# Trusted Flagger Report for the Year 2025

This report reflects the activity of Greece Fact Check during the year 2025 in the context of its capacity as a Trusted Flagger under the Digital Services Act, with an exclusive geographical focus on Greece. Greece Fact Check is a fully independent organization and belongs to the civil non-profit company "Observatory Against Disinformation." It is not connected in any way, in terms of funding or operations, with any online platform.

By Decision No. EETT 1134/18 of September 25, 2024 of the Hellenic Telecommunications and Post Commission (EETT), Greece Fact Check was granted the status of a trusted source for reporting illegal content (Trusted Flagger), in accordance with Article 22 of Regulation (EU) 2022/2065. This recognition was based on the organization's experience, independence, and proven ability to identify and document harmful content affecting users in Greece.

Following this decision, Greece Fact Check was officially recognized as a Trusted Flagger by a number of very large online platforms. During the first half of 2025, these platforms created dedicated communication channels and provided general guidelines regarding the reporting procedure.

Specifically, recognition was granted by Google on 24 February 2025, by Dailymotion on 1 March 2025, and by Meta on 19 March 2025, followed by AliExpress, X, eBay, TikTok and Snapchat by the end of April 2025. These recognitions enabled the substantial activation of the Trusted Flagger role during the year.
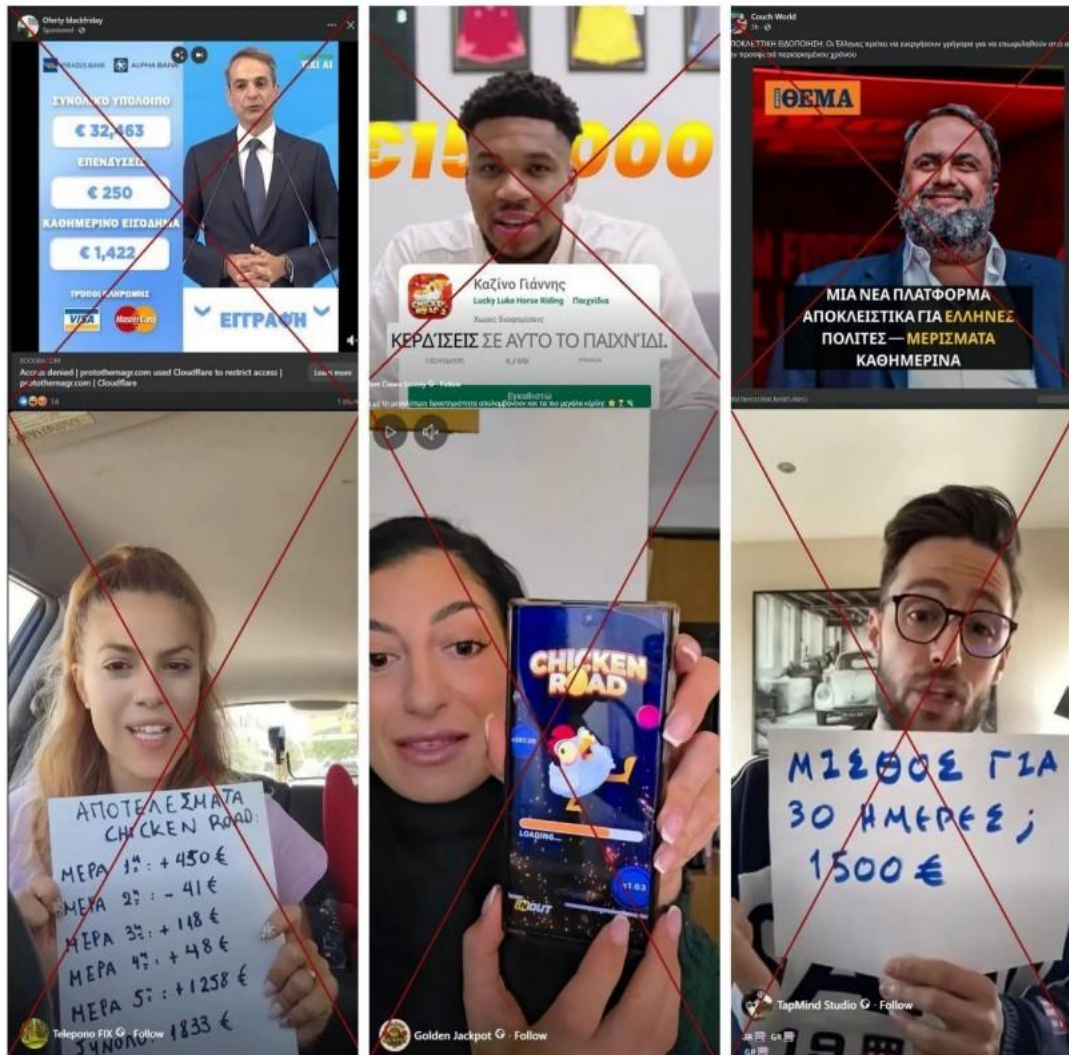
**Meta**

From March 24 to December 31, 2025, the majority of Greece Fact Check's activity as a Trusted Flagger focused on Meta's platforms, where a particularly high volume of fraudulent and harmful content was observed. During this period, a total of 778 reports were submitted to Meta. Many of these reports included more than one instance of different content, mainly from Facebook and to a lesser extent from Instagram. In total, the reports concerned 4,177 unique links, the majority of which were sponsored content, namely advertisements. This content concerned almost exclusively financial and medical scams.

