

June 29, 2023

Mr. Mark Zuckerberg  
Chairman and Chief Executive Officer  
Meta  
1 Hacker Way  
Menlo Park, CA 94025

Mr. Sam Altman  
Chief Executive Officer  
OpenAI  
3180 18th Street, Suite 100  
San Francisco, CA 94110

Mr. Sundar Pichai  
Chief Executive Officer  
Alphabet Inc. and Google LLC  
1600 Amphitheatre Parkway  
Mountain View, CA 94043

Mr. Emad Mostaque  
Chief Executive Officer  
Stability AI  
88 Notting Hill Gate  
London, U.K. W11 3HP

Dr. Dario Amodei  
Chief Executive Officer  
Anthropic  
548 Market Street  
San Francisco, CA 94104

Mr. Elon Musk  
Owner  
Twitter  
1355 Market Street, Suite 900  
San Francisco, CA 94103

Mr. Shou Zi Chew  
Chief Executive Office  
TikTok  
5800 Bristol Parkway, Suite 100  
Culver City, CA 90230

Mr. Satya Nadella  
Chief Executive Officer  
Microsoft  
One Microsoft Way  
Redmond, WA 98052

Mr. David Holz  
Chief Executive Officer  
Midjourney  
611 Gateway Boulevard, Suite 120  
San Francisco, CA 94080

Dear Mr. Zuckerberg, Mr. Musk, Mr. Altman, Mr. Chew, Mr. Pichai, Mr. Nadella, Mr. Mostaque, Mr. Holz, and Dr. Amodei:

I write with concerns about your current identification and disclosure policies for content generated by artificial intelligence (AI). Americans should know when images or videos are the

product of generative AI models, and platforms and developers have a responsibility to label such content properly. This is especially true for political communication. Fabricated images can derail stock markets, suppress voter turnout, and shake Americans' confidence in the authenticity of campaign material. Continuing to produce and disseminate AI-generated content without clear, easily comprehensible identifiers poses an unacceptable risk to public discourse and electoral integrity.

Online misinformation and disinformation are not new. But the sophistication and scale of these tools has rapidly evolved and outpaced our existing safeguards. In the past, creating plausible deepfakes required significant technical skill; today, generative AI systems have democratized the ability—opening the floodgates to anyone who wants to use or abuse the technology.

We have already seen evidence of generative AI being used to create and share false images.<sup>1</sup> In some instances, these have been relatively benign—such as Pope Francis depicted wearing a large white down jacket.<sup>2</sup> Others are more disturbing. In May, an AI-generated image of a purported explosion at the Pentagon went viral, causing a dip in major stock indices.<sup>3</sup> Fake news accounts recirculated these images alongside real outlets, including *RT*, a Russian state-backed media organization.

The proliferation of AI-generated content poses a particular problem for political communication. In his recent testimony before the Senate Judiciary Committee, OpenAI CEO Sam Altman identified the ability of AI models to provide “one-on-one interactive disinformation” as an area of greatest concern.<sup>4</sup> We have entered the beginning of this era.

In June, the official rapid response Twitter account of Florida Governor Ron DeSantis, a candidate for the 2024 Republican nomination for president, shared images that experts say appear to be AI-generated.<sup>5</sup> Both official and unaffiliated accounts supporting former President Trump have posted AI-generated content targeting his political rivals.<sup>6</sup>

---

<sup>1</sup> Tiffany Hsu and Steven Lee Myers, “[Can We No Longer Believe Anything We See?](#)” *New York Times*, April 8, 2023; Clare Duffy, “[Puffer coat Pope. Musk on a date with GM CEO. Fake AI ‘news’ images are fooling social media users.](#)” CNN, April 2, 2023.

<sup>2</sup> Daysia Tolentino, “[AI-generated images of Pope Francis in puffer jacket fool the internet.](#)” NBC News, March 27, 2023.

<sup>3</sup> Shannon Bond, “[Fake viral images of an explosion at the Pentagon were probably created by AI.](#)” NPR, May 22, 2023.

<sup>4</sup> Cat Zakrzewski, Cristiano Lima and Will Oremus, “[CEO behind ChatGPT warns Congress AI could cause ‘harm to the world.’](#)” *Washington Post*, May 16, 2023.

<sup>5</sup> Benjy Sarlin and Shelby Talcott, “[DeSantis campaign shares fake Trump/Fauci images, prompting new AI fears.](#)” *Semafor*, June 8, 2023; Bill McCarthy, “[Ron DeSantis ad uses AI-generated photos of Trump, Fauci.](#)” *Agence France-Presse*, June 7, 2023.

<sup>6</sup> Aditi Bharade, “[Someone made a hyper-realistic deepfake of Ron DeSantis as Michael Scott from 'The Office' wearing women's clothes. It's the latest instance of AI being weaponized to take DeSantis down.](#)” *Business Insider*, May 29, 2023; Matthew Loh, “[The Trump-DeSantis showdown is now official, and artificial intelligence is right in the middle of it.](#)” *Business Insider*, May 25, 2023.

In May, I joined colleagues to introduce the REAL Political Ads Act, which would require a disclaimer on political ads for federal campaigns that use content generated by AI.<sup>7</sup> However, as political media increasingly shifts from regulated television, print, and radio advertising to the free-for-all of social media, broader disclosure requirements must follow.

AI system developers and platforms will have to collaborate to combat the spread of unlabeled AI content. Developers should work to watermark video and images at the time of creation, and platforms should commit to attaching labels and disclosures at the time of distribution. A combined approach is required to deal with this singular threat.

