

[Cloud AutoML Vision](#)

# Making batch predictions

After you have created (trained) a model, you can make an asynchronous prediction request for a batch of images using the `batchPredict`

(<https://cloud.google.com/automl/docs/reference/rest/v1/projects.locations.models/batchPredict>)

method. The `batchPredict` method applies labels to your image based on the primary object of the image that your model predicts.

Batch prediction often offers a lower cost per inference and higher throughput than synchronous (online) prediction. However, batch prediction produces a long-running operation (LRO), meaning that results are only available once the LRO has completed.

The maximum lifespan for a custom model is 18 months as of the GA release. You must create and train a new model to continue annotating content after that time.

## Batch prediction

You can request annotations (predictions) for images by using the `batchPredict` command. The `batchPredict` command takes, as input, a CSV file stored in your Google Cloud Storage bucket that contains the paths to the images to annotate. Each line specifies a separate path to an image in Google Cloud Storage.

**batch\_prediction.csv:**

```
gs://my-cloud-storage-bucket/prediction_files/image1.jpg
gs://my-cloud-storage-bucket/prediction_files/image2.jpg
gs://my-cloud-storage-bucket/prediction_files/image3.jpg
gs://my-cloud-storage-bucket/prediction_files/image4.jpg
gs://my-cloud-storage-bucket/prediction_files/image5.jpg
gs://my-cloud-storage-bucket/prediction_files/image6.png
```

If your filenames have spaces in them use quotation marks around the Cloud Storage location. For example:

- "gs://my-cloud-storage-bucket/prediction\_files/image filename with spaces.jpg"

Depending on the number of images that you specified in your CSV file, the batch predict task can take some time to complete. Even on a small number of images batch prediction will take *at minimum* 30 minutes to complete.

REST &amp; CMD LINE

C#

JAVA

MORE ▾

Before using any of the request data below, make the following replacements:

- **project-id**: your GCP project ID.
- **location-id**: A valid location identifier. Currently you must use the following value:
  - `us-central1`
- **model-id**: the ID of your model, from the response when you created the model. The ID is the last element of the name of your model. For example:
  - model name: `projects/project-id/locations/location-id/models/I0D4412217016962778756`
  - model id: **I0D4412217016962778756**
- **input-storage-path**: the path to a CSV file stored on Google Cloud Storage. The requesting user must have at least read permission to the bucket.
- **output-storage-bucket**: a Google Cloud Storage bucket/directory to save output files to, expressed in the following form: `gs://bucket/directory/`. The requesting user must have write permission to the bucket.

#### Field-specific considerations:

- `params.score_threshold` - A value between 0.0 and 1.0. Only results with scores greater or equal to this value will be returned.

HTTP method and URL:

POST `https://automl.googleapis.com/v1/projects/project-id/locations/location-id/models/I0D4412217016962778756`

Request JSON body:

```
{
  "inputConfig": {
    "gcsSource": {
      "inputUris": [ "input-storage-path" ]
    }
  },
  "outputConfig": {
    "gcsDestination": {
```

```

    "outputUriPrefix": "output-storage-bucket"
  }
},
"params": {
  "score_threshold": "0.0"
}
}

```

To send your request, choose one of these options:

**CURL**

POWERSHELL

**Note:** Ensure you have set the [GOOGLE\\_APPLICATION\\_CREDENTIALS](https://cloud.google.com/docs/authentication/production) (<https://cloud.google.com/docs/authentication/production>) environment variable to your service account private key file path.

Save the request body in a file called `request.json`, and execute the following command:

```

curl -X POST \
-H "Authorization: Bearer "$(gcloud auth application-default print-access-token)
-H "Content-Type: application/json; charset=utf-8" \
-d @request.json \
https://automl.googleapis.com/v1/projects/project-id/locations/location-id/model

```

## Response:

You should see output similar to the following:

```

{
  "name": "projects/project-id/locations/location-id/operations/ICN9266156233314795"
  "metadata": {
    "@type": "type.googleapis.com/google.cloud.automl.v1.OperationMetadata",
    "createTime": "2019-06-19T21:28:35.302067Z",
    "updateTime": "2019-06-19T21:28:35.302067Z",
    "batchPredictDetails": {
      "inputConfig": {
        "gcsSource": {
          "inputUris": [
            "input-storage-path"
          ]
        }
      }
    }
  }
}

```

```

    }
  }
}

```

You can use the operation ID (ICN926615623331479552, in this case) to get the status of the task. For an example, see [Getting the status of an operation](#) (#get-operation).

Depending on the number of images that you specified in your CSV file, the batch predict task can take some time to complete. Even on a small number of images batch prediction will take *at minimum* 30 minutes to complete.