In addition to individual scam incidents, Greece Fact Check submitted an additional 2,444 unique links from the advertisement libraries of the pages promoting this content. Thus, the total number of links submitted to Meta amounted to 6,621. In most cases, the pages involved displayed multiple advertisements, often on a large scale and over an extended period of time. In this way, Greece Fact Check requested that Meta examine not only individual advertisements, but also the pages themselves and their advertising activity as a whole.

This practice was based on the experience that the removal of individual fraudulent advertisements containing illegal content did not always lead to a meaningful resolution of the problem. Through the systematic reporting of advertisement libraries, it became possible to review and, where deemed necessary, remove pages managing dozens, hundreds, or even thousands of fraudulent advertisements, significantly increasing the impact of the intervention.

According to Meta's responses, 697 of the 778 submitted reports led to content removal, a rate of 89.58%, while 81 reports, that is 10.41%, did not result in removal. These figures demonstrate the effectiveness of targeted and well-documented reporting in cases of systematic abuse.
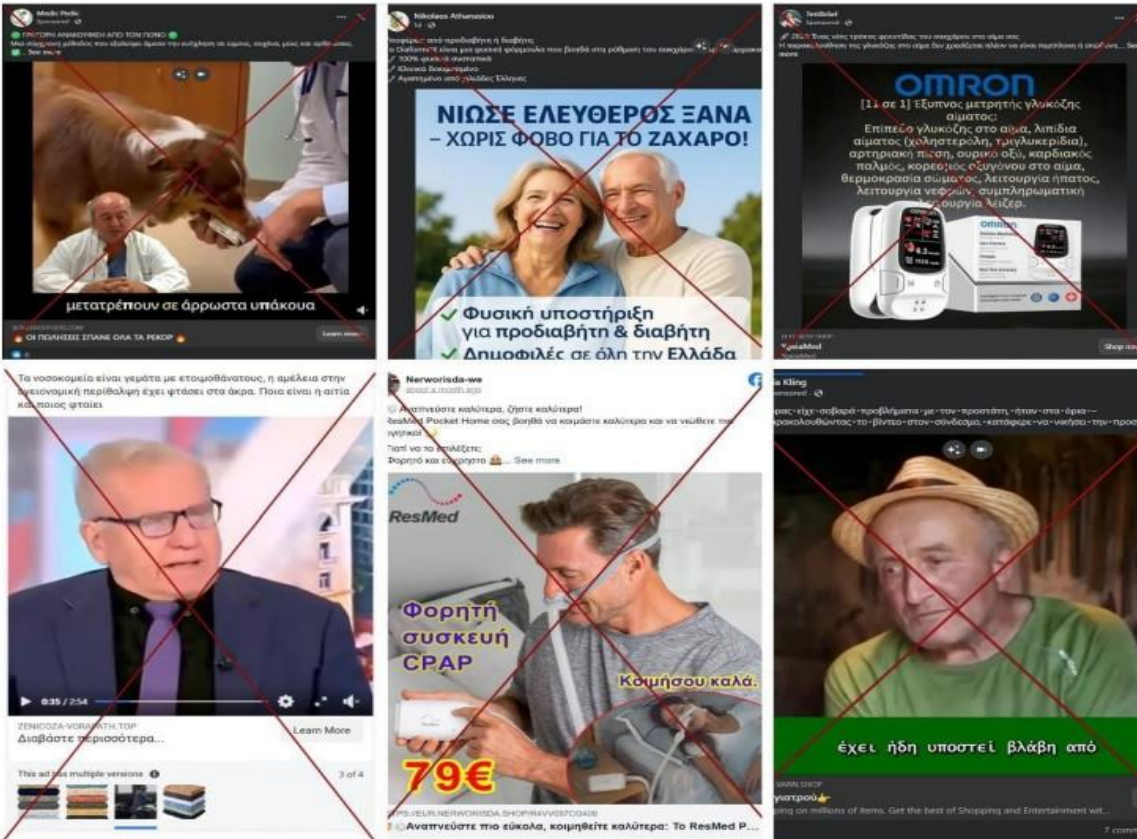
At the same time, the non-removal of content does not necessarily imply an incorrect assessment on the part of the platform. In some cases, there may be a different evaluation as to whether a specific case constitutes illegal or harmful content between the reporting entity and the competent body examining it. This differentiation may stem from a different weighing of the facts or from the evaluation of the available evidence and does not negate the role of substantiated reporting within the Trusted Flagger mechanism.
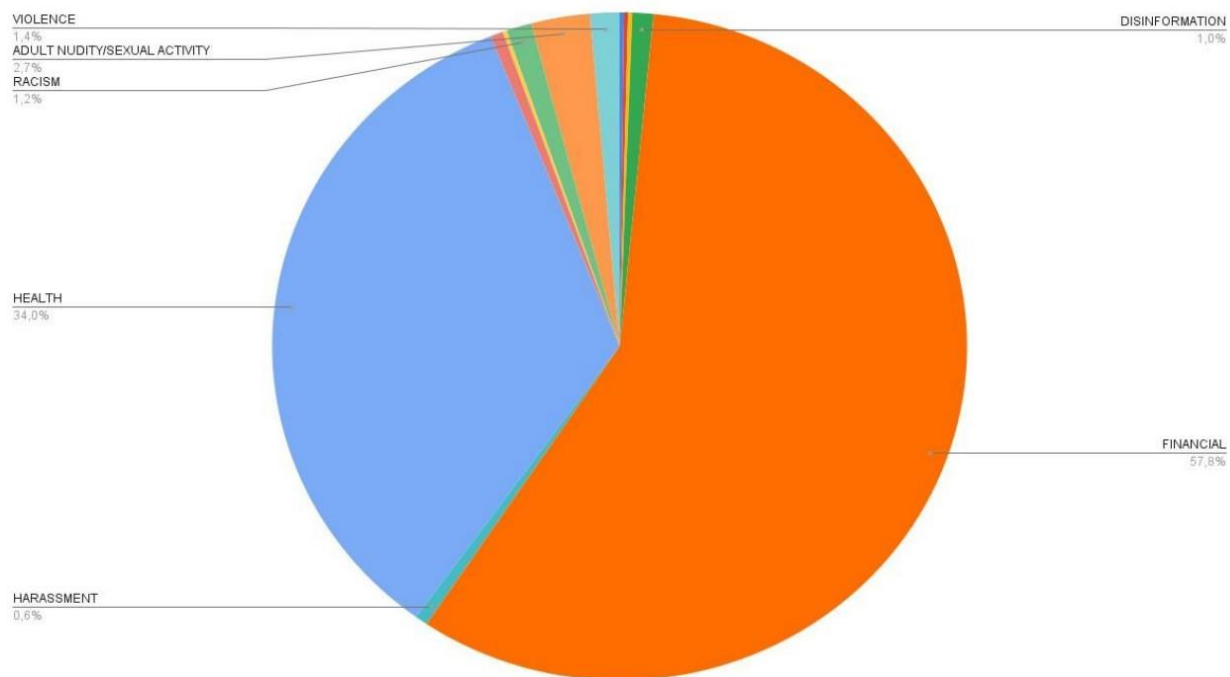
Financial scams were the most frequent category of content identified. A dominant role was played by deepfake videos and manipulated material featuring well-known individuals, such as politicians, government officials, bankers, investors and various celebrities, which were used to promote fraudulent investment schemes. There was frequent use of the names and logos of well-known mainstream media outlets, for example "Proto Thema," businesses such as the supermarket chain "Sklavenitis," banks such as National Bank, Piraeus Bank, Alpha Bank, and organizations such as OPAP, in cases of illegal betting and gambling.

