
MITIGATING VIEWER IMPACT FROM DISTURBING IMAGERY USING AI FILTERS: A USER-STUDY

Ioannis Sarridis

Centre for Research and Technology Hellas
Thessaloniki, Greece
gsarridis@iti.gr

Jochen Spangenberg

Deutsche Welle
Berlin, Germany
jochen.spangenberg@dw.com

Olga Papadopoulou

Centre for Research and Technology Hellas
Thessaloniki, Greece
olgapapa@iti.gr

Symeon Papadopoulos

Centre for Research and Technology Hellas
Thessaloniki, Greece
papadop@iti.gr

ABSTRACT

Exposure to disturbing imagery can significantly impact individuals, especially professionals who encounter such content as part of their work. This paper presents a user study, involving 107 participants, predominantly journalists and human rights investigators, that explores the capability of Artificial Intelligence (AI)-based image filters to potentially mitigate the emotional impact of viewing such disturbing content. We tested five different filter styles, both traditional (Blurring and Partial Blurring) and AI-based (Drawing, Colored Drawing, and Painting), and measured their effectiveness in terms of conveying image information while reducing emotional distress. Our findings suggest that the AI-based Drawing style filter demonstrates the best performance, offering a promising solution for reducing negative feelings (-30.38%) while preserving the interpretability of the image (97.19%). Despite the requirement for many professionals to eventually inspect the original images, participants suggested potential strategies for integrating AI filters into their workflow, such as using AI filters as an initial, preparatory step before viewing the original image. Overall, this paper contributes to the development of a more ethically considerate and effective visual environment for professionals routinely engaging with potentially disturbing imagery.

Keywords disturbing content · image style transfer · journalists · human rights investigators · mental health · artificial intelligence · gruesome imagery

1 Introduction

In the era of digital communication there is an exponential increase in media content, with potentially disturbing and traumatizing images becoming increasingly prevalent [1, 2, 3]. This issue holds particular significance for professions such as journalism and human rights investigation, where interactions with distressing visual content are occupational inevitabilities [1]. Such graphic visuals frequently encapsulate scenes of violence, harm, and suffering, provoking emotions of worry, concern, or anxiety that can lead to secondary or vicarious trauma [4, 5, 6]. For instance, professionals may be required to inspect footage from conflict zones such as the war in Ukraine [7], scenes from natural disasters, or horrific accidents. Thus, it is crucial to develop and employ solutions that can effectively mitigate the viewer's impact from such disturbing imagery.

Conventional solutions to this issue have predominantly focused on the application of traditional image filters, such as blurring [8, 9]. However, these traditional filters come with significant drawbacks. If applied too heavily, blurring can render an image virtually unrecognizable, stripping away essential details and making the content impossible to interpret [10]. If not enough distortion is applied, however, disturbing elements are not sufficiently masked, thus

failing to mitigate the negative impact on the viewer. This creates a challenging trade-off between the preservation of information and the protection of the viewer.

The rapid advancements of Artificial Intelligence (AI) in recent years have enabled its integration into numerous fields with a wide range of applications [11, 12, 13, 14]. Among the areas where AI systems have exhibited notable effectiveness is the neural image style transfer [15, 16, 17, 18], i.e., the process of altering digital images to adopt the appearance or visual style of another image.

In this paper, we investigate the potential of AI style-transfer filters to mitigate the distressing impact of graphic imagery, thereby addressing the inherent limitations of conventional blurring techniques. To this end, we have adopted three distinct styles/filters, i.e., Drawing, Colored Drawing, and Painting (see Figure 1), and conducted a user study to compare their effectiveness with that of traditional filters. It is important to stress that this study focuses on images containing explicit scenes of violence, injury, etc. It does not aim to detect or address every potentially traumatizing or distressing content due to the diverse range of triggers that different individuals may have. For instance, an image featuring a sorrowful child could potentially evoke distress, yet such images cannot easily be identified as potential distress triggers.

The 107 participants of this study (details about study set-up in Section 3) are individuals from professional fields that often entail regular engagement with distressing digital content (e.g., journalists, investigators, etc.). The conducted evaluation is primarily based on two key axes - the intensity of the negative emotional responses triggered while viewing the filtered images and the degree of information retained within these images. The latter is of high importance in the relevant professional contexts where detail identification is essential.