Companies have started taking steps to better identify AI-generated content for users. For example, non-profit organizations, like the Partnership on AI, have released suggested guidelines.<sup>8</sup> Microsoft has committed to watermark AI-generated content, and Google will begin attaching a written disclosure on Google Images.<sup>9</sup> OpenAI's DALL-E 2 adds a watermark to images it generates, and Stable Diffusion embeds watermarks into its content by default.<sup>10</sup> Midjourney, Shutterstock, and Google have committed to embedding metadata indicators in AI-generated content.<sup>11</sup>

However, these policies remain easily bypassed or alarmingly reliant on voluntary compliance. Google's process for labeling AI-generated images from third-party systems depends on self-disclosure.<sup>12</sup> Stable Diffusion's open source structure allows users to circumvent the watermarking code.<sup>13</sup> DALL-E 2's watermarks are inconspicuous and easily removed.<sup>14</sup> And while some platforms –including Meta,<sup>15</sup> Twitter,<sup>16</sup> and TikTok<sup>17</sup>– have existing policies for AI-generated images and video, such content continues to appear on users' feeds.

---

<sup>7</sup> Office of Senator Michael Bennet, "[Bennet, Klobuchar, Booker Push to Regulate AI-Generated Content in Political Ads](#)," May 15, 2023.

<sup>8</sup> Partnership on AI, *[PAI's Responsible Practices for Synthetic Media: A Framework for Collective Action](#)*, February 27, 2023.

<sup>9</sup> Kyle Wiggers, "[Microsoft pledges to watermark AI-generated images and videos](#)," *TechCrunch*, May 23, 2023; Cory Dunton, "[Get helpful context with About this image](#)," Google, May 10, 2023.

<sup>10</sup> Benj Edwards, "[AI image generation tech can now create life-wrecking deepfakes with ease](#)," *Ars Technica*, December 9, 2022.

<sup>11</sup> IPTC, "[Midjourney and Shutterstock AI sign up to use of IPTC Digital Source Type to signal generated AI content](#)," May 11, 2023.

<sup>12</sup> Dunton, "[Get helpful context with About this image](#)."

<sup>13</sup> Hany Farid, "[ChatGPT and Dall-E Should Watermark Their Results](#)," *Gizmodo*, April 2, 2023.

<sup>14</sup> OpenAI, "[How should I credit DALL·E in my work?](#)" accessed June 28, 2023.

<sup>15</sup> Meta, "[Manipulated Media](#)," accessed June 28, 2023.

<sup>16</sup> Twitter, "[Synthetic and manipulated media policy](#)," April 2023.

<sup>17</sup> TikTok, "[Integrity and Authenticity](#)," last updated March 2023.

Platforms must update their policies for a world where everyone has access to generative AI tools. They should require clear, conspicuous labels for AI-generated video and images, and where users fail to comply, should label AI-generated content themselves. Platforms should consider particular rules for official political accounts, and should release regular reports detailing their efforts to identify, label, or remove AI-generated content.

Similarly, generative AI system developers must scrutinize whether their models can be used to manipulate and misinform, and should conduct public risk assessments and create action plans to identify and mitigate these vulnerabilities. We cannot expect users to dive into the metadata of every image in their feeds, nor should platforms force them to guess the authenticity of content shared by political candidates, parties, and their supporters.

Continued inaction endangers our democracy. Generative AI can support new creative endeavors and produce astonishing content, but these benefits cannot come at the cost of corrupting our shared reality.

To that end, I request answers to the following questions by July 31, 2023:

For generative AI developers:

- What technical standards, features, or requirements do you currently employ to watermark or otherwise identify content created using your systems?
  - When were these standards, features, or requirements developed?
  - When were these standards, features, or requirements last updated?
  - What auditing processes, if any, does your organization have in place to evaluate the effectiveness of these standards, features, or requirements?
- What policies do you currently have in place for users that repeatedly violate a watermarking or identifying requirement, either by removing the identifier or avoiding it in some other way?
  - How many accounts, if any, have you suspended or removed for violating a watermarking or identifying requirement?
- What tracking system do you currently have in place, if any, to monitor the distribution of content created using your systems?
- Before deploying a model, what tests or evaluations do you use to estimate potential capabilities relating to misinformation, disinformation, persuasion, and manipulation?

- What processes do you use to estimate risks associated with misinformation, disinformation, persuasion, and manipulation? Under what circumstances would you delay or restrict access to a generative AI system due to concerns about these risks?
- What interoperable standards currently offer the highest degree of provenance assurance?

For social media platforms and search engines:

- Will you commit to removing AI-generated content designed to mislead users?
- What technical processes are currently in place to identify AI-generated content?
- How many pieces of AI-generated content did you identify in 2022 and the first quarter of 2023?
  - Of those identified, how many were removed for violating a policy?
  - If removed, what policy did they violate?
  - If not removed, was a label or other clear identifier affixed?
  - If not labeled, please provide a rationale for declining to do so.
- Do you have specific policies in place for AI-generated content posted by an official political campaign account?
  - If so, what are they?
  - If not, describe why not.
- Do you have specific policies in place for AI-generated content related to campaigns and elections?
  - If so, what are they?
  - If not, describe why not.

I appreciate your attention to this important matter and look forward to your response.

Sincerely,

A handwritten signature in blue ink that reads "Michael F. Bennet". The signature is fluid and cursive, with the first name "Michael" and last name "Bennet" clearly legible.

U.S. Senator Michael F. Bennet