# Infosys Responsible AI Toolkit – Safety tenet

## API usage Instructions

## Contents

# Introduction

AI models can sometimes generate harmful content, such as profanity, toxic language, explicit images, or sexually suggestive text. To ensure safe and responsible AI, it's crucial to implement measures that filter and prevent the generation of such content. This involves using advanced techniques to detect and mitigate harmful outputs, protecting users from exposure to inappropriate material, and maintaining a positive and inclusive online environment.

Once API swagger page is populated as per instructions given in the github repository Readme file, click on 'try it out' to use required endpoints. Details of endpoints associated with Safety tenet are outlined below.

# Analyze

**Endpoint –** /api/v1/safety/profanity/analyze
Using this API, we can check if the text contains any profane words or not and we can get the toxicity score for the same.

**Input :**

Replace the 'inputText' with the prompt we want to check for any profane words.



**Response :**

Server response

| Code | Details |
|------|---------|
| 200 | Response body |

```json
{
  "profanity": [
    {
      "profaneWord": "dummy",
      "beginOffset": 9,
      "endOffset": 14
    }
  ],
  "profanityScoreList": [
    {
      "metricName": "toxicity",
      "metricScore": 0.986
    },
    {
      "metricName": "severe_toxicity",
      "metricScore": 0
    },
    {
      "metricName": "obscene",
      "metricScore": 0.002
    },
    {
      "metricName": "threat",
      "metricScore": 0
    },
```

```json
    },
    {
      "metricName": "insult",
      "metricScore": 0.979
    },
    {
      "metricName": "identity_attack",
      "metricScore": 0.001
    },
    {
      "metricName": "sexual_explicit",
      "metricScore": 0
    }
  ]
}
```

# Censor

**Endpoint –** /api/v1/safety/profanity/censor
Using this API, we can censor any profane words identified in the text.

**Input :**

Replace the 'inputText' with the prompt we want to censor for any profane words.

POST /api/v1/safety/profanity/censor  Censor

Parameters

No parameters

Request body required

```
{
  "inputText": "You are a dummy",
  "user": "string",
  "lotNumber": "string"
}
```

Execute

**Response :**

Server response

| Code | Details |
| --- | --- |
| 200 | Response body |

```
{
    "outputText": "You are a ****"
}
```

# Image Analyze

**Endpoint –** /api/v1/safety/profanity/imageanalyze
Using this API, we can check if the image falls under any of the following labels – drawings, hentai, porn, neutral or sexy.

**Input :**
Upload the image to be analyzed.

**POST** /api/v1/safety/profanity/imageanalyze  Imageanalyze

Parameters

No parameters

Request body required

portfolio

| portfolio |

☑ Send empty value

account

| account |

☑ Send empty value

image * required
string($binary)

