Your response – Additional terms of service duties

Questions 1 – 5: Terms of service and policy statements

For all respondents

Question 1: What can providers of online services do to enhance the clarity and accessibility of terms of service and public policy statements?

Please submit evidence about what features make terms or policies clear and accessible.

Response:

Trustpilot welcomes the opportunity to respond to this consultation.

Trustpilot plays an important role in establishing trust between consumers and businesses. We help consumers make confident, informed decisions about where to buy. We help businesses to build trust, grow, and improve what they offer.

As an online reviews platform which engages both consumers and businesses, a key part of our work is ensuring that those using and engaging with Trustpilot understand how we work, and that includes understanding our terms of service and public policy statements.

In order to do this, we take a layered approach, utilising easily findable, readable and digestible information to assist those using our service. We would highlight these as key factors to assist in delivering clarity and accessibility.

For more information, please see our response to Question 3.

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: No

Question 2: How do you think service providers can help users to understand whether action taken by the provider against content (including taking it down or restricting access to it) or action taken to ban or suspend a user would be justified under the terms of service?

In your response to this question please consider and provide any evidence related to the level of detail provided in the terms of service themselves, whether services should provide user support materials to help users understand the terms of service and, if so, what kinds of user support materials they can or should provide.

Response:

We would strongly advocate for service providers using an educational and flexible approach when seeking to help users to understand action taken by the provider in this context.

In order to help users understand action taken, it is vital that the rationale and the basis for the decision is clearly explained to them in a manner that they will understand and learn from.

In order to achieve this, it is important to be cognisant of the diversity of users and to cater for this when pursuing an educational approach. What resonates with one user, may not deliver the same result for another.

As such, for materials to be educational, they should be provided in a range of formats so they can cater to different learning styles, levels of understanding and appetite for detail. For example, for some users, short, instructional videos are the best way to provide this information, whereas for others - written information is more likely to resonate and help them to understand a decision.

As a part of this, it can also be beneficial to repeat and summarise key messages for users within educational materials, and to use ways of highlighting key information - be it through bullet points, headers or other visual aids.

Any guidance or rules for the provision of such information by services should recognise that this is a developing area. It is therefore important that guidance or rules have flexibility built in to ensure that different users are catered for, and allow the capacity for change as innovation continues. A future-proofed and flexible approach is critical to continue to build informed users.

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: No

For providers of online services

Question 3: How do you ensure users understand the provisions in your terms of service about taking down content, restricting access to content, or suspending or banning a user from accessing the service and the actions you might take in response to violations of those terms of service?

In your response to this question, please provide information relating to (a) - (d) where relevant.

Response:

(a) how you ensure your terms of service enable users to understand both what is and is not allowed on your service, and how you will respond to user violations of these rules;

Response:

To help Trustpilot's users to understand these provisions, we provide key information about what is and isn't allowed on our service in clear and plain language. We believe that not only is guiding and educating users to help them follow our rules important, it is also a more effective approach than relying on simply correcting or penalising users after the fact.

As mentioned in response to Questions 1 and 2, our approach is based on educating users and we present this information via several layers, comprising multiple formats which are made available at different useful touchpoints (see the examples provided at the end of this response).

Trustpilot takes this approach because we recognise that we have a very wide range of users, who have different levels of knowledge and experience with online services. Added to this, such information - which is often very legalistic in its nature - can be complex to deliver. It is therefore clear that there can be no one-size-fits-all approach in this area.

We are always developing our approach as we seek to assist users in understanding this important information. [CONFIDENTIAL]

Examples

To illustrate our approach, we highlight the following examples.

Firstly, our Terms of Use for Consumers are written in simple, plain language to help users understand and digest the information. We also include short explanations to lead each section which help to sign post and guide consumers. For example, for the sections on terminating access to user accounts or deleting user accounts, we write: "We explain below when your access to our platform can be terminated or suspended." This is then followed by further details.

Secondly, our <u>Guidelines for Reviewers</u> outline how people should and shouldn't use our platform. Forming part of our terms, this page breaks down key information regarding what is and is not permitted on our service, providing it in bite-sized and easy to understand pieces of content that users can quickly skim read. The format taken is an "accordion" approach where people can click to expand certain sections if they want more detail, or simply skim the headings to see the critical points. This aims to assist with usability and accessibility for users to get the key information, and the level of detail they require, without being overloaded by content.

These are complemented by a policy entitled "Action We Take". This sets out the kinds of misuse and misbehaviour we encounter on the platform, the action we take to prevent it, and the consequences for anyone who breaches our Guidelines. Again, this deploys clear headings to signpost information to readers, and uses clear, short sections to explain important details.

In addition to this, we also provide this key information in a different format via our Help Center articles. Our Help Center hosts over 50 searchable articles which are designed to deliver critical, user information in different ways. Many of these take a "frequently asked questions" format posing questions which we provide short answers to - or as tips to guide users in the right direction. Some of the articles also include short videos as different content format to assist consumers.

(b) any relevant considerations about the risk of bad actors taking advantage of transparency

around your terms of service and how they are enforced;	
Response:	

[CONFIDENTIAL]

(c) details about any user support materials or functionalities you provide to assist users to better understand or navigate your terms of service or related products;

Response:

Please see our response to Question 3 part (a).

(d) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – please keep our response to (b) confidential.

Part of our answer to (a) is also redacted as confidential.

Question 4: Please describe the processes you have in place to measure user engagement with and comprehension of your terms of service and how you make improvements when required.

In your response to this request, please provide information relating to (a) – (f) where relevant.

Response:

(a) how you measure user engagement with/comprehension of your terms of service and the metrics you collect;

Response:

Trustpilot takes a holistic approach when measuring and understanding users' engagement with - and comprehension of - our terms of service. Our approach utilises a number of different sources to gather insights and feedback, these include:

A. User feedback

• Reviews of Trustpilot on our platform.

As a reviews platform, it is possible to review Trustpilot on our platform too. This provides a rich source of information and feedback which includes insights in relation to how users are engaging with, and comprehending, our terms of service.

We strongly believe that all businesses, ourselves included, can learn a lot from feedback. As such, we monitor the feedback provided in reviews of Trustpilot and it is used by teams across the company to make improvements, including in this area.

- [CONFIDENTIAL]
- [CONFIDENTIAL]

B. User actions

User complaints and whistleblower reports

Responding to, and analysing, complaints from users can be a helpful tool in building a picture as to how users are understanding our terms of service and how they are applied. For example, if complaints are based on misunderstandings of our terms and how our platform works, this can be helpful information to demonstrate if – and how - existing resources require refinement.

The same is true of whistleblower reports. These also provide insight as to how users understand our terms of service. [CONFIDENTIAL]

• User flagging of reviews

Users have the ability to flag reviews on Trustpilot if they believe the review breaches our rules. This is a further evidence point for whether users are understanding our rules based on whether reviews are being correctly flagged, or not.

