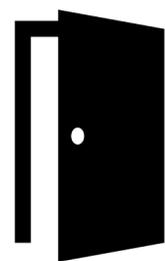


This is the ASPARC framework.

It serves to break down the user-generated content journey into 7 distinct phases.



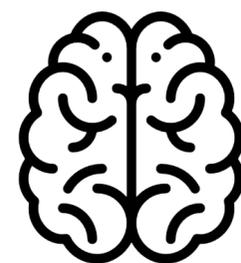
**Access**



**Sign-on**



**Participate**



**Analyse**



**Respond**



**Comply**

The models in this document were built around the ASPARC framework.

For each phase, from Sign-on to Comply, there is a functional model that provides a detailed view of the functional and architectural processes that platforms take to achieve the objectives called out below. 'Access', the process of accessing the internet, was out of scope for our investigation but is included below for completeness.

The models are generically representative of all online platforms, however significant areas of local diversity of implementation are highlighted. For more detail, please refer to the report.

Click on the icons to view the models.



## Access

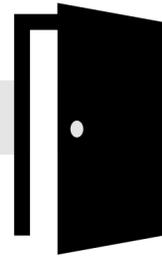
Select device

Select operating system

Select Network access

Select on-device safety features

Deploy On-Network safety features



## Sign-on

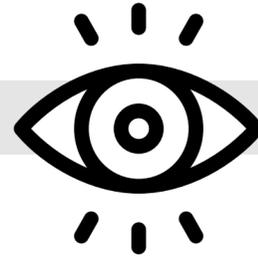
Registration

Establish user profile

Age verification

User sign in

Connect with other users



## Participate

Create content

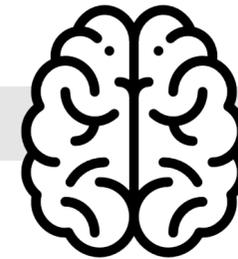
Curate content

Post static content

Livestream content

View content

Report content



## Analyse

Train classifiers

Auto-detect content

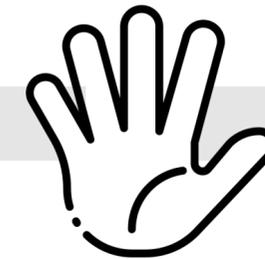
Receive user signal

Operate intelligence desk

Auto-assess content against single classifiers

Auto-assess content against combination classifiers

Human moderation



## Respond

Relegate content (reduce virality)

Remove content

Sanction user

Alert safety partners

Alert law enforcement

Invoke content incident protocol



## Comply

Manage appeals process

Maintain MI database of harmful content and bad actors

Share safety data with 3rd party partners

Publish transparency report

Audit performance

# Phase 1 Sign-on

The 'Sign-on' model represents the process by which a user signs up and establishes their account. It covers user registration, age verification, sign in and the generation of connection requests for other users.