Once the operation has completed, the `state` shows as `DONE` and your results are written to the Google Cloud Storage file you specified:

```

{
  "name": "projects/project-id/locations/location-id/operations/ICN9266156233314795",
  "metadata": {
    "@type": "type.googleapis.com/google.cloud.automl.v1.OperationMetadata",
    "createTime": "2019-06-19T21:28:35.302067Z",
    "updateTime": "2019-06-19T21:57:18.310033Z",
    "batchPredictDetails": {
      "inputConfig": {
        "gcsSource": {
          "inputUris": [
            "input-storage-path"
          ]
        }
      },
      "outputInfo": {
        "gcsOutputDirectory": "gs://storage-bucket-vcml/subdirectory/prediction-8370"
      }
    }
  },
  "done": true,
  "response": {
    "@type": "type.googleapis.com/google.cloud.automl.v1.BatchPredictResult"
  }
}

```

See the [Output JSONL files](#) (#output-jsonl) section below for a sample output file.

## Output JSONL files

When the batch predict task is complete, the output of the prediction is stored in the Google Cloud Storage bucket that you specified in your command.

In your output bucket (if applicable, in your specified directory) the files `image_classification_1.jsonl`, `image_classification_2.jsonl`, ..., `image_classification_N.jsonl` will be created, where N may be 1, and depends on the total number of the successfully predicted images and annotations.

A single image will be listed only once with all its annotations, and its annotations will never be split across files.

You can specify a minimum score in the request to limit the results returned.

Each JSONL file will contain, per line, a JSON representation of a proto that wraps image's "ID" : "<id\_value>" followed by a list of zero or more AnnotationPayload protos (called annotations), which have classification detail populated.

If prediction for any image failed (partially or completely), then an additional `errors_1.jsonl`, `errors_2.jsonl`, ..., `errors_N.jsonl` files will be created (N depends on total number of failed predictions). These files will have a JSON representation of a proto that wraps the same "ID" : "<id\_value>" but here followed by exactly one `google.rpc.Status` containing only `code` and `message` fields.

### Example JSONL file (a single .jsonl file with 4 lines/image file annotations):

`image_image_classification_0.jsonl`

#### Line 1 (daisy1.jpgannotation JSON)

```
{
  "ID": "gs://storage-bucket-vcm/img/daisy1.jpg",
  "annotations": [
    {
      "annotation_spec_id": "daisy",
      "classification": {
        "score": 0.99906391
      },
      "display_name": "daisy"
    },
    {
      "annotation_spec_id": "dandelion",
```



```

    "classification": {
      "score": 0.00085875636
    },
    "display_name": "dandelion"
  },
  {
    "annotation_spec_id": "roses",
    "classification": {
      "score": 0.000018997729
    },
    "display_name": "roses"
  },
  {
    "annotation_spec_id": "sunflowers",
    "classification": {
      "score": 0.0000041045291
    },
    "display_name": "sunflowers"
  },
  {
    "annotation_spec_id": "tulips",
    "classification": {
      "score": 0.000039702507
    },
    "display_name": "tulips"
  },
  {
    "annotation_spec_id": "--other--",
    "classification": {
      "score": 0.000014527803
    },
    "display_name": "--other--"
  }
]
}

```

### Line 2 (daisy2.jpgannotation JSON)

```

{
  "ID": "gs://storage-bucket-vcn/img/daisy2.jpg",
  "annotations": [
    {
      "annotation_spec_id": "daisy",
      "classification": {

```

```
    "score": 0.99953115
  },
  "display_name": "daisy"
},
{
  "annotation_spec_id": "dandelion",
  "classification": {
    "score": 0.00014155755
  },
  "display_name": "dandelion"
},
{
  "annotation_spec_id": "roses",
  "classification": {
    "score": 0.000011171558
  },
  "display_name": "roses"
},
{
  "annotation_spec_id": "sunflowers",
  "classification": {
    "score": 0.00030725187
  },
  "display_name": "sunflowers"
},
{
  "annotation_spec_id": "tulips",
  "classification": {
    "score": 7.7882828e-7
  },
  "display_name": "tulips"
},
{
  "annotation_spec_id": "--other--",
  "classification": {
    "score": 0.0000081920462
  },
  "display_name": "--other--"
}
]
}
```