• Engagement with content

As set out in the response to Question 3 (a), Trustpilot's terms of service are communicated in a number of ways. [CONFIDENTIAL]

Overall, these elements can create a very useful picture and source of information for how users are engaging with, and understanding, Trustpilot's terms of service. This information is drawn on by a range of teams when considering how to enhance and develop the information we provide to help users understand our terms of service and their enforcement.

(b) any behavioural research you undertake to better understand engagement with and/or comprehension of your terms of service (including any research into reasons why users do not engage with terms of service);

Response:

(c) any measures you have taken to improve engagement with and/or comprehension of your terms of service, including (but not limited to) how the findings of any behavioural research influenced these measures and/or any design changes (e.g. prompts to remind users to read the terms of the service, changes to the structure of the terms of service or changes to how users access the terms of service etc.);

Response:

[CONFIDENTIAL]

(d) costs of these processes (including the design, implementation and continued use of these processes or updated versions of these processes);

Response:

Changes to design, information, wording and other such aspects of presenting information and engaging with terms of service can be quite extensive, involving a range of different teams to draw on cross-disciplinary resources.

Making changes, particularly when it relates to the product, can require testing and careful consideration to ensure that unintended consequences don't surface elsewhere. [CONFIDENTIAL]

(e) how you evaluate the effectiveness of measures designed to improve engagement with and/or comprehension of your terms of service;

Response:

Please see response to part (a) of this question.

(f) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – part (C) is to be kept confidential. Some parts of our answer to (a) and (d) are also redacted as confidential.

Question 5: Please describe any evidence you have about the effectiveness of using different types of mechanisms to promote compliance with terms of service or change user behaviour in the event of a violation, or potential violation, of terms of service.

In your response to this request, please provide information relating to (a) – (d) where relevant.

Response:

(a) any evidence about the effectiveness of enforcement measures such as taking down content, restricting access to content, or suspending or banning user accounts in relation to encouraging users to comply with specific aspects of terms of service in the future

Response:

[CONFIDENTIAL]

(b) any evidence about how effective non-enforcement mechanisms are at reducing violations of the terms of service or repeated violations, including the type of non-enforcement mechanism and how it is implemented (e.g. prompts for users to consider the appropriateness of their content before posting it to the service (with or without links to specific provisions within the terms of service), or prompts for users to review certain provisions within the terms of service when their content is found to violate these provisions)

Response:

(c) any information and/or evidence on the costs of designing and implementing different types of enforcement or non-enforcement mechanisms (including costs of the research behind the design, implementation and continued assessment/study of these mechanisms)

Response:

(d) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – the response to (a)

Questions 6 – 8: Reporting and complaints processes

For all respondents

Question 6: What can providers of online services do to enhance the transparency, accessibility, ease of use and users' awareness of their reporting and complaints mechanisms?

In your response to this question, please provide evidence about what features make user reporting and complaints systems effective.

In your response to this question, please provide information relating to (a) – (h) where relevant.

Response:

To enhance this element of reporting or complaints mechanisms, we would advise that clear routes for these actions are provided for users to utilise when required. In Trustpilot's case, we use flag icons which are well recognised by users as symbols which enable reporting of reviews. These are placed in prominent and visible locations on each reviews, so that users can easily access and utilise them, should they need to.

Further to this, it is important that the language deployed by such mechanisms is clear and simple to aid comprehension, engagement and usability. Adopting an educational approach to help users understand can be very beneficial in interpreting potentially complex information. For example, setting out simply and clearly what information will and won't be removed can help bring rules to life.

As noted in response to earlier questions, there is also benefit to the information being relayed in different forms and ways at different touchpoints to cater for different user needs.

For the appeals process, we would advise that this is clearly explained, as well as referenced at appropriate points of the communications.

(a) reporting or complaints routes for registered users, non-registered users and potential complainants (being affected persons who are not users of the service)

Response:

(b) how to ensure that reporting and complaints mechanisms are not misused

Response:

(c) the key choices and factors involved in designing these mechanisms

Response:

(d) how users can or should be supported to report/complain about specific concerns (e.g., other users, certain types of content or, appeal content takedowns or account bans)

Response:

(e) how to ensure they are user-friendly and accessible to all users (e.g., disabled users, children)

Response:

(f) whether users are informed that their reports are anonymous (e.g., other users will not be informed about who has reported their content or account);

Response:

(g) any user support materials that explain how to use the reporting and complaints process and what will happen when users engage with these systems

Response:

(h) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: No.

For providers of online services

Question 7: Can you provide any evidence or information about the best practices for effective reporting and/or complaints mechanisms, and how these processes are designed and maintained?

In your response to this question, please provide evidence relating to (a) - (j) where relevant.

Response:

Trustpilot is a double sided platform, serving both businesses and consumers. As such, we have different, tailored reporting and complaints mechanisms for these two audiences.

Added to this, given the type of platform we are, both logged-in and non-logged-in users can read the reviews we host. We therefore provide multiple approaches to reporting content. In order to use our flagging process, a user requires an account. For those without an account, a report can be made by submitting a whistleblower complaint.

We would note that designing and maintaining reporting and complaints mechanisms requires cross-disciplinary collaboration [CONFIDENTIAL]

[CONFIDENTIAL]

As raised in previous consultation responses, our data shows a disparity in review flagging accuracy between consumers and businesses. Our 2022 transparency report noted that consumers had a 16% flagging accuracy rate, as deemed by our Content Integrity team, an improvement from 12.4% the previous year. Meanwhile, businesses had a flagging accuracy rate of 77.2%, up from 62.8% the previous year.

When considering the handling of reports from trusted flaggers – if selection of such flaggers is rigorous and they are limited to a specialized group of suitably qualified experts, then prioritising their reports could be a sensible approach.

[CONFIDENTIAL]

(a) how users report harmful content on your service(s) (including the mechanisms' location and prominence for users, and any screenshots you can provide);

Response:

(b) whether there are separate or different reporting or complaints mechanisms or processes for different types of content and/or for different types of users, including children;

Response:

(c) how users appeal against content takedowns, content restrictions or account suspensions or bans;

Response:

(d) what type of content or conduct users and non-users may make a complaint about / report, including any specific lists or categories;

Response:

(e) whether users need to create accounts to access reporting and complaints mechanisms (if there are multiple mechanisms, please provide information for each mechanism);

Response:

- (f) whether reporting and complaints mechanisms are effective, in terms of:
 - (i) enabling users to easily report content they consider to be potentially the types of content specified in the relevant terms of service, and how to determine effectiveness;

Response:

(ii) enabling, supporting or improving the accuracy of user reporting in relation to identifying the types of content specified in the relevant terms of service, and how to determine effectiveness;

Response:

(iii) enabling, supporting or improving the provider's ability to detect and take timely enforcement action against content or users as specified in the relevant terms of service, and how to determine effectiveness;

Response:

(g) whether there are any reporting or complaints mechanisms you consider to be less effective in terms of identifying certain types of content and how you determine this;

Response:

(h) the use of trusted flaggers (and if reports from trusted flaggers should be prioritised over reports or complaints from users);

Response:

(i) the cost involved in designing and maintaining reporting and/or complaints mechanisms, including any relevant issues, difficulties or considerations relating to scalability; and

Response:

(j) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – parts of paragraph 3, paragraph 4 and the final paragraph of our answer to question 7 are confidential.

