
ON THE HORIZON: INTERACTIVE AND COMPOSITIONAL DEEPPAKES

Eric Horvitz
Microsoft
Redmond, Washington
horvitz@microsoft.com

ABSTRACT

Over a five-year period, computing methods for generating high-fidelity, fictional depictions of people and events moved from exotic demonstrations by computer science research teams into ongoing use as a tool of disinformation. The methods, referred to with the portmanteau of “deepfakes,” have been used to create compelling audiovisual content. Here, I share challenges ahead with malevolent uses of two classes of deepfakes that we can expect to come into practice with costly implications for society: *interactive* and *compositional* deepfakes. Interactive deepfakes have the capability to impersonate people with realistic interactive behaviors, taking advantage of advances in multimodal interaction. Compositional deepfakes leverage synthetic content in larger disinformation *plans* that integrate sets of deepfakes over time with observed, expected, and engineered world events to create persuasive *synthetic histories*. Synthetic histories can be constructed manually but may one day be guided by *adversarial generative explanation* (AGE) techniques. In the absence of mitigations, interactive and compositional deepfakes threaten to move us closer to a post-epistemic world, where fact cannot be distinguished from fiction. I shall describe interactive and compositional deepfakes and reflect about cautions and potential mitigations to defend against them.

Keywords synthetic media, deepfakes, digital content provenance, disinformation, causal models, multimodal interaction, multimodal neural models, adversarial generative explanation

1 Introduction

Democracy depends on an informed and engaged citizenry. Democracy and civil liberties are coming under threat from new forms of disinformation—the distribution of false information with the *intention* of swaying public opinion [Horvitz(2021a)]. Fabrications of falsehoods aimed at manipulating masses have a long and tragic history. Highly successful disinformation campaigns have relied on the expertise of propagandists to formulate persuasive false narratives and to propagate verbal and visual information in support of them. Disinformation efforts have grown in sophistication over time, riding on waves of technical advances, from the printing press, to photography, radio and television, and on to internet-based social media, computer graphics, and machine learning.

We are at an inflection point with the rising capabilities of discriminative and generative AI methods. The advances in machine learning are enabling new forms of content generation and new powers of multimodal interaction. The advances are providing unprecedented tools that can be used by state and non-state actors to create and distribute persuasive disinformation. The rising capabilities frame concerns that our children and grandchildren could find themselves in a post-epistemic world where it is difficult or impossible to distinguish fact from fiction. The speed of these technical developments and new possibilities for their abuse places responsibility in the hands of computer scientists to envision technical futures, likely abuses, and potential mitigations—and to engage across disciplines, organizations, and agencies to raise awareness and collaborate on practices, policies, and regulations.

Impressive developments in neural models for producing language, visual, and audiovisual content can be exploited for influence operations. Generative neural language models can be directed with ease to synthesize persuasive writings as well as to power compelling dialog aimed at achieving specific goals. Multimodal neural models, which leverage

representations constructed from text, images, audio, and video data, can generate realistic visual content based on natural language prompts. Advances with multimodal neural models include interactive techniques that provide creators with the ability to modify or refine the generated images, using such methods as in-painting, which enables the extension or manipulation of specific regions of images, and prompt engineering, the design of language inputs to produce desired outputs. While innovations with language-centric and multimodal neural models open up new possibilities for creative expression, imagination, and education, they can serve as potent weapons of persuasion and disinformation.

AI-generated content about people and events are now being employed in fraud, impersonation, and larger cyber influence programs. Compelling synthetic media, referred to commonly as “deepfakes,” were exotic research projects just a few years ago, first appearing as startling demonstration videos linked with conference papers. Today, open source toolkits are available for producing deepfakes, lowering the bar on required expertise to generate and then distribute them at lightning speed across social media.

We can expect deepfakes to become difficult to discriminate from reality. While numerous methods can be used to generate deepfakes, the challenge with distinguishing fact from fiction is easy to see for the generative adversarial networks (GAN) methodology [Goodfellow et al.(2014)]. GANs are an iterative technique where the machine learning and inference employed to generate synthetic content is pitted against systems that attempt to discriminate generated fictions from fact. Both the generator and the detector become increasingly better in the process, with the generator learning over time how best to fool the detector. With this process at the foundation of deepfakes, neither pattern recognition techniques nor humans will be able to reliably recognize deepfakes.

Turning to the focus of this paper, to date, deepfakes have been constructed and distributed as one-off, stand-alone creations. We can expect to see the rise of new forms of persuasive deepfakes that move beyond fixed, singleton productions [Horvitz(2022)]. Malevolent uses of *interactive* and *compositional* deepfakes will leverage advances in the multimodal interaction research community, and wider AI and HCI communities. I will describe these two new expected classes of synthetic media and touch on directions for defending against them.

2 Interactive Deepfakes

A constellation of advances in generative AI methods, coupled with frontier research on multimodal interaction, can enable new interactive forms of deepfakes. One of the earliest demonstrations of the use of generative methods for compelling impersonation was presented as an interactive prototype named Face2Face, published at CVPR 2016 [Thies et al.(2018)]. The project demonstrated how a commodity PC could be used to perform real-time tracking of a *source actor* to control the pose, mouth, and facial expressions of rendered *target actors*, including well-known politicians (see Figure 1). The authors reported that “the resulting synthesized model is so close to the input that it is hard to distinguish between the synthesized and the real face.” More recent generative methodologies can be coupled with real-time tracking to provide similarly rich interactive deepfakes.

