# The Illustrated Man: A brief analysis of the regulatory challenges of deepfake technologies (2,986 words)

*And yet the past, though of its nature alterable, never had been*
*altered. Whatever was true now was true from everlasting to*
*everlasting. It was quite simple. All that was needed was*
*an unending series of victories over your own memory.*

—GEORGE ORWELL, "NINETEEN EIGHTY-FOUR"

In the early 1990s, a team of visual effects artists spent six months painstakingly working on one sequence for *Forrest Gump*. This sequence was mainly comprised of archival footage of John F Kennedy—dozens of true historical shots manipulated and spliced together, so as to make it appear that the President shook Forrest Gump's (played by Tom Hanks, who was seven years old at the time of Kennedy's assassination) hand, turned to the camera, and laughed to the audience, "I believe he said he had to go pee."[1]

Fast forward thirty years and similar synthetic alterations to existing video and audio have become known as "deepfakes". Code capable of achieving the same goals as *Forrest Gump*'s visual effects team is open source and executable on hardware found in most relatively modern personal computers;[2] Tom Scott jokingly lists the system requirements for running deepfake programs as being, "A Windows 10 PC, a powerful Nvidia graphics card, moderate tech skills, and poor or no empathy".[3] In fact, face swapping technology is even freely available in apps as ubiquitous as Instagram and Snapchat, run on modern smartphones.[4] The creation of deepfakes is no longer achievable only by George Lucas' visual effects team over the course of weeks or months, but by individuals in a matter of days, if not hours.

Progressions in deepfake technology will unlock myriad wonderous possibilities—actors posthumously or remotely reprising or completing roles;[5] medical advancements, such as the opportunity for one who has lost speech abilities to reproduce their own voice through transcription; educational opportunities that may be born from the ability to truthfully recreate historical figures and events on screen. Simmering below these possibilities, however, looms the threat of deepfake technology being the cornerstone of an erosion of public discourse as we know it. A falsified video of a politician admitting to a crime which they did not commit,

---

[1] Peyton Reed "Through The Eyes of Forrest Gump: The Making of an Extraordinary Film" (2001) Paramount Pictures; Robert Zemeckis "Forrest Gump" (1994) Paramount Pictures

[2] Ivan Perov and others "DeepFaceLab: A simple, flexible and extensible face swapping framework" (12 May 2020) arXiv:2005.05535 [cs.CV]

[3] Tom Scott "Faceswapping, Unethical Videos, and Future Shock" *YouTube* (6 February 2018)

[4] Kelsey Warner "Instagram rolls out augmented reality filters after Snapchat's face swap success" *The National (*International Edition Online*,* 14 October 2019)

[5] See: J. Naruniec, L Helminger, C Schroers and R.M. Weber "High-Resolution Neural Face Swapping for Visual Effects" (29 June 2020) 39:4 Eurographics Symposium on Rendering 2020; DisneyResearchHub "High Resolution Neural Face Swapping for Visual Effects" (29 June 2020) YouTube; Alissa Wilkinson "Hollywood is replacing artists with AI. Its future is bleak." *Vox* (Online, 29 January 2020)

circulated days before a general election;[6] the synthetic production of pornography of a public figure and its exploitation for blackmail;[7] the alteration of reality on a scale previously only imaginable through the works of Philip K Dick. Regulators face the challenge of implementing measures which retain the beneficial use of this technology while preventing its malicious exploitation.

## I.     IDENTIFYING DEEPFAKES

The first hurdle faced by regulators of deepfake technology is identifying when it has been used in the first place. While it is comforting to imagine a future in which the exploitation of deepfake technology is flaunted by similarly powerful 'good' artificial intelligence, Dixon Jr foresees a future in which "continuing advances in deepfake technology may render it impossible for one AI-driven system to detect a video created or modified by another AI-driven system."[8] If regulation of deepfakes is to have effect, regulators must be able to determine when a piece of content in question is legitimate or not.

