

# Labeling Deepfake Videos Reduces Exposure But Not Persuasiveness

Amanda Chen<sup>a,1</sup>, David Hagmann<sup>a</sup>, Christie Pang<sup>b</sup>, and George Loewenstein<sup>c</sup>

Increasingly realistic AI generated deepfake videos are raising alarm about their potential use in disinformation campaigns. Viewers may be misled into believing the depicted content is real and thus get persuaded on policy-relevant beliefs. A commonly proposed solution is to require disclosure, such as via text embedded in the video that informs the audience that the content is AI generated. Across four large scale experiments (N = 7,107), we show that the concern about persuasion is justified, but that disclosure is only partially effective. Specifically, participants who watched realistic but AI generated videos in which a speaker provided arguments related to regulation of AI technology were persuaded regardless of whether it included a warning label, which did not diminish the video's persuasiveness. Viewers rated labeled deepfakes as less deceptive, but otherwise evaluated them and their senders no differently from unlabeled deepfakes. However, disclosure and labeling are effective at the extensive margin: people were less likely to watch a video that is labeled as a deepfake, and less likely to share it with others. Our findings suggest that labels are effective when users opt-in to viewing the content, or can share it with others—but not in contexts such as news media feeds that don't offer users the opportunity to selectively view, or not view, content.

Deepfakes | Misinformation | Disclosure | Persuasion

Videos and images generated by artificial intelligence have become increasingly realistic, often making it difficult to distinguish them from authentic recordings. Synthetic media that are difficult to distinguish from reality (so-called “deepfakes”) pose risks for elections, markets, and national security because they can be deployed at scale to mislead audiences and erode trust in information ecosystems (1–3). Recent evidence underscores both their realism and their psychological potency: AI-synthesized faces can be indistinguishable from real faces and even perceived as more trustworthy (4), and synthetic political videos can alter judgments about truth and news credibility (5).<sup>\*</sup> These concerns have motivated policymakers, platforms, and news organizations to seek policies that can avoid the deception of even highly discerning consumers.

One potential solution relies on disclosure regimes (textual or on-screen labels indicating that content was generated by AI) as a potential remedy. Disclosure is a familiar policy tool. In domains ranging from consumer finance to medicine and journalism, disclosure is frequently mandated on the intuition that informing people will enable them to discount biased or low-quality information (6). When recipients are informed about problems related to the source of information, they may become more skeptical, recognize strategic incentives, and subsequently ignore the message and resist attempts at persuasion. Disclosures also promote transparency and accountability norms, which may further reduce naive acceptance of deceptive content. Some jurisdictions and institutions are extending this logic to AI generated content by requiring or encouraging labels that identify deepfakes (for example the European Union's AI Act, which requires watermarking for AI generated or manipulated video, 7).

However, a large literature in psychology, economics, and law cautions that mandated disclosure often has limited effects and can sometimes backfire: disclosures may be overlooked due to limited attention or cognitive load, misunderstood, or even license more biased behavior and greater compliance with conflicted advice—for example, when advisees feel obligated to heed disclosed advice or advisers feel morally “absolved” after disclosure (8–11). Relatedly, corrections and warnings often fail to fully neutralize misinformation due to continued-influence effects (12). In a related context, informing people that they have been microtargeted has no

<sup>\*</sup>Generative AI technology is advancing at a rapid pace. In the referenced work, viewers became skeptical of real content rather than being misled. But the ability to generate plausible artificial video content has rapidly advanced over the past few years and, we suspect, will continue to get harder to distinguish from real content.

## Significance Statement

Realistic deepfake videos are concerning to policymakers, journalists, and academics who worry about the spread of misinformation. One solution implemented in some jurisdictions and by some news organizations is to provide disclosure, affixing a clear label to such videos. While one may think and hope that this reduces the impact of deepfakes on viewers, we find that this is only partially true: deepfake videos labeled as such are no less persuasive than those without disclosure. However, disclosure reduces both demand for and supply of deepfake videos (willingness to watch and willingness to share, respectively). This suggests that labeling may be effective in contexts like social media, but not when aired by trusted news sources as part of their programming.

Author affiliations: <sup>a</sup>The Hong Kong University of Science and Technology, Department of Management, Hong Kong; <sup>b</sup>The Hong Kong University of Science and Technology, Division of Emerging Interdisciplinary Areas, Hong Kong; <sup>c</sup>Carnegie Mellon University, Department of Social and Decision Sciences, Pittsburgh, 15213

GL initiated the research question and all authors contributed to the design of the studies. CP generated the deepfake videos. AC programmed the survey experiments and conducted the analyses. DH replicated the analyses. AC and DH wrote the first draft of the manuscript and all authors edited and approved the final version.

The authors declare no conflicts of interest.

<sup>1</sup>To whom correspondence should be addressed. E-mail: zchengj@connect.ust.hk

125 impact on the persuasiveness of the message (13). Thus, while  
126 labels may increase awareness that content is AI generated,  
127 it remains unclear whether they meaningfully reduce the  
128 persuasive impact of what that content says.<sup>†</sup>

129 Disclosure of how content was created can reduce persuasion  
130 to the extent that it changes how people evaluate the message.  
131 That is, they may come to understand that the video itself is  
132 not real and then ignore the position that is being presented.  
133 However, it is possible that people are nonetheless persuaded  
134 because they may still attend to the underlying argument,  
135 which is independent of whether the speaker actually held  
136 the view inherent in what they are shown as saying. Indeed,  
137 persuasion can occur even when audiences know content is  
138 fictional: narratives routinely shift beliefs and intentions,  
139 including when they are not presented as originating from a  
140 fictitious communicator (15). In one notable example, access  
141 to a Brazilian telenovela that featured smaller families reduced  
142 fertility, consistent with viewers updating expectations and  
143 social norms from fictional portrayals (16). Thus, raising  
144 awareness about the AI generated nature of a video may  
145 make people view it as fictional, yet the content may remain  
146 persuasive. We test whether clearly identified deepfake  
147 content remains persuasive in our first study.

