# Moderation

Learn how to build moderation into your AI applications.

## Overview

The moderations endpoint is a tool you can use to check whether text is potentially harmful. Developers can use it to identify content that might be harmful and take action, for instance by filtering it.

The models classifies the following categories:

| CATEGORY | DESCRIPTION |
|---|---|
| hate | Content that expresses, incites, or promotes hate based on race, gender, ethnicity, religion, nationality, sexual orientation, disability status, or caste. Hateful content aimed at non-protected groups (e.g., chess players) is harassment. |
| hate/threatening | Hateful content that also includes violence or serious harm towards the targeted group based on race, gender, ethnicity, religion, nationality, sexual orientation, disability status, or caste. |
| harassment | Content that expresses, incites, or promotes harassing language towards any target. |
| harassment/threatening | Harassment content that also includes violence or serious harm towards any target. |
| self-harm | Content that promotes, encourages, or depicts acts of self-harm, such as suicide, cutting, and eating disorders. |
| self-harm/intent | Content where the speaker expresses that they are engaging or intend to engage in acts of self-harm, such as suicide, cutting, and eating disorders. |
| self-harm/instructions | Content that encourages performing acts of self-harm, such as suicide, cutting, and eating disorders, or that gives instructions or advice on how to commit such acts. |
| sexual | Content meant to arouse sexual excitement, such as the description of sexual activity, or that promotes sexual services (excluding sex education and wellness). |
| sexual/minors | Sexual content that includes an individual who is under 18 years old. |

| CATEGORY | DESCRIPTION |
|---|---|
| violence | Content that depicts death, violence, or physical injury. |
| violence/graphic | Content that depicts death, violence, or physical injury in graphic detail. |

The moderation endpoint is free to use for most developers. For higher accuracy, try splitting long pieces of text into smaller chunks each less than 2,000 characters.

ⓘ  We are continuously working to improve the accuracy of our classifier. Our support for non-English languages is currently limited.

## Quickstart

To obtain a classification for a piece of text, make a request to the moderation endpoint as demonstrated in the following code snippets:

```curl
Example: Getting moderations                                    curl ⌄  ⧉

1  curl https://api.openai.com/v1/moderations \
2    -X POST \
3    -H "Content-Type: application/json" \
4    -H "Authorization: Bearer $OPENAI_API_KEY" \
5    -d '{"input": "Sample text goes here"}'
```

Below is an example output of the endpoint. It returns the following fields:

- `flagged` : Set to `true` if the model classifies the content as potentially harmful, `false` otherwise.

- `categories` : Contains a dictionary of per-category violation flags. For each category, the value is `true` if the model flags the corresponding category as violated, `false` otherwise.

- `category_scores` : Contains a dictionary of per-category raw scores output by the model, denoting the model's confidence that the input violates the OpenAI's policy for the category. The value is between 0 and 1, where higher values denote higher confidence. The scores should not be interpreted as probabilities.

```
1  {
2      "id": "modr-XXXXX",
3      "model": "text-moderation-007",
4      "results": [
```

```
 5          {
 6              "flagged": true,
 7              "categories": {
 8                  "sexual": false,
 9                  "hate": false,
10                  "harassment": false,
11                  "self-harm": false,
12                  "sexual/minors": false,
13                  "hate/threatening": false,
14                  "violence/graphic": false,
15                  "self-harm/intent": false,
16                  "self-harm/instructions": false,
17                  "harassment/threatening": true,
18                  "violence": true
19              },
20              "category_scores": {
21                  "sexual": 1.2282071e-6,
22                  "hate": 0.010696256,
23                  "harassment": 0.29842457,
24                  "self-harm": 1.5236925e-8,
25                  "sexual/minors": 5.7246268e-8,
26                  "hate/threatening": 0.0060676364,
27                  "violence/graphic": 4.435014e-6,
28                  "self-harm/intent": 8.098441e-10,
29                  "self-harm/instructions": 2.8498655e-11,
30                  "harassment/threatening": 0.63055265,
31                  "violence": 0.99011886
32              }
33          }
34      ]
35 }
```

ⓘ We plan to continuously upgrade the moderation endpoint's underlying model. Therefore, custom policies that rely on `category_scores` may need recalibration over time.