The findings of this user study confirmed the potential of AI filters to protect the mental well-being of such professionals. In particular, it was observed that compared to the conventional Blurring filters, the Drawing filter was more effective in reducing the negative emotional impact of viewing distressing images, as evidenced by the lower mean ratings used to measure negative feelings (i.e., -34.14%). In addition, Drawing maintained a significant amount of image detail (97.19%) necessary for various professional purposes, which is not the case for the Blurring filter (6.54%). It is worth noting that the absence of color and the regional consistency of the Drawing filter were the two major advantages compared to the other filters. Furthermore, feedback from participants indicated a broad acknowledgment of the potential utility of AI filters in their professional contexts. They highlighted specific stages in their workflow where such filters could be beneficially incorporated, proposed additional enhancements that could facilitate this integration, and noted potential limitations. The main contributions of this paper are the following:

- Exploring the application of AI style transfer filters as a valuable tool for mitigating the emotional impact caused by disturbing digital content, with a focus on professions such as journalism and human rights investigation, where exposure to distressing imagery is a routine occurrence.
- A comprehensive user study comparing the effectiveness of AI-based filters against traditional blurring techniques. The results indicate the promising performance of the AI style transfer filters.
- Presenting user feedback, detailing potential workflow integration points, potential improvements, and limitations.

The remainder of this paper is organized as follows: Section 2 provides an overview of related work, Section 3 presents the methodology followed for the user study conducted, and the results of the performed analysis are detailed in Section 4. Finally, we conclude with Section 5, summarizing our findings, outlining the study's limitations, and suggesting directions for future research.

2 Related Works

Professions Associated with Exposure to Disturbing Digital Imagery. The exposure to disturbing user-generated content (UGC) has been recognized as a significant issue across multiple professions, including journalism, human rights investigations, content moderation, and criminal justice, among others. Zeng et al. [2] delve into the ethical responsibilities of news organizations towards journalists processing UGC, emphasizing the risk of secondary trauma and Post-Traumatic Stress Disorder (PTSD) symptoms. Similarly, the authors of a study conducted for Eyewitness Media Hub¹ highlight that journalists engaged in the verification and editing of traumatic UGC can suffer 'secondary trauma' and symptoms associated with PTSD [1]. Feinstein et al. [19] and Reid [3] also align with this viewpoint, suggesting that the frequency or duration of exposure to graphic imagery escalates the likelihood of vicarious trauma. These studies recommend protective measures such as staff rotation, peer support, and preemptive hiring warnings.

¹<http://eyewitnessmediahub.com/>



Figure 1: Examples of AI-based and conventional filters.

Hill et al. [20] and Baker et al. [21] further emphasize the emotional impact of reporting on traumatic events and reviewing graphic war crime imagery for journalists and human rights investigators, respectively. Both studies highlight the risk of secondary trauma and the need for strategies to mitigate this risk. Pearson et al. [22] provide insight into the harms experienced by online extremism and terrorism researchers due to their exposure to distressing content. Furthermore, psychological traumas on content moderators are highlighted in several studies [23, 24]. Finally, in-depth interviews with human content moderators exposed to child sexual abuse material (CSAM), focusing on the individual and organizational coping strategies, are presented in [25]. In particular, this study highlights the importance of social support, role validation, and work-life separation, revealing a preference for mandatory, specialized therapy.

Mitigation Strategies and Approaches. Employing image blurring to decrease the exposure of moderators to harmful data is studied in [8]. However, blurring often compromises the conveyance of crucial image information, hampering a moderator’s comprehension of the depicted content. Furthermore, the potential of grayscale and blurring filters to minimize the emotional impact on content moderation workers is explored in [10]. However, similar usability concerns, such as the obscurity of image content and eye strain are highlighted. Consequently, achieving a balance between preserving essential information and protecting viewers remains an ongoing challenge.

In addition to the professional contexts, image blurring has been employed by social media platforms such as Instagram² to protect users from potentially disturbing content (content warning screens) [26, 27]. However, several studies underscore the limitations of this method. A comprehensive analysis [28] related to Instagram’s sensitive content screens found their efficacy in deterring users from accessing negative content to be low, even among individuals presenting with mental health issues. An effort of addressing the limitations of content warning screens suggests that providing additional information along with content warnings can reduce user engagement [29]. The underlying idea is that being informed about the content of an image can deter users from viewing the original image. Given these insights, the aim of this paper is to contribute to this field by exploring the utilization of AI filters for mitigating the effects of viewing disturbing imagery. By comparing these advanced AI approaches with conventional methods, we aim to deepen our understanding of this field and help devise more effective solutions to protect the mental well-being of individuals professionally required to interact with disturbing digital content.

3 Methodology

A detailed description of the methodology, including the technical details, study format, and study distribution, is outlined in this Section.

3.1 Image Style Transfer Algorithm

Central to our methodology is the use of the Progressive Attentional Manifold Alignment (PAMA) [18] style transfer algorithm, which operates on the premise of aligning the content manifold to the style manifold. This is a sophisticated, three-staged process that involves a channel alignment module, an attention module, and a spatial interpolation module. Each module serves a distinct yet interconnected function. The channel alignment module focuses on related content and style semantics, the attention module is responsible for establishing correspondence between features, and the spatial interpolation module then adaptively aligns the manifolds. One of the key characteristics of PAMA is its capacity to alleviate the often-encountered style degradation problem, thus generating stylization outcomes that achieve state-of-the-art quality. In particular, PAMA offers regional consistency, content preservation, and high style quality. The inputs into the algorithm include the style image and the image set to be transformed.