Significant progress on methods for recognition, generation, and interaction pave the way to creating interactive deepfakes. Interactive deepfakes can raise the bar on the persuasiveness of impersonation, bringing new forms of “presence” and engagement via audio, visual, and audiovisual channels. Advances in enabling technologies include work on speech recognition, speech production, and appropriate visual renderings of expressions associated with voice activity that enable source actors to have their voice re-rendered as that of the rendered target actors. On voice impersonation, over a decade ago, methods for real-time *voice conversion* were demonstrated, enabling a source speaker to render utterances in another person’s voice. More recent advances include the use of deep neural models for generating natural-sounding voices from text [Arik et al.(2017)] and the efficient cloning of a target voice from just a few samples of speaking [Arik et al.(2018)].

An important challenge in multimodal interaction with rendered avatars is generating natural expressions associated with voice activity. Recent advances in methods for *neural voice puppetry* [Thies et al.(2020)] enable compelling real-time generation of appropriate expressions along with synthesized voice. As of 2020, the end-to-end pipeline for mapping audio features of source utterances to a person-specific expression *and* the generation of a photo-realistic rendering took 5ms on an Nvidia 1080Ti [Thies et al.(2018)]. Beyond human voices, audio capture and re-generation tools, such as WaveNet [Oord et al.(2016)], enable the generation of arbitrary background sounds in a deepfake. We are seeing advances in multimodal capabilities for recognition and generation that will soon enable audio, graphical, and language technologies to be woven into toolkits that enable the interactive control of fictional renderings of targeted personalities in teleconferences via the pose, expressions, and voice of a controlling actor.

Scale can be achieved with moving beyond manual control of multimodal impersonation technologies. Various forms of automation can be used to enable a system to convince viewers of the presence of an individual in audio calls and audiovisual conferencing. For example, ambient patterns of motion and attention have long been used in automated

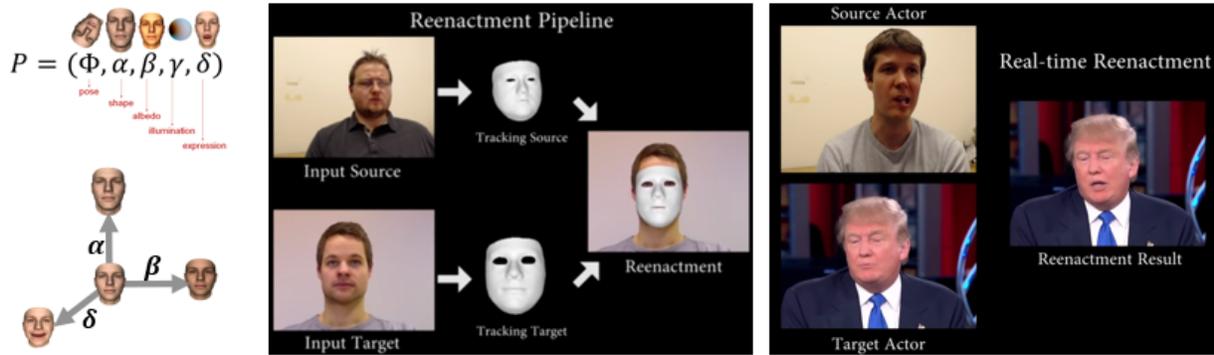


Figure 1: Early demonstration of end-to-end pipeline for real-time visual tracking and synthesis of face pose and expression drawn from sets of video frames (derived from [Thies et al.(2018)]).

visual agents. Expressions, gestures, and other behaviors, such as nods, laughter, applause, and patterns of attention, as captured by the pose and gaze of attendees in a multiparty teleconference call, could be mirrored by an automated avatar based on following and reflecting signals drawn from automated analyses of the identity, speech activity, speech recognition, pose, and expressions of attendees.

Simple, compelling strategies for projecting the presence of impersonating avatars in audiovisual conference calls include generating the usual greetings and goodbyes of teleconferences. More sophisticated directions with automation can harness rudimentary dialog capabilities that are available today to enable a rendered agent to respond appropriately to specific triggers in the flow of a conversation—such as the call for a vote, agreement or disagreement, or for an opinion which could flow from an agent on command.

Automated interactive deepfakes could be endowed with basic understandings of the status of flow of a conversation to inform decisions about if and when to interject, leveraging prior work in the multimodal interaction research community on predictive models of when the floor is yielded by a speaker to other participants in a multiparty setting [Bohus and Horvitz(2011)]. The prior work demonstrated the importance of sensing, prediction, and decision-making to guide the fine structure of the timing of turns in multiparty settings. Automated agents must be able to predict the source and target of utterances, and, more generally, the state and dynamics of shifts of the floor in multiparty conversations. Advances in multiparty understanding and interaction will enable new forms of automated impersonation.

Beyond fully manual and automated control of an impersonation, there are opportunities for mixed-initiative approaches, where a source actor can be on standby to take over basic automation of natural patterns of pose and expression as needed when alerted via signals about engagement, complexity, or confusion. With background matching and cloning, a rendering could be slipped into a live conversation without a participant recognizing the substitution, perhaps done briefly enough for an important intervention, such as agreeing with or voting on a proposal.

Interactive deepfakes can be enhanced in numerous ways with auxiliary audiovisual content, for instance, with introducing simulated or real events occurring in the background that are synchronized with appropriate responses of the impersonating avatar, such as reactions to nearby explosions. The methods could be employed to create persuasive fabricated outcomes, such as the real-time injury of a leader during a call or to convince the public that a leader who has perished is alive and in command.

3 Compositional Deepfakes

I use *compositional deepfakes* to refer to a concerning, feasible direction with malevolent uses of synthetic media in influence campaigns on a larger scale: the integration of multiple coordinated deepfakes and/or fabricated events with real-world occurrences to build fictional explanations or *synthetic histories*. Compositional deepfake *plans* include near real-time, as well as the pre- and post-dating one or more deepfakes that project the synthetic history into the past and the future, respectively. Figure 3 depicts a canonical compositional deepfake plan and resulting synthetic history, where a sequence of two fabricated “past” deepfake media pieces are injected between two world occurrences and time-stamped as happening at appropriate times between the two events. Moving into the future in this canonical synthetic history, an in-world event is fabricated to complete the persuasive storyline. As highlighted in the figure, the response to such plans can be monitored and updated with the creation and injection of additional past and future fictional or fabricated events.

