

Kaltura's AI-powered moderation service

Last Modified on 06/07/2026 2:22 am IDT

 This article is designated for administrators.

About

Kaltura's AI-powered moderation service helps you review video content more easily by automatically checking it for potentially inappropriate or sensitive material.

Instead of watching every video from start to finish, you get a moderation report that shows what might need attention and where it appears in the video.

The moderation service is available through [Kaltura Publishing Workflow](#).

How moderation works

The moderation service reviews video content automatically and highlights parts that may need attention. It checks both what appears on screen (visual moderation) and what is being said in the video (verbal moderation).

Visual moderation uses computer vision to analyze images and scenes in the video, while **verbal moderation** evaluates the video transcript using an LLM.

Content is evaluated against moderation policies, either Kaltura's default policies or your organization's customized guidelines. Visual and verbal moderation are evaluated separately, and each requires its own policy.




Creating customized policies requires Professional Services hours.

Moderation policy examples

When configuring the **Run content moderation** action, you can click **Show policy preview** to see what the selected policy checks for.

Run content moderation

The moderation agent will evaluate whether the media complies with the selected policy.

 Media will not be published if it fails the content moderation regardless of publish settings

Profile 

Select policy

[Show policy preview](#) 

Cancel

Below is a preview of a default **Corporate Content Integrity & Compliance** policy.

Run content moderation

Corporate Content Integrity & Compliance

Ensures content complies with corporate guidelines and prevents the upload of inappropriate material.

#	Rule	Weight
#1	Hate Speech & Discrimination Prohibits content that promotes violence, hatred, or discrimination based on race, gender, religion, nationality, disability, or sexual orientation.	12.5%
#2	Explicit & Sexual Content Restricts sexually explicit, pornographic, or suggestive content, including nudity or inappropriate depictions of individuals.	12.5%
#3	Violence & Gore Prevents the upload of content containing graphic violence, self-harm, or extreme cruelty.	12.5%
#4	Profanity & Inappropriate Language Filters and flags content containing excessive profanity, offensive language, or derogatory remarks.	12.5%
#5	Harassment & Cyberbullying Blocks content that includes threats, personal attacks, or any form of intimidation toward individuals or groups.	12.5%
#6	Illegal Activities & Dangerous Behavior	12.5%

Cancel

Below is a preview of a default **Corporate Visual Content Integrity & Compliance** policy.

Run content moderation

Corporate Visual Content Integrity & Compliance
Ensures content complies with corporate guidelines and prevents inappropriate visual content.

#	Rule	Weight
#1	Explicit Nudity Frames depicting sexual activity or full nudity.	14.3%
#2	Violence Scenes containing physical harm, injury, or gore.	14.3%
#3	Hate Symbols Flags logos, gestures, banners linked to hate groups.	14.3%
#4	Drugs Depictions of drug use or drug paraphernalia.	14.3%
#5	Self-harm / Suicide Visuals suggesting self-harm behavior.	14.3%
#6	Graphic Medical Content Gore shown for educational purposes should be tagged.	14.3%
#7	Non-Explicit Nudity Non-explicit exposure of intimate parts, kissing, and implied sexual acts	14.3%

Cancel Done

Reviewing results

After a video is reviewed, the service generates a moderation report.

The report includes:

- Up to 10 findings per policy rule
- Timestamps showing where each issue starts in the video
- Results for visual moderation, verbal moderation, or both

This helps reviewers jump directly to the relevant parts of the video.

Based on the moderation report, you can decide whether to approve or block the content manually or apply an automated workflow if one is configured.



- Moderation results are generated automatically by AI and may contain errors. The service highlights parts of the video that may need review, but a human reviewer always makes the final decision.
- Responsibility for approving or blocking content rests with the licensee.

How to use the moderation service

The moderation service is run through [Kaltura Publishing Workflow](#).

You can:

- add '[Run content moderation](#)' as a workflow action
- choose the relevant moderation policy
- review results after the workflow runs

For step-by-step instructions, see [Create a workflow](#).
