# Volume 2: The causes and impacts of online harm

Ofcom's Register of Risks

| Question 1: |
| --- |
| i)     Do you have any comments on Ofcom's assessment of the causes and impacts of online harms? |
| Response: **While risk assessment is a valuable practice, it often falls short in adequately addressing the complexities of emerging technologies and their potential harms. For instance, merely identifying that planes falling out of the sky is a risk does not suffice. What's crucial is establishing binding requirements to ensure such catastrophic events are prevented.** |
| ii)     Do you think we have missed anything important in our analysis? Please provide evidence to support your answer. |
| Response: **Both research findings and real-world incidents of public harm underscore the need to transition from high-level theoretical frameworks to a more focused examination of the technical intricacies underlying these harms. It's imperative to adapt risk assessment methodologies by integrating compliance mechanisms and precisely defining what constitutes risk, enabling us to effectively mitigate potential harms.** |
| iii)    Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: **No** |

| Question 2: |
| --- |
| i)     Do you have any views about our interpretation of the links between risk factors and different kinds of illegal harm? Please provide evidence to support your answer. |
| Response: |
| ii)     Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

# Volume 3: How should services assess the risk of online harms?

## Governance and accountability

| Question 3: |
| --- |
| i)      Do you agree with our proposals in relation to governance and accountability measures in the illegal content Codes of Practice? |
| Response: **No, making external audit requirements is the most effective step towards enforcement and implementation.** |
| ii)      Do you think we have missed anything important in our analysis? Please provide evidence to support your answer. |
| Response: **Audits and oversight mechanisms, proven successful in various industries including technology, is undeniable. In democratic contexts, products entering the market undergo rigorous scrutiny, encompassing external oversight to assess the implementation of risk mitigation measures. Implementing independent third-party audits ensures robust inspection of systems, safeguarding individuals throughout technological processes.** |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: **No** |

| Question 4: |
| --- |
| i)      Do you agree with the types of services that we propose the governance and accountability measures should apply to? |
| Response: |
| ii)      Please explain your answer. |
| Response: |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

## Question 5:

| | | |
|---|---|---|
| i) | Are you aware of any additional evidence of the efficacy, costs and risks associated with a potential future measure to requiring services to have measures to mitigate and manage illegal content risks audited by an independent third-party? | |

Response: **External audits offer superior effectiveness, despite the associated risks and costs, mirroring the success seen in other sectors where mandatory regulations, such as seatbelt laws and vaccine clinical trials, have vastly improved safety standards. Similarly, AI and online systems benefit from rigorous safety testing and auditing to enhance user protection.**

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

Response: No

## Question 6:

| | | |
|---|---|---|
| i) | Are you aware of any additional evidence of the efficacy, costs and risks associated with a potential future measure to tie remuneration for senior managers to positive online safety outcomes? | |

Response:

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

Response:

## Service's risk assessment

## Question 7:

| | | |
|---|---|---|
| i) | Do you agree with our proposals? | |

Response:

| | | |
|---|---|---|
| ii) | Please provide the underlying arguments and evidence that support your views. | |

Response:

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

Response:

*Specifically, we would also appreciate evidence from regulated services on the following:*

| Question 8: |
| --- |

| i) | Do you think the four-step risk assessment process and the Risk Profiles are useful models to help services navigate and comply with their wider obligations under the Act? |
| --- | --- |

Response: **Risk-assessment is a necessary step but remains a limited exercise. Regulators must shift focus towards establishing thresholds and specific metrics to evaluate system performance comprehensively. Algorithmic auditing emerges as a critical tool, bridging the gap in ensuring compliance with evolving regulatory standards, especially considering the highly contextual nature of online system impacts.**

**Hence, while the four-step risk assessment process provides valuable insights, it's imperative to address the question of implementation and enforcement. Industries seek clear compliance thresholds from regulators to effectively navigate the regulatory landscape. Until then, regulatory requirements remain primarily compliance-driven.**