Our study adopts three specific styles for transformation: grayscale drawing, slightly colored drawing, and painting. The grayscale drawing imparts a monochrome filter to the imagery, simplifying the visual content while preserving essential structural information. The slightly colored drawing adds a minimal amount of color, providing additional visual clues. The painting style transforms the image into a Renaissance rendition, further distancing the viewer from the graphic reality of the content. Finally, as regards the conventional filters, we employed two blurring filters. The first one applies the blur across the entirety of the image, whereas the second one selectively blurs only the portion of the image that contains the disturbing content.

3.2 Study Design

The first user-study segment was dedicated to profiling the participants. This preliminary part of the study comprised typical demographic questions and two profile-building questions. The latter aimed to indicate the participants’ frequency of exposure to potentially disturbing UGC and their level of comfort or discomfort when exposed to graphic imagery.

The second segment constitutes the core of the study. It was carefully structured to include two main phases, the first of which involved five transformed images—one for each filter under consideration (i.e., three AI-based and two traditional filters). In this phase, participants were asked to rate a select subset of emotions from the Positive and Negative Affect Schedule (PANAS) scale. These emotions were carefully chosen for their relevance to the disturbing nature of the images - Distressed, Upset, Scared, Irritable, Nervous, Jittery, and Afraid. By rating these specific feelings after viewing each transformed image, participants were able to provide an empirical measure of their affective response, thus giving us an understanding of the emotional impact each filter had. The same procedure was followed for the original images to establish the baseline emotional reactions. In addition to this emotion rating, participants were asked to engage in an interpretative exercise. They were prompted to provide a free text description of what they believed each image depicted. This exercise allowed us to determine how successfully each filter retained the necessary information. In the second phase of image filter evaluation, participants were shown four more image sets. Each set contained five variations of a single image, showcasing the effects of each filter. Instead of focusing on specific emotions as in the first phase, participants were asked to rate the overall level of disturbance triggered by each (filtered) image. This aspect of

²<https://www.instagram.com/>

Table 1: Age distribution.

Age group	Percentage
18-30	22.43%
30-45	44.86%
45-60	29.91%
>60	2.80%

Table 2: Gender distribution.

Gender	Percentage
Male	54.21%
Female	41.12%
Non-binary	3.74%
Prefer not to say	0.93%

the study was aimed at understanding the overall effectiveness of each transformation in mitigating the negative impact of the original images.

The final segment of the study was designed to utilize the collective expertise and insights of the participants. A general feedback question was posed to participants, inviting them to share their thoughts on how AI technology, and specifically the AI style transformation approaches, can contribute to protecting users from the negative impact of being exposed to graphic imagery. The intent was to gain insights that would aid in further refining our approach, bridging gaps, and possibly revealing new research directions.

The questionnaire used in this study is available as supplementary material. Note that it contains disturbing content.

3.3 Distribution of the Study

The study was distributed to a diverse array of professionals whose roles often necessitate engagement with potentially disturbing content. To this end, personalized emails were sent to carefully chosen, targeted individuals, including researchers, journalists, investigators, fact-checkers, documentalists, editors, political scientists, and producers. This initiative led to responses from more than 42 organizations, amounting to a total of 86 participants. In addition to the focused outreach to specific professionals, the study was also made accessible through an open call for participation. This initiative garnered an additional 21 responses from various professions, including forensic analysts, operations managers, sociologists, technologists, post-production supervisors, and systems engineers. This combination of targeted and open-call distribution strategies aimed to diversify the sample population, ensuring a comprehensive evaluation of the proposed approach’s efficacy across different contexts and levels of exposure to disturbing content.

4 Results

4.1 Demographics and profiling questions

Beginning with the demographics of the participants, there was a diverse group in terms of age distribution, as evidenced by Table 1. The largest proportion of respondents, 44.86%, fell into the 30-45 age bracket, reflecting a participant pool primarily composed of mid-career professionals. This was followed by the 45-60 age group, representing almost a third of the sample at 29.91%. Younger participants, aged 18-30, constituted 22.43% of the sample, while those aged over 60 were least represented at 2.80%.

Looking at gender diversity, as outlined in Table 2, the distribution was predominantly binary. Male participants accounted for over half of the total at 54.21%, while females constituted 41.12%. Non-binary individuals represented a smaller proportion at 3.74%, and a minimal percentage of 0.93% opted not to disclose their gender.

Regarding the frequency of exposure to potentially disturbing UGC, as reported in Table 3, it was found that the largest portion of participants, namely 34.58%, encountered such content multiple times a week. Those who reported daily exposure constituted 22.43%, closely trailed by respondents who encounter disturbing material several times a month (21.50%). A lesser proportion, 16.82%, came across such content several times a year, while 4.67% of the participants almost never encountered disturbing material online.