In the United States, a legislative solution to this problem was put forward in the DEEP FAKES Accountability Act.[9] Although it failed to pass beyond the House of Representatives, its first section dealt with identifying deepfakes: it proposed a legislative requirement that deepfake content be embedded with a digital watermark and that it be accompanied by aural and written statements explaining the synthetically altered nature of that material.[10] The approach of requiring creators to warn consumers through disclaimers is in some areas used by certain platforms—Twitter's Parody, newsfeed, commentary, and fan account policy, for example, requires that non-affiliated parody accounts should "clearly indicate" that they are not affiliated with the 'legitimate' account in both the user's bio and their username.[11] This may go some way in ameliorating the problem of identification, but it does not completely discharge it; if disputes arise as to the authenticity of content, regulators will need some way to solve it that does not rely on the honesty of creators.

David Doermann, Former Project Manager for the United States Defense Advanced Research Project Agency, implies that there are already processes in place to be used if artificial intelligence is not a viable solution:[12]

> [J]ust like a court of law you're going to have one side saying one thing and another side saying the other thing and there's going to be cases where there's nothing definitive.

---

[6] Christiano Lima "'Nightmarish': Lawmakers brace for swarm of 2020 deepfakes" *Politico* (Online, 13 June 2019)

[7] Drew Harwell "Fake-porn videos are being weaponized to harass and humiliate women: 'Everybody is a potential target'" *The Washington Post* (Online, 30 December 2018)

[8] Herbert B Dixon Jr "Deepfakes: More Frightening Than Photoshop on Steroids" (Summer 2019) The Judges' Journal; Chicago Vol. 58, Iss. 3, 35-37 at 2.

[9] Defending Each and Every Person from False Appearances by Keeping Exploitation Subject to Accountability Act of 2019 H.R.3230

[10] Defending Each and Every Person from False Appearances by Keeping Exploitation Subject to Accountability Act of 2019 above n 9 at §1041(a)-(c)

[11] Twitter "Parody, newsfeed, commentary, and fan account policy"

[12] "House Intelligence Committee Hearing on 'Deepfake' Videos" *C-Span* (13 June 2019).

Perhaps, then, the best available solution exists in the form of the current court system. It does have its flaws—which I later discuss—but until a perfect artificial intelligence becomes available, the court may be the best remaining option. The New Zealand Court has so far accepted the adoption of such an approach; although it has not yet been fully tested. In *R v Iyer*[13] it was stated that, "If the defendant alleges that the material has been fraudulently created, or misrepresents the nature of the communication, then that would be a matter to be tested by evidence as part of the defence case."[14]

II.      NEW ZEALAND'S EXISTING INTERNET LEGISLATION

Once deepfake content has been recognised, regulation will need to maintain a careful balance which restricts its malicious use while allowing its use for beneficial industries and purposes. There are two main pieces of legislation which New Zealand currently has in place that could be applied to malicious uses of deepfake technology: the Film, Videos, and Publications Classifications Act 1993 (FVPCA) and the Harmful Digital Communications Act 2015 (HDCA).

The FVPCA is the legislation that was used by the New Zealand Classification Office to prohibit the distribution of both the livestreamed video and the 'manifesto' of the Christchurch Mosque Attacker in 2019, where they were deemed "objectionable" under s 23.[15] The FVPCA therefore has a use in prohibiting the widespread online sharing of harmful material; however, its uses are likely limited to a few narrow scenarios. In *Living Word Distributors v Human Rights Action Group (Wellington)*,[16] the Court found that the FVPCA was subject to a "subject matter gateway"[17] of sex, horror, crime, cruelty, or violence "rather than to the expression of opinion or attitude".[18] This means that while the FVPCA could potentially be used to prevent the sharing of synthetic pornography created via deepfake technology, it would likely not apply to such technology's use for nefarious political purposes.