148 Another way labels could affect persuasion is by affecting  
149 exposure to the content. If people understand that a video  
150 is a deepfake, they may want to avoid watching it (17).  
151 Thus, even if a deepfake is as persuasive as a real video  
152 conditional on exposure, people may be less likely to view it.  
153 For example, while false claims about COVID-19 vaccines  
154 (flagged as misinformation) nonetheless reduced vaccination  
155 intention when viewed, they were viewed less frequently than  
156 unflagged content (18). However, exposure may have been  
157 influenced by Facebook’s downranking policy. Labels which  
158 simply inform participants that a social media image was  
159 generated using AI have only limited impact on participants’  
160 willingness to engage with the posts (14). On the other hand,  
161 people have been found to be more receptive to opposing  
162 views when expressed by AI (19) and may be curious to see  
163 realistic AI generated content due to its novelty and thus  
164 may be more likely to watch something labeled as a deepfake  
165 (20, 21). This could then lead to more persuasion overall. We  
166 test this experimentally in the context of AI generated video  
167 content in Study 2.

168 Relatedly, people may be less willing to share a deepfake with  
169 others, thus also reducing how likely others are to have the  
170 chance to view it. Prior work shows that simple prompts  
171 that shift attention to accuracy can substantially reduce  
172 intentions to share misinformation, and that false content  
173 spreads differently from true content in social networks (22,  
174 23). Thus, it may be that people second-guess their decision  
175 to share a video if they are informed that it is a deepfake.  
176 Moreover, people may also hesitate more to share content  
177 that they know will be identified as a deepfake to downstream  
178 viewers. We test this in Study 3, and recruit a separate sample  
179 to evaluate perceptions of those who sent AI generated content  
180 (Study 3b).

181 Our predictions point to two margins on which labels could  
182 matter. On the *intensive* margin, labels could reduce the

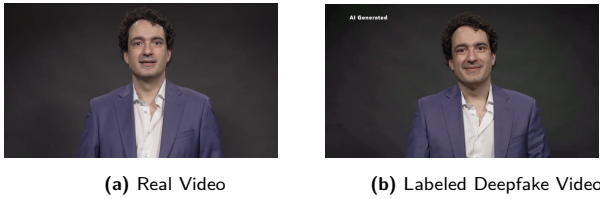
183 <sup>†</sup>There are multiple ways to label AI generated content; for example distinguishing between  
184 process-based labels, which reveal how content is created, and harm-based labels, which  
185 highlight its potential to mislead (14). We focus on labels that discuss how the content was  
186 created to parallel existing and proposed disclosure requirements.

187 credibility and trustworthiness of the source and thus make  
188 the content less persuasive. The prior literature on disclosure,  
189 however, suggests that this is unlikely. On the *extensive*  
190 margin, labels could reduce the number of people who are  
191 exposed to the content via the demand side (willingness to  
192 watch the video) and via the supply side (willingness to share  
193 the video). This second pathway is relevant in the many cases  
194 in which people choose what content to consume and share,  
195 such as on social media.

196 Across four large-scale experiments (N = 7,107), we test  
197 whether disclosing how content was created shapes the  
198 persuasive impact of AI generated media, the core proposition  
199 underlying AI-labeling policies. In Study 1, participants  
200 watched a video seeking to persuade them for or against  
201 the need for regulating AI technology. They were randomly  
202 assigned to watch a real video, an AI generated deepfake, or  
203 the identical deepfake but with a label displayed with the  
204 video. In Study 2, participants chose whether to watch a  
205 video in favor of or against the need for AI regulation. In one  
206 treatment, we informed them that one of the videos was a  
207 deepfake and our measure is which video they elected to watch.  
208 In Study 3, participants viewed both videos (in favor of and  
209 opposed to the need for AI regulation). One of the videos was  
210 real and the other was a deepfake, which was either unlabeled  
211 or labeled. Our main outcome measure is which of the videos  
212 they shared with another participant. Finally, in Study 3b,  
213 we asked participants to watch the shared videos and asked  
214 them to evaluate the videos as well as their perceptions of  
215 the participant who shared them. Across our studies, we find  
216 that labels are effective at the extensive margin—reducing  
217 voluntary exposure and the willingness to share—but not at  
218 the intensive margin—persuasion conditional on exposure.  
219 Sharing a labeled as opposed to an unlabeled deepfake video  
220 also did not change perceptions of the sender, but people  
221 thought a labeled deepfake video was less deceptive. These  
222 results suggest that labels may help prevent the spread of false  
223 beliefs when people can opt in or out of viewing and sharing  
224 content, but are unlikely to reduce persuasion when labeled  
225 deepfakes are aired within news programming or otherwise  
226 passively encountered.

## 227 Experiments

228 We generated two scripts on the topic of AI regulation using  
229 ChatGPT, with one arguing in favor of AI regulation, while  
230 the other is opposed. The scripts are of approximately equal  
231 length (see the OSF repository for the full texts). We then  
232 recorded a professor reading each script with the help of  
233 a teleprompter, with each video lasting a little over four  
234 minutes. Then, we used publicly available recordings of the  
235 same professor to train an algorithm to generate two deepfake  
236 videos using LatentSync. Specifically, we took the video  
237 arguing in favor of AI regulation and generated a deepfake  
238 such that the speaker instead delivered the speech opposed  
239 to AI regulation, and vice versa (“deepfake” videos). Thus,  
240 we have for each argument (for and against regulation) a  
241 real video and a paired AI generated deepfake with the same  
242 content. We then generated a “labeled deepfake” by adding  
243 text identifying the video as “AI Generated” (see Figure 1).  
244 We used these materials throughout our experiments.  
245 In Study 1, we examine if adding a disclosure label for AI  
246 generated deepfake videos reduces persuasion among viewers.  
247  
248



**Figure 1.** Screenshots displaying the real video (left side) and the watermark labeling used in the deepfake video (right side).

We recruited participants from Prolific ( $n = 1,800$ ;  $M_{\text{Age}} = 42.5$ , 47.8% Male) and first measured their attitudes toward AI regulation using a six-item index, such as “Government regulation is necessary to ensure that generative artificial intelligence is used ethically,” and “Government regulation will stifle innovation in generative artificial intelligence (reverse coded)” using sliders from -100 to 100, with endpoints labeled as “Strongly disagree” and “Strongly agree” and the midpoint labeled “Neither agree nor disagree.” We averaged the responses across the items (three items reverse-coded,  $\alpha = 0.86$ ). We then randomly assigned them to one of six treatments, varying both the version of the video and the content of the speech. Participants viewed either the real recording of the argument, a deepfake in which the speaker was altered to present the script for the other position, or a labeled deepfake in which they were informed in advance that the video was created using artificial intelligence and the video contained a watermark. They were further randomly assigned to view either arguments for or against regulation. After watching the video, we asked participants the same six items to measure their position on AI regulation, with their previous responses pre-selected. We calculated a measure of persuasion that is the change in the index, signed in the direction of the argument in the video. We further asked participants how persuasive they thought the video was in shaping their views (5-point Likert scale ranging from “not at all” to “extremely”) and whether they thought it was generated using AI (“Yes” or “No”). We incentivized participants to pay attention by informing them prior to watching the video that they would later be asked three questions about the content and could earn a bonus of 10 cents for each question they answer correctly (on average, they answered 1.97 questions correct). The survey ended with basic demographic questions.