Table 3: Frequency of exposure to potentially disturbing UGC.

Frequency	Percentage
Almost never	4.67%
Several times a year	16.82%
Several times a month	21.50%
Several times a week	34.58%
Daily	22.43%

Table 4: Self-perceived reactions to exposure of potentially graphic imagery.

Response	Percentage
Graphic imagery does not affect me negatively	4.67%
I rarely react negatively	36.45%
I sometimes react negatively	39.25%
I often react negatively	14.95%
I almost always react negatively	1.87%
Other responses	2.79%

Table 4 presents the distribution of self-perceived reactions to exposure to graphic imagery. The 39.25% of the participants indicated they sometimes react negatively to such content, while a slightly smaller proportion, 36.45% reported rarely reacting negatively. Those who regularly had negative reactions comprised 14.95% of the sample. A small fraction of 4.67% indicated that graphic imagery does not affect them negatively. Only 1.87% of participants claimed they almost always react negatively to such imagery, with a few respondents, i.e., 2.79%, providing other responses.

4.2 Trade-off between conveyed information and mitigation of negative feelings

As regards emotion alleviation, the Painting style filter illustrated a promising performance with an average negative feeling mitigation of 38.03% as presented in Table 5. The strongest mitigation effect was observed on feelings of being upset and distressed (i.e., the feelings that demonstrated the highest values w.r.t. the original image), registering a significant decrease of 49.20% and 44.63%, respectively. The least affected was the feeling triggered less when viewing the original image (i.e., irritability), with a mitigation rate of 27.09%. Although this reduction spectrum suggests the potential of the Painting filter in diminishing the overall emotional distress incited by graphic images, it was not without its drawbacks. While 87 out of 107 participants (i.e., 81.31%) were able to describe the content of the image, the provided responses revealed that the inherent abstraction of the Painting filter occasionally added an extra layer of distress, while some participants compared it to a piece of disturbing art. For instance, one of the responses was: ‘*An injured person (though it is very unclear, and that’s what makes it a bit disturbing)*’. This unintended consequence indicates that while the Painting style filter has a definite potential in reducing negative emotional reactions, it may unintentionally introduce certain elements of unease.

Furthermore, Table 6 presents the results for the Colored Drawing filter, indicating an average emotional mitigation of 17.96%. The highest mitigation was observed for feelings of distress, 30.25%, while feeling of fear saw the least

Table 5: Painting Style: Feelings while watching the image. The rating scale ranges from 1 (low) to 5 (high).

Feeling	Filtered	Original	Mitigation
Distressed	1.729 ± 0.907	3.122 ± 1.178	44.63%
Upset	1.439 ± 0.826	2.833 ± 1.295	49.20%
Scared	1.458 ± 0.872	2.061 ± 1.299	29.26%
Irritable	1.421 ± 0.847	1.949 ± 1.205	27.09%
Nervous	1.402 ± 0.775	2.122 ± 1.310	33.93%
Jittery	1.262 ± 0.649	2.163 ± 1.298	41.65%
Afraid	1.364 ± 0.719	1.990 ± 1.343	31.46%
mean	1.439 ± 0.649	2.322 ± 1.065	38.03%

Table 6: Colored Drawing Style: Feelings while watching the image. The rating scale ranges from 1 (low) to 5 (high).

Feeling	Filtered	Original	Mitigation
Distressed	1.374 ± 0.694	1.970 ± 1.096	30.25%
Upset	1.346 ± 0.754	1.680 ± 1.034	19.88%
Scared	1.308 ± 0.679	1.530 ± 0.948	14.51%
Irritable	1.252 ± 0.616	1.500 ± 0.948	16.53%
Nervous	1.327 ± 0.684	1.490 ± 0.959	10.93%
Jittery	1.243 ± 0.564	1.520 ± 0.948	18.22%
Afraid	1.355 ± 0.743	1.510 ± 0.987	10.26%
mean	1.315 ± 0.570	1.603 ± 0.891	17.96%

Table 7: Drawing Style: Feelings while watching the image. The rating scale ranges from 1 (low) to 5 (high).

Feeling	Filtered	Original	Mitigation
Distressed	1.748 ± 0.912	2.804 ± 1.213	37.66%
Upset	1.626 ± 0.906	2.649 ± 1.267	38.62%
Scared	1.439 ± 0.815	1.897 ± 1.262	24.14%
Irritable	1.449 ± 0.849	1.876 ± 1.235	22.76%
Nervous	1.421 ± 0.790	1.990 ± 1.311	28.59%
Jittery	1.430 ± 0.766	2.031 ± 1.311	29.59%
Afraid	1.393 ± 0.844	1.844 ± 1.292	24.46%
mean	1.501 ± 0.766	2.156 ± 1.132	30.38%

mitigation, i.e., 10.26%. It is worth noting that the original image was less disturbing (i.e., approximately 1.6 on the 1-5 rating scale) among the images involved in this study, which justifies the relatively low emotional mitigation (i.e., 17.96%). Although a total of 88 participants (i.e., 82.24%) could comprehend the image, some participants reported difficulties in identifying specific objects or elements within the image.