Choose File  profane.jpg

Execute

**Response :**

Server response

| Code | Details |
| --- | --- |
| 200 | Response body |

```
{
  "analyze": {
    "drawings": 0.002193321008235216,
    "hentai": 0.004320410545915365,
    "neutral": 0.9761048555374146,
    "porn": 0.01596055552363957,
    "sexy": 0.0014208250213414443
  },
  "ORIGINAL": "iVBORw0KGgoAAAANSUhEUgAAAKMAAACUCAIAAABwcQ7CAACK00lEQVR4nMz92bMlaXIfiLn7t8Rytrv1vblUZmWtXb0CDaBBLISBwlCaEYdjGuOj9CSTTDZG0LveZfoL9KIHmUz5g2S5ZjRGG41mSJBDYkgOBiDY6K26urasqlzveu7ZYvkWd9dDnHsrqzqruqsBA/nZybCTccGNEXH+fb78/Oce+J//+dvwBQMRP/fGAIIoGWe+cP3piw4hpAlQFEgRQIOgKFgAUKTnvqYAMB22y472NzSQ4epikREAiAlASRGMwcwREa1FzRxjb1C99wwqioB0QI6ZYGcIIfKAgIDgBAwYiSbQGAZ2xUV+QcPjc..."
  YvsKXAT79tv2Fd/Yz8tYX7P8Fv4QIQAAICoIASvDvkCh/2SEICKAAQNt7oKrKogBgjLWWwCtnVEBVQhJRVRZAIiqKwlc+9WvYHgMAiJT+hm+D/5W/9zm5/vxy/6KhQAACQFsxXw16ftYMq1kFAADNL3lK/xYHKhAiAaoqsVqDjARIKuzJKhkFigKSM7NGQSFEIAIZZjwAIYC5OtTfzPgySf+8OBER8NM1+Yu1NxjB4aIIAbZ/ifC5v9pe7bBHQf8dXvQ0acYQg85Yp8IgThwAGFUJUUMCo0TWkyELiTWlFCSPCsdIV3Il1102MJoWvpOSuTuArD2u+5HeuDaf+imtaEBQ+u0ZRh6t9/s8IZPv97VO+qVP/Gxo4n8Rdnzkpa KDpA ICkbA05Gc4sKSGqc9Yb24deUmBMYJ33hfM2WcxM/OkiIVIg3R5YAX55yzusAfqlv//8nPhqaxoGpwy/7As/dwQSvD45QUBAJQUEBBRQApTrm2gA+N8pMO9mGGCGAAAVUurJOkIkUQV0RE4RYt6rximlLseQg2h5cwSsscgCAK5IqAAKRrfW5kB/+TV67SX8 8uPaNfikHtn2gl+4/8VDkBRACLsUAIAUcBCSXS/q7fh3WW/Dd4K3BGAR2MASZSArC8lNL7ts17De0aNraw1ql2KIffgvboC0AogABEQgKCKUYCVIua/4vhqHtm1Xv31PbLBrgxrerg2UgEUsSWqDJ7a8D9UIv2bu/hfbRAagJRFgSIq25AUziLzumsuTk8 unhzf2NvfP9wdV44AkUVR0ZsmiaIDtdvSDmJUBQW/YiiCX8WPed6o218c9FxFFdu3wxvVQcbPv9meynOTgBSartnfv9FuVgAyqsu22wOnsixJIYRQ1JUKLlbroqrq8WS5XFrnXuh861cJOuGXmIK/7HGIAGQbW2znIwzi7rqOUMu64pSOnz599Mlhi4vST7///w+988xt/90/9IZjcL880X3qpR3k2gIwARi00DZ0MLVEJEX7Nyuh1/2ehXIhwmf+smv6i8Yv8tR0NhmRJFIGybFlwwk8p0+Mc7HrCKEejUdVnQFT5AgvuOZ/mwPli/wlVVaVt4A5x2azbrsNEY3iH49t3bn7w3s92pvaN1+7v3v3v3/7v3v3v"

# Generate Image

**Endpoint – /api/v1/safety/profanity/imageGenerate**
Using this API, we can generate images based on the prompt and can check under which
label(drawings, hentai, porn, neutral or sexy) they would fall.