At the same time, extensive use of content created with artificial intelligence was recorded, as well as recurring videos of individuals claiming to have won large sums of money through gambling applications.

Medical scams mainly concerned products unlawfully advertised as treatments for chronic conditions or as "medical" devices without scientific substantiation. Deepfakes or artificial intelligence-generated material featuring well-known Greek doctors or other famous personalities were often used in order to enhance the credibility of false claims.

To a much lesser extent, cases of pornographic content, racism, harassment and violence were also identified and reported. At the same time, incidents of disinformation were recorded and reported but did not lead to action by Meta. These incidents, including cases of state-sponsored propaganda, showed particularly high frequency and repetition during the examined period. However, the submission of such reports on our part was gradually reduced, as no meaningful response or systematic action was observed from the platform following these reports. In the majority of cases, disinformation incidents were reviewed by fact-checkers at a later stage within the framework of the 3PFC program, through systematic monitoring and analysis of content circulating on the platform.

Pie chart showing categories: DISINFORMATION 1,0%; FINANCIAL 57,8%; HARASSMENT 0,6%; HEALTH 34,0%; RACISM 1,2%; ADULT NUDITY/SEXUAL ACTIVITY 2,7%; VIOLENCE 1,4%

**Geographical origin of scams**

Although no systematic recording of the countries of origin of the scams was conducted during the reporting period, overall experience from monitoring and analyzing the incidents allows for the safe conclusion that the overwhelming majority of administrators of the pages and accounts involved are not based in countries of the European Union. This picture emerged from recurring patterns, account details and indicators accompanying advertising activity.
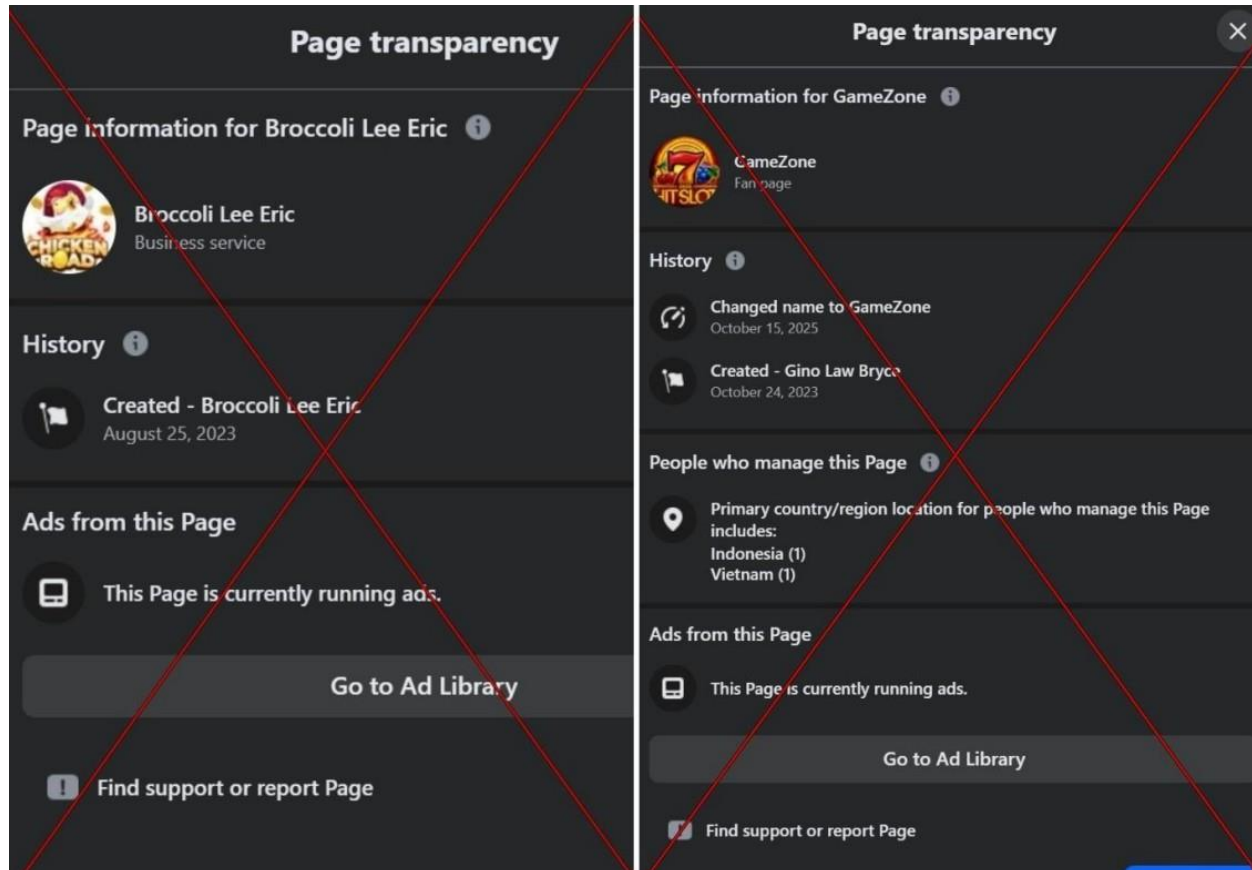
**Page hijacking**