Table 7 shows that the Drawing style filter particularly excelled in preserving the interpretability of the image and mitigating the negative feelings. A majority of participants (i.e., 97.19% or 104 out of 107) successfully identified several details, suggesting that this style maintained a high level of clarity. For example, one of the responses was the following: ‘A dead man lying on the floor in front of two other people, one in Flipflops (so no soldiers, but private people)’. The average reduction in negative emotions was significant, averaging 30.38%. It is worth noting that the feelings most profoundly triggered by the original images, such as being upset and distressed, experienced the highest reduction, with 38.62% and 37.66%, respectively.

As regards the Partially Blurring filter, a significant majority of 103 participants (i.e., 96.26%) could interpret the image but primarily relied on the unblurred regions. In addition, Table 8 reports emotional mitigation results, the Partial Blurring style filter had a mean mitigation score of 25.54%. Similarly to the previous filters, it was most effective on feelings of distress and being upset, with reductions of 31.96% and 30.78%, respectively. The least impacted emotion was fear, with a mitigation of only 16.10%.

Table 8: Partial Blurring Style: Feelings while watching the image. The rating scale ranges from 1 (low) to 5 (high).

Feeling	Filtered	Original	Mitigation
Distressed	2.299 ± 1.143	3.379 ± 1.178	31.96%
Upset	2.215 ± 1.182	3.200 ± 1.260	30.78%
Scared	1.692 ± 1.032	2.096 ± 1.329	19.27%
Irritable	1.626 ± 0.995	2.232 ± 1.364	27.15%
Nervous	1.822 ± 1.156	2.295 ± 1.487	20.61%
Jittery	1.757 ± 1.071	2.404 ± 1.483	26.91%
Afraid	1.766 ± 1.194	2.105 ± 1.403	16.10%
mean	1.883 ± 1.000	2.529 ± 1.169	25.54%

Table 9: Blurring Style: Feelings while watching the image. The rating scale ranges from 1 (low) to 5 (high).

Feeling	Filtered	Original	Mitigation
Distressed	1.710 ± 0.858	2.773 ± 1.342	38.33%
Upset	1.607 ± 0.939	2.814 ± 1.294	42.89%
Scared	1.486 ± 0.817	1.701 ± 1.165	12.63%
Irritable	1.355 ± 0.717	1.835 ± 1.304	26.16%
Nervous	1.551 ± 0.838	1.794 ± 1.258	13.54%
Jittery	1.477 ± 0.872	1.794 ± 1.241	17.67%
Afraid	1.439 ± 0.815	1.670 ± 1.143	13.83%
mean	1.518 ± 0.734	2.054 ± 1.111	26.09%

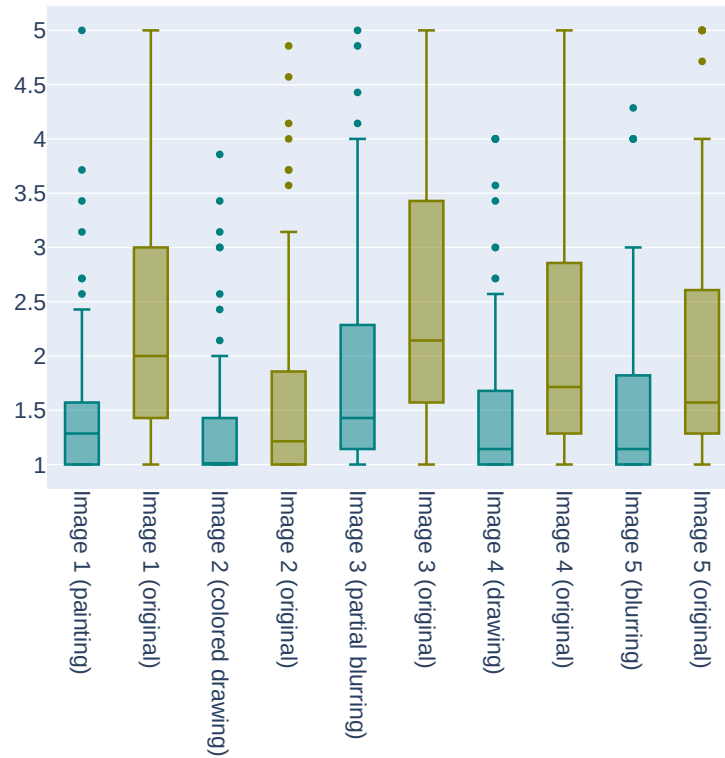


Figure 2: Mean negative feelings for all filtered and original images.