**In conclusion, augmenting existing compliance frameworks with a focus on impact assessment represents a crucial step forward. By incentivizing product teams to prioritize safety metrics from the inception of system development, we can foster a culture of responsible innovation and mitigate inherent risks effectively.**

| ii) | Please provide the underlying arguments and evidence that support your views. |
| --- | --- |

Response:

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

Response:

| Question 9: | |
|---|---|
| i) | Are the Risk Profiles sufficiently clear? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Do you think the information provided on risk factors will help you understand the risks on your service? |
| Response: | |
| iv) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| v) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

## Record keeping and review guidance

| Question 10: | |
|---|---|
| i) | Do you have any comments on our draft record keeping and review guidance? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

| Question 11: | |
|---|---|
| i) | Do you agree with our proposal not to exercise our power to exempt specified descriptions of services from the record keeping and review duty for the moment? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

# Volume 4: What should services do to mitigate the risk of online harms

Our approach to the Illegal content Codes of Practice

| Question 12: | |
|---|---|
| i) | Do you have any comments on our overarching approach to developing our illegal content Codes of Practice? |
| Response: | |
| ii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

| Question 13: | |
|---|---|
| i) | Do you agree that in general we should apply the most onerous measures in our Codes only to services which are large and/or medium or high risk? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

| Question 14: | |
|---|---|
| i) | Do you agree with our definition of large services? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

| Question 15: | |
|---|---|
| i) Do you agree with our definition of multi-risk services? | |
| Response: | |
| ii) Please provide the underlying arguments and evidence that support your views. | |
| Response: | |
| iii) Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | |

| Question 16: | |
|---|---|
| i) Do you have any comments on the draft Codes of Practice themselves? | |
| Response: | |
| ii) Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | |

| Question 17: | |
|---|---|
| i) Do you have any comments on the costs assumptions set out in Annex 14, which we used for calculating the costs of various measures? | |
| Response: | |
| ii) Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | |

## Content moderation (User to User)

| Question 18: | |
|---|---|
| i) Do you agree with our proposals? | |
| Response: | |
| ii) Please provide the underlying arguments and evidence that support your views. | |
| Response: | |
| iii) Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | |

## Content moderation (Search)

| Question 19: |
| --- |
| i)          Do you agree with our proposals? |
| Response: |
| ii)         Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)       Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

## Automated content moderation (User to User)

| Question 20: |
| --- |
| i)          Do you agree with our proposals? |
| Response: |
| ii)         Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)       Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 21: |
| --- |
| i)          Do you have any comments on the draft guidance set out in Annex 9 regarding whether content is communicated 'publicly' or 'privately'? |
| Response: |
| ii)         Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

***Do you have any relevant evidence on:***

| Question 22: |
| --- |
| i)          Accuracy of perceptual hash matching and the costs of applying CSAM hash matching to smaller services; |
| Response: |
| ii)         Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)       Is this response confidential? (if yes, please specify which part(s) are confidential) |

| | |
|---|---|
| Response: | |

**Question 23:**

| | | |
|---|---|---|
| | i) | Ability of services in scope of the CSAM hash matching measure to access hash databases/services, with respect to access criteria or requirements set by database and/or hash matching service providers; |
| Response: | | |
| | ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | | |
| | iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | | |

**Question 24:**

| | | |
|---|---|---|
| | i) | Costs of applying our CSAM URL detection measure to smaller services, and the effectiveness of fuzzy matching for CSAM URL detection;; |
| Response: | | |
| | ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | | |
| | iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | | |

**Question 25:**

| | | |
|---|---|---|
| | i) | Costs of applying our articles for use in frauds (standard keyword detection) measure, including for smaller services; |
| Response: | | |
| | ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | | |
| | iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | | |

| Question 26: | | |
|---|---|---|
| i) | An effective application of hash matching and/or URL detection for terrorism content, including how such measures could address concerns around 'context' and freedom of expression, and any information you have on the costs and efficacy of applying hash matching and URL detection for terrorism content to a range of services. | |
| Response: | | |
| ii) | Please provide the underlying arguments and evidence that support your views. | |
| Response: | | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