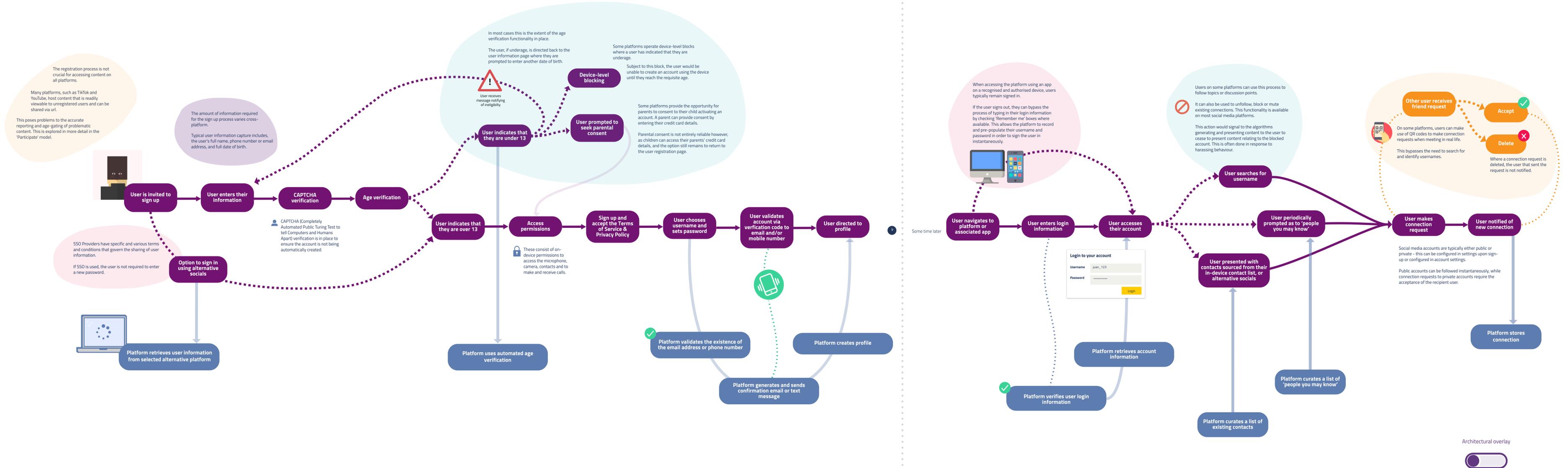
## Key



Step in user journey



Platform activities



# Phase 1 Sign-on

The 'Sign-on' model represents the process by which a user signs up and establishes their account. It covers user registration, age verification, sign in and the generation of connection requests for other users.

## Key

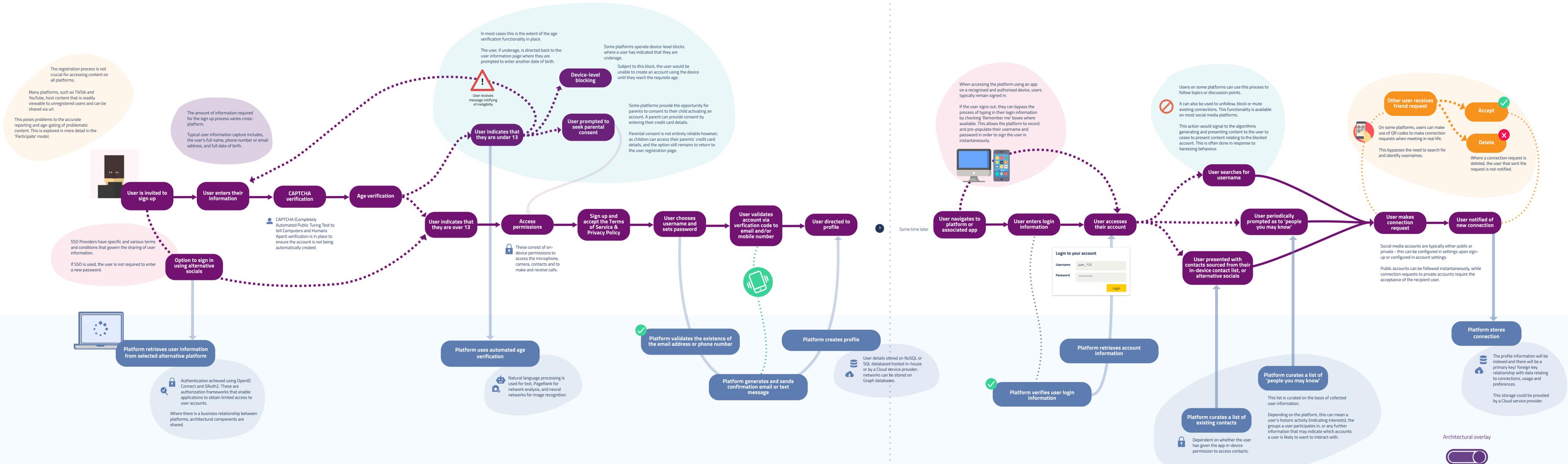


Step in user journey



Platform activities

- security
- harmful content
- machine learning
- profile/registration
- database/persistence
- image recognition
- cloud
- management information capture
- external service
- search/analysis