In contrast to the other styles, the Blurring filter drastically affected image interpretability. Only a small fraction of participants (i.e., 6.54%) could recognize the subject matter of the image, suggesting a high potential for information loss with this filter. Regardless of its impact on interpretability, the Blurring style filter achieved a mean mitigation of 26.09%. However, the significant information loss renders this filter less suitable for professionals where comprehending image details is crucial. Additionally, it is important to underscore that all discrepancies between the filtered and original images are statistically significant, affirmed by extremely small p-values ($< 1e-10$).

For further highlighting the discrepancies between filtered and original images, Figure 2 visualizes the mean negative feelings values reported in Tables 5-9. Overall, these findings endorse the Drawing style as the most effective filter in terms of maintaining a balance between interpretability and negative feelings mitigation. Its exceptional performance in preserving the image content while also significantly reducing negative feelings positions it as an optimal choice for professionals needing to interpret graphic images without undue emotional distress.

Table 10: Question: "How disturbing do you consider the following images?". Mean value across the 4 images. The rating scale ranges from 1 (not disturbing) to 5 (highly disturbing).

Style	Mean	Std
Drawing Style	1.977	0.733
Colored Drawing Style	2.439	0.799
Painting Style	2.692	0.804
Partially blurred	3.371	0.926
Blurred	3.002	0.873

Table 11: Question: "If the system you use in the scope of your work would provide the option to inspect images using this filter, to what extent would you use this option?". The rating scale ranges from 1 (low) to 5 (high).

Style	Mean \pm Std
Blurring	2.523 \pm 1.231
Partial Blurring	3.486 \pm 1.231
Painting	2.542 \pm 1.276
Colored Drawing	2.505 \pm 1.239
Drawing	3.000 \pm 1.259

4.3 Direct Filters Comparison

To further explore the relative efficacy of different filters in mitigating the emotional impact of disturbing images, we involved four additional images in the study. Each image was subjected to all five filtering styles, and participants were asked to rate how disturbing they found the filtered images on a scale from 1 (not disturbing) to 5 (highly disturbing). The results of this investigation provide a direct comparison between the styles and allow us to examine the potential of AI filters in comparison to conventional approaches (i.e., blurring filters).

As presented in Table 10, based on mean disturbance ratings, the Drawing style filter was found to be the least disturbing with a mean score of 1.977 with a standard deviation equal to 0.733. The Colored Drawing and Painting style filters followed with scores of 2.439 ± 0.799 and 2.692 ± 0.804 , respectively. The Partially Blurring and Blurring styles were perceived as the most disturbing, with mean scores of 3.371 ± 0.926 and 3.002 ± 0.873 , respectively. These findings showcase high discrepancies among the filtering styles, with the AI-based Drawing style filter outperforming both the rest AI-based filters and the conventional blurring techniques.

Overall, the results above underscore the advantage of AI-based filters over traditional filters. However, as mentioned in Section 4.2, they should be considered in combination with the image interpretability, where the Drawing style offered the optimal trade-off across the evaluated filters.

4.4 Practical Use and General Feedback

To assess the practical usability of each filter, we asked the participants, "If the system you use in the scope of your work would provide the option to inspect images using this filter, to what extent would you use this option?". The responses, which ranged from 1 (would not use) to 5 (would use extensively), are compiled in Table 11.

With an average rating of 3.486, the Partial Blurring filter exhibits the greatest adoption rate among all filters. This is primarily attributed to its ability to blur only the distressing regions of an image, preserving crucial details. This aspect appears particularly beneficial to professionals who require comprehensive analysis of images in their work. The full Blurring filter, on the other hand, garnered a lower mean score of 2.523, due to its tendency to conceal most of the image information. Regarding the AI-based filters, the Drawing outperformed the Colored Drawing and Painting styles with a mean rating of 3.0 compared to 2.505 and 2.542, respectively. The preference for the Drawing filter can be attributed to its capacity to preserve the visual structure of the image while simultaneously distancing the viewer from the original scene. It is also worth noting that the black-and-white nature of the Drawing filter helps to distance viewers from the reality of the content, which is not the case for the Colored Drawing filter. From a technical standpoint, a standard framework incorporating such filters would encompass two AI models. The first one would be tasked with differentiating between potentially disturbing and safe content [5], while the second model [18] would then apply the proposed filters to the content classified as potentially disturbing by the first model.

Many participants highlighted that significant limitations of such filters (both AI-based and conventional) exist – as it is often essential for professionals to view and investigate every minor detail in an original image, there were several suggestions on potential strategies to incorporate the proposed AI filters into their routine workflows. The idea of using filters as a preparatory step before viewing the original image was brought up by several participants. By first viewing a filtered version, the viewer can prepare themselves emotionally for the impact of the real image, thus potentially reducing distress. This approach may be particularly effective in contexts where exposure to the original image is ultimately unavoidable, such as investigative journalism or forensics. Furthermore, color was mentioned as a significant factor in perceiving images as disturbing. This aligns with previous psychological research suggesting that certain colors can evoke strong emotional responses [30]. Thus, adjusting the color palette of an image accordingly could be an effective way to reduce its emotional impact. Moreover, applying AI filters only to the regions of an image that depict disturbing content (as in the Partially Blurring filter) was another interesting suggestion. This targeted approach could maintain much of the image’s original context and detail, while still protecting the viewer from the most distressing elements.

