# Microsoft Cognitive Services: Content Moderator

## UNDERSTANDING THE CORE ASPECTS OF THE CONTENT MODERATOR APIS

**Eduardo Freitas**
DATA CAPTURE SPECIALIST

https://edfreitas.me

# Overview

What Are Microsoft Cognitive Services?

Why Moderate Content?

Ways to Moderate Content

Content Moderator APIs

Accessing the Image, Text and Video Moderation APIs

# What Are Microsoft Cognitive Services?

**PAAS**

**Set of APIs that perform specific AI features**

**Hosted on Microsoft Azure**

**Enable AI application development**
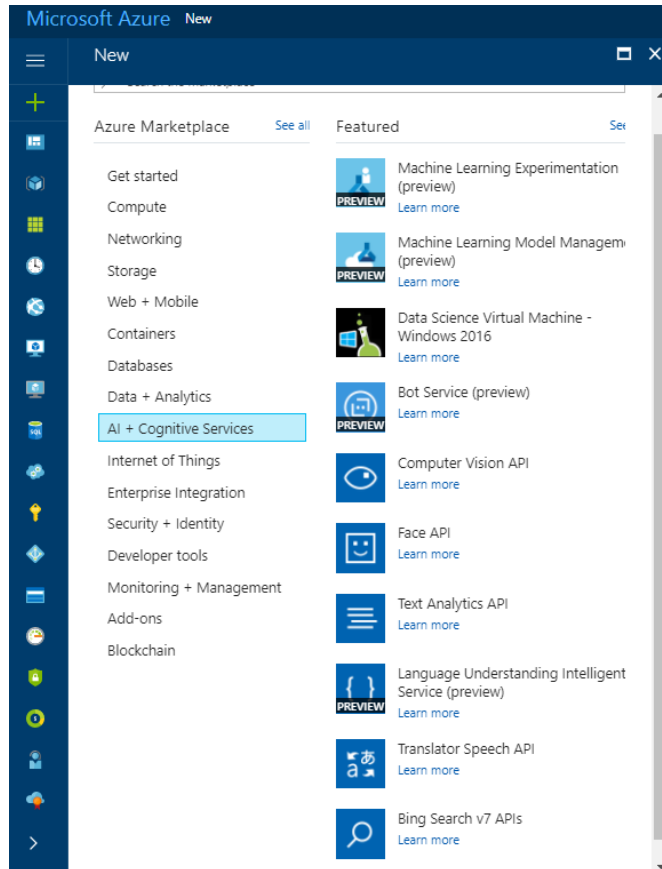
# Core AI Features of Cognitive Services

# Microsoft Azure Cognitive Services: Computer Vision API

By Eduardo Freitas

https://www.pluralsight.com/courses/microsoft-azure-cognitive-services-computer-vision-api

# Getting Started with Azure



Azure subscription

API endpoint and subscription key

Send and receive JSON

SDKs for various programming languages

# Why Moderate Content?

Flag and filter out unwanted content that creates risk.

# Content Moderation Verticals

## Online

Content generated from messaging, social media and online platforms

## Enterprise

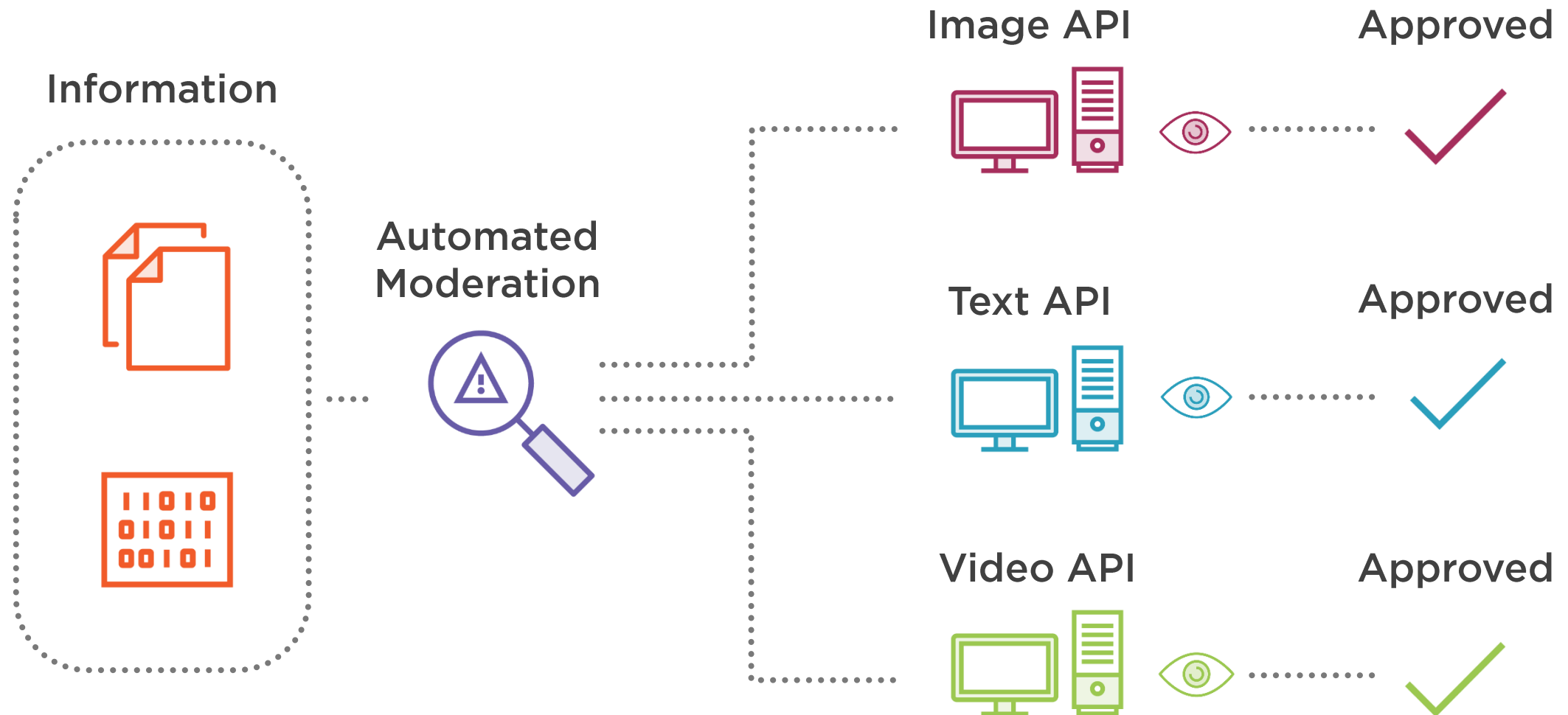Content generated from enterprise systems and platforms