The HDCA's purpose is to "deter, prevent, and mitigate harm caused to individuals by digital communications".[19] Under s 22, a person commits an offence if:

(a)    the person posts a digital communication with the intention that it cause harm to a victim; and

(b)    posting the communication would cause harm to an ordinary reasonable person in the position of the victim; and

(c)    posting the communication causes harm to the victim.

Like the FVPCA, this provision appears to have some use in combatting instances of synthetic pornography, although it too has its limits; particularly in its requirements for intention to cause harm and the

---

[13] *R v Iyer* [2017] DCR 82

[14] *R v Iyer,* above n 13 at [42]

[15] Classification Office to Chief Censor "Notice of decision under section 38(1)" (18 March 2019) OFLC Ref 1900148.000

[16] *Living Word Distributors v Human Rights Action Group (Wellington)* [2000] 3 NZLR 570

[17] *Living Word Distributors,* above n 16 at [29]

[18] *Living Word Distributors,* above n 16 at [28]

[19] Harmful Digital Communications Act 2015 at s 3(a)

necessity to actually cause harm to the victim. These requirements are discussed in detail in *R v Iyer*,[20] where the defendant posted semi-nude photos of his ex-wife to Facebook. They, too, can be frustratingly high for victims:[21]

> It was clear from the inclusion of the word "serious" [in the s 4 definition of "harm" being "serious emotional distress"] that the intended harm must be more than trivial. Being merely upset or annoyed as a consequence of a digital communication would not be sufficient to invoke the sanction of criminal law.

Such a finding presents major issues for the HDCA's useful applicability to many possible instances of malicious deepfake technology use; as the HDCA's purpose is targeted at individuals, the HDCA is useful where there is a clear individual victim—as would likely be the case in the production of synthetic pornography—but it fails to provide useful recourse when no such individual victim exists. Robert Chensey and Danielle Citron hypothesise "a fake video of a white police officer shouting racial slurs or a Black Lives Matter activist calling for violence".[22] In these types of scenarios, the harm to individuals pales in comparison to the greater harms suffered to groups, such as ethnic groups, or even an entire population who benefits from being able to vote based on reliable information. While there are devastating effects, those effects are unlikely to be felt by any single person to the extent of meeting the "serious harm" element. Chensey and Citron further note the possible issue of a "liar's dividend"[23] affecting political discourse; if regulators fail to adequately address the malicious use of deepfakes, it becomes possible for public figures to negate their actual wrongdoings by rejecting legitimate evidence which sits contrary to their own interests—as Donald Trump briefly attempted to claim regarding his infamous Access Hollywood tape[24]—affecting the function of a state's functioning democratic system. These wider harms are not considered by New Zealand's existing legislation. In fact, the Law Commission's report notes that:[25]

> Criticisms of some new media publishers focused on the lack of adherence to any ethical code, the publication of unsubstantiated information, including damaging allegations, the publication of information suppressed by the courts, and a failure to adequately differentiate between opinion and fact.

In the pursuit of ensuring New Zealand's "strong public interest in ensuring there are effective mechanisms for holding the media to account for the exercise of their power and for remedying harms arising from any breaches of ethical and professional standards,"[26] an overhaul of existing legislation is needed, if not new legislation with an entirely new approach to these wider issues.

The ways in which the FVPCA and the HDCA operate can be contrasted against Twitter's "synthetic and manipulated media policy"

---

[20] *R v Iyer*, above n 14 at [43]—[63]

[21] *R v Iyer*, above n 14 at p 82 and at [57]

[22] Robert Chensey and Danielle Citron "Deepfakes and the New Disinformation War: Deepfakes and the New Disinformation War: The Coming Age of Post-Truth Geopolitics" 98 Foreign Aff. 147 (2019) at 151

[23] Robert Chensey and Danielle Citron above n 22 at 151—152

[24] Billy Bush "(Opinion): Yes, Donald Trump, you said that" *The New York Times* (Online, 3 December 2017)