In addition, the importance of variety and flexibility in filter options was emphasized by several participants. As user responses to different filters can vary widely based on individual sensitivities, having a range of filter styles to choose from could cover all individual requirements. Some participants also highlighted that AI filters could also prove particularly useful when dealing with large volumes of images. Finally, the use of AI filters for repeated viewings of an image was noted. After the initial viewing of the original, filters can be applied in subsequent viewings to prevent the repeated experience of negative emotions.

The following quotes are direct transcriptions from a subset of participants:

- *‘Ultimately, in order to do an investigation, I will always eventually have to look at the original. With a technology as the one proposed, you advance a step from completely blurring (or overlaying) an image to giving the user some idea of what the image (original) may depict.’*
- *‘Those tests were really interesting and showed (to me) how much changing the color (especially the color red) makes an impact. So it would definitely help in my job (journalist/fact-checker) to have the possibility to use such filters by default. Sometimes we WILL have to look at the original picture, of course, if we need to investigate it further, but having a default filter making these less violent would be awesome. We would then only be forced to see the ones we need to investigate further.’*
- *‘I think the most important thing in limiting distress, for me personally, is that the photo allows me to have a symbolic understanding of what is happening without providing too many distinguishing characteristics. The black-and-white line drawing method in particular seems excellent. To that extent, I would be happy to use AI filters for researching gruesome topics if they allowed me to better understand information without suffering too many negative emotional effects.’*
- *‘I think it would be great to give these different options to journalists who are facing disturbing images and to let them choose the style they want to use (depending on the document they are looking at and depending on the way they react to those types of images - they need to be able to see original images if necessary of course). I think the best one is the drawing option in black and white, but maybe other styles would work better for other people. I would rather suggest only masking the zones which are graphic such as blood, wounds, and signs of starvation, instead of applying a new style to the whole image because it often suppresses any reality. Some filters on the whole image cartoonize it, making it look more as a contemporaneous artwork than some masked reality.’*
- *‘In some cases filters improve the content as they romanticize it in a special way. While in other cases they make the situation worse as they remove information and make you imagine whatever you want.’*

5 Conclusion

In this paper, we introduce a user study that investigates the potential of AI-based filters for mitigating the emotional impact caused by disturbing imagery, aiming to support professionals who regularly encounter such content in the context of their work or related activities. The comprehensive study provided valuable insights into the effectiveness of different filter styles, with the Drawing style filter emerging as a particularly effective solution that maintains image interpretability while significantly reducing negative emotions. Although limitations certainly exist, most notably the necessity for professionals to inspect every detail in the original images, the participants proposed potential strategies for integrating these AI filters into their workflows, such as utilizing AI filters as an initial, preparatory step to viewing the full image. Future studies can refine these filter techniques, test new ones, and experiment with the proposed integration methods to further optimize the balance between necessary exposure to critical content and the mitigation

of its emotional impact. To conclude, there is a clear need for more research and activities in this domain. We hope that with our work we can contribute to reducing secondary or vicarious trauma of investigators, supporting the mental well-being of those who, because of the nature of their work and activities online, are exposed to graphic and potentially damaging imagery.

Ethics

All participants were informed why the research is being conducted, whether or not anonymity is assured, and how the data they are collecting is being stored. We confirm that all the subjects have provided appropriate informed consent via the Google Forms platform. Finally, the ethics committee of the Centre for Research and Technology Hellas has granted ethical approval for this study.

Disclosure statement

The authors report there are no competing interests to declare.

Acknowledgment

This work was supported by the EU H2020 project MediaVerse under Grant Agreement 957252.