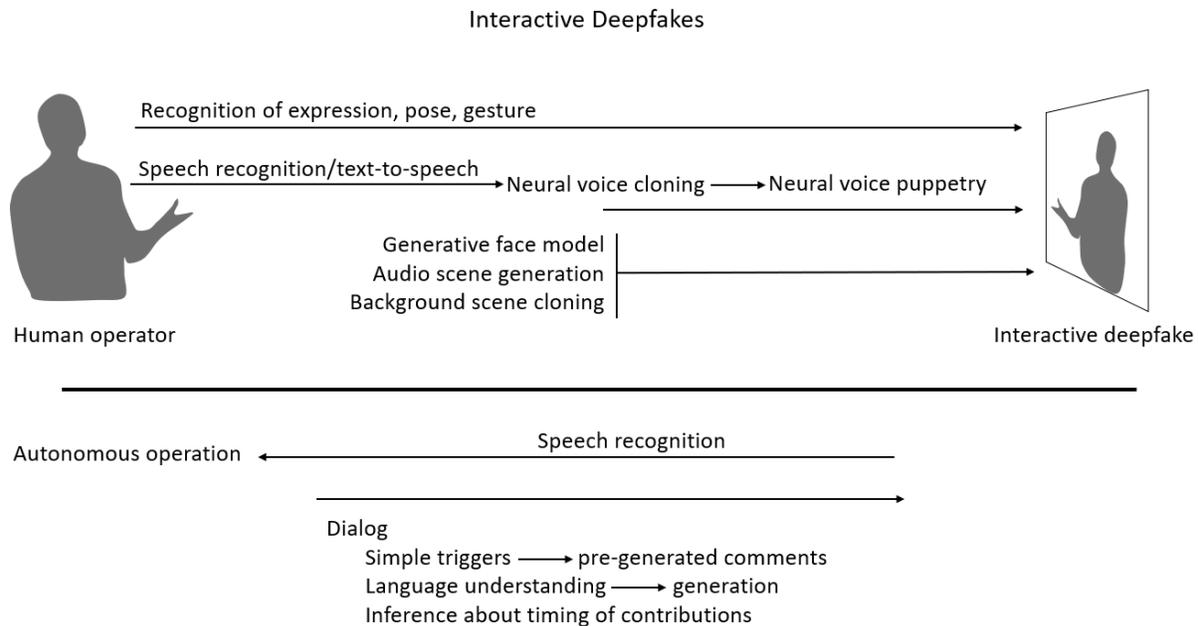


Figure 2: Interactive deepfakes. Top: Set of multimodal interaction technologies required for building an interactive deepfake system. Methods include use of high-fidelity renderings with generative AI methods, recognition of pose, expression, and gesture, speech recognition, neural voice cloning, and neural puppetry. Bottom: Opportunities for automation and mixed-initiative control span a range of sophistication, from simple trigger detection and issuance of pre-generated responses to richer models of language understanding and generation.

Compositional deepfakes can be designed to create fictional narratives that are persuasive in their ability to tie together and provide powerful explanations of sets of events in the world to citizens and government leaders. It is not hard to imagine how the explanatory power of custom-tailored synthetic histories could out-compete the explanatory power of the truthful narratives that describe the actual motivations, decisions, and causes of events that are observed in the world. State and non-state actors who seek to generate, execute, and monitor and refine the influences of persuasive compositional deepfakes can harness qualitative models and understandings about how people compose and come to be committed to beliefs about specific explanations. Such understandings include insights and practices developed over many decades to persuade populations, including long-honed nation-state propaganda and more modern commercial sales and marketing campaigns, and recent efforts to optimize engagement and clickthrough with online services that sit at the heart of modern commerce.

More fundamentally, compositional deepfake plans can leverage insights about human psychology, including studies of rich sets of biases identified and studied within the cognitive psychology area of judgment and decision making [Kahneman et al.(1982)]. Social and cognitive psychologists have studied how people weave sets of events in the world into convincing narratives. This research includes efforts to understand the formation and commitment of groups to *conspiracy beliefs*, defined as explanations for events or situations attributed to one or more actors working secretly to achieve goals that are unlawful, unfair, or malicious [Zonis and Joseph(1994)].

Several investigators have suggested that people are more predisposed to believe in conspiracies when they perceive that a society faces a crisis. They cite times of rapid change, such as the second industrial revolution at the start of the Twentieth Century and the period before World War II, with the internationalization of the economy, new forms of automation, and then the worldwide depression, as examples of this [Uscinski and Parent(2014)]. At such times, populations may feel particularly anxious and insecure, especially citizens who feel that they have little power or voice in decisions [Noble(1966)]. These periods of time are ripe for the growth of conspiracy theories that rise to challenge existing political leadership, behavioral norms, and acceptance of groups who may be viewed as outsiders [Pipes(1999)].