**Input :**
Enter the prompt to generate the image

| POST | /api/v1/safety/profanity/imageGenerate | Imageanalyze |
|------|----------------------------------------|--------------|

## Parameters

No parameters

### Request body required

portfolio

> portfolio

☑ Send empty value

account

> account

☑ Send empty value

prompt * required
string

> Generate a beach bikini image

**Execute**

## Response :

Server response

| Code | Details |
|------|---------|
| 200 | |

Response body

{
  "analyze": {
    "drawings": 0.0001658086485593677,
    "hentai": 0.00036795716327404,
    "neutral": 0.014667540788650513,
    "porn": 0.015493784099817276,
    "sexy": 0.969304919242859
  },
  "ORIGINAL": "iVBORw0KGgoAAAANSUhEUgAAAgAAAAIACAIAAAB7Gk0tAAEAAElEQVR4nFT9SZBtWZYdhq2197n3vsbdfxM/uozIpiqrQaE5VQAIwGg5SQgQjUZKJomSgUOQMpPJNJBME8010YTimKaBJpJJQ85gmgpGACZKFIBCV4WqykJlH5Hx43f+3f01996z99Lgn0
    gYCQlCSJTcHv8YAhGCR8ggIzth4//RKp9b4CpBNl+S7avQAqEIKqtCIwUsv05JEASygTYP+JxiSSatQtify+xXXRbA0ig5BJSkkYp5aRAMkIgARnUvwggyQ1IpUhjSk6JApgBmZmSgNMT8nYZ7V6QBBIyEkqQmQINSAhG67eO0DdrRABmgvoi5jCKQv/OgmhQGvs2AAS1Ows1
    tj+rGX7Yfmt//7/7Cf//B9999Q//yecvrg/V353q6/ul8uaLn/34v/f73yq7svV1ovuoGuGQluXrL9/Cn+uf/ejlv/nl4acPtXzg9v4Rmvt78+fuH9/fr+X6+ee7f++D6+y9u/vRnr87D8Pufbv/j/+tFHwri3RuW/Fd/9OVmd12P53P91pDf1pm7/7b/7/2a2Tm1/7/vmFX301fP46b+
    bEZDaHg5Z0zYEYCE4E4KZoe08ICW7BJkUCLYj2A9v22NSixBtu7erU8KNVErtELUz294TBElZjobkJay1l7fT146FEwrbmeprmEl6RAA0SgEHRWOLBZfT088s29/DgOhHR6T1pZdEYW1cVkcOJtTCo5BeTmMkrL+Yta0mhZ5Z5RCoIDkRckKCwYlhbIEG0UGpYIAyANB52cL9
    bLF08bWVxKEwLZyUIvYm2KSbdFJgMoWudlwQl2io2/39vkpth3avyEv+a0Fd0nZU6Qyr+El/QFKSdlia18lRKi/IHusDpHeFicEUEyEgGzn16BAploguCtedrLawc+Mvt/7W1+SsXjSEF6+ElNG8A3UtIVopPUgIntH/lAkaATE9t052U+y2oqzf54o9tvRjiSagENQ2XBAT2
    80cdHDc+uNztyIAnURRs9++6znez5b3/6e8yT70c/ffX//10f3Nv+ul98cy054YJ7Dm/eHv7Kt55wmFf7P/y1+f7V6DNf//D9L78+EvY/SxbvDZ8fPFrmznv2KZVApvlhsbpm+wdZqCP4wNTUg/yI+f7V6N0f//D9L78+EvY/SxbvDZ8fPFrmzVKZVApvlhsbpm+wdQwCKUjXqMX5l2rHT5eTqRst
