to react more strongly to emotional content, making it easier to spread misinformation.</i>
<i>Fabricated Content</i>	<i>Creating entirely false news stories, images, or videos about a party or its members. This could include deepfakes, doctored videos, or fake news articles designed to discredit the targeted party.</i>
<i>Echo Chambers</i>	<i>Amplifying misinformation within closed or like-minded groups on social media, creating an environment where false information is repeatedly shared and reinforced, making it appear more credible.</i>
<i>Astroturfing</i>	<i>Creating the illusion of grassroots support against a party by using bots or paid individuals to spread misinformation or negative opinions, making the campaign appear as coming from the public rather than its actual orchestrated origins.</i>
<i>Misleading Headlines and Clickbait</i>	<i>Using sensational or misleading headlines to entice people to engage with content that is biased, partially true, or completely false. The actual content often does not support the headline, but the headline itself influences perceptions</i>

Commonly, the goal of these campaigns is to damage the credibility and reputation of targeted party leaders. This can be through spreading fabricated stories about their personal lives, questioning their competence or integrity, or attributing false statements to them.

Another common tactic is to suppress voter turnout by creating confusion or apathy around the election process. This can involve spreading false information about the date, time, or location of the election, or suggesting that the outcome has already been predetermined.

For the purposes of this report, we will be using these two tactics.

PART 2: Generating Text

In 2016, Facebook was flooded with fake news regarding Hilary Clinton, Bill Clinton and Donald Trump, [penned by Macedonian teenagers](#) in order to make advertising revenue. Writing 5 to 10 articles daily, these sites could make over \$1,000 per day.

Since 2020, and the advent of GPT-2, large language models have been able to output coherent, but false, news articles. This means that humans are no longer the limiting factor for the word count of fake news - a language model can easily output hundreds of articles in a matter of minutes.

Most widely available tools will initially refuse to generate scripts for manipulative news, due to fine-tuning and guidance from the model creators themselves.

Using a technique known as 'prompt-hacking' we can easily trick the AI into producing coherent, convincing fake news.



With roughly 200 characters of 'prompt-hacking', we were able to make a popular free AI chat tool output defamatory scripts about Keir Starmer and Rishi Sunak, in addition to speech from each of them encouraging voters to stay at home.

Running total: 5 minutes, £0 spent

Part 3: Generating Supportive Imagery

Over the past year, the quality and consistency of AI-generated imagery has improved substantially. Image generation models can now produce realistic photos with coherent objects, textures and lighting.

However, synthetic images often still contain subtle flaws that include misshapen anatomy, inconsistent details, or an artificially smooth and polished appearance.

Due to this, it is difficult to generate the type of images which will support a misinformation narrative. Outlandish creations are easily debunked.



Due to this, we expect that **repurposing** or **editing** authentic imagery to be a bigger threat for people intending to manipulate the election.

In our example, we have used videos of party leaders from other news reports to lend credibility to our 'scandal'.



Tactics such as deceptively editing, mislabeling or mischaracterizing authentic imagery are a significant threat. This is referred to as 'shallow-faking' and may include:

- Claiming it shows them at a controversial location or with a controversial figure.
- Selectively cropping an image to remove context.
- Editing signs or symbols in a photo
- Presenting an old image as a recent event to make false claims.

Running total: 10 minutes, £0 spent

Part 3: Generating Audio

To generate convincing audio with widely-available tools, only 30 seconds of source audio is needed. An 'instant voice clone' can be created immediately after uploading an audio or video clip.



Platforms which offer this service openly allowed us to create clones of newsreaders and celebrities. These platforms let you create up to 10 voice clones and 30 minutes of audio for \$1.