The design palette for compositional deepfakes includes methods for targeting specific individuals or groups, including the creation of multiple narratives, each custom-tailored to different populations. Borrowing from advances in e-commerce, other tools facilitate automated experimentation and refinement on subpopulations in advance of broader

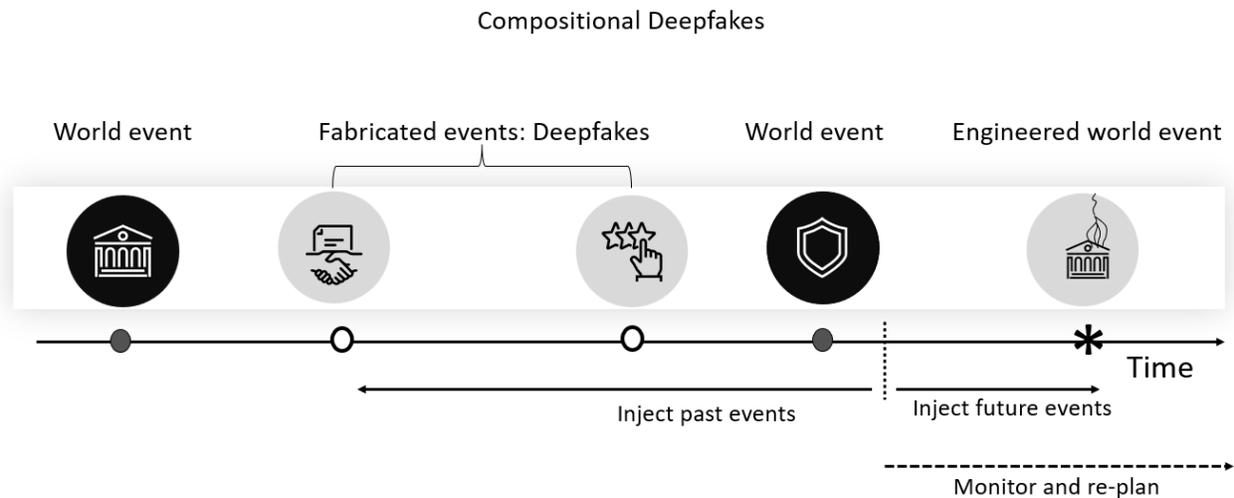


Figure 3: Compositional deepfake plans. Persuasive synthetic histories can be created via compositional deepfake strategies that interleave sequences of deepfakes and engineered in-world events with event occurrences in the world. In advance of general releases, testing and refinement can be performed on subpopulations. Follow-on monitoring can allow for fine-grained control of synthetic histories in response to population reactions via modification or extension of synthetic histories.

releases via testing of the finalized persuasive explanations. As indicated in the Figure 3, ongoing monitoring and refinement can be executed via the updating or “editing” of synthetic histories, where in-world events can be fabricated and deepfakes can be added anywhere on the timeline.

Other approaches include creating and pre-positioning one or more deepfakes or larger alternate compositional plans into obscure online holding areas at specific dates and following up with deletions and amplifications of content as needed, conditioned on the flow of real-world events. Multiple variants of deepfakes can be generated and quietly prepositioned and specific subsets can be brought forward for amplification or deleted depending on events in the world. Forms of pre-positioning have been described as a known disinformation strategy.

As an example cited as pre-positioning of disinformation [Microsoft(2022)], a false claim [Cercone(2022)] of U.S.-funded biolabs in Ukraine being connected to bioweapons development was positioned in a relatively obscure place on YouTube in November 2021 as part of a Russian television series. The story was lifted into prominence at the start of the invasion of Ukraine on February 24th, 2022, when it was simultaneously referred to as a known finding “from last year” by ten Russian-controlled or highly influenced news sites and then amplified on social media.

Beyond a goal of long-term persuasion, compositional deepfakes can be constructed for use in time-limited operations so as to achieve local goals regardless of whether the fabrications will eventually come to light. For example, the plans can incorporate the temporary disruption of electric power, sensing, or communication aimed at diminishing the ability of people to do real-time fact-checking or discovering evidence of manipulation. Benign explanations for sensing or communication outages can suppress suspicion of a link between an outage and the larger compositional deepfakes.

4 Adversarial Generative Explanation

Tragic situations and outcomes for humanity over the course of history are testimony to the ability of expert propagandists to construct and execute persuasive disinformation based on their intuitions and experience. Propagandists have demonstrated how they can weave together real-world observations, fabricated events, and fictional media to create powerful narratives, create and stoke conspiracy beliefs, and achieve goals of moving populations to action—or to acquiescence and inaction [Horz and Kocak(2022)]. Despite the demonstrated capabilities of people to manually author compositional disinformation operations, new forms of assistance and automation is feasible.

Sketches and more complete synthetic histories, along with recommended sets of disinformation actions, may one day be provided by “persuasion toolkits.” Recent advances in machine learning and inference can power new engines of persuasion in the form of advisory tools or automated services that can generate and persuasive narratives that run