Question 8: What actions do or should services take in response to reports or complaints about content that is potentially prohibited or accounts engaging in potentially prohibited activity?

In your response to this question, please include information relating to (a) – (g) where relevant.

Response:

The different actions we take to combat misuse and misbehaviour are summarised in an easy-to-read <u>Action We Take policy</u>, but we also provide more bite-sized details during the user journey, in the product, and on our Help Center.

As part of providing transparent information about our processes, we set expectations by describing what users can expect to happen in the event that they have flagged content or made a complaint, or if their review has been flagged. For example, we state what is likely to happen after a business flags reviews (we explain this via a Help Center article), and what happens where a consumer's review is flagged (see our Help Center article), or after a dispute ticket is raised (How can I dispute a decision made by the Content Integrity Team?).

[CONFIDENTIAL]

- (a) what proportion of reports are reviewed, and what proportion result in action taken including;
 - (i) any potential variation in the number and actionability (i.e., the proportion that result in a takedown or other action) of reports or complaints in relation to different provisions within your terms of service;

Response:

(ii) any differences for cases involving multiple reports/complaints about a single piece of content or user;

Response:

(iii) the costs associated with reviewing reports;

Response:

- (b) whether any reports or complaints are expedited or directed to specialist teams, including:
 - (i) the criteria for this;

Response:

(ii) the cost involved in facilitating this;

Response:

- (c) the extent to which relevant individuals (content creators, users, and non-registered or logged-out users) are informed about the progress of their report or complaint, including:
 - (i) if they are not, the reasons why;

Response:

(ii) if they are, what is included when users are informed about the progress of their report (e.g. receipt of the report, the progress of the report through the service's review process, and/or the outcome of the report);

(iii) the technical mechanisms/process to inform any relevant individuals about the progress of their report (e.g., whether non-registered users are provided an opportunity to provide an email address);

Response:

(iv) any differences in responses to different types of reports (e.g., reports about content or an account a user believes violates the terms of service, about the provider not operating in line with its terms of service, or about the accessibility, clarity or comprehensibility of those terms of service);

Response:

(v) the costs associated with responding to reports;

Response:

(d) what happens to the content while it is being assessed/processed (e.g., if and how it may still be found or viewed by other users);

Response:

(e) any internal or external timeframes or key performance indicators (KPIs) for reviewing and/or acting on reports or complaints;

Response:

(f) any user support materials that are used or should be used to support users understand the service's responses to reports, or how users can appeal moderation decisions about their content or accounts, or about decisions taken in response to reports they have submitted about other users' content or accounts;

Response:

(g) any other information.

Response:

[CONFIDENTIAL]

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – the last part of the response to question 8 should be confidential, as well as all of (g).

Questions 9 – 15: Moderation

For all respondents

Question 9: Could improvements be made to content moderation to deliver more consistent enforcement of terms of service, without unduly restricting user activity? If so, what improvements could be made?

In your response to this question, please provide information relating to (a) –(c) where relevant.
Response:
[CONFIDENTIAL]
(a) improvements in terms of user safety and user rights (e.g., freedom of expression), as well as any relevant considerations around potential costs or cost drivers;
Response:
(b) evidence of the effectiveness of existing moderation systems including any relevant examples of the accuracy, bias and or effectiveness of specific moderation processes;
Response:
(c) any other information.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response: Yes – the answer to this question is to be confidential.
For providers of online services
Question 10: Please describe circumstances where you have taken or would take enforcement action against content or users outside of what is set out publicly in your terms of service and the reasons for taking this action.
In your response to this question, please provide information relating to (a) – (e) where relevant.
Response:
Trustpilot has been in operation since 2007 and over time, we have developed relatively

comprehensive terms for our business and consumer users.

[CONFIDENTIAL] As part of our efforts to maintain trust with users, we also aim to be as transparent as possible for both business and consumer users on our service. Where we make significant changes to how we enforce our terms or the actions we take, we typically also amend our terms or Help Center information to reflect and communicate the new or amended policy.

(a) the types of action taken, and frequency of these actions (including per type of action);

(b) how relevant content or users were or would be brought to your attention;

[CONFIDENTIAL]

Response:

(c) any policies, approaches or processes you have used or would use to guide moderation decisions in these cases;

Response:

- (d) whether new policies are or would be written in response to these cases, and if so:
 - (i) whether and when these new policies are written before enforcement action is taken or after:

Response:

(ii) when and how these new policies would be added to or included in your publicly available terms of service;

Response:

(e) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes - selected paragraphs in our answer to question 10 are confidential.

Question 11: If you are made aware of content or an account that potentially violates your terms of service, please describe any relevant circumstances which might not result in enforcement action, immediately or at all.

In your response to this question, please provide describe (with examples) any relevant circumstances relating to (a) - (e).

Response:

As stated above, what type of action we take in each case will depend on the gravity of the problem or the activity involved, and our action can range from education and warnings through to formal legal action where we have exhausted other methods to resolve repeated misuse.

We take the view that user education can play a vital role in improving how people use online services, whether as a precursor to, or in tandem with other actions to enforce our rules. For example, for breaches of our terms that involve users posting inappropriate content in their reviews, where we can keep the overall rating and review but take an educational approach that helps the consumer understand how to amend or reframe their offending text to ensure it fits within the rules, that is what we do.

Over time, we have engaged in a number of projects to enhance the clarity of our flagging reasons and processes, educate users on what types of content we will or will not remove, and to help users understand how we balance different competing rights and priorities. [CONFIDENTIAL] This includes letting users know via the flagging process that we will not remove content simply because readers of reviews don't like it or agree with it, or because they personally may find it distasteful. Unfortunately, we still receive a proportion of reported reviews that fall within these latter two categories, and many of these will be assessed by our team as "invalid" reports where we do not remove content. The overall ratio of invalid to valid reports has improved over time, but we continue to innovate in this area to encourage quality user-generated content, and accurate and useful flagging.