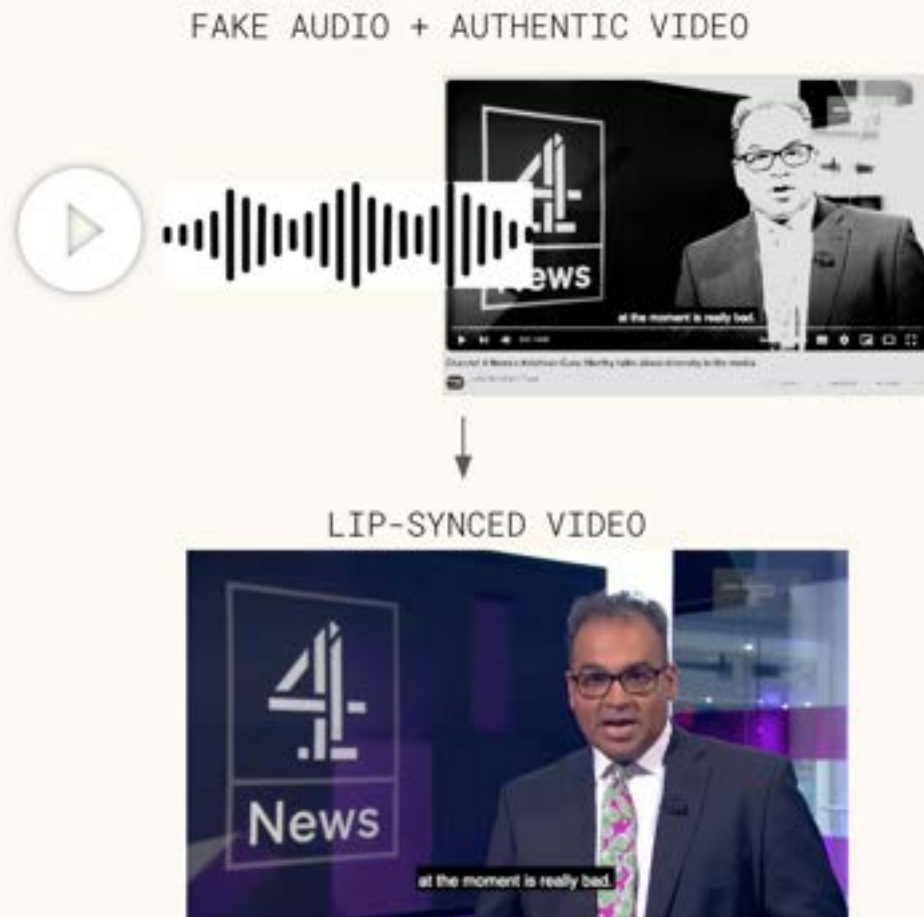
Some of the largest platforms did not allow us to create voice-clones of Prime Minister Rishi Sunak and opposition leader Keir Starmer. However, smaller start-up platforms using the same technology permitted these voice clones on their free trial.

Running total: 20 minutes, £0.80 spent

Part 4: Generating Video 'Deep Fakes'

Once a voice-clone is generated, an existing video of a newsreader or celebrity can be edited with AI to lip-sync the audio with a video.

30 seconds of audio and standard-definition video cost 16p to produce and was generated in 1 minute and 20 seconds.



Currently, the AI model used to generate these lip-syncs still contains artefacts such as: blurry mouth, occasional de-syncs and cuts mid-sentence. However, these are easily covered up with editing and by adding more misleading imagery.

Running total: 21 minutes, £0.96 spent

OPERATION OUTPUTS

Key examples below show what was possible with widely available tools. Each of these clips were created for less than £1

Click on each image to view the video ([Google Drive](#)).

Party Leaders



Prime Minister
Rishi Sunak
'Stay At Home' Suppression



Opposition Leader
Keir Starmer
'Stay At Home' Suppression

TV News



BBC News
Rishi Sunak 'Tampering' Scandal



ITV News
Rishi Sunak 'Tampering' Scandal



Sky News
Keir Starmer 'Tampering' Scandal



Channel 4 News
Keir Starmer 'Tampering' Scandal

Radio / Podcast News



The News Agents
Keir Starmer 'Tampering' Scandal



Times Radio
Keir Starmer 'Tampering' Scandal



LBC
Rishi Sunak 'Tampering' Scandal

HOW TO SPOT A DEEP-FAKE

While the technology is advancing rapidly, there are several signs that a video has been generated with 'deep-fake' and voice cloning AI technology.