    dagZo8I1rLLN2hwDXaoAUhISiOU7USBym3SSKQQALISSOXEUNmzN/r3CHa8g4b5sp/89t0KwEgIAwWLRrjEqQSswcd2AC/nESZmiKl2zDsqznZrpKjmhg0hJ86prG0ZaZTBpSQ8E/fnTcm0aanQTQYZAGdISN00h2fNsPF5oMGICdg0yKLRFTuFz90e63NcUjGGaPQq2KUI/ElS
    u+yezbUN7XAC6ZO8+oNrfU4B5yxcI0NvVt83/oogTBrgPGYLFJmtPiKIXhugf3qrfNqNBbPlo+xLJgC09nrK+ir2y2wlTUtjPF+I7X5ySeYLoP+VqqM4JOOJSZ7YNFElFFL/ZenLL3CijlRcfD9flL/1Pf/mT//LTzfJw5maaIuZlj1/OdhqTwzJGNDs6f//+/e3zj06fDPs
    +ef/vRnn3226eb50/VB2ZI7H893DQbErw2A+jbHWNY9fj++XoU71+nuY/oKGHQTRCZhdMPolCqD62TGCRKphIqORdeBHZgbGyqWkEomzfmY959xvgPoDvKQlwq17czHQHM5fn0nt55sbPe+bzII8va3JtISaum7VR9IwADKRV10ZbSEgLD+rTIb5nrESUAQ0YM82x4CUj2Kdh
    6vpzvDisbxkxF+8r9HcwAwuCREZDJYEglqJAyUgkTrBidJXygdfDoFImgnAkzteVst+VCC2QL+t/UTAJBU0J6jOTtBreVyIbce8S/wI1eTQCgf4MM95t9dEI1aDZEI+FS9XwTLEjATOqlMRq9ggsqx6X6U69m+4bL55g5OshJwXjhlWBopBN5SesQKctvPrU1IFLKy+Sq0et
    pekQIDJfgyUoHVoxX4k2ze8FIyNthOtvZxoh//CV8EMmQmZiR5pyJGD/eZF+fGLq8Ppdu8407GW1hv++P060PbT+qHFAHfYayVRj+vLZ+XtX/ng36l8mN+/f4iH2yeL/urV//JbP/jjf3+3055nz27//PZZzfDfjsr1Pv4fZXSSvuVuk0JfT3Wri8nd/cJwAUEq0qdKRNhwmqjMfZTYS4u3ndT2z
    oGGmhFlwOyehHMXzIEL6UoAMAbmQLj4y8ez3IDKYmEWn1nHT2zbzcQRCgft3qXXOAOVCM/kwYj2VioxyNLs84FtevPx0N3KdzlpmwJMCmCftnbCShl7bh3m3ICFX27FXZeBLaPRMd4bQ826MXsPGyPZkwoeEEK2KVUzMes225V+xklM/Utj05j08+rGhPBbIkdgGiiSQEgYT08rSrG
    CNofLMIZIdA/RMFKYmLttMFFV+AYw88BshwiZCN/MrLHWjwocGetmwkyDIOLKOwBJAt+sNoVZK75NxhhkmXosHBTIIwmHqga3ddXTIwXpa4nb7Libf6vQXadv1tIazTk+1Y0R4lDvS1UE8sjW6llKHVe0oFS8EKWilv7GlfmIb+equ0itm4AF60XYr0tu3ZiXmBLYmBbas1ZZ
    78PnC8cnh6TAUPCzHlw/Hr94+je3nv/NXrn7z/Bs3X7+4B8Xa94y7+oX2xy2voyxqxjHVPNb1aivb70/j+HR3+8Py/Pfyfvpvw/DEE3s4HlnGQ3rHY98fxTj0+eDbzsn/7rebljeHCb16qziuq9WVPNb1a1zv70/j+HR3+8Py/Pfypw/DEE3s4HlnGQ3rHY89DnMVYfSl5
    Ebz2e2zJNkm0V4KVGaPe4Dw3rx0a7PZGxcmhCKzrRlwlDdrJNDZX1Q3VhnwzKZKInMFBUtgpX/ZJ8wPp7Kjsao134InRNhmStF2zktDEBLPLLSP66SDbCtIfHLcce+dfDQMufjxmhodYEQbO20dvJ6edP2aGbjb/XbAtDY0K0rjUmUnD3q+0n/60vnv2TVw9/+HsT59P87a1s
    aaGwnnRi+9980cv3/vV33ny609uXj8/fX16//WPf7l+96N1NnA7DHOpyUoUbULMOADPX6svbl86t3f6k7H/dN9vo3GVm+uW2GJnyhq3NmWmebPW44qxvA8M5SZEvgHeATaQ11ttNEY9/N1m+m9WNuLQ41yG50Cxneavxn31kv0MyQj+ogSMi+0ctNWYi2x5M0Wl4go1FiI37al0