References

- [1] Sam Dubberley, Elizabeth Griffin, and Haluk Mert Bal. Making secondary trauma a primary issue: A study of eyewitness media and vicarious trauma on the digital frontline. *Eyewitness Media Hub*, pages 1–69, 2015.
- [2] Yuan Zeng. Danger, trauma, and verification: eyewitnesses and the journalists who view their material. *Media Asia*, 45(1-2):55–59, 2018.
- [3] Alastair Reid. How are journalists at risk of vicarious trauma from ugc. *Journalism. co. uk*, 13, 2014.
- [4] Markos Zampoglou, Symeon Papadopoulos, Yiannis Kompatsiaris, and Jochen Spangenberg. A web-based service for disturbing image detection. In *International Conference on Multimedia Modeling*, pages 438–441. Springer, 2016.
- [5] Ioannis Sarridis, Christos Koutlis, Olga Papadopoulou, and Symeon Papadopoulos. Leveraging large-scale multimedia datasets to refine content moderation models. In *2022 IEEE Eighth International Conference on Multimedia Big Data (BigMM)*, pages 125–132. IEEE, 2022.
- [6] Dhruvi Shah. What’s it like when your job involves wading through others’ suffering? i was left weeping and hopeless. *The Guardian*, 2023. (Accessed: July 16th, 2023).
- [7] Jochen Spangenberg. How war videos on social media can trigger secondary trauma. *Deutsche Welle*, 2022. (Accessed: July 16th, 2023).
- [8] Anubrata Das, Brandon Dang, and Matthew Lease. Fast, accurate, and healthier: Interactive blurring helps moderators reduce exposure to harmful content. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 8, pages 33–42, 2020.
- [9] Giancarlo Fiorella. How to maintain mental hygiene as an open source researcher. *Bellingcat*, 2022. (Accessed: July 16th, 2023).
- [10] Sowmya Karunakaran and Rashmi Ramakrishan. Testing stylistic interventions to reduce emotional impact of content moderation workers. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7, pages 50–58, 2019.
- [11] Yaniv Taigman, Ming Yang, Marc’ Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1701–1708, 2014.
- [12] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A Bharath. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1):53–65, 2018.
- [13] Jesús Bobadilla, Fernando Ortega, Antonio Hernandez, and Abraham Gutiérrez. Recommender systems survey. *Knowledge-based systems*, 46:109–132, 2013.
- [14] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10781–10790, 2020.
- [15] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016.

- [16] Yongcheng Jing, Yezhou Yang, Zunlei Feng, Jingwen Ye, Yizhou Yu, and Mingli Song. Neural style transfer: A review. *IEEE transactions on visualization and computer graphics*, 26(11):3365–3385, 2019.
- [17] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11326–11336, 2022.
- [18] Xuan Luo, Zhen Han, and Linkang Yang. Progressive attentional manifold alignment for arbitrary style transfer. In *Proceedings of the Asian Conference on Computer Vision*, pages 3206–3222, 2022.
- [19] Anthony Feinstein, Blair Audet, and Elizabeth Waknine. Witnessing images of extreme violence: a psychological study of journalists in the newsroom. *JRSM open*, 5(8):2054270414533323, 2014.
- [20] Desiree Hill, Catherine A Luther, and Phyllis Slocum. Preparing future journalists for trauma on the job. *Journalism & mass communication educator*, 75(1):64–68, 2020.
- [21] Elise Baker, Eric Stover, Rohini Haar, Andrea Lampros, and Alexa Koenig. Safer viewing: A study of secondary trauma mitigation techniques in open source investigations. *Health and Human Rights*, 22(1):293, 2020.
- [22] Elizabeth Pearson, Joe Whittaker, Till Baaken, Sara Zeiger, Farangiz Atamuradova, and Maura Conway. Online extremism and terrorism researchers’ security, safety, and resilience: findings from the field. 2023.
- [23] Andrew Arsht and Daniel Etcovitch. The human cost of online content moderation. *Harvard Journal of Law and Technology*, 2018.
- [24] Miriah Steiger, Timir J Bharucha, Sukrit Venkatagiri, Martin J Riedl, and Matthew Lease. The psychological well-being of content moderators: the emotional labor of commercial moderation and avenues for improving support. In *Proceedings of the 2021 CHI conference on human factors in computing systems*, pages 1–14, 2021.
- [25] Ruth Spence, Amy Harrison, Paula Bradbury, Paul Bleakley, Elena Martellozzo, and Jeffrey DeMarco. Content moderators’ strategies for coping with the stress of moderating content online. *Journal of Online Trust and Safety*, 1(5), 2023.
- [26] Victoria Bridgland, Payton J Jones, and Benjamin W Bellet. A meta-analysis of the effects of trigger warnings, content warnings, and content notes. 2022.
- [27] Victoria ME Bridgland and Melanie KT Takarangi. Something distressing this way comes: The effects of trigger warnings on avoidance behaviors in an analogue trauma task. *Behavior Therapy*, 53(3):414–427, 2022.
- [28] Victoria ME Bridgland, Benjamin W Bellet, and Melanie KT Takarangi. Curiosity disturbed the cat: Instagram’s sensitive-content screens do not deter vulnerable users from viewing distressing content. *Clinical Psychological Science*, 11(2):290–307, 2023.
- [29] Erin T Simister, Victoria ME Bridgland, Paul Williamson, and Melanie KT Takarangi. Mind the information-gap: Instagram’s sensitive-content screens are more likely to deter people from viewing potentially distressing content when they provide information about the content. *Media Psychology*, pages 1–20, 2023.
- [30] Naz Kaya and Helen H Epps. Relationship between color and emotion: A study of college students. *College student journal*, 38(3):396–405, 2004.