[25] Law Commission *The news media meets 'new media': rights, responsibilities and regulations in the digital age* (NZLC R128, March 2013) at 93

[26] Law Commission above n 25 at 95

which takes into account three main factors: (1) whether the content is synthetic or manipulated; (2) whether the content is shared in a deceptive manner; and (3) whether the content is likely to impact public safety or cause serious harm.[27] Twitter's definition of "serious harm" varies considerably from that of the HDCA. It takes into account the safety of both individuals and groups, the risk of possible stalking or intimidation of both individuals and groups; and explicitly includes synthetic pornography (although synthetic pornography is also covered more comprehensively by the site's "non-consensual nudity policy").[28] While there are policy and philosophical discussions to be had about the extent to which such an approach might affect users' freedom of expression,[29] there is little doubt that Twitter's considerations and definitions cast a much wider net than those of the FVPCA and HDCA in the prevention of malicious deepfake technology use.

III.     WOULD ADOPTING TWITTER'S DEFINITIONS SOLVE THE ISSUE?

While the FVPCA and the HDCA are both flawed, the answers to their problems do not necessarily lie simply in amending them to plug their leaks; further issues inherently exist in the use of courts to solve deepfake issues, even if they apply ideal legislation. This is due to the fact that the use of the courts at all introduces a latency in the response to an issue of malicious deepfake technology use. In the case of the Christchurch Mosque Attack, it took three days for the Classification Office's s 23 FVPCA decision regarding the livestream to be released.[30] Even if the Classification Office were to avoid increasing this turnover rate after becoming the first port of call for malicious deepfake use, three days is an alarmingly long length of time—particularly as it has been accepted both in global discourse and in the New Zealand courts that the issues associated with deepfake technology do not rest solely in the creation of synthetic content, but in how quickly and widely that content might be shared.[31] Further, as the Law Commission notes in its report on the regulation of damaging content published by mass news media, redress through the courts remains inaccessible on practical and financial levels for a large number of potential victims:[32]

> [W]hile it is true that citizens have the right to seek redress through the courts when the published content breaches the law, the reality is that the expense of pursuing a civil action for defamation or breach of privacy means this is simply not a meaningful remedy for most private citizens.

Courts are likely to be too difficult and costly to access and too slow to be themselves an adequate response; they are the legislative equivalent of the ambulance at the bottom of the cliff when dealing with issues as fast-moving and permanently harmful as those deepfake technology can create. If the harms of malicious deepfakes are to be minimised, faster and more efficient solutions are likely to be required.

---

[27] Twitter "Synthetic and manipulated media policy"

[28] Twitter "Non-consensual nudity policy) (November 2019)

[29] New Zealand Bill of Rights Act 1990 at s 14

[30] Classification Office to Chief Censor "Notice of decision under section 38(1)" (18 March 2019) OFLC Ref 1900148.000

[31] Law Commission above n 25 at 96; R v Iyer, above n 13 at [69].

[32] Law Commission above n 25 at 96

## IV.    REGULATING CONTENT MODERATION

In a United States House Intelligence Committee Hearing, Danielle Citron put forward the suggestion that platforms be better incentivised to moderate content that is posted by users.[33] Support for such an approach appears to exist in New Zealand, too; the Law Commission states that individuals harmed by false reporting "are often reliant not just on the original publisher but also on remote parties – such as Google – to remove the damaging content."[34] Some platforms have already adopted this approach of moderating users' content—earlier this year, Facebook attached a "partially false" disclaimer to a video of Nancy Pelosi appearing to slur her speech as if she were drunk,[35] and Twitter has recently been applying similar disclaimers to misleading tweets made by President Trump.[36] While deepfake technology and social media platforms' rules on transparency have not existed for long enough to have been tested against each other such that it is possible to say whether or not this approach might be successful, there is a certain logic in having social media platforms more actively moderate their users' content. Platforms have access to huge amounts of data and are well-versed in the way that users engage with the platform. They should, in theory, be perfectly placed to implement internal systems that allow for fast detection, analysis, and removal of malicious deepfakes.