Trustpilot is a service that allows consumers and businesses to help one another, but we are independent of both. We don't always remove businesses that repeatedly misuse or abuse our platform. This is because removing a business from our platform in this way means that consumers lose visibility about a business and the way it operates. It can be more beneficial for consumers to see this information publicly, rather than preventing it from being surfaced. In fact, removal would potentially further benefit and aid those businesses who may be trying to mislead. For cases where a paying customer has repeatedly misused our services and we have elected to terminate the business relationship, we may keep the business's profile and reviews visible on our service, but place a warning banner or label on the page and downgrade the functionality of the business accounts to the bare minimum, meaning they can only respond to and report reviews but are unlikely to be able to continue any abuse. We allow this functionality because businesses still need this ability, and because we rely on our community, including businesses, to flag reviews which breach our Guidelines.

(a) circumstances that relate to issues or challenges within your content moderation system (e.g. moderator error, language or local knowledge gaps, content is no longer available (e.g. livestream), nuance/context of content means it is found non-violative, further investigation needs to be done before action can be taken);

Response:

(b) circumstances that relate to issues or challenges within your terms of service and/or associated policies (e.g. new iterations of a harm falls outside the scope of internal moderation policies, individual piece of content is only of concern at scale (but itself does not violate policies);

Response:

(c) circumstances that relate to competing priorities (e.g., freedom of expression, public interest concerns);

Response:

(d) circumstances that would be understood by a user who has read the terms of service and why or why not, (e.g., the terms of service sets out exception for not removing violating content (e.g. news content), or transparency is not provided to avoid empowering bad actors);

Response:

(e) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes - part of the third paragraph in our response to this question is confidential.

Question 12: What automated systems do you have in place to enforce terms of service provisions about taking down or restricting access to content or suspending or banning accounts?

In your response to this question, please provide information relating to (a) – (d).

[CONFIDENTIAL] (a) the suitability/effectiveness of automated systems to identify content or accounts likely to violate different provisions within your terms of service, including the factors that materially impact suitability/effectiveness (e.g. language of content, type of content) including: (i) the suitability/effectiveness of automated systems to take down content, apply access restrictions or ban accounts in relation to any or certain provisions within your terms of service without further assistance from human moderation; Response: (ii) how you use your recommender systems to restrict access to certain content, and how you measure the effectiveness and any unintended consequences of using the recommender system in this way; Response: (iii) whether and how automated moderation systems differ by type of content (e.g., audio, video, text) or type of violation (of provisions within your terms of service) and any relevant information about costs of these different systems; Response: (iv) how data is used to develop, train, test or operate content moderation systems is sourced for different provisions within your terms of service; Response: (v) how performance/effectiveness/accuracy of automated systems are assessed and improvements then made, including any relevant considerations or differences for different provisions within the terms of service (e.g., tolerance level for false negatives and false positives between different provisions); Response: (vi) how and when automated systems are updated, and the trigger for this (e.g., in response to changing user behaviour or emerging harms); Response: (vii) what safeguards are employed to mitigate biases or adverse impacts of automated content moderation (e.g., on privacy and/or freedom of expression), and any relevant considerations or differences for different provisions within the terms of service; Response: (b) the range and quality of third-party content moderation system providers available in the

UK, particularly for different provisions within your terms of service;

(c) the process and costs associated with expanding use of existing automated moderation systems for additional provisions in your terms of service, and any relevant barriers or challenges in deploying these automated moderation systems or expanding or upgrading these systems to cover new or additional provisions;

Response:

(d) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – the response to this question is confidential.

Question 13: How do you use human moderators to enforce terms of service provisions about taking down or restricting access to content, or suspending or banning accounts?

In your response to this question, please provide information relating to (a) – (c).

Response:

[CONFIDENTIAL]

- (a) how you determine your services' resource requirements in relation to human moderation, and the factors (or key factors) that impact these requirements (e.g., increases in content or users, the range or types of content prohibited in your terms of service or technological advances in your automated system) including;
 - (i) which languages are covered by your moderation team and how you decide which languages to cover;

Response:

(ii) whether moderators are employed by the service or outsourced, or are volunteers/users and any differences regarding how different provisions within the terms of service are moderated;

Response:

(iii) whether and how moderators are vetted, and any relevant consideration for how moderators are assigned to different roles relating to different provisions within the terms of service;

Response:

(iv) the type of coverage (e.g., weekends or overnight, UK time) moderators provide and any relevant considerations for different provisions within the terms of service;

Response:

(b) the process and costs associated with extending the use of human moderation for new/additional provisions in your terms of service, and any relevant barriers or challenges to

adding new/additional provisions in your terms of service in relation to your human moderation resources;
Response:
(c) any other information.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response: Yes – the response to this question is confidential.
Question 14: What training and support is or should be provided to moderators, and what are the costs incurred by providing this training and support?
In your response to this question, please provide information relating to (a) – (g).
Response:
[CONFIDENTIAL]
(a) whether certain moderators are specialised in certain harms or subject material relating to different provisions in the terms of service;
Response:
(b) how services can/should/do assess the accuracy and consistency of human moderation teams;
Response:
(c) the impact of mental health or well-being support for moderators on the effectiveness of content moderation (including impacts on turn-over in moderation teams);
Response:
(d) whether training is provided and/or updated (including for emerging harms), and the frequency of these updates;
Response:
(e) the costs of creating training materials and support systems, and then the costs of updating or expanding these materials and systems (when relevant/required);
Response:
(f) how training, guidance and/or any relevant support systems and/or materials are provided to moderators including which moderators it is provided to (internal, contract, volunteer etc);
Response:

(g) any other information.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: Yes – the response to this question is confidential.

Question 15: How do human moderators and automated systems work together, and what is their relative scale in relation to each other regarding how you ensure your terms of service are enforced?

In your response to this question, please provide information relating to (a) – (e).

Response:

Our current approach is a dynamic partnership where human content moderators are assisted by tuned AI and machine learning tools. These are applied in different forms across different stages of the review journey, and can be adjusted by our teams to respond to emerging trends or challenges. For example, we have a number of automated systems that scan reviews at the outset - these include filters that highlight profanities in reviews and prevent users posting them, filters that scan reviews for particularly harmful or violent content and either remove or flag it for human moderator assessment, and bespoke automated technology to detect fake reviews. All of these operate before reviews are even posted to the platform and viewed by users of our service.

[CONFIDENTIAL]

We adopt a process of iteration and adjustment, learning from feedback provided by our Content Integrity team, Fraud & Investigation teams and those handling disputes about decisions made. Together with other feedback about Trustpilot's platform (including reviews about Trustpilot left on our service), this information helps improve the accuracy and robustness of our systems.

[CONFIDENTIAL]

On an ongoing basis, we continue to innovate and improve and adjust our processes and ways of working to improve efficiency and accuracy.