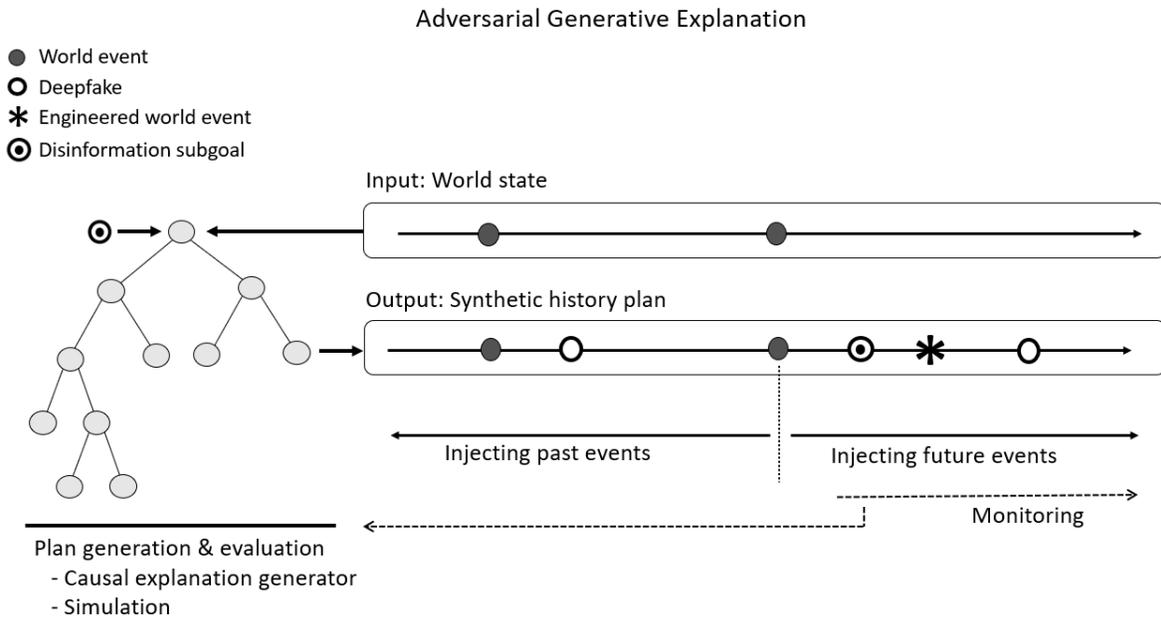


Figure 4: Adversarial generative explanation (AGE). Adversarial explanation systems employ causal inference and cognitive models to perform adversarial attacks on understandings of world events. AGE methods depend on search and optimization methods to generate persuasive explanations. The methods take as input key disinformation goals and in-world occurrences to date and generate and test persuasive synthetic histories composed of configurations of synthetic and in-world events. Objective functions for guiding the generation and search can take into consideration the persuasive power of plans and their ability to achieve key disinformation goals. Signals collected via monitoring can guide re-planning with additional content created and injected as prior or follow-on synthetic media and the engineering of additional in-world events.

counter to truthful explanations. Such systems could be employed for many goals, including political campaigns, legal disputes, and sales and marketing. They could also be aimed at guiding the creation and fielding of sequences of deepfakes and fabrications of real-world events to create persuasive compositional deepfakes as part of generating and target narratives. I refer to these hypothesized tools as *adversarial generative explanation (AGE)* systems.

AGE systems can be developed by leveraging causal reasoning and psychological models to generate narratives that run counter to truthful explanation of events and intentions. An AGE system could use search and optimization to create persuasive alternative or *contrastive* explanations for sets of observed events in the world. The generation of contrastive, false narratives could take into consideration prior understandings and biases in target populations about the desires, intentions, and actions of people. Such generations could be guided by specific disinformation goals input as sketches of false narratives or specifications about the actions of specific individuals or organizations as “must-include” events that are injected into the optimization.

Beyond building contrastive explanations based on existing or expected real-world events, AGE systems could be harnessed to generate compositional deepfakes, considering ideal sets of synthetic media and fabricated real-world events to introduce as part of generating explanations. Given a set of disinformation goals, recent sets of world events, a production budget, and time horizon, an AGE system could generate alternate synthetic histories, each potentially containing sets of recommended deepfakes and in-world events.

Figure 4 shows key components of a proposed AGE system, taking into consideration a set of known observations and then countering realistic causal understandings with the generation of adversarial narratives via a process of causal inference, psychological modeling, composition and search. The search can consider different sets of ideal injections of synthetic media and fabricated in-world events to achieve one or more goals, such as promoting a fictional history and narrative persuading populations about the actions and intentions of specific people, organizations, or nation states.

Causal plans can be generated and evaluated with the assistance of a simulator that uses psychological models in its evaluation. During execution, monitoring can provide input for ongoing revision of the plans, including accretion of new synthetic and fabricated in-world events. Beyond generating proposals for larger synthetic histories, such tools

could provide new forms of "causal in-painting," identifying ideal point-wise injections of synthesized events to overlay on real-world events so as to build a plausible story.

The feasibility of developing powerful AGE systems is supported by recent advances in machine learning, representation, and inference. The systems can draw on methods for performing causal reasoning and explanation [Pearl(2009), Spirtes et al.(2000), Halpern(2016)], including blossoming efforts for integrating causal knowledge into deep learning [Athey et al.(2021)]. Particularly relevant are recently developed techniques for constructing plausible explanations via multimodal inference that leverages visual, audio, and language capabilities to produce human-understandable causal explanations of events in the world [Gerstenberg et al.(2021)] and developments of benchmarks and methods for evaluating causal explanations [Roemmele et al.(2011)]. Such methods could be harnessed to generate and rank adversarial explanations.

Other AI advances that could provide core competencies to AGE systems include deep neural models that are trained to make commonsense "if-then" inferences about the consequences of actions in the world [Sap et al.(2019)] and about human feelings [Rashkin et al.(2018)]. Relevant advances also come via research on counterfactual reasoning aimed at enhancing explanation of automated inferences. This work spans efforts in machine learning [Verma et al.(2020), van der Waa et al.(2018)] and Bayesian networks [Koopman and Renooij(2021)]. AGE systems might also harness automated inferences about the representation of events in narratives [Chambers and Jurafsky(2008)], including methods for computing expected sequences of events [Sap et al.(2022)].

Moving from the computational to psychological realms, the generation and evaluation of candidate adversarial explanations would benefit from research on the psychological front, where researchers have developed qualitative causal models for describing the effectiveness of propaganda based in psychological models [Horz(2021), Horz and Kocak(2022), Woodward(2005), Lagnado et al.(2013), Lagnado(2021), Baker et al.(2017), Lombrozo and Carey(2006), Kirfel et al.(2021)] and models of persuasion [Wood and Eagly(1981)]. Other relevant research on the psychological front focuses on gaining understandings of language that can provoke strong emotions [Vu et al.(2014)], a factor demonstrated to play a role in the influence of disinformation [Martel et al.(2020)].