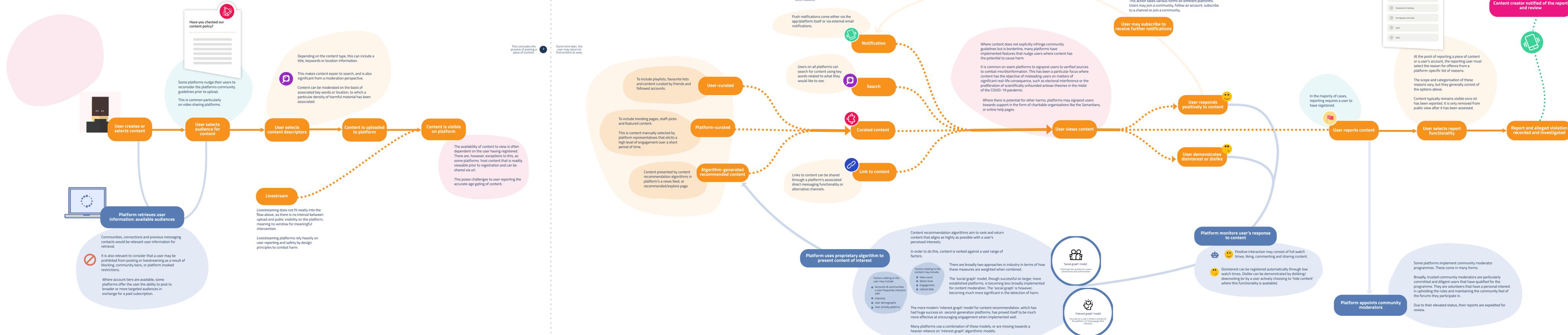
Architectural overlay



## Phase 2 Participate

The 'Participate' model details a user's engagement with content on a platform. It covers the processes of creating, uploading and curating content, the various methods of discovering content to view and interact with, and the process by which users report objectionable content.

### Key



# Phase 2 Participate

The 'Participate' model details a user's engagement with content on a platform. It covers the processes of creating, uploading and curating content, the various methods of discovering content to view and interact with, and the process by which users report objectionable content.

## Key



- security
- machine learning
- database/persistence
- cloud
- external service
- search/analysis
- harmful content
- profile/registration
- image recognition
- management information capture



## Phase 3 Analyse

The 'Analyse' model details the measures put in place to analyse content and user behaviour in order to determine whether it is harmful. These measures include the operation of an intelligence desk, the receipt of user signals, the application of single and combination classifiers and the triage process in place to prioritise content for human review.

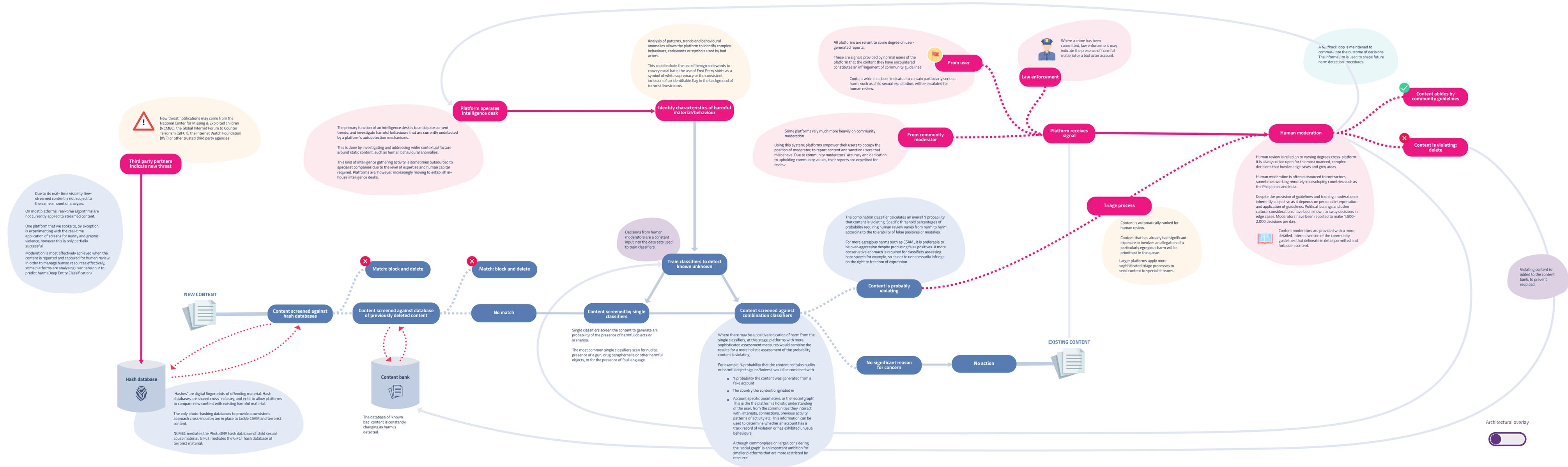
### Key



Human activity on behalf of platform



Platform activities



**Third party partners indicate new threat**

New threat notifications may come from the National Center for Missing & Exploited Children (NCMEC), the Global Internet Forum to Counter Terrorism (GIFCT), the Internet Watch Foundation (IWF) or other trusted third party agencies.