## Automated content moderation (Search)

| Question 27: | | |
|---|---|---|
| i) | Do you agree with our proposals? | |
| Response: | | |
| ii) | Please provide the underlying arguments and evidence that support your views. | |
| Response: | | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

## User reporting and complaints (U2U and search)

| Question 28: | | |
|---|---|---|
| i) | Do you agree with our proposals? | |
| Response: | | |
| ii) | Please provide the underlying arguments and evidence that support your views. | |
| Response: | | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

## Terms of service and Publicly Available Statements

| Question 29: |
| --- |
|     i)        Do you agree with our proposals? |
| Response: |
|     ii)       Please provide the underlying arguments and evidence that support your views. |
| Response: |
|     iii)     Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 30: |
| --- |
|     i)        Do you have any evidence, in particular on the use of prompts, to guide further work in this area? |
| Response: |
|     ii)       Please provide the underlying arguments and evidence that support your views. |
| Response: |
|     iii)     Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

## Default settings and user support for child users (U2U)

| Question 31: |
| --- |
|     i)        Do you agree with our proposals? |
| Response: |
|     ii)       Please provide the underlying arguments and evidence that support your views. |
| Response: |
|     iii)     Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 32: |
| --- |
|     i)        Are there functionalities outside of the ones listed in our proposals, that should explicitly inform users around changing default settings? |
| Response: |
|     ii)       Is this response confidential? (if yes, please specify which part(s) are confidential) |

| Response: |
|---|

| **Question 33:** |
|---|
| i)      Are there other points within the user journey where under 18s should be informed of the risk of illegal content? |
| Response: |
| ii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

## Recommender system testing (U2U)

| **Question 34:** |
|---|
| i)      Do you agree with our proposals? |
| Response: **In order to ensure comprehensive safety assessment, it is imperative to gather safety metrics across various stages of development. While model testing provides valuable insights, it must be complemented by real-world impact metrics. As auditors, we employ a multifaceted approach, designing tests to measure actual engagement and outcomes through experiments and user interactions. Moreover, it's essential to establish clear thresholds for these metrics. When should action be taken? What constitutes a metric falling above or below the acceptable threshold? These are the critical questions we address at Eticas.ai, drawing from our extensive experience and expertise. We are eager to share our insights and collaborate on relevant projects.** |
| ii)      Please provide the underlying arguments and evidence that support your views. |
| Response: **Relying solely on model metrics is akin to endorsing a vaccine based solely on its performance in controlled laboratory settings, without undergoing clinical trials on humans.** |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: **No** |

| **Question 35:** |
|---|
| i)      What evaluation methods might be suitable for smaller services that do not have the capacity to perform on-platform testing? |
| Response: |
| ii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

*We are aware of design features and parameters that can be used in recommender system to minimise the distribution of illegal content, e.g. ensuring content/network balance and low/neutral weightings on content labelled as sensitive.*

| Question 36: | | |
|---|---|---|
| i) | Are you aware of any other design parameters and choices that are proven to improve user safety? | |
| Response: | | |
| ii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

## Enhanced user control (U2U)

| Question 37: | | |
|---|---|---|
| i) | Do you agree with our proposals? | |
| Response: | | |
| ii) | Please provide the underlying arguments and evidence that support your views. | |
| Response: | | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

| Question 38: | | |
|---|---|---|
| i) | Do you think the first two proposed measures should include requirements for how these controls are made known to users? | |
| Response: | | |
| ii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

| Question 39: | | |
|---|---|---|
| i) | Do you think there are situations where the labelling of accounts through voluntary verification schemes has particular value or risks? | |
| Response: | | |
| ii) | Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | | |

## User access to services (U2U)

| Question 40: |
|---|

| | |
|---|---|
| i) | Do you agree with our proposals? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