Our preregistered primary outcome is the attitude change in the direction of the argument in the video. We fit an OLS regression with dummy variables for whether the video is a Deepfake and whether there is a Label present, such that the coefficient for the Label estimates the impact of the disclosure on the persuasiveness of a deepfake video. We also include a dummy variable to capture whether the video advocated for or against regulation because arguments may not be equally persuasive or participants who favor or oppose regulation may not be equally persuadable. We show these results in Table 1 and, visually, in Figure 2. Notably the real video opposing regulation was very persuasive, changing beliefs by almost 18 points on average. The deepfake video was not significantly less persuasive. And, consistent with our preregistered prediction, our key finding shows that disclosure via adding a label did not reduce the impact on beliefs relative to the unlabeled deepfake. When looking at our Likert

measure (column 2), participants who were assigned to watch the deepfake video reported it to be less persuasive, but no less so if they saw the video that was labeled. Participants who saw the pro-regulation video reported it as more persuasive, while their beliefs changed less than those who saw the anti-regulation videos.

Among participants watching the labeled deepfake video, 90.4% reported that they thought it was generated with AI, which was significantly greater than those watching the unlabeled deepfake (77.8%;  $\chi^2(1, n = 1198) = 34.51, p < .001$ ). Thus, the watermark informed some people who would otherwise not have noticed it. Notably, some participants did not trust the label and still believed it to be a real video. Conversely, among participants who watched the real video, 27.2% incorrectly thought it was AI generated.

Both real and deepfake videos shifted attitudes. Adding an on-screen “AI generated” label to deepfakes increased the share of participants who recognized the video as a deepfake, but, consistent with our prediction and prior work on the limits of disclosure, did not significantly reduce the video’s impact on beliefs. However, labels may still meaningfully mitigate downstream influence by reducing exposure and distribution in contexts where audiences opt in to viewing content, or have control over whether to forward it. That is, people may not wish to be exposed (or expose others) to content they know is fake. Conversely, participants may be *curious* about AI generated videos, particularly given their novelty, and thus seek them out even more.

Study 2 examines whether disclosure labels impact voluntary exposure to belief-incongruent deepfakes. Prior work suggests that participants prefer exposure to content that aligns with their own position (24). This study compares whether labeling a counter-attitudinal deepfake video as such further depresses viewing rates or instead increases viewing rates because participants may be curious about what the video looks like. We recruited participants from Prolific ( $n = 1,551$ ;  $M_{\text{Age}} = 45.04$ , 47.7% Male) and first measured their attitudes toward AI regulation using the same six-item index as in Study 1 ( $\alpha = 0.8$ ) and classified them as being either for regulation (a score above zero) or against regulation (below zero).<sup>‡</sup> Participants then chose one of two videos to watch: one aligned with their initial attitude (the Real video) and one opposed to it (the Deepfake). Participants were randomly assigned to a Label condition in which the belief-opposing video was disclosed as “AI generated” prior to their choice and later with an embedded watermark, or a Baseline condition in which they were not informed that it was a deepfake. After viewing the chosen video, participants again reported their attitudes on the same scale. The key outcomes were (i) whether participants chose the belief-opposing video (exposure), and (ii) the magnitude of attitude change toward the opposing view, regardless of which video they chose. Participants again reported whether they thought the video was AI generated, and we incentivized them to listen to the arguments with a bonus for correctly answering questions about the content of the video.

We first estimate how labels affect people’s willingness to select the belief-opposing video. Contrary to our expectation, curiosity did not lead people to opt into watching the

<sup>‡</sup>14 participants had an average index score of zero and were randomly assigned to one of the two positions.

AI generated video. Instead, the label effectively reduced exposure from 38.7% to 30.5% ( $\chi^2(1, n = 1551) = 11.26, p < .001$ ). We report in column 3 of Table 1 a linear probability model that also controls the direction of the argument presented in the deepfake video. Participants who were opposed to regulation (and hence whose deepfake video corresponded to a pro-regulation argument) were more likely to select the deepfake video ( $p < 0.001$ ). Notably, however, even controlling for this, the effect of the label still reduced willingness to select the deepfake video by 6.3 percentage points ( $p < 0.01$ ).

We next look at whether labeling had an impact on aggregate beliefs. For example, it may be that people who are deterred from watching the video would not have been persuaded by the argument. We therefore compare the change in beliefs across all participants in the two treatments, regardless of which video they watched. Participants who were informed that the video was AI generated prior to their choice revised their attitude by 3.2 points in the direction of their initial position—that is, they became more convinced about their initial view. Those in the Baseline condition revised their attitudes by 0.98 away from their initial position ( $t(1549) = 3.07, p = .002$ ). Labeling the deepfake thus reduced the persuasive impact of the AI generated argument by deterring people from watching it who would otherwise have been persuaded. Contrary to our prediction and concern, curiosity about AI generated content did not increase voluntary exposure.

Study 3 evaluates whether disclosure and labeling reduce downstream diffusion by lowering willingness to share deepfakes, as those spreading false information may worry about reputational concerns (25). Participants recruited from Prolific ( $n = 1,832; M_{\text{Age}} = 44.13, 45.3\%$  Male) began by completing the six-item scale measuring their attitude toward AI regulation ( $\alpha = 0.86$ ). Participants with a score above zero were classified as pro-regulation, and those with a score below zero as anti-regulation.<sup>§</sup> They then watched two videos: a real video arguing against their position, and a deepfake video with arguments consistent with their own view. Unlike in the previous study, we set the belief-congruent video to be the deepfake because participants were likely most motivated to share content that supported their own position, and we wanted to test whether knowing that this content was AI generated was sufficient to dampen that sharing motive. We then randomly assigned them to one of three conditions that differed only in the information they received about the deepfake video. In the Baseline condition, they were not informed that the video was AI generated. In the Label condition, participants were informed that the video aligned with their view was AI generated and the video had a watermark identifying it as such. Finally, in the Disclosure condition, participants were told that the video was a deepfake but there was no watermark. Participants then watched both videos and chose which video to share with another user. Importantly, in the Disclosure condition, participants knew that the recipient would not be informed that the video they shared was a deepfake.