Due to its real-time visibility, live-streamed content is not subject to the same amount of analysis.

On most platforms, real-time algorithms are not currently applied to streamed content.

One platform that we spoke to, by exception, is experimenting with the real-time application of screens for nudity and graphic violence, however this is only partially successful.

Moderation is most effectively achieved when the content is reported and captured for human review. In order to manage human resources effectively, some platforms are analysing user behaviour to predict harm (Deep Entity Classification).

**Platform operates intelligence desk**

The primary function of an intelligence desk is to anticipate content trends, and investigate harmful behaviours that are currently undetected by a platform's autodetection mechanisms.

This is done by investigating and addressing wider contextual factors around static content, such as human behavioural anomalies.

This kind of intelligence gathering activity is sometimes outsourced to specialist companies due to the level of expertise and human capital required. Platforms are, however, increasingly moving to establish in-house intelligence desks.

**Identify characteristics of harmful material/behaviour**

Analysis of patterns, trends and behavioural anomalies allows the platform to identify complex behaviours, codewords or symbols used by bad actors.

This could include the use of benign codewords to convey racial hate, the use of Fred Perry shirts as a symbol of white supremacy or the consistent inclusion of an identifiable flag in the background of terrorist livestreams.

**From user**

All platforms are reliant to some degree on user-generated reports.

These are signals provided by normal users of the platform that the content they have encountered constitutes an infringement of community guidelines.

Content which has been indicated to contain particularly serious harm, such as child sexual exploitation, will be escalated for human review.

**Law enforcement**

Where a crime has been committed, law enforcement may indicate the presence of harmful material or a bad actor account.

**From community moderator**

Some platforms rely much more heavily on community moderation.

Using this system, platforms empower their users to occupy the position of moderator, to report content and sanction users that misbehave. Due to community moderators' accuracy and dedication to upholding community values, their reports are expedited for review.

**Content is probably violating**

The combination classifier calculates an overall % probability that content is violating. Specific threshold percentages of probability requiring human review varies from harm to harm according to the tolerability of false positives or mistakes.

For more egregious harms such as CSAM, it is preferable to be over-aggressive despite producing false positives. A more conservative approach is required for classifiers assessing hate speech for example, so as not to unnecessarily infringe on the right to freedom of expression.

**Triage process**

Content is automatically ranked for human review.

Content that has already had significant exposure or involves an allegation of a particularly egregious harm will be prioritised in the queue.

Larger platforms apply more sophisticated triage processes to send content to specialist teams.

**Human moderation**

Human review is relied on to varying degrees cross-platform. It is always relied upon for the most nuanced, complex decisions that involve edge cases and grey areas.

Human moderation is often outsourced to contractors, sometimes working remotely in developing countries such as the Philippines and India.

Despite the provision of guidelines and training, moderation is inherently subjective as it depends on personal interpretation and application of guidelines. Political leanings and other cultural considerations have been known to sway decisions in edge cases. Moderators have been reported to make 1,500-2,000 decisions per day.

Content moderators are provided with a more detailed, internal version of the community guidelines that delineate in detail permitted and forbidden content.

**Feedback loop**

A feedback loop is maintained to communicate the outcome of decisions. The information is used to shape future harm detection procedures.

## Phase 3 Analyse

The 'Analyse' model details the measures put in place to analyse content and user behaviour in order to determine whether it is harmful. These measures include the operation of an intelligence desk, the receipt of user signals, the application of single and combination classifiers and the triage process in place to prioritise content for human review.