    IwhoeM5GNwE/CkNkPSRIwtoVkNgMLZEYZMhu71dY95Na0Fgfbq4NamOOl3o6Lpte2QZoRsExKFemEOYNyqJi3YNyQKbq6CbB7YPRYgF7K9Y6nifZ59JBzAerN5GLdRtItPWhhICFvEie+2QINA4YuBbv/4pF96Wgf/f2AZmq4oHpe4kdXh/p36VFbIP1XBIhWuaGzN+0tlEA/"

# Video Analyze

**Endpoint –** /api/v1/safety/profanity/videosafety

Using this API, we can check under which labels the uploaded video belongs to – drawings, hentai, neutral, porn or sexy and can mask those profane objects in the video.

**Input :**
Upload the video file to be analyzed and masked



**Response :**



# Nude Image Analyze

**Endpoint –** /api/v1/safety/profanity/nudanalyze
Using this API, we can check if the uploaded image contains any nudity and can blur the same.

**Input :**
Upload the image to be analyzed and blurred.

**POST** /api/v1/safety/profanity/nudanalyze Nudanalyze

**Parameters**

No parameters

**Request body** required

portfolio

| portfolio |

☑ Send empty value

account

| account |

☑ Send empty value

**image** * required
string($binary)

Choose File | profane.jpg

**Execute**

**Response :**

Server response

| Code | Details |
|------|---------|
| 200 | |

Response body

{
    "blurredImage": "iVBORw0KGgoAAAANSUhEUgAAAKMAAACUCAIAAABwcQ7CAACFNElEQVR4nMz92bPMmXXIfiLn7WSLiW++WNSeqrMpau3pBAw2gQSyEgcJQIsXhmMb4qHm5STY2etC73mXSPyCTmR5kMkkPkkkaaQw2Go1AkENiSA4IgAR6q66ufcnK
5e7FFstZ3F0P57s3b2vlVVc1YCCPFRUVGV98cSPCj+8/94P/xZ+9CV8wEPGpHQM1oq5fOeHq22ddQkgToCiQIoAaQVCwAKBI105TAC1X/bKr/Y0PVLh8W6QEAGICUFIEYIZ8zRERrUTPHO8hU7z2D1iKgU7sATlkg38ihsYDAACAEjBjJZnQA4DMb1s+5h6c
6w9c4GeD3J2fbmvtnP0fufcfzm/CVEAAJAU8AEUIJ/j0j5VYcgIIACoG7fgaoqiwKAMdZaAq+cUQFVCUIEVVkAiaiqKt/4NKOxhewOA1FL6G34N9iue9xRdP8/uXzQUCEAAaEvmy0HXZ63h27hUAAD6f8Zb+HQ5UIEQCVFViIQYZCZBU2JNVMgoUBSRnZo2CQo
hABF3mPAAhgLm81N/M+D3KfS6c1Aj4hCd/vvQGIIgeihBg+0uEp361fdpyREH/PWZ6KpYI3xKAz1qKwiAgHAEZVQtSQwCit9WYIQmINKQXJ48ox0iVdCXU7o0nhawmSyxv42sOaL/k7V4pTf06GeFgSFz/Ioanna6z8jk035210+0q3/DQ0sN0VVIt6wqACAgW
4NEhjNLSojqnPXGDmGQF8gTWOd95bxNFjMTP2ESIgXS7YUV4Ktr3sID9JXPvz4nvhSPQzHK8Mt0+NwVSPDq5gQBAZUUEBBQQAlQrl6iAe8/r8j8maFFDQEAKqQ0W4WESKIK6IicISS810xSSn2OIQfRp0LIWGORBQBIEVEBFIxutZNAfnUevbI5vvq4Mg2+