(a) how and when automated systems or human moderators are deployed in the moderation process;

Response:

(b) the costs of different systems or processes and of using different combinations of these systems and processes. In the absence of specific costs, please provide indication of cost drivers (e.g., moderator location) and other relevant figures (e.g., number of moderators employed, how many items the service moderates per day);

Response:

(c) how the outputs of human moderators, or appeal decisions are used to update the automated systems, and what steps are taken to mitigate bias;

(d) whether there are any relevant differences or considerations for costs or quality assurance processes for moderating different provisions within the terms of service; and
Response:
(e) any other information.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response: Yes - the second and fourth paragraphs of our answer to question 15 are confidential.

Your response – News publisher content, journalistic content and content of democratic importance

Questions 16 - 17: Identifying, defining, and categorising journalistic content, news publisher content and content of democratic importance

For all respondents

Question 16: What methods should service providers use to identify and define journalistic content and content of democratic importance, particularly at scale?

In your response to this question, please provide information relating to (a) where relevant.

Response:

(a) how journalistic content and content of democratic importance can be described in the terms of service so that users can reasonably be expected to understand what content falls into these categories.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

For providers of online services

Question 17: What, if any, methods are in place for identifying, defining or categorising content as journalistic content, content of democratic importance or news publisher content on your service?

In particular, please provide any evidence regarding the effectiveness of any existing methods.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Question 18: Moderating journalistic content, news publisher content and content of democratic importance

For providers of online services

Question 18: What considerations are taken into account when moderating journalistic content, news publisher content and content of democratic importance?

In your response to this question, please provide information relating to (a) – (e) where relevant.

Response:

(a) once identified, how journalistic content, news publisher content and content of democratic importance is actioned and what kind of action is taken; and how that differs from the moderation of other types of content

Response:

(b) the factors that are or should be considered when taking action (e.g.: downranking/removal/suspension/ban or other) regarding this content

Response:

(c) the proportion of all journalistic content, content of democratic importance and news publisher content actioned upon by you that is actioned based on algorithmic decision making

Response:

(d) the proportion of all journalistic content, content of democratic importance and news publisher content actioned upon by you that is reviewed by human moderators and on what basis content is escalated to be reviewed by human moderators

Response:

(e) any insights into the costs of moderating journalistic content and content of democratic importance, including set up and ongoing costs in terms of employee time and other material costs.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

Questions 19 – 21: Complaints and appeal processes for journalistic content, news publisher content and content of democratic importance

For all respondents

Question 19: What complaint, counter-notice or other appeal processes should be in place for users to contest any action taken by service providers regarding journalistic content and content of democratic importance?

In your response to this question, please provide information relating to (a) and (b) where relevant.

Response:
(a) examples of effective redress mechanisms that you consider would be most suited to these
content types
Response:
(b) briefings, investigations, transparency reports, media investigations and research papers that
provide more evidence
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 20: What initiatives could service providers use to create and increase awareness
about the process for users to complain and/or appeal content decisions and to minimise its'
misuse?
In your response to this question, please provide information relating to (a) and (b) where
relevant.
Response:
(a) any known impacts of over-removal or erroneous removal of news publisher content,
journalistic content or content of democratic importance
Response:
(b) briefings, investigations, transparency reports, media investigations and research papers
regarding misuse of such speech protective provisions
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
For providers of online services
Question 21: What are the current complaints, counter-notice or other appeal processes for
users to centest any action taken by you regarding journalistic centent, news publisher centent

Question 21: What are the current complaints, counter-notice or other appeal processes for users to contest any action taken by you regarding journalistic content, news publisher content and content of democratic importance on your service?

In your response to this question, please provide information relating to (a) and (b) where relevant.

Response:

(a) any initiatives taken to create and increase awareness about the process for users to complain and/or appeal content removals

(b) any measures currently in place to prevent individual or systematic misuse of any protections for news publisher content, journalistic content or content of democratic importance.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Questions 22 – 24: Other information for journalistic content, news publisher content and content of democratic importance
For providers of online services
Question 22: Do you carry out any internal impact assessments to understand the freedom of expression and privacy implications of existing policies regarding journalistic content, news publisher content and content of democratic importance?
In your response to this question, please provide information relating to (a) and (b) where relevant.
Response:
(a) explain which elements of your service design or operation they relate to and which factors they take into account
Response:
(b) provide relevant briefings, investigations, transparency reports, media investigations and research papers.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 23: What, if any, measures are in place to ensure that protection of content of democratic importance applies in the same way to a wide diversity of political opinion?
In your response to this question, please provide information relating to (a) where relevant.
Response:
(a) whether there are any additional measures/safeguards that are put in place during local or national elections.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:

For all respondents

Question 24: What, if any, measures can online service providers put in place to ensure that protection of content of democratic importance applies in the same way to a wide diversity of political opinion?

In your response to this question, please provide information relating to (a) where relevant.

Response:

(a) whether there are any additional measures/ safeguards that can be put in place during local or national elections

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

Your response - User empowerment duties

Question 25: Detecting and moderating relevant content

For providers of online services

Question 25: What processes do you use to detect relevant content and how do you moderate it?

In your response to this request, please provide information relating to (a) – (g) where relevant.

Response:

(a) what systems you use for detection

Response:

(b) further to the above, if there are any important features that you take into account to make distinctions between content, e.g. features that might identify a piece of content as promotional suicide material versus content intended to support users at risk of suicide

Response:

(c) where distinctions are made, the extent to which content is actioned automatically, by human moderation, through user reports, other methods or a combination of methods

Response:

(d) any insight into the cost of these processes, including set-up and on-going costs, in terms of employee time and any other material costs

Response:

(e) whether relevant content is allowed or prohibited on your service

(f) whether you measure the incidence of users encountering such content, and if yes, whether these systems are different to those measuring other types of content, including illegal content Response:

(g) if you offer users separate complaints procedures for moderated legal content versus illegal content, how often users report content through these channels, and what proportion of content is removed following a complaint

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

Question 26: Impact of relevant content

For all respondents

Question 26: Can you provide any evidence on whether the impact of relevant content differs between adults and children on user-to-user services?

We are interested in particular in briefings, investigations, transparency reports, media investigations and research papers that provide more evidence.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

Question 27 and 28: Experience of specific types of users

For all respondents

Question 27: Can you provide evidence around the types of adult users more likely to encounter relevant content, and the types of adult users more likely to be affected by such content?

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

For all respondents

Question 28: How do you consider the experience of users who have a protected characteristic, or those considered to be vulnerable or likely to be particularly affected by certain types of content?

In your response to this request, please provide information relating to (a) – (c) where relevant.