Regulators must be careful to not allow social media platforms to retain too much power in the moderation of deepfake content, however; users may have legitimate fears similar to those that the Law Commission discussed in regards to news media where there is the "potential for a small number of publishers (mainstream and new media) to exert undue influence on the news agenda and public opinion".[37] If regulation of deepfake content is left to the moral whims of social media moguls, there is still a great lack of transparency as to how the issue of malicious deepfake technology is actually targeted and dealt with, and no single standard by which any platform is required to abide.

Citron's suggestion would largely be implemented in the United States via an amendment of § 230(c)(1) Communications Decency Act,[38] to state that "No provider or user that engages in reasonable moderation practices of an interactive computer service shall be treated as the publisher or speaker of any information provided by another information content provider".[39] While such a change could be difficult to enact in the United States due to clashes with the First Amendment,[40] implementing equivalent legislation in New Zealand would face no such legal roadblocks (although it is imaginable that there would be intense policy discussion on the matter). New Zealand has no similar legislation, but there is precedent

---

[33] "House Intelligence Committee Hearing on 'Deepfake' Videos" above n 12

[34] Law Commission above n 25 at 96

[35] Hannah Denham "Another fake video of Pelosi goes viral on Facebook" *The Washington Post* (Online, 4 August 2020)

[36] Gilad Edelman "Twitter Finally Fact-Checked Trump. It's a Bit of a Mess" *Wired* (online, 27 May 2020)

[37] Law Commission above n 25 at 95

[38] Communications Decency Act (1996) 47 U.S.C. § 230

[39] "House Intelligence Committee Hearing on 'Deepfake' Videos" above n 12

[40] United States Constitution, amend I

for a similar approach to content moderation: in *Jensen v Clark*[41] where a printing company allowed the printing of defamatory content and where that printing company knew that the material was an "attack by students of the University on one of their Professors"[42] it was held that "To print this material without any more investigation than was possible in a telephone discussion between the secretary and the manager of the company was, in my view, irresponsible".[43] While there are clear differences between a newspaper printing company monitoring the content of what comes out of its physical presses and a social media company monitoring its potentially hundreds of millions of users' content, Citron's approach is in essence a variation of that applied in *Jensen v Clark*.

Citron's approach enjoys the advantages of incentivising platforms to be proactive in preventing the creation and sharing of malicious deepfake content, and these moderation practices may be judicially reviewed for greater transparency and regulatory oversight. While many platforms already release annual transparency reports without regulatory incentives—such as Twitter through its Transparency Center which was started in late August 2020,[44] and Reddit's annual transparency report which has since 2019 included a greater look into its moderation practices including its removal of content breaching its "Content Manipulation" rules,[45] these initiatives are voluntary and not universally followed by platforms. Citron's approach, therefore, appears to strike a greater balance of efficiency and transparency than either court-led or platform-led approaches alone.

## V. Conclusion

There are a great number of topics in this discussion which, for the sake of brevity, I have been unable to cover: the inevitable complexity in balancing effective regulation with free expression; the global nature of the internet and its effects on enforcement methods; international regulatory approaches; and the potential for defamation law to be used by regulators in this area, to list but a few. However, this demonstrates part of the challenge of regulating deepfake technology: it is an area in which there are endless complexities and considerations, to the extent that there is almost certainly no one "correct" approach to the issue.

Unfortunately, regulators do not have the luxury of time; either deepfake technology is already advanced enough that it is capable of causing widespread and devastating harm to individuals, groups, and democracies, or it very soon will be, and it is clear that current legal avenues are inadequate in dealing with the legal tests that deepfake technology threatens. It is likely that the mitigation of these issues will require discussion, analysis, and maintenance for decades to come, but in order for the outcomes of such efforts to be successful, regulators are in a race against the clock to prepare today's legal systems for tomorrow's world.