Blurry Mouth: Deepfakes often struggle to replicate the exact movements of the lips and tongue, leading to blurriness or unnatural movements. Look for inconsistencies in the mouth area. Pay attention to the sync between lip movements and spoken words.

Strange Pronunciation of Words: AI-generated voices often don't quite match the speaker's natural speech patterns. Listen for odd phrasing, mispronunciations, or unusual cadence. Compare the voice and speaking style to other known recordings of the individual. Significant deviations in tone, pace or pronunciation can be red flags.

Repetitive Movement: AI algorithms might repeat certain movements unnaturally. This can manifest as repetitive blinking, head movements, or facial expressions that don't vary as much as they would in a real person. Observe the overall fluidity of movement. If you notice a mechanical or looping quality to gestures and facial expressions, this could indicate a deep-fake.

CONCLUSION

Our 'Red Team' operation has shown that creating convincing, harmful false content is now trivial.

Cost and technical expertise are no longer barriers to generating fake news at scale. Even those with limited skills can create deep-fakes of authoritative figures saying anything they choose.

We created many deep-fakes manually within 30 minutes for less than 1 each. This process can be mechanised easily, ultimately allowing fake news to be created on 'autopilot'.

While video deep-fakes are still detectable by the human eye, misinformation actors will likely take a 'quantity over quality' approach. Flooding platforms with huge volumes of synthetic media may overwhelm our ability to discern truth from fiction.

Repurposing authentic video out of context will be a key strategy for those seeking to deceive. Familiarity with these 'shallow-fake' techniques is our best defence.

Ultimately, the outcome of our elections should be decided by an informed public, not by those spreading disinformation.

The AI-powered misinformation threat is serious, but it is not insurmountable. Protecting our democracy requires a coordinated effort from government, tech platforms, media, and citizens alike.

There is no silver bullet for countering misinformation. A variety of counter-measures, focusing on both the long and short term will yield the best results.

Investing and promoting media literacy initiatives will help people to spot manipulation techniques. Platforms improving their moderation of election-related content will have an impact, but the exact implementation of labelling, fact-checking and reach-limited has consequences on its effectiveness.

It is important to remember that AI is a tool which can be used for many purposes. It can be used to inform, entertain and enlighten, or to deceive, manipulate and divide. However, AI-powered fact-checking and analysis should be utilised during the election campaign to fight against misinformation.

Part one of our 'Red Team' operation has focused on the generation of misinformation. Part two will focus on how **AI may be used to distribute and amplify harmful misinformation** - and how platforms such as TikTok, Meta and X may already do this.

Recommendations and Next Steps

1. The UK public must learn to expect that **a large amount of AI-generated misinformation may be generated and shared during the election campaign.**
2. Familiarity with deep-fakes and other manipulation techniques is the best defence against misinformation. See: [Thinks Disinformation Intervention](#). However, there is **no silver bullet for countering misinformation.**
3. The UK should take a 'portfolio approach' with interventions via **public information, action by government institutions** and action by **technology platforms themselves.**
4. Steps to limit deep-fakes of authoritative and political figures by commercial platforms are welcomed. However, this may have limited effectiveness, given the open-source nature of most of this technology.
5. More research must be done to understand how **AI may be used to distribute and amplify harmful misinformation.**

ABOUT FENIMORE HARPER

Fenimore Harper is a media insight and research firm. It was founded in 2021 by Marcus Beard, after working as a communications adviser at HM Treasury, Cabinet Office and 10 Downing Street.

During the COVID-19 pandemic, Marcus led the UK Government's digital counter-misinformation strategy. Working internationally, he advised partners such as NATO, the G7 and the governments of Ukraine, Bulgaria and North Macedonia.

Marcus' writing and research has appeared in The Times, The Telegraph, The Independent, The Guardian, and Bloomberg. He has also provided political commentary for BBC News and Sky News.

CONTACT: marcus@fenimoreharper.com - 07809323683

ADDITIONAL EDITING:

VANESSA TRACEY

FREYA JOSEPH