During the reporting period, the phenomenon of Facebook page hijacking was repeatedly identified, through the acquisition of unauthorized access to page management credentials. These pages, once taken over by third parties, were renamed and repurposed in terms of content in order to be subsequently used for the dissemination of fraudulent advertisements.

This practice allows perpetrators to exploit already existing pages with a history of activity, followers and apparent credibility, largely bypassing the theoretical initial control mechanisms that are supposed to apply to newly created pages. However, the latter appears in practice not to apply.

The phenomenon of page hijacking reinforces the systemic dimension of the risk associated with fraudulent advertising and highlights the need for stricter controls in cases of sudden changes in ownership, name or activity of pages, especially when these are combined with mass advertising activity.

In several cases, however, Meta did not provide information regarding the country of origin of administrators, which limits the possibility of a more precise geographical mapping. This lack hinders the full understanding of the cross-border dimension of fraudulent networks and underscores the importance of greater transparency regarding the basic characteristics of entities that make use of the platform's advertising systems.

An example follows of a page where the origin of the administrators is not visible and an example of a page that was stolen to be used for scams.
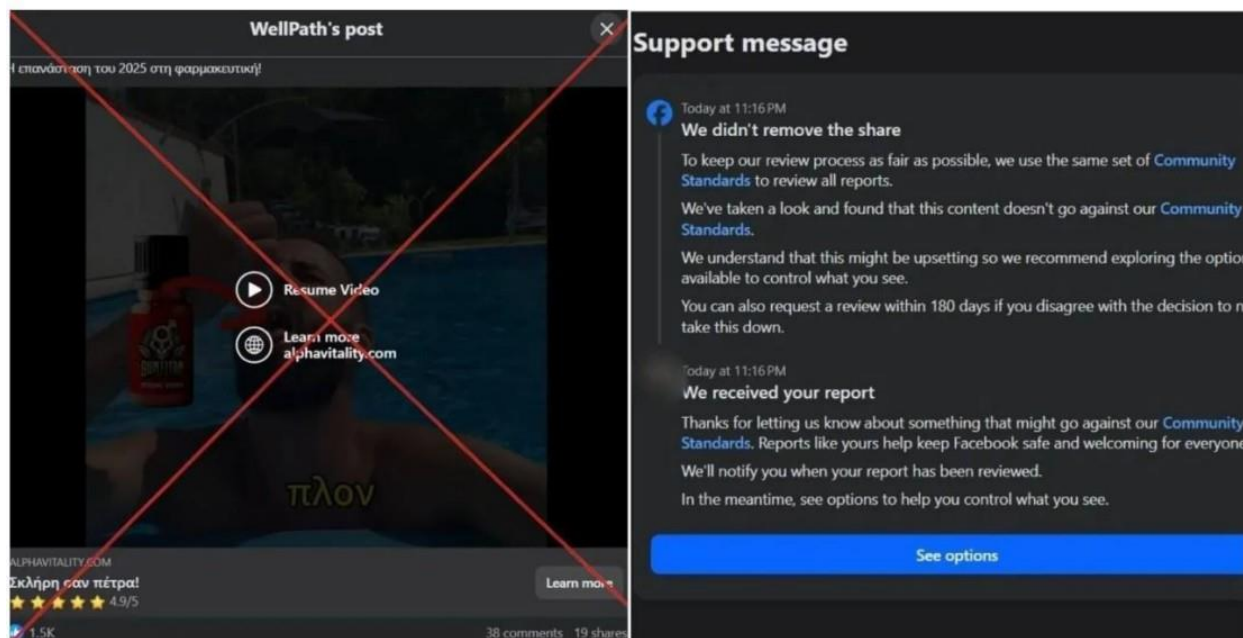


**Geographical dispersion**

Meta's advertising tools provide advertisers with the ability to tailor their content and targeting with great precision, depending on the objectives they pursue. An advertisement may appear in multiple versions and may be further customized in terms of users' age range as well as the geographical distribution of the audience to which it is addressed. This flexibility, although legitimate for lawful commercial uses, is often exploited by malicious actors.