## Peer-to-Peer

Content generated from peer communications or gaming platforms

# Ways to Moderate Content

# Three-way Automation

**Information**

**Automated Moderation**

**Image API** — **Approved**

**Text API** — **Approved**

**Video API** — **Approved**

Is it possible to achieve end-to-end content moderation automation?

# Hybrid Moderation

**Automated Content Moderation**

>=90%

**Human Verification**

<=10%

# Computer + Human-in-the-loop

**Computer**

Does most of the work

**Human**

Does the final check

Content moderation is specific to each organization.

# Content Moderator APIs

# APIs Overview

## Human-in-the-loop

## Review API: Jobs, Reviews, Workflows

| Image API | Text API | Video API |
|---|---|---|
| Adult<br>Racy<br>OCR | Profanity<br>Adult<br>Racy<br>Offensive<br>Malware | Adult<br>Racy |

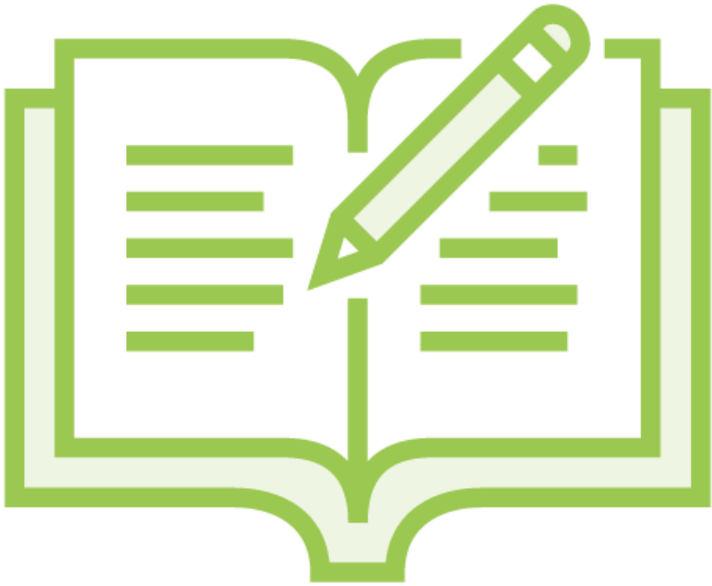# Accessing the API

Subscribe to Azure Content Moderator

Obtain the subscription keys

Choose a programming language

Invoke the API using a programming language

# Basic API Requirements

Images must have a minimum of 128 pixels

Images must not be larger than 4MB

Maximum of 1024 chars of extracted text

```
--curl -v -X POST
"https://[location].api.cognitive.microsoft.com/contentmode
rator/moderate/v1.0/ProcessImage/Evaluate?CacheImage={boole
an}"

-H "Content-Type: application/json"

-H "Ocp-Apim-Subscription-Key: {subscription key}"

--data-ascii "{body}"
```

# cURL Request Example

- location

- Ocp-Apim-Subscription-Key

- CacheImage (optional)

- {body}

**Content-Type**
image/gif
image/jpeg
image/png
image/bmp
application/json

# In C#

```csharp
public static async void MakeRequest()
{
    var client = new HttpClient();

    var queryString = HttpUtility.ParseQueryString(string.Empty);
    client.DefaultRequestHeaders.Add("Ocp-Apim-Subscription-Key", "{subscription key}");
    queryString["CacheImage"] = "{boolean}";

    var uri = "https://.../ProcessImage/Evaluate?" + queryString;
    HttpResponseMessage response;
    byte[] byteData = Encoding.UTF8.GetBytes("{body}");

    using (var content = new ByteArrayContent(byteData))
    {
        content.Headers.ContentType = new MediaTypeHeaderValue("image/png");
        response = await client.PostAsync(uri, content);
    }
}
```

```json
{

    "AdultClassificationScore": x.xxx,

    "IsImageAdultClassified": <Bool>,

    "RacyClassificationScore": x.xxx,

    "IsImageRacyClassified": <Bool>,

    "AdvancedInfo": [],

    "Result": false,

    "Status": {

        "Code": 3000,

        "Description": "OK",

        "Exception": null

    },

    "TrackingId": "<Request Tracking Id>"

}
```

◄ Adult Classification Score

◄ True, if image contains Adult Content

◄ Racy Classification Score

◄ True, if image contains Racy Content

◄ Status, which contains a Description and Exception (if applicable)

◄ Tracking ID

# Why use a library?

Demo

Accessing the Image, Text and Video Moderation APIs

# Summary

## Microsoft Cognitive Services

- Subscribe to Azure
- Accessing the Image, Text, Video Moderation APIs

## Fundamentals

- Automated moderation
- Hybrid moderation
- Human verification