Participants whose belief-congruent deepfake video was not disclosed as being AI generated shared it more than half

the time (62%,  $t(613) = 6.15, p < .001$ ). When participants were informed that the video they were about to watch was a deepfake, their likelihood of sharing it declined by 16.5 percentage points ( $t(1219) = 5.83, p < .001$ ), and by 17.6 percentage points if they were informed and the video had a watermark that would also inform the recipient ( $t(1223) = 6.24, p < .001$ ). Notably, whether the watermark was included (and hence whether the recipient would also be aware that the participant knowingly shared a deepfake video) did not reduce willingness to share significantly ( $t(1216) = 0.39, p = .696$ ). After watching both videos and making their sharing decisions, participants evaluated whether they thought each of the two videos was real or AI generated. Both disclosure and disclosure with a watermark helped participants identify that deepfake content was AI generated (83.2% and 84.8%, respectively, versus 75.1% without disclosure, both  $p < 0.001$ ).

Accuracy did not differ with or without the watermark ( $t(1216) = 0.75, p = .452$ ). Notably, we also find evidence for an implied truth effect (26), albeit one that here is not an error. In the absence of a label, 38.4% of participants thought the real video was in fact AI generated as well. The disclosure and watermark treatments lower this share to 24.4% and 22.1%, respectively (both  $p < 0.001$ ). That is, about 15% of participants who were informed that one of the videos was AI generated inferred from the lack of disclosure about the other video that it must have been real.

Finally, we can examine one pathway through which disclosure may deter sharing: it can eliminate plausible deniability. People may believe that a video is a deepfake, but given some uncertainty, they nonetheless share the video to propagate their own position. Disclosure without a watermark removes this channel, even though the recipient would not be aware that the sender knowingly shared misinformation. The watermark, finally, removes the uncertainty also for the recipient, who would then know that the sender knowingly shared fake content.

To examine whether disclosure deters sharing by eliminating plausible deniability, we look at the subset of participants across the three treatments who correctly identified the deepfake video as such (i.e., not participants who may have genuinely thought that the deepfake video was real). This analysis was not preregistered. Among participants who were not informed and hence correctly inferred this on their own, 58.1% nonetheless shared the deepfake video. Among those who were informed (but who knew that the recipient would not be informed), this dropped by 13.1 percentage points ( $t(964) = 4.13, p < .001$ ), and by 14.9 percentage points among those whose video was also watermarked ( $t(977) = 4.70, p < .001$ ); but again, the difference between the two informed treatments was not different ( $t(1021) = 0.55, p = .583$ ). This suggests that the decrease in sharing results from the reduction in plausible deniability: when people merely suspect that it was AI generated, they were nonetheless willing to share it. But when ambiguity was removed, they no longer shared the AI generated video even when recipients would remain unaware.

In order to avoid deception, we then recruited a separate sample of Prolific participants to watch the video shared by participants in Study 3 ( $N = 1,924; M_{\text{Age}} = 42.34, 43.6\%$  Male). We refer to this as Study 3B. Participants were

<sup>§</sup>22 participants had an average index score of zero and were randomly assigned to one of the two positions.

497 informed that the video they were going to watch was shared  
498 with them by another Prolific participant who chose between  
499 sending one of two videos, with one in favor and one opposed  
500 to AI regulation. They were not informed that the choice  
501 was between a real and a deepfake video. After watching  
502 the full video, participants answered three questions related  
503 to trust (in the speaker, the argument presented, and the  
504 participant who shared the video), a question about the  
505 perceived strength of the argument in the video, and a  
506 question about whether they would rely on the sender for  
507 other information (all on five point Likert scales ranging from  
508 “Not at all” to “Extremely”). We then asked participants if  
509 they thought the video they watched was real or AI generated,  
510 and asked them to rate their support for mandating the  
511 labeling of AI generated content on a seven point Likert scale  
512 (ranging from “Strongly oppose” to “Strongly support”). We  
513 report the results descriptively, but made no preregistered  
514 predictions. Notably, participants were not assigned to all  
515 conditions equally because assignment depended on the videos  
516 that were shared by other participants. Because most senders  
517 (and receivers) were in favor of regulation, more real videos  
518 opposed regulation and were thus counter-attitudinal to the  
519 receivers. Because trust and agreement are correlated (27), we  
520 report results separately for whether the participant agreed  
521 or disagreed with the argument in the video. Figure 3  
522 shows that participants trusted the speaker in the video,  
523 the argument, and the sender less when they watched a  
524 deepfake than a real video. They also thought the argument  
525 was weaker and were less willing to rely on the sender for  
526 information. Notably, however, this did not differ across  
527 labeled or unlabeled conditions, regardless of whether the  
528 argument was aligned or opposed to the participant’s own  
529 view.

530 84.4% of participants correctly identified the unlabeled  
531 deepfake video as AI generated, and this increased to 91.7%  
532 with the AI Generated watermark ( $t(1010) = 3.09, p = .002$ ).  
533 Notably, the lowest accuracy was for the real video: only  
534 67% correctly identified it as a real video. Regardless of how  
535 they identified the videos, participants found the real video  
536 less deceptive than the unlabeled deepfake (1.88 vs. 2.37,  
537  $t(1633) = 8.38, p < .001$ ) and also the labeled deepfake (2.08,  
538  $t(1199) = 2.78, p = .005$ ). Notably, adding the label led to  
539 the video being perceived as less deceptive ( $t(1010) = 3.27,$   
540  $p = .001$ ). Finally, participants expressed very high support  
541 for a mandate to label AI generated content: on a 7-point  
542 Likert scale ranging from “Strongly oppose” to “Strongly  
543 support,” 69% expressed the strongest support and only 2  
544 were opposed (average 6.46).