Response:	
having a pro	ervice collects any information about users that could be used to identify the otected characteristic, vulnerable or likely to be particularly affected by certaind, if so, what information you collect
Response:	
(c) if you co	nduct any research into the experience of the above users on your service
Response:	
Is this respo	onse confidential? (if yes, please specify which part(s) are confidential)
Response:	
Question content	ns 29 and 30: Features employed to enable greater control
For all respo	ondents
	9: What features exist to enable adult users to have greater control over the t by encounter?
content the	ey encounter?
content the	
content the In your resp Response:	ey encounter? conse to this request, please provide information relating to (a) – (d) where relating to (a) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (d) – (d)
In your resp Response: (a) features	ey encounter? conse to this request, please provide information relating to (a) – (d) where relating to (a) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (d) – (d)
content the In your resp Response: (a) features wish to see Response:	ey encounter? conse to this request, please provide information relating to (a) – (d) where relating to (a) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (a) – (d) where relating to (b) – (d) where relating to (d) – (d)
content the In your resp Response: (a) features wish to see Response:	ey encounter? conse to this request, please provide information relating to (a) – (d) where relating to (a) – (d)
content the In your resp Response: (a) features wish to see Response: (b) features Response: (c) features	ey encounter? conse to this request, please provide information relating to (a) – (d) where relating to (a) – (d)
content the In your resp Response: (a) features wish to see Response: (b) features Response: (c) features	onse to this request, please provide information relating to (a) – (d) where relating to users to reduce the likelihood of them encountering content they offered to users to alert them to the presence of certain categories of content of the users to enable them to control their interactions with different types.
content the In your resp Response: (a) features wish to see Response: (b) features Response: (c) features users (e.g., Response: (d) whether	onse to this request, please provide information relating to (a) – (d) where relating to users to reduce the likelihood of them encountering content they offered to users to alert them to the presence of certain categories of content of the users to enable them to control their interactions with different types.
content the In your resp Response: (a) features wish to see Response: (b) features Response: (c) features users (e.g., Response: (d) whether	onse to this request, please provide information relating to (a) – (d) where relations to this request, please provide information relating to (a) – (d) where relations to differed to users to reduce the likelihood of them encountering content they confered to users to alert them to the presence of certain categories of content of the control their interactions with different typon-verified)

For providers of online services

Response:

Question 30: How do you design features to enable adult users to have greater control over the

your system in deploying them? (For the purposes of our evidence base we are interested in features that enable control over a range of content, not solely relevant content).
In your response to this request, please provide information relating to (a) – (d xi) where relevant.
Response:
(a) how you measure and what evidence you can provide around the effectiveness of these features in terms of achieving their respective aims to prevent adults from encountering content that they do not want to see
Response:
(b) how you measure user engagement with these features, and any evidence you can provide around this
Response:
(c) how you ensure that these features are suitable for all adult users and that they're easy to access, including considerations for users with protected characteristics and/or vulnerable users
Response:
(d) how you decide when to offer users these features, or how to present the use of these features to users. This includes but is not limited to the following aspects, i) – xi).
Response:
i) how you develop the user need for these features, and the factors considered when determining to develop them
Response:
ii) whether these features are on by default, and in what circumstances
Response:
iii) whether these features are personalised for specific types of users
Response:
iv) when to offer users these features
Response:
v) whether, when or how often to remind users of these features - this can mean reminding users to make an initial choice, or checking if a user wants to update the initial choice later on (and if so, how frequently)
Response:
vi) where users learn about these features

vii) how to provide information about these features, including the level of detail and the words used to describe complex or technical concepts

Response:

viii) whether users have choice of controls over specific types of content

Response:

ix) how you decide whether to iterate, replace or keep such features

Response:

x) any other factors not already covered above that you take into account when considering such features

Response:

xi) any insight into the cost of these features, including set-up and on-going costs (in terms of employee time and any other material costs) as well as any intended and unintended impacts on the service more broadly (e.g., the technical feasibility of implementing filter tools, or reducing functionality based on verification status).

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

Your response - User identity verification duties

Question 31 and 32: Circumstances where user identity verification is offered and how

For all respondents

Question 31: What kind of user-to-user services currently deploy identity verification and in what circumstances?

In your response to this request, please provide information relating to (a) – (c) where relevant.

Response:

In 2021, Trustpilot introduced voluntary identity verification measures.

We are continually listening to our users and seeking to understand their views about the challenges of the online world. During this work, we found that there was a search amongst our users for greater context and a willingness to help one another – a sense of altruism within our community. These insights led us to introduce the optional ability for users to verify their identity, enabling the user to prove that they are not only a real person, but they are who they claim to be with their account.

The process requires users to safely and securely share a copy of their government-issued photo ID, as well as to take a selfie. This is the same technology deployed by banks, healthcare providers and educational institutions.

Where a user chooses to do this and their identity is verified, this is then recognised on their profile with their reviews being marked as coming from a 'verified reviewer'. This is an additional 'trust signal' for other users to see.

That said, Trustpilot has introduced this as a voluntary measure and, as a reviews platform, we believe it is critical that this remains a voluntary initiative for our users. This is on the basis that our users may not be willing or even able to verify their identity using digital tools, or they may not have access to certain types of government-issued identity documents. We do not believe that such users should be excluded from our service.

(a) the ways in which these identity verification methods are beneficial, both to the user and to the service

Response:

(b) what documentation you understand to be necessary for different types, or levels, of identity verification on user-to-user services

Response:

(c) whether you believe there are there any other circumstances where identity verification should be offered on user-to-user services.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: No.

For providers of user-to-user services that provide some types of identity verification for individual adult users

Question 32: In respect of the identity verification method(s) used on your service, please share any information explaining:

(a) in what circumstances identity verification is offered on your service and why, and to which category/categories of users

Response:

(b) what evidence and steps are taken to verify the identity of a user, e.g., which attributes are checked, what aspects of verified users are known only to the provider and what aspects are made available for other users to see, including whether processes regarding adult users are different to those regarding children

Response: See response to Question 31

(c) whether the process is, or can be, tailored to users in different geographical areas, such as the UK

Response:

(d) whether you engage third party providers to provide all or part of this identity verification process and, if so, which providers

Response: We use Veriff as our partner to provide this service.

e) once a user has their identity verified, what this allows them to do on your service, and if relevant, what activities this enables on another service

Response: Once a user's identity has been verified, their reviews are marked as coming from a 'verified reviewer'.

f) how your identity verification policies have been developed, including any research that you can share

Response:

g) any steps you take to ensure that identity verification is available to all adult users, including users who may not be able to access certain types of identity verification

Response:

h) any consideration around users who may be vulnerable participating in the identity verification method

Response:

i) how you manage the identity verification of users who have multiple accounts

Response:

j) how you manage different identity verification methods operating simultaneously on your service, such as forms of age verification that require ID to complete the process, monetised schemes and notable user schemes, and how you consider user perceptions of these different methods

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response: No.

Question 33: Cost and effectiveness of these methods

For all respondents

Question 33: Please share any information about the costs and the effectiveness of identity verification methods

In your response to this request, please provide information relating to:

- (a) (d) where relevant for all respondents, and
- f) and g) where relevant for providers of user-to-user services that provide some types of identity verification for individual adult users.