pkW2feBnHn/2ECQFENyyAgCQAhaS69Ws3o5/n+U2bOejgBIBeDQGkEQJyPqKQGPXtes1TObe2MZao9qnGPIA3qurAK0AAhABAQiqGAX40mT+K46vZ5FdydWvbpEVvVJ4ujwbqQCK2VJViqVW/oVKpH9zD/+LDQIAkLqqUEQlG4TK0WRe9+3Z8dHZg8c39vb
3D3cnjSNAZFFU9KZNouhA7Xa+gxhVQcGv6Yrg17fjrit1+/OdnkuvYrtbdlQLja/vbG/l2iQghbZv9/dvdJsVgIxHdddvgHNd16QQQqhGjQouVuuqaUaT6XXX5tm490/jWr+N0wleYg1/10kQAsVUZtvMRCrn7vifUetRwSo8fPrz/8YeLs/Of/sUPv/vtb/
3+f/B7YPKWPD18/vk86dHFGY33FBGIUIkgoyIAE4hBg9dDC5dDRJ5xcMsPX/V5FeFKZn5Vmv6i8fMsNZ1PxySJlEFy7NhwQgAZeuNc7HtcGI0n42aUAVNICM945n+XA+WL7CVW9o23gDnHdrPu+g0RT5aTO8/deu+dn+3M7GuV3Nu7sdODPc6+3qknEXEAz
IoIoIQIiMqqCvA1vEpS4F/0UX5BSj/Tv3qGVwagOXZ96xAIVVJoPBmkEHrnHWrK0UKj1tqUMjM75x5eMZc//7e+4h3+lUd5TgHVYk9dfeGcyzGKlJzRyWyac3x8dPTRxx+06BVydX529MjyBuL26Pzkxq26mhhLwsAqgqQoSqokIoL0bLH6zPn1V7FWvzal
n+1kf8EJCCIh5KHd3dnJKXTD2rkpSBrW58i5tsY4G2NMgogGyRqilMK/EzfrC2YGImzt3W2I71LqiIiAimYCaqp672B//3D/9G5KnPqVPTt9aHm8U8n+4c3hwtbGV34WAAiAyx8zIMoKOfB99fvBq/v4+s9lv7pz9szf5SPMT75VaJyRbBqLD490V4vTqX8
updQuzz/95KOXXvmGfFathsyKrqmRXEZp5/7fr3uHF7XzEQFh+wEABFRQEoQ+dKO6QW/7bnO6vKi9e/wN1959d7650f/8f/fv9nP9ofv+whm9gefXR6oHZ3emj8EqKgCgoDAyko1Cp+Tzfja/jT8CRw8YVzNHyOrk+dgIgGVONgMA/rSQfvvI3mcX37+dk
p0Yrj/wQfPP/eCs2QQHBprfGKJMXr378byfjYPKV7JF/lM18/QGNYMckCAxqgFsKaqmr35vd/+27/1TiXh4vHy+OHYAgx88fjR7Par4JzxXggUISsDZQLEZyurLw5RfMH9f5HFenWdX9wC+mqqWvLQG83I8vz+h++Hdj3yr1stjx7cn83HdeVUVVWNMCycs
jj3bFH2NzDwqw95RK2qKuTUx1A19Wxvx1V+3bWPTx5fLM9fe+01f/gf/QPgTMCxXe+MmvPHj/PQa2aLZIwB0KyZgcH8lVTv1xq/IXWfPePK/ETBEutDQVDvaNrUCCl0K9BEGh88+PjNn/74pVdemuzMk0pUzppj6oX78cQBKqqgCsCXfODGVvGpIwIgiPq1
Ps/8Q4rASIx2695AJs0ImUAyRyJqmsZZn1Ia+jgaje69+LIau2y72d7uaDT6pV/6JRI1rP1yIUNEzhbUAoIqiKAqgHzRdPrF6PIl4+dL76dV71Nfq25fsSqoKCii6pNIC9feOtR3fvrjw72ZN/KHP/yzf/JP//83f//vvvvCN1xJQzAkrJ6AE2TpMYWlo60f
Ipfmj28THE7YDAVXR8qcJlJRKGlEFUYnIGMMpX3sGgUs39Lp0/OyjPUOMMlLG5gGsJqPZgmCZW0A55KpqCG0IGUSbamxBIucr1erwIt2PFvaDdR9wy9a70vbDzZ39/el0Y+w6RmDyBEDWKiEXobYNUNy+0c/w+VPfwteiP0TLI/9V/elrN85IiKiIe0WGok