## 545 Discussion

546 Widespread availability of tools that can create realistic, but  
547 false, video content is raising concerns about the impact of  
548 such content on public opinion. People may believe that  
549 events that never happened are real, and may generally  
550 become more skeptical of news footage even when it is  
551 real. One approach to addressing this concern is to mandate  
552 labeling via text embedded in the video informing viewers  
553 of the artificial nature of the content. We show that this  
554 labeling may be partially effective: people are less likely to  
555 watch and share content that they know is AI generated.  
556 However, we also show that this labeling has no effect on

559 the persuasiveness of a message among those who do view it.  
560 That is, people’s opinions about a policy-relevant belief (here:  
561 support or opposition to government regulation of AI) are  
562 impacted by deepfake videos, regardless of whether they are  
563 labeled as such. This suggests that labels may be effective  
564 in contexts such as social media, in which people can choose  
565 what content to view and share, but not in settings where  
566 people are passively exposed to content, such as streaming  
567 news channels.

568 Another limitation of disclosure is that it requires cooperation  
569 from the content creator. In organized disinformation  
570 campaigns, for example, an actor seeking to deceive viewers  
571 (e.g., to influence the outcome of an election) would not abide  
572 by a disclosure requirement. While regulation could mandate  
573 embedded watermarks by video generation tools, this would  
574 limit their use in legitimate content creation (e.g., for the  
575 purpose of entertainment or education) and would, moreover,  
576 be easily circumvented by relying on open source tools.

577 Generative artificial intelligence is likely to have a growing  
578 impact on belief formation and public opinion. For our  
579 research, we used ChatGPT to write persuasive arguments  
580 for and against government regulation of AI. These arguments  
581 were persuasive and impacted people’s views. Moreover, our  
582 video format gave participants plenty of opportunities to  
583 recognize it as a deepfake: because they featured a single  
584 speaker for four minutes, with today’s technology, it was not  
585 possible to create flawless deepfakes. In the future, advances  
586 in the technology could make it more difficult or impossible to  
587 detect them. Even so, many participants thought the video  
588 was real—and perhaps also concerningly, many participants  
589 thought the real video was a deepfake. Notably, this is not  
590 the only way in which AI can influence public opinion. As  
591 people increasingly chat with large language models and use  
592 them to get advice and as tutors, actors could strategically  
593 insert content that shifts the model’s views on a topic—a  
594 bias that would then be passed on to the users. Moreover,  
595 custom instructions provided to models could more directly  
596 influence how they respond to questions. As far as we know,  
597 there are no regulations guiding such custom instructions,  
598 nor a requirement that they be disclosed.

## 600 Extended Methods

601 The studies in this paper have been approved by the human  
602 subjects review board at the Hong Kong University of Science  
603 and Technology.

604 **Study 1.** Participants first answered six questions related  
605 to their views on AI regulation: “Government regulation  
606 is necessary to ensure that generative artificial intelligence  
607 is used ethically,” “I believe that without government over-  
608 sight, generative artificial intelligence could be harmful to  
609 society,” “The government should impose strict regulations  
610 on the development and deployment of generative artificial  
611 intelligence,” “Government regulation will stifle innovation in  
612 generative artificial intelligence,” “Government interference  
613 in generative artificial intelligence should be minimal to  
614 allow for technological advancement,” and “Regulation of  
615 generative artificial intelligence could hinder its potential  
616 positive impact.” All questions were asked on sliders ranging  
617 from -100 (Strongly disagree) to +100 (Strongly agree), with  
618 a neutral 0 point labeled as “Neither agree nor disagree.”  
619 To calculate our measure of support for AI regulation, we  
620

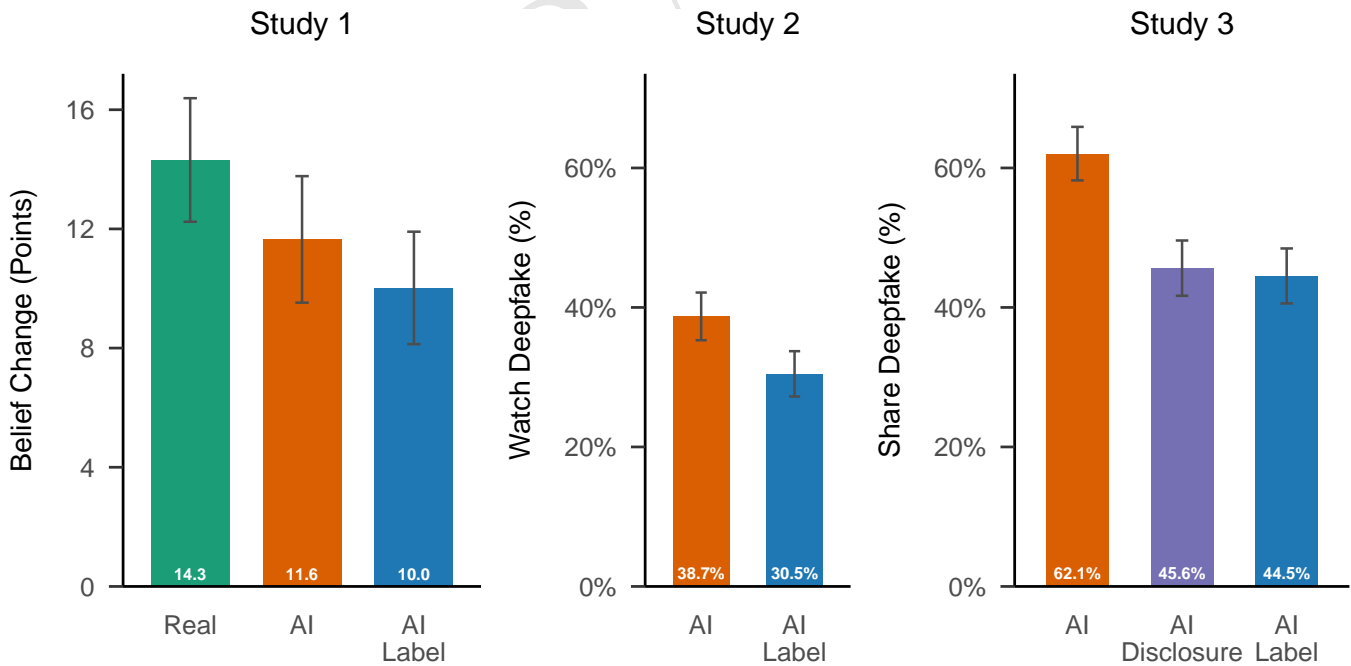
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682