Response:

(a) any insight into the cost of identity verification methods, including set-up and on-going costs, in terms of employee time and any other material costs, as well as any intended and unintended impacts on services more broadly

(b) how effective these identity verification methods are in verifying the identity of a user for the particular purpose for which verification is carried out
Response:
(c) any other benefits or unintended consequences from these schemes existing
Response:
(d) the safeguards necessary to ensure users' privacy is protected
Response:
For providers of user-to-user services that provide some types of identity verification for individual adult users
(e) any unintended consequences of implementing identity verification, such as the impact this may have on your site's ecosystem
Response:
(f) how you envisage your service operating in the digital identity market, bearing in mind moves towards cross-industry and federated identity schemes
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 34 and 35: User attitudes and demand for identity verification of user-to-user services
For all respondents
Question 34: What are user attitudes and demand for identity verification on user-to-user services?
In your response to this request, please provide information relating to (a) – (d) where relevant.
Response:
As noted in response to Question 31, we introduced voluntary identity verification in response to

feedback and sentiment from our users.

verification to signify authenticity

Response:

Response:

Response:

engage with them

(a) whether they value verification being offered on a service

(b) whether verification influences user behaviour, such as whether they perceive identity

(c) attitudes towards non-verified, anonymous or pseudonymous users and the willingness to

(d) who you deem to be 'vulnerable' in terms of verifying their identity online – for example, whether this includes users unable to access or less likely to hold identification documentation,
and those who may become vulnerable by displaying their identity to other users.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
For providers of user-to-user services that provide some types of identity verification for individually users
Question 35: How do you measure engagement with your identity verification methods?
In your response to this request, please provide information relating to (a) and (b) where relevant
Response:
Engagement with our optional identity verification method can be measured by uptake, as well a in relation to user feedback with regards to how they perceive and view reviews from those deemed 'verified users' in contrast with those who do not choose to verify their identity,
(a) take-up of identity verification by your users
Response:
(b) any insight into whether identity verification has any other effect on user behaviour, such as the content that users post and the amount that they engage with your service.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Resnanse:

Your response - Fraudulent advertising

Questions 36 - 42: Overarching considerations

For all respondents

Question 36: Please provide evidence of the following:

(a) The most prevalent kinds of fraudulent advertising activity on user-to-user and search services (e.g. illegal financial promotions, misleading statements, malvertising)

Response:

(b) The harms associated with different kinds of fraudulent advertisements, the severity of such harms, and, if relevant, how this varies by user group

(c) The key challenges to successfully detecting different types of fraudulent paid-for advertising, and how these challenges can be minimised or resolved
Response:
(d) The prioritisation of suspected fraudulent advertising within all categories of harmful advertising queues, e.g. account verification, user reports, appeals
Response:
(e) The proportion of fraudulent advertisements that are currently estimated to remain undetected by services' systems.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 37: What technological developments aiding the prevention/detection of fraudulent advertisements do you anticipate in the coming years, and how costly and effective do you expect them to be? What are the challenges/barriers to their development?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 38: If you have information/evidence/suggested mitigations to share which may be useful in the preparation of codes of practice, which is not covered by the questions above, please include these under 'Overarching considerations'.
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
For providers of online services
Question 39: What proportion of all paid-for advertising on your service is identified as fraudulent advertising?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:

Question 40: Does your service take any steps to warn users of the risk of encountering fraudulent advertising or to educate them about how to identify potentially fraudulent advertising?

Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 41: Please provide information regarding the proportion of successfully identified fraudulent advertisements that are identified via:
(a) automated systems
Response:
(b) human processes
Response:
(c) user reports
Response:
(d) other (please provide further detail).
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 42: What is the average and/or median time taken between the identification of a fraudulent advertisement and its removal/other actions taken? (If other actions taken, please specify what they are).
Response:
Response: Is this response confidential? (if yes, please specify which part(s) are confidential)
Is this response confidential? (if yes, please specify which part(s) are confidential)
Is this response confidential? (if yes, please specify which part(s) are confidential)
Is this response confidential? (if yes, please specify which part(s) are confidential) Response:
Is this response confidential? (if yes, please specify which part(s) are confidential) Response: Question 43: Proactive technology
Is this response confidential? (if yes, please specify which part(s) are confidential) Response: Question 43: Proactive technology For all respondents Question 43: Please provide any evidence you have regarding proactive technologies which
Is this response confidential? (if yes, please specify which part(s) are confidential) Response: Question 43: Proactive technology For all respondents Question 43: Please provide any evidence you have regarding proactive technologies which could be used to identify fraudulent advertising activity.
Is this response confidential? (if yes, please specify which part(s) are confidential) Response: Question 43: Proactive technology For all respondents Question 43: Please provide any evidence you have regarding proactive technologies which could be used to identify fraudulent advertising activity. In particular, we are interested in information related to the following points: (a) The kinds of proactive technology which are/could be applied to identify or prevent
Is this response confidential? (if yes, please specify which part(s) are confidential) Response: Question 43: Proactive technology For all respondents Question 43: Please provide any evidence you have regarding proactive technologies which could be used to identify fraudulent advertising activity. In particular, we are interested in information related to the following points: (a) The kinds of proactive technology which are/could be applied to identify or prevent fraudulent advertising

(c) The effectiveness, accuracy and lack of bias of such technology (including compared to alternative proactive and non-proactive methods) in relation to detecting fraudulent advertising and accounts which post fraudulent advertising material
Response:
(d) How proactive technologies are maintained and kept up to date
Response:
e) Information related to the associated time and/or costs for set-up, operation, and human review
Response:
f) The cost of integrating such technologies: (a) for the first time; and (b) when updating these technologies over time
Response:
g) Whether there are cost savings associated with these technologies
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 44: Advertising onboarding and verification For all respondents
Question 44: Please provide any evidence you have regarding the processes for advertiser onboarding and verification related to protections against fraudulent advertising. In your response, please indicate whether these processes are currently implemented in respect of services which are in scope of the Act or whether they stem from another sector In particular, we are interested in information related to the following points:
(a) The criteria which advertisers are verified against, including documentation/evidence used to support verification, and what advertisers are required to declare
Response:
(b) The role of (a) automated processing and (b) human processing in the verification process, and how they interact
Response:
(c) The costs associated with advertiser verification and how those costs vary as scale increases
Response:
nesponse.
(d) The percentage of advertiser accounts that are verified