**Do you have any supporting information and evidence to inform any recommendations we may make on blocking sharers of CSAM content? Specifically:**

| Question 41: | |
|---|---|
| i) | What are the options available to block and prevent a user from returning to a service (e.g. blocking by username, email or IP address, or a combination of factors)? |
| Response: | |
| ii) | What are the advantages and disadvantages of the different options, including any potential impact on other users? |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

| Question 42: | |
|---|---|
| i) | How long should a user be blocked for sharing known CSAM, and should the period vary depending on the nature of the offence committed? |
| Response: | |
| ii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

**There is a risk that lawful content is erroneously classified as CSAM by automated systems, which may impact on the rights of law-abiding users.**

| Question 43: | |
|---|---|
| i) | What steps can services take to manage this risk? For example, are there alternative options to immediate blocking (such as a strikes system) that might help mitigate some of the risks and impacts on user rights? |
| Response: | |
| ii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

## Service design and user support (Search)

| Question 44: |
| --- |
| i)       Do you agree with our proposals? |
| Response: |
| ii)     Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)    Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

## Cumulative Assessment

| Question 45: |
| --- |
| i)       Do you agree that the overall burden of our measures on low risk small and micro businesses is proportionate? |
| Response: |
| ii)     Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)    Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 46: |
| --- |
| i)       Do you agree that the overall burden is proportionate for those small and micro businesses that find they have significant risks of illegal content and for whom we propose to recommend more measures? |
| Response: |
| ii)     Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)    Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 47: |
| --- |
| i)       We are applying more measures to large services. Do you agree that the overall burden on large services proportionate? |
| Response: |

| | |
|---|---|
| ii)      Please provide the underlying arguments and evidence that support your views. | |
| Response: | |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | |

## Statutory Tests

| Question 48: | |
|---|---|
| i)      Do you agree that Ofcom's proposed recommendations for the Codes are appropriate in the light of the matters to which Ofcom must have regard? | |
| Response: | |
| ii)      Please provide the underlying arguments and evidence that support your views. | |
| Response: | |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) | |
| Response: | |

# Volume 5: How to judge whether content is illegal or not?

The Illegal Content Judgements Guidance (ICJG)

| Question 49: |
| --- |
| i)      Do you agree with our proposals, including the detail of the drafting? |
| Response: |
| ii)      What are the underlying arguments and evidence that inform your view? |
| Response: |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 50: |
| --- |
| i)      Do you consider the guidance to be sufficiently accessible, particularly for services with limited access to legal expertise? |
| Response: |
| ii)      Please provide the underlying arguments and evidence that support your views. |
| Response: |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

| Question 51: |
| --- |
| i)      What do you think of our assessment of what information is reasonably available and relevant to illegal content judgements? |
| Response: |
| ii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |

# Volume 6: Information gathering and enforcement powers, and approach to supervision.

## Information powers

| Question 52: | |
|---|---|
| i) | Do you have any comments on our proposed approach to information gathering powers under the Online Safety Act? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

## Enforcement powers

| Question 53: | |
|---|---|
| i) | Do you have any comments on our draft Online Safety Enforcement Guidance? |
| Response: | |
| ii) | Please provide the underlying arguments and evidence that support your views. |
| Response: | |
| iii) | Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: | |

# Annex 13: Impact Assessments

| Question 54: |
|---|
| i)      Do you agree that our proposals as set out in Chapter 16 (reporting and complaints), and Chapter 10 and Annex 6 (record keeping) are likely to have positive, or more positive impacts on opportunities to use Welsh and treating Welsh no less favourably than English? |
| Response: |
| ii)      If you disagree, please explain why, including how you consider these proposals could be revised to have positive effects or more positive effects, or no adverse effects or fewer adverse effects on opportunities to use Welsh and treating Welsh no less favourably than English. |
| Response: |
| iii)      Is this response confidential? (if yes, please specify which part(s) are confidential) |
| Response: |