683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744

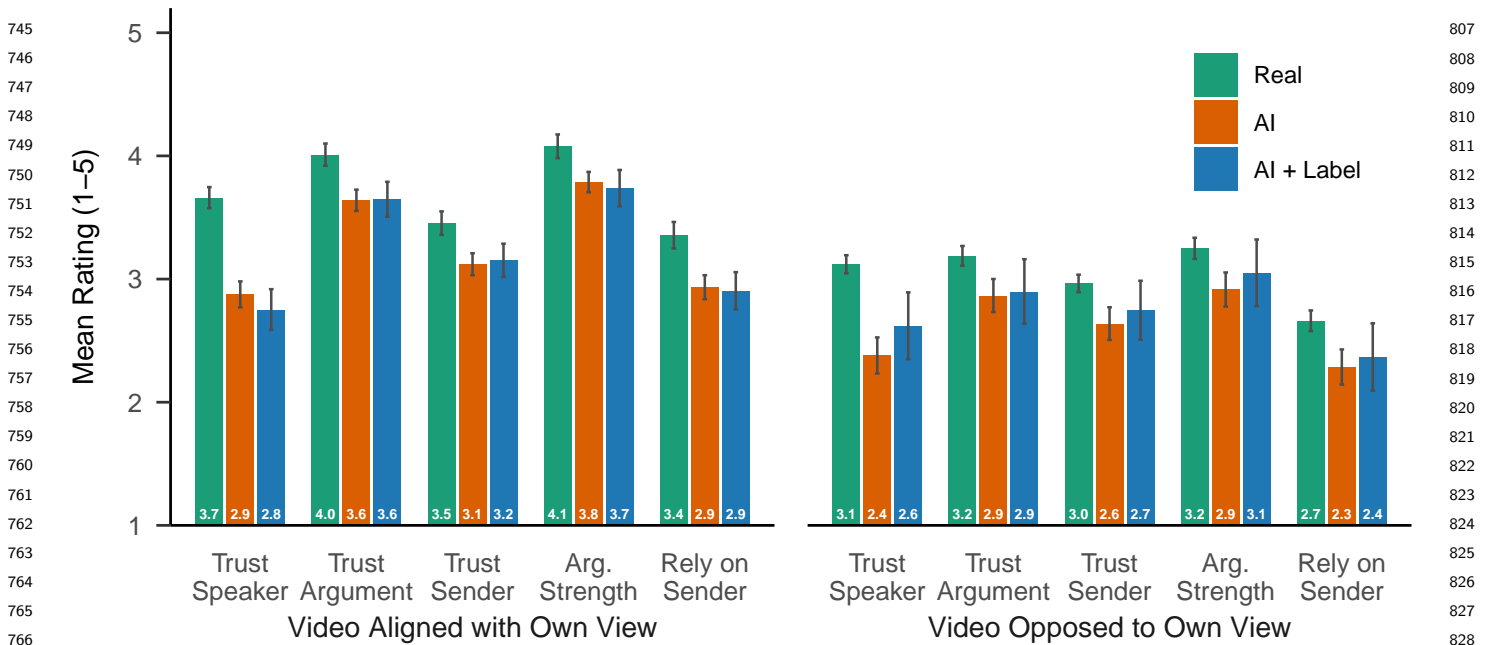
**Table 1. OLS regression results from all three studies examining the effect of labeling on persuasion, voluntary exposure, and willingness to share deepfake content. Column 1 shows the change in beliefs about the need for AI regulation in the direction of the position advocated in the video. Column 2 shows how persuasive participants thought the video was. Column 3 shows the likelihood of watching the belief-opposing deepfake video, and Column 4 shows the degree of persuasion away from the originally held position. Finally, Column 5 shows the likelihood of sharing a belief-aligned deepfake video.**

	Study 1		Study 2		Study 3
	Belief Change	Persuasiveness	Exposure	Belief Change	Share Deepfake
Deepfake (vs. Real)	-2.714+	-0.319***			
	(1.456)	(0.069)			
Label Shown	-1.542	-0.095	-0.062**	-4.178**	-0.175***
	(1.450)	(0.068)	(0.023)	(1.360)	(0.028)
Disclosure Only					-0.164***
					(0.028)
Pro-Regulation Video	-6.461***	0.324***	0.268***		
	(1.186)	(0.056)	(0.027)		
Constant	17.791***	3.040***	0.310***	0.981	0.621***
	(1.205)	(0.057)	(0.018)	(0.959)	(0.020)
N	1800	1800	1551	1551	1832

+ p < 0.1, \* p < 0.05, \*\* p < 0.01, \*\*\* p < 0.001



**Figure 2. Main results from three experiments on disclosure of AI generated content. Adding a disclosure label did not reduce the persuasiveness of the video's arguments relative to an unlabeled deepfake (Study 1). Participants were less likely to watch a belief-opposing video when they were informed that it was a deepfake (Study 2). Finally, informing participants that a belief-aligned video was a deepfake reduced their willingness to share it, even if the recipients would not be informed (Study 3). Error bars show 95% confidence intervals.**



**Figure 3.** Evaluations of speaker, arguments, and sender as rated by participants in Study 3B. Deepfakes are rated as less trustworthy across dimensions, and participants are less willing to rely on the sender for other information than after watching a real video. Notably, however, we observe no difference between the unlabeled and labeled deepfake videos.

reverse-coded the last three items and averaged participants' responses.

After the baseline attitude measure, participants were informed that they would watch a video in which a professor of AI discussed whether the use of artificial intelligence should be regulated by the government. They were informed that the video was approximately four minutes long, and that they would be able to earn a bonus of up to 30 cents for answering three multiple-choice questions about the content of the video.

We created six videos presenting arguments for and against AI regulation as discussed in the main section of the paper. Participants were randomly assigned to watch one of the videos in a 3 × 2 design. On one dimension, we varied whether the video they saw was Real, a Deepfake, or a Labeled Deepfake (which included a watermark showing "AI Generated"); and on the other dimension, the video either made arguments in favor of or against government regulation of generative artificial intelligence. In the Labeled Deepfake treatment, participants were additionally informed prior to watching the video that while the professor was real, the video itself was created by artificial intelligence.

