# Content Moderation

By :Atheer Alzhrani
Mais Alshahri

# Table of *contents*

**01**

What is Content Moderation

**02**

Content Moderation Solutions

**03**

How Does it Work

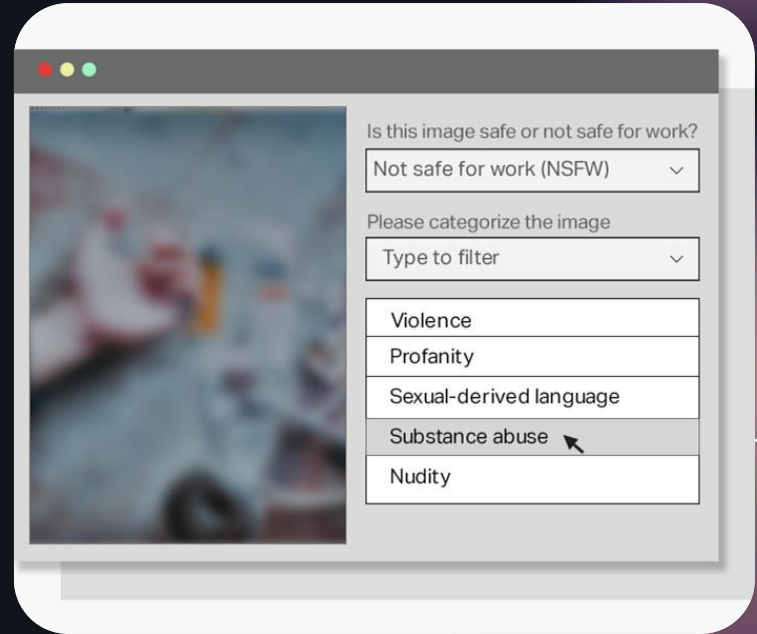**04**

Different types of content moderation

**05**

AWS Content Moderation Solution

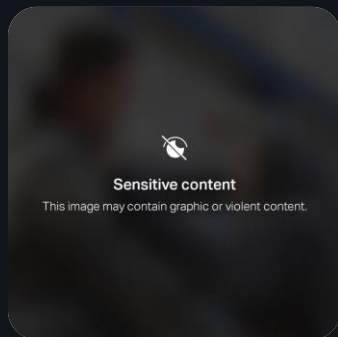**06**

Content moderation API demo

# What is Content Moderation?

- Content Moderation provide accurate monitoring of pictures, video, text and other multimedia content.

- Not only does Content Moderation help users to reduce adult, violence, terrorism, drugs and other illegal or inappropriate content, but can also minimize spam advertising and other user experience pain points.
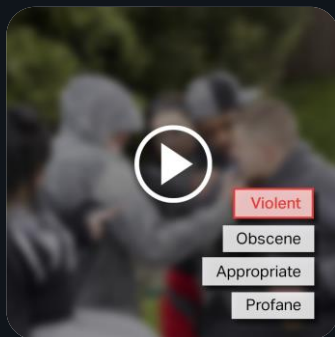
Is this image safe or not safe for work?
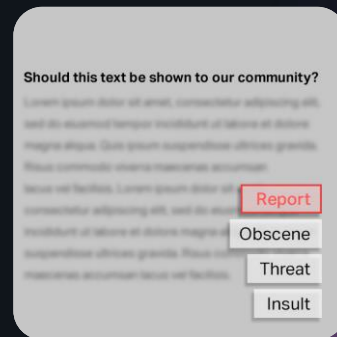
Not safe for work (NSFW) ∨

Please categorize the image

Type to filter ∨

Violence

Profanity

Sexual-derived language

Substance abuse

Nudity

# Content Moderation Solutions

## IMAGE MODERATION

Sensitive content
This image may contain graphic or violent content.

## VIDEO MODERATION

Violent
Obscene
Appropriate
Profane

## TEXT MODERATION

Should this text be shown to our community?