ht8Fz4QwWn1b5ecupPT47svvvDCK69kw1SYKERxe/+YCRAtbmPL29z2Vr6JiKqIIAAYJKIyCUzOWa/5JSKIWBgaPqMAf3EZyQgKRFAEbSZVVFCQ2jtWTZFFlBAl5cQsnJpR1fD8+9++CHBw2A09zK5ZrW0MIAEaQQuo1qBowaACX94oIhZyXuI1H888+LXGX
30KYqvIiH8LJEaJOz/+oDYIOV+cHH/w3gff+uZ3nn/uBVEjYB5MAuK2U3TpwDzhL70ikvfeGFO+JAWkY8tQf4b4hcImCqCoemITkf4icXUCQcjXDxoBBck4r6o5byVHzjnloJIn31zORm/y4QcvPP9cSimDdjGOdmf2QOZbs04QUUmWVtnP6wY7rlx/88c2f
f0zqukhE1dxtLo4eGsJjyq4XS+fca6+9juRY5dQoWEWnaAANKAFQFpBnjesPrK0ikI1KW3A1sZ75dv7KQwCUQAnkup0vACLCCKpojCGiYleOq6oGefje2yNLtw4OjLVisBXevX07Ifhlp7htpERfI6y4Rtq/lsf5a6P05003RDSAKEnChvu2cXazXLzx2us
3btzshixqslLoG2JQBLRQQBxAiAm0/wWmBQApXDMIZQQkBEVSLMOcrCw5VFESUrZklAyp8lpt/MZ5GyAiScDKACoAgcEJIURGdc2gNohpLlbPjyg9nJycfvv/5nVuGwFZ+YImIk8Mb0ZlsCycjgiEwpJbUGnhCzus7T9H+r8LQJR3+16ant4pTVVWRtvrGEF
QWIQUJPeQELN94701DDiALowIBkQJi8VqVWRXXN0wCrrQpWQQOErlxcAY1xVeViJAW5VVgfAFQNABjz14BSIlAEAdBLAQ6wtekos7jKkbHC2aNWhhwI9+3Re+9Q6G7vv8DDnlXWIbjZDtRNIkgEfPVs5gD26prPVNJ4LYfOJSr8y8dV7svCF7nuT87FJ1v4D
8bs+jDGiHjWl8QpEtUYhJS61cXefIYiB/v7Nw9utJuNq+Z9RgACJURCUATZEhBRcZulx+3rLtdHIlTVnGJKiUGttc45s5V/QIjFasctSO3Jc12GD+kXCAgSMIDQpYZVJAFSIHJknFXuNAQidt7w0C+0Hzx6960XD3YgBZY0xLyO6YVXn+sJ2BCTEVJQREEE
Q8Kw5VQpL/iSnE+Z3589+Iwh4p/zbH/N0lsvx9WRvu+Xy+WoaZxzd24eemc2y4VDMCokjCBG5Tq780UQZmbWzJIZckLOmnLqu83i4uz0+PT40CXZyXp5XjnjHXIrnClb651x31x/Ib9YxBcACmQAQVD1CklR7tPZyhjDnNPQY84VKg/t8uTR4vj+y3cOu9U
FAAwpDWzzm7eDQCLiLUry6joE+lkP9udZZF/37q9T1z7Tj4RLaVye69qWnvX3tsLzisDw5EZJjNkM6b2PPn799Tfyh5/88T/7Z6/+0q/dfewb69PVuJ5o7SILc0YSRGDOVVUlzsqQUih8XiCCqfHJ8qsyqt2c35+bm43ZLGuR8YbzvrCi88f3rglwjlG5z
wBxBCJbNd1ADCZTq31wzDEO8hjvK0QUU5YmYi894jI/OxkICqACGDx3wu0SEvEIAvH2N7Y2z/drGsCq3lC8Cc/+ouXDnYxD9P3+GSligrVdL5/+85xr4QF8X75ei+FKjMzflkq8ik/my6V21PHv2ga4CXU4evqafmifGrRkVttDVpeZcdnYk++8hu/9aknW
99/+NBXtt/0R/c//LOs3/z+76zzMHQRrENg4WxIR0OTYlQRZ22FBpUdmhRD6r5H9++LMDN3oeuGvgtd4gwA1tqPP3hn07brdVtV1UuvvPb66688f3ro1qr211plz5pTVm3MejUaz6azrhusmzxW54ksGldkLBgAUCdUgGk5cNKOTx492z55,Zy/wi,ghw
6CATcujbyHnVxf0XX0i1TWW4wA1QoCAkUa4AmV/095/5zU8pGa+rS/8aIyd69QYVtEjgmLPUo+/85t/+5NHp/YfHL919fuS5Xy0eHF354e1b9e7hdLI7SGxTEARrrCpbYyIzD1EBIOcQuV0vF6eni7OzzmDvdVi3i2A39Yr0YQgBDq4sferuvY5,ysyl.u,

💾 Download