After watching the assigned video, participants were then presented with the identical six-item scale from the beginning of the survey, with the sliders placed at the locations of their initial response. We again calculated an average score after reverse coding three of the items, and our main measure of belief change was the difference in the scale in the direction of the argument made in the video (e.g.,  $Posterior\_Score - Prior\_Score$  for those viewing an argument in favor of regulation). We then asked them "how persuasive do you think the video was in shaping your views on artificial intelligence regulation," on a five-point Likert scale ranging from "Not at all persuasive" to "Extremely persuasive."

Participants then answered three multiple-choice questions, each with three options, about the content of the video. The questions were identical for the Real, Deepfake, and Labeled Deepfake conditions, but differed for videos in favor of regulation and those opposed to regulation, as those videos contained different arguments. Participants were then asked if they thought the video was in fact generated using artificial intelligence (yes or no), and the survey concluded with basic demographic questions (age, education, gender, ethnicity, race, and political affiliation) and an open-ended box for comments for the researchers.

**Repeat Participants.** The study we report here is a replication of a design in which we asked the six-item scale only after watching the video. We collected the data presented here after running Study 2. Due to an error setting the eligibility criteria on Prolific, 351 participants who had previously completed a related study were able to enter again. Because they might have been assigned to different treatments previously, we excluded all duplicate respondents and increased our sample to arrive at the preregistered number of first-time respondents. The full data with duplicates is available on OSF and all results replicate with the full sample.

**Study 2.** Participants first reported their beliefs on AI regulation using the identical six item scale as in Study 1. Based on their responses, we classified them as Pro Regulation (an average above zero) or Anti Regulation (below zero), with those who scored exactly zero randomly allocated to one side. They were then given the choice to watch one of two videos: one explaining why artificial intelligence should be regulated, or one explaining why it should **not** be regulated. For all participants, the video opposing their initial view was the deepfake video. In the Baseline condition, participants were not aware of this when they made their choice and they saw no indication that the video was AI generated. In the

869 Label condition, participants were informed which of the two  
870 videos was a deepfake. We also disclosed how we generated  
871 the video: by taking the recording arguing for the opposite  
872 position and using specialized software to replace the words.  
873 Our main outcome measure was which of the two videos  
874 they elected to watch: the real, belief-congruent video or  
875 the deepfake that opposed their initial view. After making  
876 their choice, we informed participants that they would earn  
877 a bonus for answering three questions about the content  
878 of the videos at the end of the survey to incentivize them  
879 to pay attention. They then watched their chosen video,  
880 before again answering the six-item scale of support for AI  
881 regulation with their initial responses selected by default.  
882 We calculate the persuasiveness of the videos, conditional  
883 on exposure, by looking at the change in the scale signed in  
884 the direction of the argument in the video. In addition, we  
885 looked at the extent to which a deepfake “campaign” would  
886 have been effective, examining persuasion in the direction of  
887 the fake video regardless of the video the participants elected  
888 to watch.

889 The survey then concluded in the same way as Study 1:  
890 participants answered the identical three content questions,  
891 whether they thought the video they had watched was  
892 generated by artificial intelligence, provided demographic  
893 information, as well as optional comments for the researchers.  
894 **Study 3.** Paralleling our previous two studies, participants  
895 began by completing the six-item scale on AI regulation  
896 attitude. As in Study 2, we assigned participants to be Pro  
897 (Anti) Regulation if their average response was above (below)  
898 zero, and randomly assigned participants to one side if their  
899 score was precisely zero. We then informed them that they  
900 would watch two videos on AI regulation, one making an  
901 argument in favor and one making an argument against more  
902 regulation. For all participants, the video aligned with their  
903 initial position was the deepfake.

904 In the “Baseline” condition, participants got no further infor-  
905 mation. In the Disclosure and Label conditions, participants  
906 were informed that the two videos further differed in how they  
907 were created: one video was a real recording of a professor,  
908 while the other video was created by replacing the text in  
909 the video. In the Label condition, the video also contained  
910 watermark text identifying it as “AI Generated.” We informed  
911 them that they would be asked three questions about the  
912 content of the two videos, and participants watched the two  
913 videos in random order. They then selected one of the two  
914 videos they wanted to share with a participant we would  
915

- 916 1. B Chesney, D Citron, Deep Fakes: A Looming Challenge for Privacy, Democracy, and  
917 National Security. *California Law Rev.* **107**, 1753–1820 (2019).
- 918 2. B Paris, J Donovan, Deepfakes and Cheap Fakes, Technical report (2019).
- 919 3. K Hackenburg, et al., The levers of political persuasion with conversational artificial  
920 intelligence. *Science* **390**, eaea3884 (2025).
- 921 4. SJ Nightingale, H Farid, AI-synthesized faces are indistinguishable from real faces and  
922 more trustworthy. *Proc. Natl. Acad. Sci.* **119**, e2120481119 (2022).
- 923 5. C Vaccari, A Chadwick, Deepfakes and Disinformation: Exploring the Impact of  
924 Synthetic Political Video on Deception, Uncertainty, and Trust in News. *Soc. Media +*  
925 *Soc.* **6**, 2056305120903408 (2020).
- 926 6. D Weil, M Graham, A Fung, Targeting Transparency. *Science* **340**, 1410–1411 (2013).
- 927 7. MJ Block, A Critical Evaluation of Deepfake Regulation through the AI Act in the  
928 European Union. *J. Eur. Consumer Mark. Law* **13**, 184–192 (2024).
- 929 8. O Ben-Shahar, CE Schneider, *More Than You Wanted to Know: The Failure of*  
930 *Mandated Disclosure*. (Princeton University Press), (2014).
- 931 9. DM Cain, G Loewenstein, DA Moore, The Dirt on Coming Clean: Perverse Effects of  
932 Disclosing Conflicts of Interest. *The J. Leg. Stud.* **34**, 1–25 (2005).

931 subsequently recruit. Participants saw a screenshot from  
932 each of the two videos and selected whether to share the one  
933 in favor or the one against regulation. In the disclosure and  
934 label conditions, they were again reminded that the video  
935 aligned with their position was a deepfake, and in the label  
936 condition the screenshot also showed the “AI Generated”  
937 watermark.

938 Finally, participants answered the three content questions,  
939 indicated whether they thought each of the two videos was  
940 real or AI generated, and answered the identical demographic  
941 questions and optional open-ended comment as in the  
942 previous two studies.