Report
Obscene
Threat
Insult

# CONTENT MODERATION
## CASE STUDIES

- **Client Profile:** Leading e-commerce site

- **Client data type:** Customer reviews

- **Challenge:** Content moderation and approval of user-generated content for the site

- **Outcome:** iMerit Content Moderation team reviewed all users being onboarded, with service level agreements, and the task was completed as required with all submissions moderated accurately.
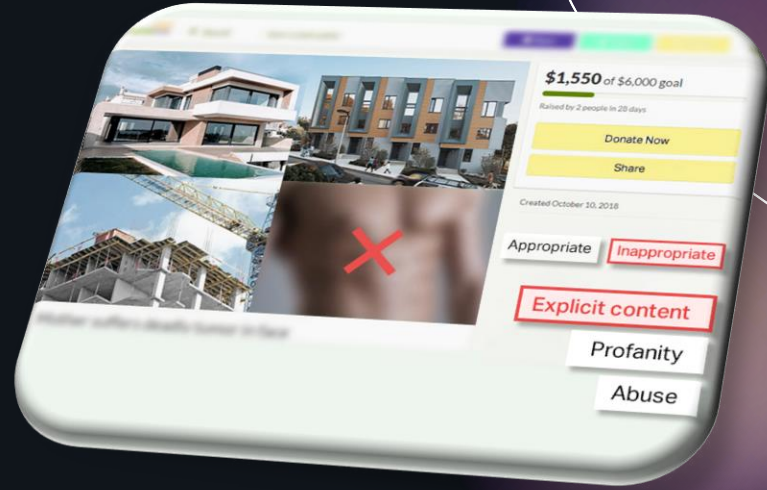
# CONTENT MODERATION
## CASE STUDIES



**Client Profile:** Donation record platform

**Client data type:** Campaign images submitted by users

**Challenge:** Moderation and approval of user-submitted campaign content for donation campaigns

**Outcome:** iMerit's image moderation team helped the client interpret subjective guidelines for disturbing and explicit imagery content and flag inappropriate material not adhering to the guidelines, along with actual abuse.

# COMMON CONTENT MDERATTION METHODS

- **Automated:**

  Automated moderation uses artificial intelligence (AI) and machine learning algorithms to review content.

- **Pre moderation:**

  Pre-moderation involves reviewing content before it becomes publicly visible on a platform.

- **Post moderation**

  Allows content to be published immediately but is reviewed after it goes live.

- **Reactive moderation:**

  Reactive moderation relies on user reports to identify and remove inappropriate content.

- **Distributed moderation:**

  Distributed moderation involves users in the moderation process, allowing them to vote on content to determine its appropriateness.

# Content Moderation Demo

- IMAGE MODERATION DEMO

- Text Moderation Demo

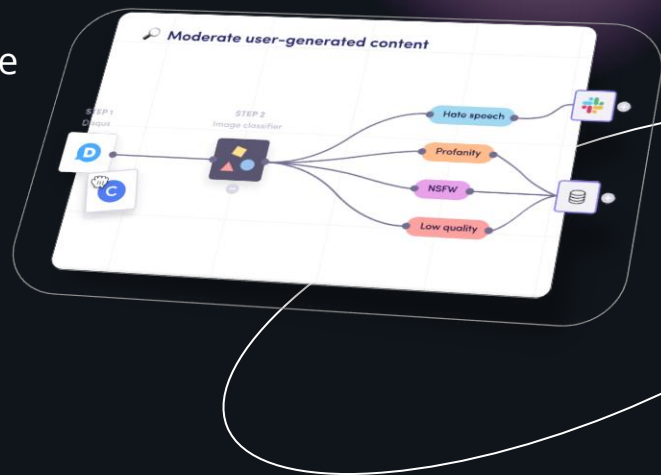# AWS Content Moderation solution

# Conclusion

- Content moderation is important for a safe online environment, but there is no one-size-fits-all solution.

- The best strategy depends on the platform, user base, and needs, and understanding different methods helps platform owners make informed decisions.

# Resources

- https://imerit.net/content-moderation/

- https://www.clarifai.com/solutions/content-moderation

# Thanks!