### Key

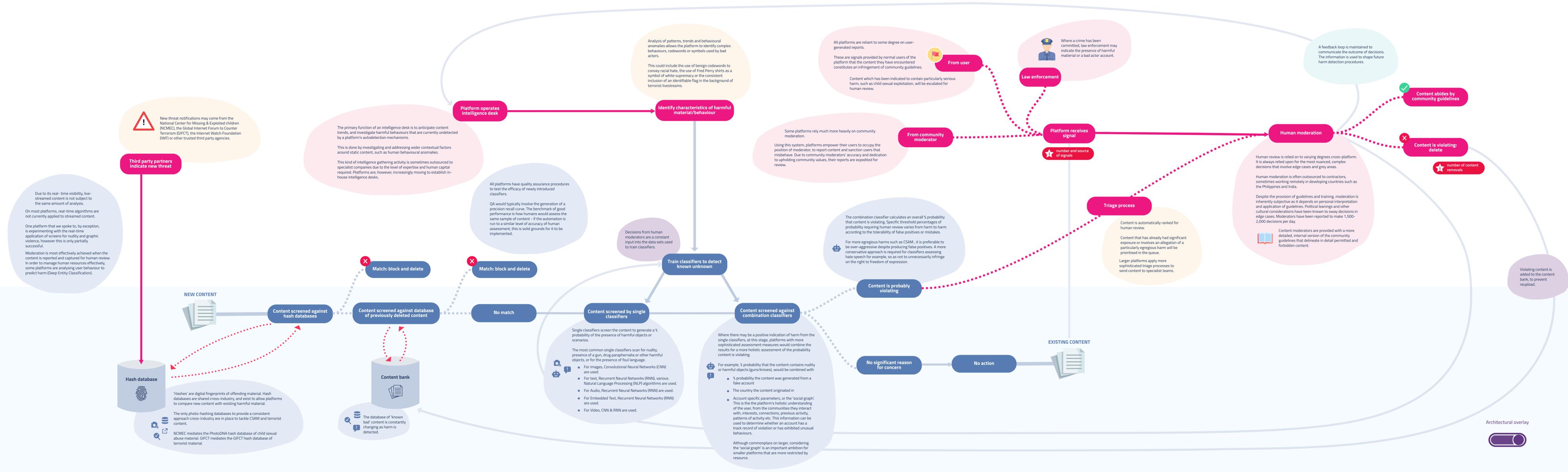


Human activity on behalf of platform



Platform activities

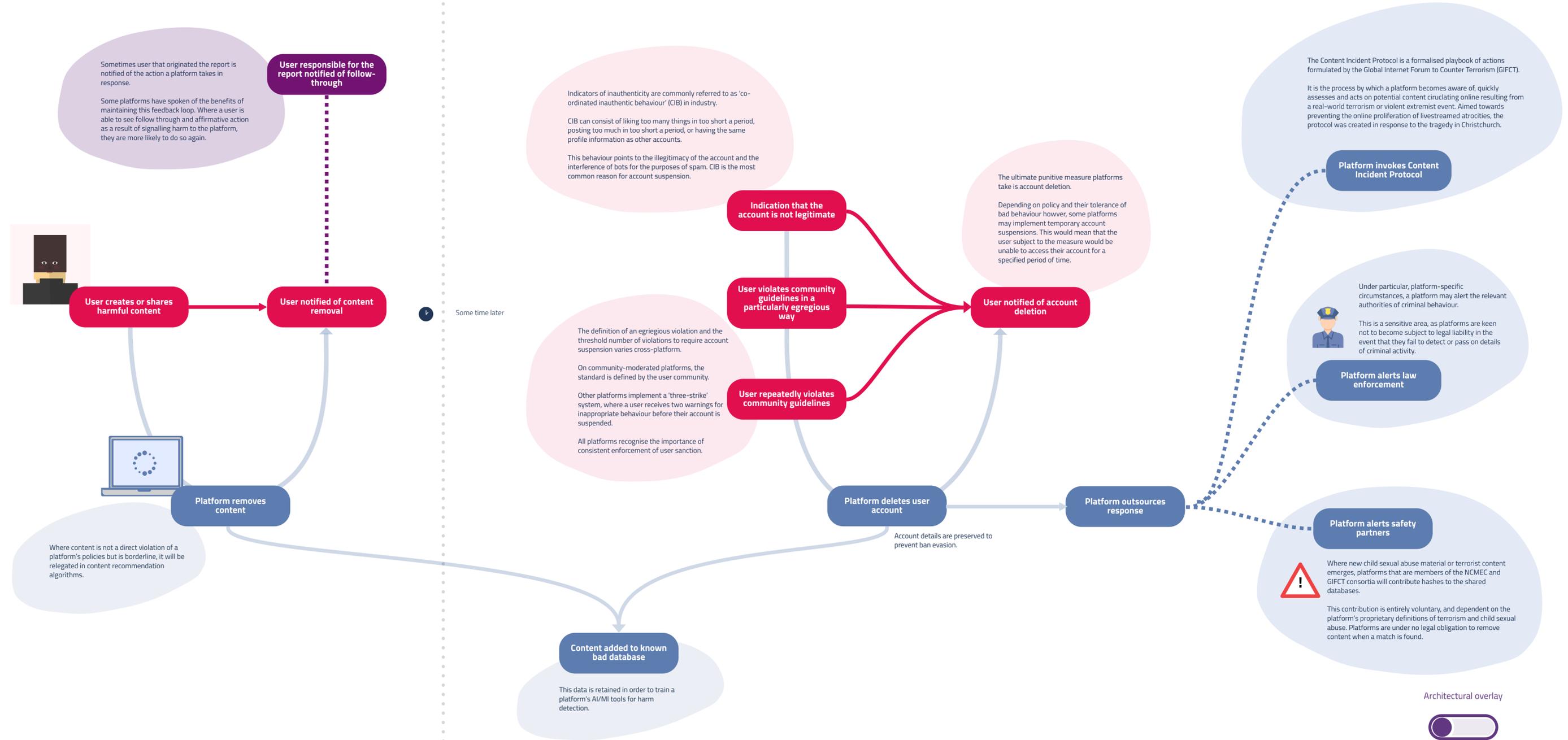
- security
- machine learning
- database/persistence
- cloud
- external service
- search/analysis
- harmful content
- profile/registration
- image recognition
- management information capture