5 Preparing for Advances in Persuasion and Disinformation

What might be done to defend against the expected development of integrative and compositional deepfakes? Directions ahead span efforts and innovations in the realms of technology, policy, and practice.

Journalism and reporting. We need to nurture high-quality journalism and reporting, including the support of local and international news organizations. Efforts include ensuring that trusted reporters are on the ground to observe and record events. The rise of new technologies for impersonation, generating persuasive narratives, and manufacturing synthetic histories will raise the bar on expectations and requirements around reporting. Beyond professional reporting, we need to continue to explore opportunities to engage (and protect) citizen journalists, who can provide multiple, independent signals about world events, including the capture and sharing of photos and audiovisual content from multiple cameras, each with certifiable metadata, such as the location and time of content capture (see the efforts of Witness [Witness(2022)]).

Media literacy. We will need to foster media literacy and to raise awareness about new forms of manipulation and their growing power to impersonate, fabricate, and persuade. There is evidence that education, including special programs of "inoculation" and pre-bunking can help can raise alertness and resistance to various disinformation tactics [Roosenbeek et al.(2022)]. Educational programs will have to keep pace with the technical prowess of cyber influence operations, with special attention to raising awareness about the nature and operation of interactive and compositional deepfakes as these methods come into practice. Education and awareness needs to include efforts to disparage truthful explanations that we can expect to come with the rise of powerful disinformation methods: A world of pervasive, persuasive disinformation is conducive to the discrediting of actual happenings—a phenomenon referred to as the "liar's dividend." Such an approach to discrediting real-world events has been leveraged by malevolent actors who take advantage of common understandings about the ease with which photos can be doctored [Leibowicz(2021)].

Authenticity protocols. We will need to stay alert to new forms of generative content and to continue to pursue means for detecting and thwarting inappropriate uses in influence and deception. For example, special attention will be needed to assuring and asserting the identity of participants in critical private and public meetings. New authenticity confirming protocols, such as real-time *authenticity challenges*, may need to be introduced to identify interactive deepfakes via required tests of competency and knowledge. New practices of multifactor authentication of identity may become necessary for admittance into online meetings or appearances in videos.

Content provenance. We will need to rely increasingly on formal cryptographic pipelines and standards for authenticating the provenance of digital content. Methods and tools for certifying content provenance are recent developments. The methods employ systems and protocols that certify the source and history of edits made to photos or audiovisual content [England et al.(2021)]. Digital content provenance methods can raise the level of trust in digital content by helping consumers to understand the organizational source of the content, such as a trusted journalism organization. The methods employ a tamper-proof manifest that travels along with the content which includes the origin of the content and sequence of modifications that may have been made since the publication of the material [England et al.(2021), Aythora et al.(2020), Horvitz(2021b), C2PA(2022)]. The methods certify that media has not been modified beyond what is indicated in the manifest and that the holder of the cryptographic key created or modified the manifest. Editing and processing without a compliant content provenance pipeline invalidates the manifest.

Efforts in digital content provenance include work to push authentication closer to reality, with the goal of cryptographic “glass-to-glass” certification—that is, ensuring that the photons hitting the light-sensitive surface of cameras are faithfully rendered as photons on displays. This work includes efforts to build special phones that cryptographically encode time and location along with audiovisual content [Truepic(2022)]. There is work to be done on extending provenance to certifying reality itself. Such work will need to make use of intensive red-teaming aimed at ensuring that the chain of authentication cannot be broken nor gamed. For interactive deepfakes, it will be important to field real-time versions of digital content authentication.

Digital content provenance solutions have leveraged private and public distributed digital ledger technologies, including the common blockchain form of ledger. Distributed digital ledgers can provide indelible histories of sequences of media posts. The wide use of public digital ledgers for media will make it difficult to change history with false time stamps for newly posted audiovisual content or other attempts to rewrite the course of events.

Watermarks and fingerprints. Related to core efforts with digital content provenance are methods that embed in digital content an indelible watermark that withstands well-intentioned and adversarial edits and modifications. Potential watermarks include the use of encoded urls or other coding that point to stored versions of the original content and accessible via privately stored keys. This direction of effort includes the use of soft-hash fingerprints of the content itself that is stored in a database that includes information about the content and its creation. Work on indelible watermarks and robust fingerprinting could be employed to mark fabricated content as synthetic media. Such watermarking and fingerprinting could be useful for ensuring that synthetic media created for such uses as satire, envisioning, and art are not misinterpreted as capturing real-world events. Indelible watermarking should support signing by creators as well as enable for anonymous publication to protect the authors from retribution by abusive organizations and governments.

Detection. Research will be needed on the detection and disruption of compositional deepfake campaigns. Detection of compositional plans and synthetic histories will benefit from the development of new tracking tools and ongoing vigilance. Directions include careful monitoring of nation states and organizations with a history of developing disinformation for such behaviors as placing and removing pre-positioned content and sequencing of deepfakes. Adversarial generative modeling tools may be useful for generating, interpreting, and “pre-bunking” or reacting to compositional deepfakes.

Regulation and self-regulation. Moving to policy and law, nations and localities will need to consider balanced actions in the regulatory realm aimed at squelching the creation and influence of deepfakes for impersonation and other forms of disinformation, while enabling and protecting free speech. National and international conventions norms, regulations, and laws with stiff penalties may help to stem the tide of new forms of synthetic media aimed at disinformation. However, laws will have to be inspected carefully and be subjected to open debate and ongoing refinement. In the nearer-term, self-regulation of corporate and academic research labs may be helpful in keeping the most powerful tools and packages out of the hands of malevolent actors. Nonetheless, we must assume that innovations in multimodal interaction, generative AI, and causal modeling and explanation will spread worldwide at lightning speeds and will be harnessed with both beneficent and malevolent intentions.