During the reporting period, the phenomenon of malicious pages simultaneously targeting users in more than one country with the same or similar advertising content was repeatedly observed. In several cases, the same fraudulent material was identified as being promoted in parallel in multiple Member States of the European Union, including Greece. This content, that is, the pages, were reported and, in several cases, removed following review by the platform.

This practice demonstrates that fraudulent advertising campaigns are not confined within national borders but often operate with a cross-border logic, making use of geographical targeting tools to maximize reach and financial benefit. This geographical dispersion reinforces the systemic dimension of the risk and highlights the need for coordinated action both at the level of platforms and at the level of supervisory authorities.
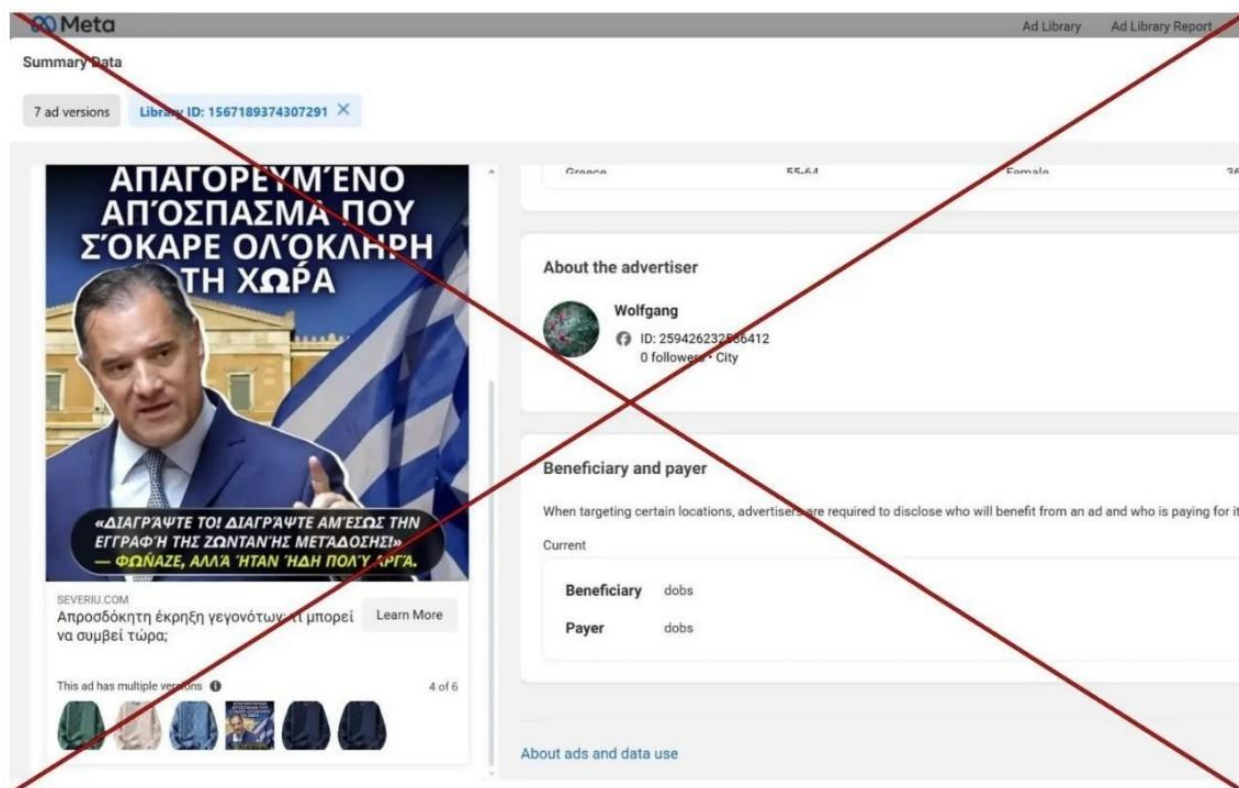
**Reports from users**