Phase 4  
**Respond**

The 'Respond' model provides detail on how platforms respond to the finding of harmful content. It covers the process of content removal, the application of user sanctions and the procedures in place to ensure that the relevant authorities are alerted.

Key



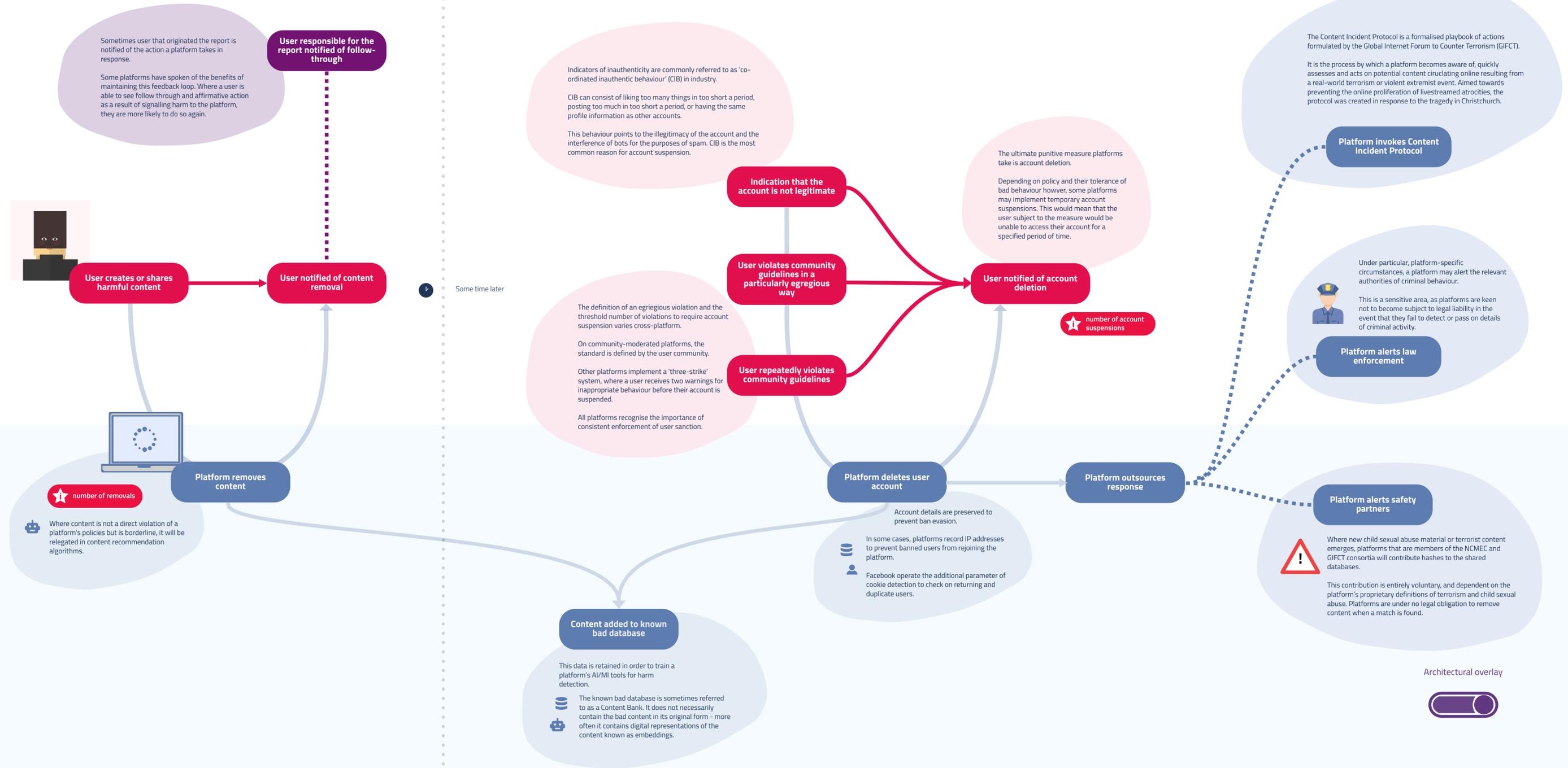
# Phase 4 Respond

The 'Respond' model provides detail on how platforms respond to the finding of harmful content. It covers the process of content removal, the application of user sanctions and the procedures in place to ensure that the relevant authorities are alerted.

### Key



- security
- machine learning
- database/persistence
- cloud
- external service
- search/analysis
- harmful content
- profile/registration
- image recognition
- management information capture