---

[41] *Jensen v Clark* [1982] 2 NZLR 268

[42] *Jensen v Clark* above n 41 at 276

[43] *Jensen v Clark* above n 41 at 276

[44] Twitter Inc. "Introducing the new Twitter Transparency Center" blog.twitter.com (19 August 2020)

[45] Reddit Inc "Reddit Content Policy"

Bibliography

I.   CASES

   *Jensen v Clark* [1982] 2 NZLR 268
   *Living Word Distributors v Human Rights Action Group (Wellington)*
[2000] 3 NZLR 570
   *R v Iyer* [2017] DCR 82

II.   LEGISLATION

New Zealand

   Film, Video, and Publications Classifications Act 1993
   Harmful Digital Communications Act 2015
   New Zealand Bill of Rights Act 1990

United States

   Communications Decency Act (1996) 47 U.S.C. § 230
   Defending Each and Every Person from False Appearances by
Keeping Exploitation Subject to Accountability Act of 2019 H.R.3230 §1041
   United States Constitution, amend I

III.   BOOKS

   Danielle Citron "Hate crimes in cyberspace" (Harvard University
Press, 2014)
   George Orwell "Nineteen Eighty-Four" (Penguin Press, United
Kingdom, 2011)
   Ray Bradbury "The Illustrated Man" (Bantam Books, United States,
1967)

IV.   JOURNAL ARTICLES

   Bobby Chesney & Danielle Citron "Deep Fakes: A Looming
Challenge for Privacy, Democracy, and National Security" 107 Calif. L. Rev.
1753 (2019)
   Herbert B Dixon Jr "Deepfakes: More Frightening Than Photoshop
on Steroids" (Summer 2019) The Judges' Journal; Chicago Vol. 58, Iss. 3, 35-
37
   Ian Reilly "F for Fake: Propaganda! Hoaxing! Hacking! Partisanship!
and Activism! in the Fake News Ecology" (14 January 2018) The Journal of
American Culture 41:2, 139-152
   Ivan Perov and others "DeepFaceLab: A simple, flexible and
extensible face swapping framework" (12 May 2020) arXiv:2005.05535
[cs.CV]
   J. Naruniec, L Helminger, C Schroers and R.M. Weber "High-
Resolution Neural Face Swapping for Visual Effects" (29 June 2020) 39:4
Eurographics Symposium on Rendering 2020
   Jessica Ice "Defamatory Political Deepfakes and the First
Amendment" 70 Case W. Res. L. Rev. 417 (2019)

Robert Chensey and Danielle Citron "Deepfakes and the New Disinformation War: The Coming Age of Post-Truth Geopolitics" 98 Foreign Aff. 147 (2019)

Joshua S. Sellers "Legislating against Lying in Campaigns and Elections" 71 Okla. L. Rev. 141 (2018)

V.     PARLIAMENTARY AND GOVERNMENT MATERIALS

New Zealand

Classification Office to Chief Censor "Notice of decision under section 38(1)" (18 March 2019) OFLC Ref 1900148.000

United States

"House Intelligence Committee Hearing on 'Deepfake' Videos" *C-Span* (13 June 2019). <https://www.c-span.org/video/?461679-1/house-intelligence-committee-hearing-deepfake-videos>

Energy and Commerce Subcommittee "Hearing on 'Americans at Risk: Manipulation and Deception in the Digital Age'" (8 January 2020) <https://energycommerce.house.gov/committee-activity/hearings/hearing-on-americans-at-risk-manipulation-and-deception-in-the-digital>

VI.     REPORTS

Law Commission *The news media meets 'new media': rights, responsibilities and regulations in the digital age* (NZLC R128, March 2013)