Red-teaming and continuous monitoring. Whether in the realms of technology or policy, threat models and impact assessments will be needed and proposed solutions and mitigations must be carefully “red-teamed” to ensure robustness of the methods to creative attacks. Well-intentioned innovations aimed at addressing concerns can introduce new avenues for attacks on specific systems and end-to-end operations. As examples, we need to understand how content reported to come via the phones of multiple citizens might be fabricated and how manifests or watermarks on content might be removed and replaced, and how these and other attacks can lead to confusion or to the bolstering of disinformation.

6 Conclusion

We can expect that advances in machine learning and interaction will be leveraged in new forms of persuasion and disinformation. I touched on three concerning directions, including interactive deepfakes, compositional deepfakes, and adversarial generative explanation. Computer scientists innovating at the frontiers are uniquely qualified to anticipate scientific innovations and to kick off envisioning across multiple disciplines about how the advances may be harnessed, both for the greater good and for malevolent ends. Vigilance will be needed on potential uses of our models, data, results, and technologies for disinformation. As we progress at the frontier of technological possibilities, we must continue to envision potential abuses of the technologies that we create and work to develop threat models, controls, and safeguards—and to engage across multiple sectors on rising concerns, acceptable uses, best practices, mitigations, and regulations.

Acknowledgments

I thank Neil Coles, Paul England, Andrew Jenks, Ram Shankar Siva Kumar, Sarah McGee, Georgianna Shea, Allison Stanger, Subramaniam Vincent, Rand Waltzman, and Ben Zorn for their feedback on an earlier version of the manuscript.

References

- [Arik et al.(2018)] Sercan Arik, Jitong Chen, Kainan Peng, Wei Ping, and Yanqi Zhou. 2018. Neural voice cloning with a few samples. *Advances in neural information processing systems* 31 (2018).
- [Arik et al.(2017)] Sercan Arik, Mike Chrzanowski, Adam Coates, Gregory Diamos, Andrew Gibiansky, Yongguo Kang, Xian Li, John Miller, Andrew Ng, Jonathan Raiman, et al. 2017. Deep voice: Real-time neural text-to-speech. In *International Conference on Machine Learning*. PMLR, 195–204.
- [Athey et al.(2021)] Susan Athey, Yoshua Bengio, Eric Horvitz, and Judea Pearl. 2021. *Panel: Challenges and opportunities of causality*. <https://www.microsoft.com/en-us/research/video/panel-challenges-and-opportunities-of-causality>
- [Aythora et al.(2020)] Jatin Aythora, Rebecca Burke-Agüero, Amaury Chamayou, Sylvan Clebsch, M. Costa, Julianna Deutscher, N. Earnshaw, Lauren Ellis, Paul England, Cédric Fournet, Monique Anne Gaylor, Christina Halford, Eric Horvitz, Andrew Jenks, Kevin Kane, Matthew Lavalley, Stark Lowenstein, B. MacCormack, Henrique S. Malvar, Sean O’Brien, Judy Parnall, Elissa M. Redmiles, Alex Shamis, Isha Sharma, Jack W. Stokes, Sam Wenker, and Anika Zaman. 2020. Multi-stakeholder media provenance management to counter synthetic media risks in news publishing.
- [Baker et al.(2017)] Chris L Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B Tenenbaum. 2017. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour* 1, 4 (2017), 1–10.
- [Bohus and Horvitz(2011)] Dan Bohus and Eric Horvitz. 2011. Decisions about turns in multiparty conversation: From perception to action. In *Proceedings of the 13th international conference on multimodal interfaces*. 153–160.
- [C2PA(2022)] C2PA. 2022. *Overview of the Coalition for Content Provenance and Authenticity*. <https://c2pa.org/>
- [Cercone(2022)] Jeff Cercone. 2022. *There are no US-run biolabs in Ukraine, contrary to social media posts*. <https://www.politifact.com/factchecks/2022/feb/25/tweets/there-are-no-us-run-biolabs-ukraine-contrary-social-media-posts>
- [Chambers and Jurafsky(2008)] Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised Learning of Narrative Event Chains. In *Proceedings of ACL-08: HLT*. Association for Computational Linguistics, 789–797.
- [England et al.(2021)] Paul England, Henrique S Malvar, Eric Horvitz, Jack W Stokes, Cédric Fournet, Rebecca Burke-Agüero, Amaury Chamayou, Sylvan Clebsch, Manuel Costa, John Deutscher, et al. 2021. Amp: Authentication of media via provenance. In *Proceedings of the 12th ACM Multimedia Systems Conference*. 108–121.
- [Gerstenberg et al.(2021)] Tobias Gerstenberg, Max Siegel, and Joshua Tenenbaum. 2021. What happened? Reconstructing the past through vision and sound. (2021).
- [Goodfellow et al.(2014)] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in neural information processing systems* 27 (2014).
- [Halpern(2016)] Joseph Y Halpern. 2016. *Actual causality*. MIT Press.