**Line 3 (dandelion1.jpgannotation JSON)**

```
{
  "ID": "gs://storage-bucket-vcn/img/dandelion1.jpg",
  "annotations": [
    {
      "annotation_spec_id": "daisy",
      "classification": {
        "score": 0.0000041204139
      },
      "display_name": "daisy"
    },
    {
      "annotation_spec_id": "dandelion",
      "classification": {
        "score": 0.99971503
      },
      "display_name": "dandelion"
    },
    {
      "annotation_spec_id": "roses",
      "classification": {
        "score": 4.9584577e-7
      },
      "display_name": "roses"
    },
    {
      "annotation_spec_id": "sunflowers",
      "classification": {
        "score": 0.00027974427
      },
      "display_name": "sunflowers"
    },
    {
      "annotation_spec_id": "tulips",
      "classification": {
        "score": 3.8392983e-7
      },
      "display_name": "tulips"
    },
    {
      "annotation_spec_id": "--other--",
      "classification": {
        "score": 2.6729541e-7
      },
      "display_name": "--other--"
    }
  ]
}
```



```
]
}
```

#### Line 4 (dandelion2.jpgannotation JSON)

```
{
  "ID": "gs://automl-batch-iod-vcm/img/dandelion2.jpg",
  "annotations": [
    {
      "annotation_spec_id": "daisy",
      "classification": {
        "score": 0.00023957422
      },
      "display_name": "daisy"
    },
    {
      "annotation_spec_id": "dandelion",
      "classification": {
        "score": 0.99976045
      },
      "display_name": "dandelion"
    },
    {
      "annotation_spec_id": "roses",
      "classification": {
        "score": 1.7562879e-8
      },
      "display_name": "roses"
    },
    {
      "annotation_spec_id": "sunflowers",
      "classification": {
        "score": 3.2643279e-9
      },
      "display_name": "sunflowers"
    },
    {
      "annotation_spec_id": "tulips",
      "classification": {
        "score": 1.3378423e-8
      },
      "display_name": "tulips"
    },
    {
```

```
"annotation_spec_id": "--other--",
"classification": {
  "score": 4.6433613e-9
},
"display_name": "--other--"
}
]
}
```

## Getting the status of an operation

REST & CMD LINE

C#

GO

MORE ▾

Before using any of the request data below, make the following replacements:

- **project-id**: your GCP project ID.
- **operation-id**: the ID of your operation. The ID is the last element of the name of your operation. For example:
  - operation name: `projects/project-id/locations/location-id/operations/IOD5281059901324392598`
  - operation id: `IOD5281059901324392598`

HTTP method and URL:

```
GET https://automl.googleapis.com/v1/projects/project-id/locations/us-central1/ope
```

To send your request, choose one of these options:

CURL

POWERSHELL

**Note:** Ensure you have set the [GOOGLE\\_APPLICATION\\_CREDENTIALS](https://cloud.google.com/docs/authentication/production) (<https://cloud.google.com/docs/authentication/production>) environment variable to your service account private key file path.

Execute the following command:

```
curl -X GET \  
-H "Authorization: Bearer "$(gcloud auth application-default print-access-token)  
https://automl.googleapis.com/v1/projects/project-id/locations/us-central1/opera
```

You should see output similar to the following for a completed **import operation**:

```
{
  "name": "projects/project-id/locations/us-central1/operations/operation-id",
  "metadata": {
    "@type": "type.googleapis.com/google.cloud.automl.v1.OperationMetadata",
    "createTime": "2018-10-29T15:56:29.176485Z",
    "updateTime": "2018-10-29T16:10:41.326614Z",
    "importDataDetails": {}
  },
  "done": true,
  "response": {
    "@type": "type.googleapis.com/google.protobuf.Empty"
  }
}
```

You should see output similar to the following for a completed **create model operation**:

```
{
  "name": "projects/project-id/locations/us-central1/operations/operation-id",
  "metadata": {
    "@type": "type.googleapis.com/google.cloud.automl.v1.OperationMetadata",
    "createTime": "2019-07-22T18:35:06.881193Z",
    "updateTime": "2019-07-22T19:58:44.972235Z",
    "createModelDetails": {}
  },
  "done": true,
  "response": {
    "@type": "type.googleapis.com/google.cloud.automl.v1.Model",
    "name": "projects/project-id/locations/us-central1/models/model-id"
  }
}
```

Except as otherwise noted, the content of this page is licensed under the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) (https://creativecommons.org/licenses/by/4.0/), and code samples are licensed under the [Apache 2.0 License](https://www.apache.org/licenses/LICENSE-2.0) (https://www.apache.org/licenses/LICENSE-2.0). For details, see our [Site Policies](https://developers.google.com/terms/site-policies) (https://developers.google.com/terms/site-policies). Java is a registered trademark of Oracle and/or its affiliates.

Last updated January 22, 2020.