e) Whether advertisers are permitted to publish advertisements on the service while the verification process is ongoing
Response:
f) Whether there are additional/specific verification checks for advertisers placing adverts of certain kinds or targeting certain audiences, such as about specific products or services, or targeting users under the age of 18
Response:
g) Whether the verification of an advertiser account expires after a certain amount of time or certain activity, such as when advertisers make changes to their account or profile
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 45: Service review of submitted advertisements/sponsored search results For all respondents
Question 45: Please provide any evidence you have regarding the processes that services in
scope of the Act have in place to review submitted paid-for advertisements and identify
fraudulent advertising material.
fraudulent advertising material. In particular, we are interested in information related to the following points:
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii)
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact Response: (c) The red flags which trigger advertisement review processes both (i) prior to and (ii) after
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact Response: (c) The red flags which trigger advertisement review processes both (i) prior to and (ii) after publication and the basis on which those red flags are selected
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact Response: (c) The red flags which trigger advertisement review processes both (i) prior to and (ii) after publication and the basis on which those red flags are selected Response:
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact Response: (c) The red flags which trigger advertisement review processes both (i) prior to and (ii) after publication and the basis on which those red flags are selected Response: (d) The timescales for review
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact Response: (c) The red flags which trigger advertisement review processes both (i) prior to and (ii) after publication and the basis on which those red flags are selected Response: (d) The timescales for review Response: (e) What happens to the advertisement's visibility and reach, if it is flagged as suspected as
In particular, we are interested in information related to the following points: (a) The percentage of submitted advertisements which are reviewed both (i) prior to and (ii) after publication Response: (b) The role (i) automated processing and (ii) human processing play in the review process and how they interact Response: (c) The red flags which trigger advertisement review processes both (i) prior to and (ii) after publication and the basis on which those red flags are selected Response: (d) The timescales for review Response: (e) What happens to the advertisement's visibility and reach, if it is flagged as suspected as being fraudulent (either by a user or automated system)

(g) Whether trusted flagger reporting is employed to inform services' review processes. If it is, how is it applied, what guidelines / criteria does it follow, and who are those trusted flaggers?

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:

Question 46: Advertiser appeals of verification/review decisions

For all respondents

Question 46: Please provide any evidence you have regarding advertiser appeals of verification/review decisions relating to fraudulent advertising on services in scope of the Act. In particular, we are interested in information related to the following points: (a) The role of (i) automated processing and (ii) human processing in the appeals process, and how they interact; Response: (b) The level of proof required for an appeal to be accepted; Response: (c) The most frequent bases for appeals against sanctions decisions on fraudulent advertising content Response: (d) The ratio of decisions that are appealed against Response: (e) The costs associated with appeals Response: (f) The proportion of appealed decisions which are upheld and overturned Response: Is this response confidential? (if yes, please specify which part(s) are confidential) Response:

Question 47: User reporting mechanisms

For all respondents

Question 47: Please provide any evidence you have regarding user reporting mechanisms for fraudulent advertising on services in scope of the Act.

In particular, we are interested in information related to the following points:

(a) What user reporting tools there are for paid-for advertisements, and how these tools differ from those for user-generated content and/or search results and other search functionalities that are not paid-for advertising Response: (b) What percentage of user reports of advertisements relate to suspected fraudulent content, and the processes for taking action in relation to such reports Response: (c) Any statistics you can share on (i) the number of user reports of suspected fraudulent advertising received and resolved over a specific period and (b) the number of initial decisions appealed by users who made the report Response: (d) The criteria used to classify and prioritise user reports Response: (e) The median and/or average time it takes to respond to a user report, and any measures that are in place to ensure timely and accurate responses to user reports Response: (f) Any measures taken to make user reporting tools accessible, easy to use and easy to find for users Response: (g) How transparency and communication is maintained with users who have submitted reports

Question 48: Use/involvement of third parties

Is this response confidential? (if yes, please specify which part(s) are confidential)

For all respondents

Response:

Response:

Question 48: Please provide any evidence relevant to fraudulent advertising that you have, regarding the involvement and role of third parties in the provision of paid-for advertisements on services in scope of the Act.

In line with the proportionality criteria under sections 38(5) and 39(5) of the Act, we welcome information related to how the involvement of third parties impacts the degree of control that services have over fraudulent advertising content.

We also welcome information regarding contractual arrangements and how those arrangements are enforced.

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)

Response:
Question 49: Generative AI and deepfakes
For all respondents
Question 49: Please provide any evidence you have regarding the impact of generative AI developments and deepfakes on the incidence and detection of fraudulent advertisements on services in scope of the Act.
In particular, we are interested in information related to the following points:
(a) The frequency of deepfake fraudulent advertisements' occurrence, in absolute terms and/or as a proportion of all fraudulent advertisements, and how you expect this to evolve in the future
Response:
(b) What methodologies/technologies are currently employed to detect fraudulent advertisements which include deepfake or otherwise AI-generated content, and the effectiveness of these tools
Response:
(c) Whether detection technologies are developed in-house or acquired from a third-party, and how long it takes to develop and/or integrate those tools into wider systems
Response:
(d) The accuracy of detection methods, including true positive and false positive rates
Response:
(e) The costs associated with the development/acquisition and deployment of these detection mechanisms
Response:
(f) The types of deepfake or AI-generated content (in terms of either media type or subject) in fraudulent advertisements that are most difficult to detect i) via automated processes, ii) by human moderators, iii) by service users
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:

Your response – Access to information about a deceased child's use of a service

Questions 50 – 55: Processes for requesting information about a deceased child's use of a service

For all respondents

Question 50: What kinds of information might parents want to see about their child's use of the service?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 51: How long should it take to receive information in response to a request?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 52: What mechanisms could, or should services provide for parents to find out what they need to do to obtain information and updates in these circumstances?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 53: What support or information do parents need to guide them through the process of making a request?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
For providers of online services

Question 54: What kinds of information do you provide and how do you provide this information?

In your response to this request, please provide information relating to (a) where relevant.

example whether there are technological, cost or privacy factors that mean certain kinds of information may not be feasible to provide
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Question 55: How long does it typically take you to provide information in response to a request?
In your response to this request, please provide information relating to (a) where relevant.
Response:
a) How long should it reasonably take services to provide information in these circumstances?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
Questions 56 and 57: Complaints systems
For all respondents
Question 56: What can providers of online services do to ensure the transparency, accessibility, ease of use and users' awareness of complaints mechanisms in relation to deceased user information request processes?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:
For providers of online services
Question 57: Can you provide any evidence or information about the best practices for effective complaints mechanisms which could inform an approach to complaints about information request processes pertaining to a deceased user?
Response:
Is this response confidential? (if yes, please specify which part(s) are confidential)
Response:

a) If there are certain types of information you cannot provide, please explain why, for

Question 58: Evidence

For providers of online services

Response:

Question 58: What kinds of evidence do you require about the identity of the person making the request and their relationship to the deceased user?

In your response to this request, please provide information relating to (a) and (b) where relevant.

Response:

(a) Do you, or would you, require different kinds of evidence in the event that the deceased user is a child?

Response:

(b) What evidence do, or would, you require that a user is deceased?

Response:

Is this response confidential? (if yes, please specify which part(s) are confidential)