It is also important to note that, in cases where the same or similar content had previously been reported to Meta through the usual reporting mechanisms available to individual users, its removal from the platform was not observed. Below is a case of pornographic content submitted by a user, reviewed and found not to violate the Community Standards:
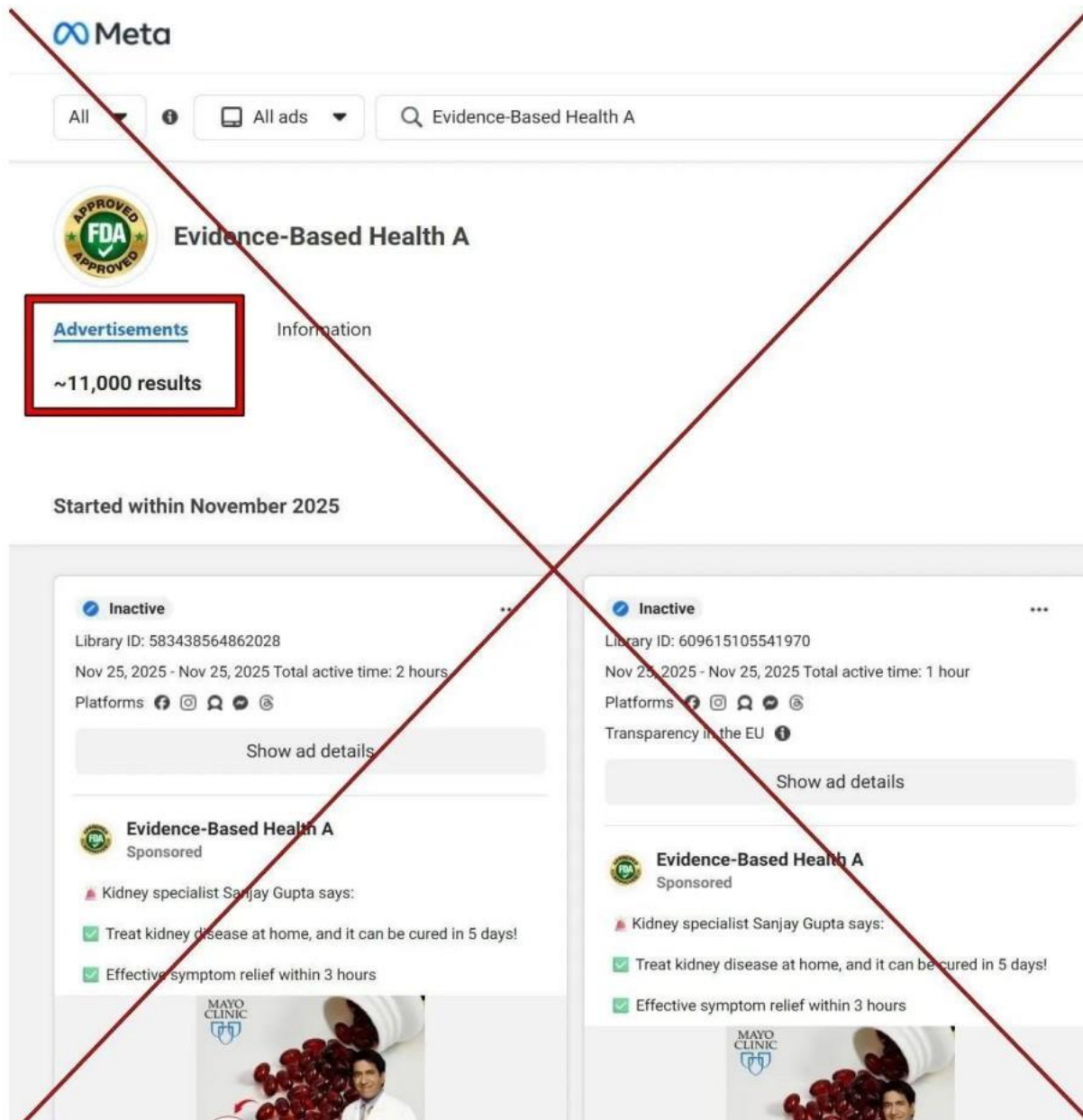


On the contrary, removal took place only following submission of a report in the capacity of Trusted Flaggers, because in addition to pornographic material, it contained advertising for an unlicensed medical product. This finding highlights on the one hand the value of the Trusted Flagger mechanism, and on the other hand the significant limitations in the effectiveness of Meta's existing content management mechanisms when these are activated exclusively through individual reports. Particularly in cases of sponsored content and organized or repeated scam practices, the experience of Greece Fact Check indicates that the lack of prioritization and systematic evaluation may lead to prolonged exposure

of users to harmful content. From the perspective of the Digital Services Act, this observation is directly linked to the need to strengthen Meta's mechanisms for assessing and mitigating systemic risks, particularly with regard to the management of advertising content and the operation of reporting channels.

With regard to financial transparency in advertisements, namely who pays and who benefits, the situation is usually unclear, because the information collected and disclosed by Meta does not provide transparency. In the example below, we see that the beneficiary and the payer do not disclose their identity and are simply listed as "dobs."



The overall picture that emerged demonstrates that fraudulent advertising constitutes a systemic risk within the framework of the Digital Services Act, as well as national legislation and supervision. In many cases, pages were created and on the same day began mass promotion of fraudulent advertisements, even dozens simultaneously. This fact highlights structural weaknesses in control mechanisms and allows for the rapid exploitation of advertising systems. Below is an example of a page that created 11 thousand advertisements within the span of one week:

In some cases, it was found that Meta had already removed individual advertisements, with justifications such as that the content did not comply with the platform's Advertising Standards or that the advertisement had been displayed by an account or page that was subsequently disabled due to violation of those standards. However, despite these individual actions, it was observed that the same pages were able to continue publishing new advertisements of a similar nature.

This finding suggests that the removal of individual advertisements or even the disabling of accounts does not necessarily entail a meaningful restriction of recurring fraudulent activity. From an operational perspective, the absence of an automatic or immediate link between previous violations and the ability

to continue advertising activity allows the recycling of the same practices and undermines the effectiveness of enforcement measures.



**The role of artificial intelligence (AI)**

Of particular importance is the growing role of artificial intelligence in the production and scaling of fraudulent content. During the period under review, systematic use of artificial intelligence tools was observed for the creation of false advertising materials, fake images and videos, as well as for the mass production of variations of the same narrative, with the aim of bypassing platform detection mechanisms. This use enables the rapid reproduction of content, its adaptation to different audiences and languages, and the enhancement of the apparent credibility of scams, thereby intensifying systemic risk for users and making effective enforcement at the level of individual reports more difficult. The use of artificial intelligence to promote illegal content may constitute the most serious threat to the effective protection of users in the digital environment. In this light, the exploitation of artificial intelligence tools for illegal purposes does not merely represent a technical and technological evolution of the phenomenon, but a factor that qualitatively and substantially alters the level of systemic risk.
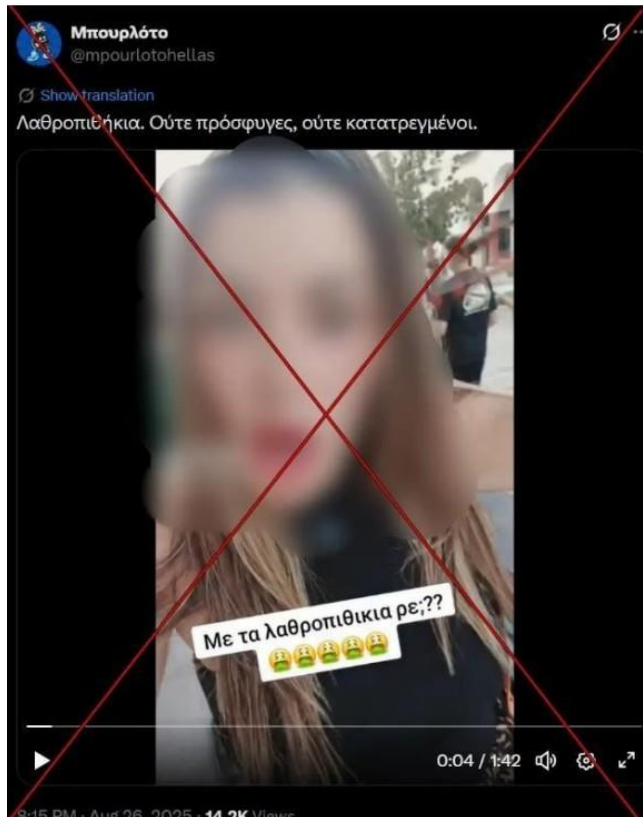
**Summary**