943 **Study 3B.** We then recruited a new sample of participants  
944 to watch the videos that had been selected for sharing. All  
945 participants first reported their attitudes toward AI regulation

946 using the same six item scale as in the previous studies. They  
947 were informed that another Prolific participant had previously  
948 watched two videos, one arguing in favor of AI regulation  
949 and one arguing against it, and had chosen to share one  
950 of them. Depending on the condition and the decision of  
951 the original participant, the shared video was either a real  
952 recording, a deepfake without a label, or a deepfake with  
953 a label. After viewing the video, participants rated the  
954 trustworthiness of the speaker, the argument, and the person  
955 who shared the video, as well as the perceived strength of the  
956 argument and how much they would rely on the sender for  
957 other information. Participants also reported how deceptive  
958 they found the video. All measures were assessed on five-point  
959 scales. Finally, participants indicated whether they believed  
960 the video was real or AI generated and reported their support  
961 for requiring labeling AI generated content on a seven point  
962 scale ranging from “Strongly oppose” to “Strongly support”.

## 965 Data Availability Statement

966 Screenshots of all materials, response data, code to reproduce  
967 all figures and analyses, the videos with the real and deepfake  
968 arguments for and against AI regulation, and the AsPredicted  
969 preregistration reports are all available via OSF: [https://osf.io/jfk8/overview?view\\_only=03c99b09c519413b8626210f0c5d1251](https://osf.io/jfk8/overview?view_only=03c99b09c519413b8626210f0c5d1251)

## 974 References

- 975 10. S Sah, G Loewenstein, DM Cain, The burden of disclosure: Increased compliance with  
976 trusted advice. *J. Pers. Soc. Psychol.* **104**, 289–304 (2013).
- 977 11. G Loewenstein, CR Sunstein, R Golman, Disclosure: Psychology Changes Everything.  
978 *Annu. Rev. Econ.* **6**, 391–419 (2014).
- 979 12. S Lewandowsky, UKH Ecker, CM Seifert, N Schwarz, J Cook, Misinformation and Its  
980 Correction: Continued Influence and Successful Debiasing. *Psychol. Sci. Public*  
981 *Interest* **13**, 106–131 (2012).
- 982 13. F Carrella, A Simchon, M Edwards, S Lewandowsky, Warning people that they are  
983 being microtargeted fails to eliminate persuasive advantage. *Commun. Psychol.* **3**, 15  
984 (2025).
- 985 14. C Wittenberg, Z Epstein, G Péloquin-Skulski, AJ Berinsky, DG Rand, Labeling  
986 AI-generated media online. *PNAS Nexus* **4**, pgaf170 (2025).
- 987 15. M Green, TC Brock, The role of transportation in the persuasiveness of public  
988 narratives. *J. personality social psychology* **79**, 701–721 (2000).
- 989 16. E La Ferrara, A Chong, S Duryea, Soap Operas and Fertility: Evidence from Brazil.  
990 *Am. Econ. Journal: Appl. Econ.* **4**, 1–31 (2012).
- 991 17. S Carney, I Riveros, S Tully, Made with AI: Consumer Engagement with Media  
992 Containing AI Disclosures (2024).

993	18. J Allen, DJ Watts, DG Rand, Quantifying the impact of misinformation and vaccine-skeptical content on Facebook. <i>Science</i> <b>384</b> , eadk3451 (2024).	1055
994	19. L Lu, ZL Tormala, A Duhachek, How AI sources can increase openness to opposing views. <i>Sci. Reports</i> <b>15</b> , 17170 (2025).	1056
995		1057
996	20. G Loewenstein, The psychology of curiosity: A review and reinterpretation. <i>Psychol. Bull.</i> <b>116</b> , 75–98 (1994).	1058
997	21. DE Berlyne, <i>Conflict, Arousal, and Curiosity</i> , Conflict, Arousal, and Curiosity. (McGraw-Hill Book Company, New York, NY, US), pp. xii, 350 (1960).	1059
998		1060
999	22. G Pennycook, et al., Shifting attention to accuracy can reduce misinformation online. <i>Nature</i> <b>592</b> , 590–595 (2021).	1061
1000	23. S Vosoughi, D Roy, S Aral, The spread of true and false news online. <i>Science</i> <b>359</b> , 1146–1151 (2018).	1062
1001		1063
1002	24. W Hart, et al., Feeling validated versus being correct: A meta-analysis of selective exposure to information. <i>Psychol. Bull.</i> <b>135</b> , 555–588 (2009).	1064
1003	25. S Altay, AS Hacquin, H Mercier, Why do so few people share fake news? It hurts their reputation. <i>New Media &amp; Soc.</i> <b>24</b> , 1303–1324 (2022).	1065
1004		1066
1005	26. G Pennycook, A Bear, ET Collins, DG Rand, The Implied Truth Effect: Attaching Warnings to a Subset of Fake News Headlines Increases Perceived Accuracy of Headlines Without Warnings. <i>Manag. Sci.</i> <b>66</b> , 4944–4957 (2020).	1067
1006		1068
1007	27. D Hagmann, JA Minson, CH Tinsley, Personal narratives build trust across ideological divides. <i>J. Appl. Psychol.</i> <b>109</b> , 1693–1715 (2024).	1069
1008		1070
1009		1071

1009	<b>ACKNOWLEDGMENTS.</b> We are most grateful to a colleague	1071
1010	for recording the arguments in favor and against regulation and	1072
1011	allowing us to generate deepfake videos based on those recordings.	1073
1012	The arguments were generated by ChatGPT and do not necessarily	1074
1013	reflect his views. This research is partially funded by the Hong	1075
1014	Kong Research Grants Council Early Career Scheme under Award	1076
1015	Number 26506323, the HKUST Institute for Emerging Market	1077
1016	Studies with support from EY, and internal funding from Carnegie	1078
1017	Mellon University.	1079

DRAFT

1017		1079
1018		1080
1019		1081
1020		1082
1021		1083
1022		1084
1023		1085
1024		1086
1025		1087
1026		1088
1027		1089
1028		1090
1029		1091
1030		1092
1031		1093
1032		1094
1033		1095
1034		1096
1035		1097
1036		1098
1037		1099
1038		1100
1039		1101
1040		1102
1041		1103
1042		1104
1043		1105
1044		1106
1045		1107
1046		1108
1047		1109
1048		1110
1049		1111
1050		1112
1051		1113
1052		1114
1053		1115
1054		1116