# Nude Video Safety

**Endpoint –** /api/v1/safety/profanity/nudanalyze
Using this API, we can check if the uploaded video contains any nudity and can blur the same.

**Input :**

Upload the video to be analyzed and blurred.

**POST** `/api/v1/safety/profanity/nudvideosafety` Videoanalyzenud

## Parameters

No parameters

**Request body** required

**video** * required
string($binary)

[Choose File] MicrosoftTeams-video.mp4

**Execute**

## Response :

Server response

| Code | Details |
|------|---------|
| 200  | Response body |

{
  "nudanalyze": {},
  "BLURRED": "AAAAHGZ0eXBpc29tAAAACAGlzb21pc28ybXA0MQAAAAhmcmVlAWeMKG1kYXQAAAGzABAHAAAABthAAwhM2nHCcPAYlBiFIDfh4D58AOZDoAlFtL8PAfE/8woXITLaVSDaYGKgYElEeAZaB4iAPBhiDAthe2PAVlKGwd8GAQoJHiwuwch2HnwVy4Mcftg8JAI+
c1ieKhzqtKnbJI0LDQ48mDr/w788yAHR7AYrLNgPlQA7qnx39oPAfZ/+SghfES o4HoPk//ZccMfAsHwGA7BiZAHiYA8DZYBYGPw/U2P1TNoGAMJhFoEAfHgCT7Y/+q/z2/TtNIUBlKWMzlKNjz6VpMBHwGwYjULTLxew1rTRfPL/HAaQBZwYmpUwnBgJA3AYKKjLXX8NaWm6G
Tpo3vdW2L5pM1qQrESTu/HKQXjq/BjQPjQ8aOBWiD+79T7Q8zA7bpLoKgGJ1WgMJ2U1ETwKwGDQGJpC+W+LBz8PU3gfRgCVPfTpwYCweJAeL/+fg+dADtOkHJZ7/2wK/BEB3wZCLenTph4HgPlQA/X2uDIAeF/6xwDwkAeoWgQB4T/zBiddOk8DVpw8B5D/5B9H/7GPgPxoO9
LBapYD9gch7jQe0HfBKgEUGGWNh54PAY2D6UAOg2Dw8Aewgr/g+jADq9ld1r3mMQsgw28D53/2MqzWc/qssUDkCwFN8D80NBm217NnKuxaOIHi7CslLwMpy0cNArwKFgFgYElmsiBoFUWyFaa+KqDCK0yLFkulgMbBg6GyEEUHx4AcwHoEQeDgF/pwYEsYJh82+Zo28BkHx/A
AD75YDgVAPJQA4C1ko+Sh18FgxUDwV/mnBirwPEwB/gYjyYCAPDf+JYmBiUFSD50ADOH4esy2lxsCyTufA+Z/50PGh6DIQeK/8wYElgbB4L/j+DxMAeD7H/yR/ocAK0tdWx5rw52lkNwDEB5b/58jhz/8z/GGP+AwDIg8aB8j/5VTpm2weFgE5xoHhf/MCOXKNLvVZsDdDZg+
gDfwPgxUD4CAehZaiRMDLgifAwddoKRA0Dg8BX86LwWA3EwMHZYOAYmHLQKUsESNg+X/9mZ9tD5YmBjpp0nqgPAM7Rtg2Pmi8SalD3reYD8x8YUVAfV6DgU6oqDmlvgfP/+4fz7Z35cPUoPCQB6ZPMw8CmcLGmWmmy0Dgh2OAfR/+xLL/i1MOAZBwG8aHB6wB4DxP/inAh8A
HyoAd/58//3mxuBdXPAWdwVzYPkQA9tmgYra+rbb82C5LPbv/q8+gLAd4GPi4PByDFfiwDYMAifqgeE/84DFPwomWg28cfpwZCej6aWAb8BWL4aHIPnf/ato54b+Ag5JwckAt4DFgfR/+8JcWYTJ0/gYoB8qAHjFf5uFob4DCsdgxK0m+DSv/2NEgh4y2Ct7GwY0wAYnZq5Nc>T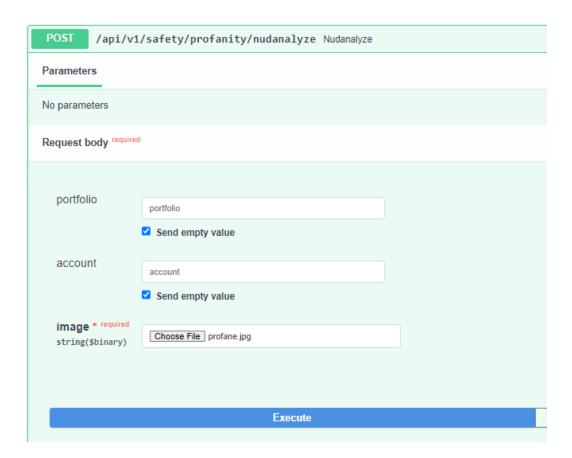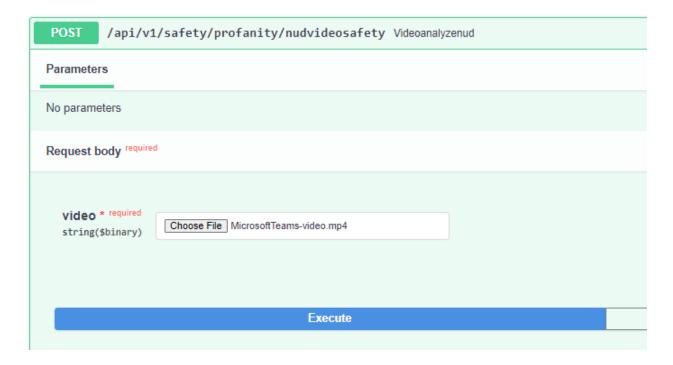