Based on the findings, fraudulent advertising on Meta's platforms constitutes a systemic risk within the meaning of the Digital Services Act, as it does not concern isolated incidents but repeated, organized and often cross-border practices that exploit structural features of the platform's advertising and management systems. The ability to instantly create or hijack pages, the rapid initiation of mass advertising activity, targeted geographical and demographic dissemination, limited transparency regarding the origin of administrators, and the absence of a meaningful link between previous violations and the ability to continue advertising activity create conditions that favor the recycling of fraudulent practices. At the same time, the low effectiveness of standard reporting mechanisms by individual users and the fragmented removal of individual advertisements without systematic review of the pages involved allow prolonged exposure of users to serious, mainly financial and health-related risks. In light of the Digital Services Act, the above demonstrate the need for more coherent and preventive risk mitigation measures that address the structural dimension of the phenomenon and not only its individual manifestations.

**X/Twitter**

A total of eleven reports were submitted during the period under review, which form the basis of the present analysis. These reports led to different decisions regarding the adoption of measures by the platform. In several cases, it was determined that the reported content did not meet the legal thresholds for intervention, including multiple decisions concerning illegal or harmful content, one case of disinformation within the framework of the Digital Services Act, and one case examining potential negative effects on public discourse or electoral processes.

At the same time, in certain cases X applied measures specific to the European Union, proceeding with the withholding of the reported content within the EU on grounds of illegal or harmful content, including one case where this decision was taken following a reconsideration of the report. One report concerning content related to violence also led to withholding within the EU, while in another case the content was found to violate X's internal rules and policies.

Below is an example of a post containing racist content that does not appear in the EU following our report. Our initial report had been rejected, but following an appeal X/Twitter decided that the content should not appear to users within the European Union, on the grounds of illegal or harmful speech.

Μπουρλότο
@mpourlotohellas

Show translation

Λαθροπιθήκια. Ούτε πρόσφυγες, ούτε κατατρεγμένοι.

Με τα λαθροπιθικια ρε;??
🤮🤮🤮🤮🤮

0:04 / 1:42

8:15 PM · Aug 26, 2025 · **14.2K** Views

Hello,

Thank you for your patience as we reviewed your appeal related to your report about the account @mpourlotohellas, regarding the following content:

https://x.com/mpourlotohellas/
status/1960390656343372208
report key: RAAEbNLTqVFdxsABPABZ-EBgvV5ABf____g
We have re-reviewed the reported content and any information shared. In accordance with applicable law, X is now withholding the reported content in the EU, specifically for the following legal grounds: Illegal or Harmful Speech.

For more information on our Country Withheld Content policy, please see this page: https://support.x.com/articles/20169222.

Please note that the user who published the reported content may challenge the above decision by filing an appeal to X, in an out-of-court dispute settlement proceeding or by filing a lawsuit with a competent court.

Sincerely,
X

Based on the activity recorded on platform X during the reporting period, it is observed that, although the volume of reports was clearly lower compared to Meta, indications of systemic risk related to the insufficient prevention and handling of harmful content remain. The cases reported, although few in number, do not negate the fact that the dissemination of such content can be achieved very rapidly through reposting mechanisms and algorithmic amplification. These data indicate that, even in environments with a lower volume of reports, there is a need for enhanced risk assessment and mitigation measures, in order to prevent the rapid escalation of harmful content and to ensure effective protection of users.

The type of report, response time and removal rates are shown in the table below:

| Επιγραμμική πλατφόρμα / Κατηγορίες παράνομου περιεχομένου | Αριθμός ειδοποιήσεων που υποβλήθηκαν σε ΜΕΤΑ | Αριθμός ειδοποιήσεων που έγιναν removed από ΜΕΤΑ | Μέσος Χρόνος Απόκρισης Meta | Αριθμός ειδοποιήσεων που υποβλήθηκαν σε X | Αριθμός ειδοποιήσεων που έγιναν removed από X | Μέσος Χρόνος Απόκρισης X | Σύνολο ειδοποιήσεων που υποβλήθηκαν | Σύνολο ειδοποιήσεων που έγιναν removed |
|---|---|---|---|---|---|---|---|---|
| K25 Χειραγώγηση πληροφοριών και παρέμβαση από τρίτες χώρες (Foreign information manipulation and interference) | 6 | 0 | 2 days - 3 months (≈26,75 days average) | 1 | 0 | 1 | 7 | 1 |
| K26 Παραποίηση πληροφοριών με στόχο να επηρεασθεί το αποτέλεσμα εκλογών (Information manipulation aimed at affecting sincerity/outcome of elections) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K29 Μη συναινετικά στοιχεία που περιέχουν βίντεο ή εικόνες με πλαστούς χαρακτήρες (deepfake) ή παρόμοια τεχνολογία που χρησιμοποιεί χαρακτηριστικά τρίτου μέρους (Nonconsensual items containing deepfake or similar technology using a third party's features) | 3 | 1 | 1-35 days (≈18 days average) | 0 | 0 | 0 | 4 | 1 |
| K46 Κίνδυνος για τη δημόσια υγεία (Risk for public health) | 264 | 235 | 1-5 days (≈2,4 days average) | 0 | 0 | 0 | 264 | 235 |
| K53 Ηλεκτρονικό "ψάρεμα" (Phishing) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K54 Επιχειρηματικό σχήμα πυραμίδας (Pyramid schemes) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K11 Διακρίσεις (Discrimination) | 0 | 0 | 0 | 2 | 2 (1 after appeal) | 2 | 2 | 2 |
| K12 Ρητορική μίσους (Hate speech) | 9 | 9 | 1-7 days (≈2,8 days average) | 3 | 2 | 1-2 | 12 | 11 |
| K13 Απειλές βίας (όπως απειλές θανάτου) (Threats of violence (such as death threats)) | 7 | 0 | 2-6 days (≈4 days average) | 1 | 1 | 1 day | 8 | 1 |
| K14 Παραποίηση ιστορικής καταγραφής, απολογία εγκλήματος κατά της ανθρωπότητας ή άρνηση εγκλημάτων πολέμου (Historical negationism, apology of crime against humanity or war crimes denialism) | 0 | 0 | 0 | 3 | 1 | 2 days | 3 | 1 |
| K30 Δημόσια παροχή πληροφοριών προσωπικής ταυτοποίησης για ένα άτομο (Doxing: publicly providing personally identifiable information about an individual) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K32 Παρακολούθηση (Stalking) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K33 Σεξουαλική παρενόχληση (sexual harassment) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K43 Πρόκληση ή υποκίνηση σε τέλεση αδικήματος επικίνδυνου για την δημόσια ασφάλεια (Provocation or incitement to commit an offense dangerous to public safety) | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |
| K72 Οργανωμένη Βία (Coordinated harm) | 7 | 6 | 1-3 days (≈1,75 days average) | 0 | 0 | 0 | 7 | 6 |
| Άλλες ειδοποιήσεις παράνομου περιεχομένου στο πλαίσιο του άρθρου 16 | 481 | 450 | 1-10 days (2,2 average) | 0 | 0 | 0 | 481 | 450 |

**Limitations**

The exercise of this role takes place within specific institutional and operational constraints. As an independent organization submitting content reports, we do not have enforcement authority nor access to internal data or systems of the platform, while the final assessment and decision on whether to remove content belongs exclusively to them. The exercise of this role is also affected by constraints related to the organization's available capacity, regardless of platform. The submission of reports and the systematic monitoring of content on platforms such as X and Meta are carried out within the framework of limited human, time and financial resources, which does not allow full and continuous coverage of the entirety of problematic activity. As a result, interventions are based on prioritization processes, with criteria such as the severity of the content, its potential reach and its relevance to issues of public interest.

This specific activity is not accompanied by financial compensation and is implemented through the reallocation of human, time and financial resources from other actions of the organization. As a result, the ability for systematic engagement in reporting and content monitoring procedures on platforms such as X and Meta depends directly on the overall financial sustainability of the organization. These conditions inevitably limit the volume and frequency of interventions, without negating the importance of this role for transparency, accountability and the evaluation of the implementation of the EU framework and national legislation by very large online platforms.

**Conclusions**

The experience of Greece Fact Check during 2025, in the context of its capacity as a Trusted Flagger, demonstrates that fraudulent advertising and organized harmful content do not constitute isolated incidents, but phenomena with a clear systemic dimension. The findings highlight structural weaknesses in advertising and content management systems. Practical experience shows that without targeted and coherent interventions at the level of pages and advertising infrastructures, enforcement remains fragmented and insufficient for the meaningful protection of users.

At the same time, the significant difference in effectiveness between standard user reporting mechanisms and Trusted Flagger reports underlines the added value of the institution, but also the need to strengthen obligations for the assessment and mitigation of systemic risks. This report demonstrates that the effective addressing of the phenomenon requires greater transparency, a preventive approach and the use of specialized expertise, so that Trusted Flaggers operate in a complementary and not compensatory manner to weaknesses in content moderation mechanisms.

Throughout the reporting period, Greece Fact Check applied safeguards to ensure that the exercise of the Trusted Flagger status was carried out in a proportional and responsible manner. Reports were based on national legislation and European directives incorporated into it, such as, indicatively, those concerning illegal online gambling and the circulation of unlicensed pharmaceutical products and devices. Content was not reported merely because of a generally controversial nature, but only when there was clear evidence of systematic abuse. Furthermore, there were no changes in the organizational, administrative or operational functioning of Greece Fact Check that would affect its independence in exercising its capacity as a trusted source for reporting illegal content. There were no changes in its legal or operational structure, in its composition or in its relationship with other entities, including online platforms, nor in its funding sources that could create a conflict of interest or a relationship of dependence. The organization continued to apply internal procedures ensuring that all Trusted Flagger activities are carried out with full operational independence from platforms, without intervention or guidance regarding the content, timing or outcome of reports, with decisions taken exclusively on the basis of internal evaluation and documentation.