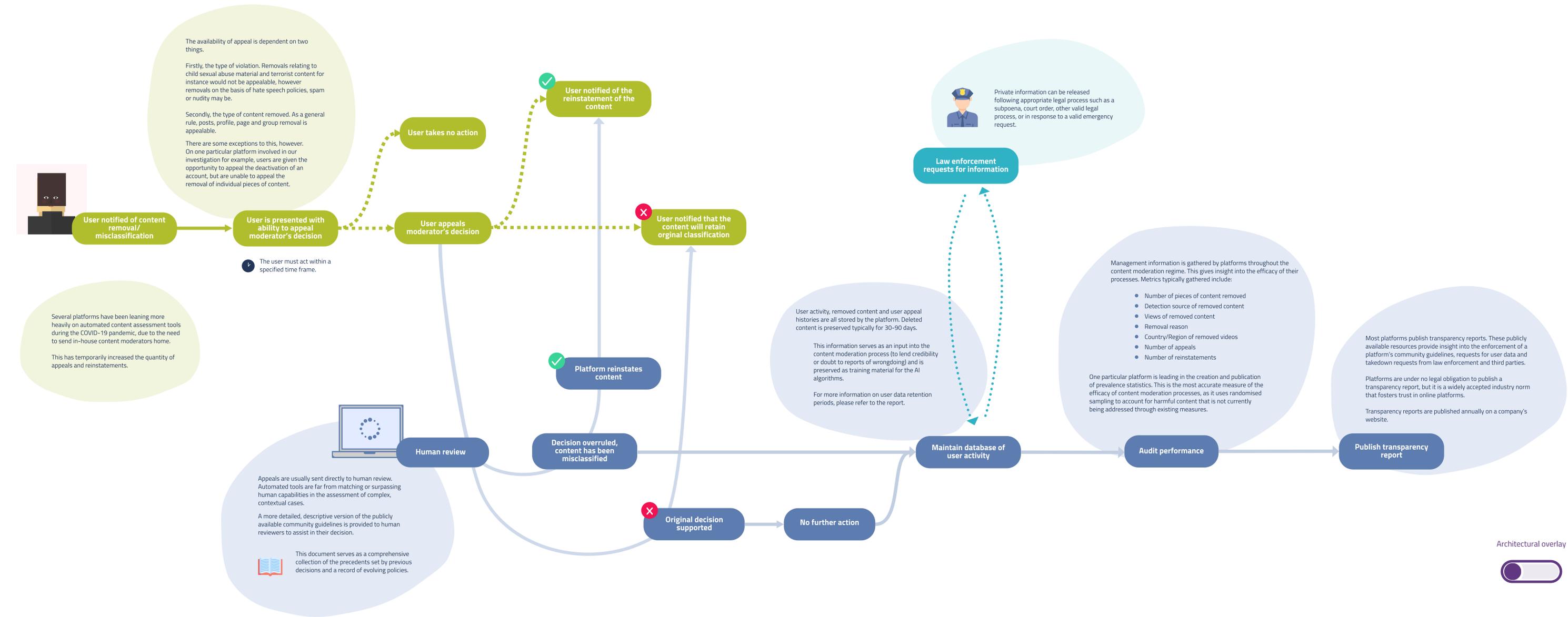
Phase 5  
**Comply**

The 'Comply' phase involves the processes the platform must undertake to fulfil its auditing and reporting responsibilities. It covers the management of the appeals process, the maintenance of a database of harmful content and activity, management information data capture and the publication of the transparency report.

**Key**

Step in user journey

Platform activities



Architectural overlay



Phase 5  
**Comply**

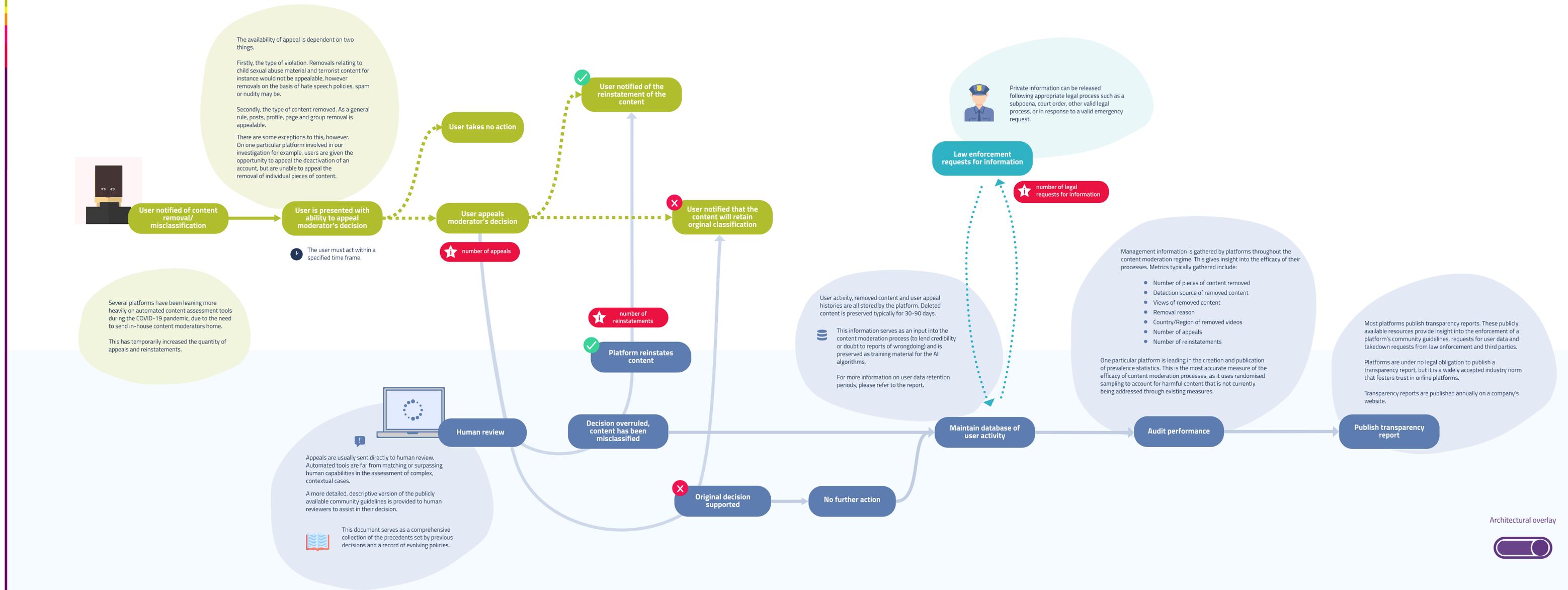
The 'Comply' phase involves the processes the platform must undertake to fulfil its auditing and reporting responsibilities. It covers the management of the appeals process, the maintenance of a database of harmful content and activity, management information data capture and the publication of the transparency report.

**Key**

Step in user journey

**Platform activities**

- security
- harmful content
- machine learning
- profile/registration
- database/persistence
- image recognition
- cloud
- management information capture
- external service
- search/analysis



Architectural overlay