VII.     INTERNET RESOURCES

Alissa Wilkinson "Hollywood is replacing artists with AI. Its future is bleak." *Vox* (online, 29 January 2020) <https://www.vox.com/culture/2020/1/29/21058521/hollywood-ai-deepfake-black-mirror-gemini-irishman-cinelytic>

Billy Bush "(Opinion): Yes, Donald Trump, you said that" *The New York Times* (Online, 3 December 2017)

Christiano Lima "'Nightmarish': Lawmakers brace for swarm of 2020 deepfakes" *Politico* (Online, 13 June 2019) <https://www.politico.com/story/2019/06/13/facebook-deep-fakes-2020-1527268>

DisneyResearchHub "High Resolution Neural Face Swapping for Visual Effects" (29 June 2020) YouTube <https://youtu.be/yji0t6KS7Qo>

Drew Harwell "Fake-porn videos are being weaponized to harass and humiliate women: 'Everybody is a potential target'" *The Washington Post* (Online, 30 December 2018)

Drew Harwell "Top AI researchers race to detect 'deepfake' videos: 'We are outgunned'" *The Washington Post* (Online, 12 June 2019)

Electronic Frontier Foundation "EFF and Coalition Partners Push Tech Companies To Be More Transparent and Accountable About Censoring User Content" (7 May 2018) <https://www.eff.org/press/releases/eff-and-coalition-partners-push-tech-companies-be-more-transparent-and-accountable>

Electronic Frontier Foundation "The Santa Clara Principles on Transparency and Accountability in Content Moderation" <https://santaclaraprinciples.org/>

Gilad Edelman "Twitter Finally Fact-Checked Trump. It's a Bit of a Mess" *Wired* (online, 27 May 2020) <https://www.wired.com/story/twitter-fact-checked-trump-tweets-mail-in-ballots/>

Hannah Denham "Another fake video of Pelosi goes viral on Facebook" *The Washington Post* (Online, 4 August 2020) <https://www.washingtonpost.com/technology/2020/08/03/nancy-pelosi-fake-video-facebook/>

Kelsey Warner "Instagram rolls out augmented reality filters after Snapchat's face swap success" *The National (*International Edition Online*,* 14 October 2019) <https://www.thenational.ae/business/technology/instagram-rolls-out-augmented-reality-filters-after-snapchat-s-face-swap-success-1.898200>

Makena Kelly "Congress grapples with how to regulate deepfakes and changes to Section 230 might be coming" *The Verge* (Online, 13 June 2019) <https://www.theverge.com/2019/6/13/18677847/deep-fakes-regulation-facebook-adam-schiff-congress-artificial-intelligence>

Reddit Inc "Reddit Content Policy" <https://www.redditinc.com/policies/content-policy>

Tom Scott "Faceswapping, Unethical Videos, and Future Shock" *YouTube* (6 February 2018) <https://youtu.be/OCLaeBAkFAY>

Twitter "Non-consensual nudity policy) (November 2019) <https://help.twitter.com/en/rules-and-policies/intimate-media>

Twitter "Parody, newsfeed, commentary, and fan account policy" <https://help.twitter.com/en/rules-and-policies/parody-account-policy>

Twitter "Synthetic and manipulated media policy" (n.d.) <https://help.twitter.com/en/rules-and-policies/manipulated-media>

Twitter Inc. "Introducing the new Twitter Transparency Center" blog.twitter.com (19 August 2020) <https://blog.twitter.com/en_us/topics/company/2020/new-transparency-center.html>

VIII.    OTHER RESOURCES

Peyton Reed "Through The Eyes of Forrest Gump: The Making of an Extraordinary Film" (2001) Paramount Pictures

Robert Zemeckis "Forrest Gump" (1994) Paramount Pictures

Eriq Gardner "Can Hollywood Adapt to the Deepfakes Era?" *The Hollywood Reporter* (10 July 2019) 52-54