- [Horvitz(2021a)] Eric Horvitz. 2021a. *Critical attention required on the increasing sophistication and scope of disinformation*. Microsoft. https://www.erichorvitz.com/Disinformation_chapter_intro_MDDR_2021.pdf
- [Horvitz(2021b)] Eric Horvitz. 2021b. *A promising step forward on disinformation*. <https://blogs.microsoft.com/on-the-issues/2021/02/22/deepfakes-disinformation-c2pa-origin-cai/>
- [Horvitz(2022)] Eric Horvitz. 2022. Artificial Intelligence and Cybersecurity: Rising Challenges and Promising Directions. In *Hearing on Artificial Intelligence Applications to Operations in Cyberspace before the Subcommittee on Cybersecurity, of the Senate Armed Services Committee, 117th Congress (May 3, 2022) (testimony of Eric Horvitz)*. <https://aka.ms/AAhee56>
- [Horz and Kocak(2022)] Carlo Horz and Korhan Kocak. 2022. How To Keep Citizens Disengaged: Propaganda and Causal Misperceptions. (2022).
- [Horz(2021)] Carlo M Horz. 2021. Propaganda and skepticism. *American Journal of Political Science* 65, 3 (2021), 717–732.
- [Kahneman et al.(1982)] Daniel Kahneman, Stewart Paul Slovic, Paul Slovic, and Amos Tversky. 1982. *Judgment under uncertainty: Heuristics and biases*. Cambridge university press.
- [Kirfel et al.(2021)] Lara Kirfel, Thomas Icard, and Tobias Gerstenberg. 2021. Inference from explanation. *Journal of Experimental Psychology: General* (2021).
- [Koopman and Renooij(2021)] Tara Koopman and Silja Renooij. 2021. Persuasive contrastive explanations for Bayesian networks. In *European Conference on Symbolic and Quantitative Approaches with Uncertainty*. Springer, 229–242.
- [Lagnado(2021)] David A Lagnado. 2021. *Explaining the evidence: How the mind investigates the world*. Cambridge University Press.
- [Lagnado et al.(2013)] David A Lagnado, Tobias Gerstenberg, and Ro'i Zultan. 2013. Causal responsibility and counterfactuals. *Cognitive science* 37, 6 (2013), 1036–1073.
- [Leibowicz(2021)] Claire Leibowicz. 2021. Preparing for a World of Holocaust Deepfakes. *Tablet Magazine* (2021). <https://www.tabletmag.com/sections/news/articles/holocaust-denial-deepfakes-misinformation-claire-leibowicz>
- [Lombrozo and Carey(2006)] Tania Lombrozo and Susan Carey. 2006. Functional explanation and the function of explanation. *Cognition* 99, 2 (2006), 167–204.
- [Martel et al.(2020)] Cameron Martel, Gordon Pennycook, and David G Rand. 2020. Reliance on emotion promotes belief in fake news. *Cognitive research: principles and implications* 5, 1 (2020), 1–20.
- [Microsoft(2022)] Microsoft. 2022. *Defending Ukraine: Early Lessons from the Cyber War*. <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RE50K0K>
- [Noble(1966)] David Noble. 1966. The Paranoid Style in American Politics and Other Essays by Richard Hofstadter. *The Canadian Historical Review* 47, 4 (1966), 373–375.
- [Oord et al.(2016)] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499* (2016).
- [Pearl(2009)] Judea Pearl. 2009. *Causality*. Cambridge university press.
- [Pipes(1999)] Daniel Pipes. 1999. *Conspiracy: How the paranoid style flourishes and where it comes from*. Simon and Schuster.
- [Rashkin et al.(2018)] Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A. Smith, and Yejin Choi. 2018. Event2Mind: Commonsense Inference on Events, Intents, and Reactions. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 463–473.
- [Roemmele et al.(2011)] Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of Plausible Alternatives: An Evaluation of Commonsense Causal Reasoning. In *AAAI spring symposium: Logical formalizations of commonsense reasoning*. 90–95.
- [Roozenbeek et al.(2022)] Jon Roozenbeek, Sander van der Linden, Beth Goldberg, Steve Rathje, and Stephan Lewandowsky. 2022. Psychological inoculation improves resilience against misinformation on social media. *Science Advances* 8, 34 (2022).

- [Sap et al.(2022)] Maarten Sap, Anna Jafarpour, Yejin Choi, Noah A. Smith, James W. Pennebaker, and Eric Horvitz. 2022. Imagined versus Remembered Stories: Quantifying Differences in Narrative Flow. (2022). <https://arxiv.org/abs/2201.02662>
- [Sap et al.(2019)] Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 3027–3035.
- [Spirtes et al.(2000)] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. 2000. *Causation, prediction, and search*. MIT press.
- [Thies et al.(2020)] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. 2020. Neural voice puppetry: Audio-driven facial reenactment. In *European conference on computer vision*. Springer, 716–731.
- [Thies et al.(2018)] Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2018. Face2Face: Real-Time Face Capture and Reenactment of RGB Videos. *Commun. ACM* 62, 1 (dec 2018), 96–104. <https://doi.org/10.1145/3292039>
- [Truepic(2022)] Truepic. 2022. *Truepic corporation homepage*. <https://truepic.com/>
- [Uscinski and Parent(2014)] Joseph E Uscinski and Joseph M Parent. 2014. *American conspiracy theories*. Oxford University Press.
- [van der Waa et al.(2018)] Jasper van der Waa, Marcel Robeer, Jurriaan van Diggelen, Matthieu Brinkhuis, and Mark Neerincx. 2018. Contrastive explanations with local foil trees. *arXiv preprint arXiv:1806.07470* (2018).
- [Verma et al.(2020)] Sahil Verma, John Dickerson, and Keegan Hines. 2020. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596* (2020).
- [Vu et al.(2014)] Hoa Trong Vu, Graham Neubig, Sakriani Sakti, Tomoki Toda, and Satoshi Nakamura. 2014. Acquiring a dictionary of emotion-provoking events. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*. 128–132.
- [Witness(2022)] Witness. 2022. *Witness organization homepage*. <https://www.witness.org/>
- [Wood and Eagly(1981)] Wendy Wood and Alice H Eagly. 1981. Stages in the analysis of persuasive messages: The role of causal attributions and message comprehension. *Journal of personality and Social Psychology* 40, 2 (1981), 246.
- [Woodward(2005)] James Woodward. 2005. *Making things happen: A theory of causal explanation*. Oxford university press.
- [Zonis and Joseph(1994)] Marvin Zonis and Craig M Joseph. 1994. Conspiracy thinking in the Middle East. *Political Psychology* (1994), 443–459.