UPLOAD

# A Quicker Way to Upload and Share

Anyone collecting camera trap photos can upload and share them with the global conservation community. Photos are stored online so you can access them from anywhere, from any device or computer, even out in the field.

Get Started
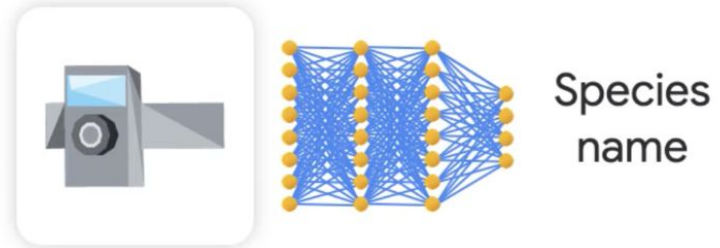
# Monitor wildlife health via image classification



https://youtu.be/hUzODH3uGg0



Open notebook in Colab

# Overview



#1 Training: builds an ML model

#2 Prediction: classifies images

# WCS Camera Traps dataset

- Approximately 1.4 million images
- Around 675 species from 12 countries
- More than 560 GB of images total
- Very unbalanced
  - Some species have tens of thousands of images
  - Many species have only a couple images
- Approximately 50% of images are empty
- Image files live in Azure Storage



http://lila.science/datasets/wcscameratraps

# Creating the images database

**beam**

| | |
|---|---|
| **Create :** *metadata_file_url* | Singleton metadata URL from LILA science |
| **FlatMap :** *(category, file_name)* | Parses the metadata file and yields (category, file_name) pairs |
| **Filter :** *(category, file_name)* | Filters out images with invalid categories |
| **Write to BigQuery** | Writes the valid labeled image file names into BigQuery |

create_images_metadata_table.py

# Creating the images database -- job graph



Dataflow

JOB GRAPH     EXECUTION DETAILS     JOB METRICS     RECOMMENDATIONS

**Job steps view**
Graph view ⌄

CLEAR SELECTION

✓ **Create None** ⌄
Succeeded
0 sec
1 of 1 stage succeeded

✓ **Get images info**
Succeeded
29 sec
1 of 1 stage succeeded

✓ **Filter invalid rows**
Succeeded
4 sec
1 of 1 stage succeeded

✓ **Write images database** ⌄
Succeeded
36 sec
11 of 11 stages succeeded

## Job info

| | |
|---|---|
| Job name | wildlife-images-database-2030e2 |
| Job ID | 2021-08-04_04_51_51-14223758464960217795 |
| Job type | Batch |
| Job status | ✓ Succeeded |
| SDK version | Apache Beam Python 3.8 SDK 2.30.0 |
| Job region ❓ | us-central1 |
| Worker location ❓ | us-central1-f |
| Current workers ❓ | 0 |
| Latest worker status | Worker pool stopped. |
| Start time | August 4, 2021 at 4:51:52 AM GMT-7 |
| Elapsed time | 6 min 58 sec |
| Encryption type | Google-managed key |

### Resource metrics ⌄

| | |
|---|---|
| Current vCPUs ❓ | 2 |
| Total vCPU time ❓ | 0.127 vCPU hr |
| Current memory ❓ | 7.5 GB |
| Total memory time ❓ | 0.478 GB hr |
| Current HDD PD ❓ | 25 GB |
| Total HDD PD time ❓ | 1.592 GB hr |
| Current SSD PD ❓ | 0 B |
| Total SSD PD time ❓ | 0 GB hr |
| Total Shuffle data | 558 B |

Logs     ≡ SHOW

Show debug panel

# Preview of the images metadata table

**BigQuery**

| category | file_name |
| --- | --- |
| tapirus bairdii | animals/0597/0707.jpg |
| equus quagga | animals/0377/1882.jpg |
| papio anubis | animals/0036/1687.jpg |
| dicerorhinus sumatrensis | animals/0329/0830.jpg |
| cephalophus nigrifrons | animals/0331/1215.jpg |
| tayassu pecari | animals/0174/0182.jpg |
| cephalophus nigrifrons | animals/0682/1295.jpg |
| giraffa camelopardalis | animals/0320/1392.jpg |
| panthera onca | animals/0564/0604.jpg |
| leopardus pardalis | animals/0576/0243.jpg |

# Training the model *(part 1) -- balancing the dataset*

| | | |
|---|---|---|
| **images :** *(category, gcs_path)* | **=** **Read from BigQuery :** *{category : Str, file_name : Str}* | Reads the images metadata information from BigQuery |

**Map :** *(category, file_name)*

Make (category, file_name) pairs

**Sample.PerKey :** *(category, List file_name)*

Get random samples of at most **max_images_per_class**

**Filter :** *(category, List file_name)*

Discard samples smaller than **min_images_per_class**

**FlatMap :** *(category, file_name)*

Flatten into (category, file_name) pairs

**FlatMap :** *(category, gcs_path)*

Lazily download *only* the images used, skip any errors

# Training the model *(part 2) -- preparing for AutoML*

**beam**

```
Create :
csv_file
```
Define the file name for a CSV file

```
images :
(category, gcs_path)
```

```
Map :
csv_file
```
Write the labeled images into the CSV file

```
Map :
(automl_dataset_path, csv_file)
```
Create an AutoML dataset

```
Map :
automl_dataset_path
```
Import the CSV entries into the AutoML dataset

```
Map :
automl_training_job_id
```
Start an AutoML training job

# Training the model -- job graph

# Training the model -- job metrics

# Create an AutoML dataset

Vertex AI

# Training the model -- precision

**Vertex AI**

← wildlife_classifier_20210629_182330    ⊞ VIEW DATASET

**EVALUATE**    DEPLOY & TEST    BATCH PREDICTIONS    MODEL PROPERTIES

≡ Filter    Filter labels

Confidence threshold ❓ ━━━━●━━━━ 0.5

| Label | Value |
|---|---|
| All labels | 0 |
| ortalis guttata | 1 |
| mazama pandora | 1 |
| sapajus apella | 1 |
| penelope purpurascens | 1 |
| crypturellus cinereus | 1 |
| unknown frog | 1 |
| tinamus sp | 1 |
| galictis vittata | 1 |
| mitu tomentosa | 0.99242 |
| grallaria andicolus | 0.98888 |
| ceratotherium simum | 0.97048 |
| xerus erythropus | 0.96231 |
| conepatus semistriatus | 0.95987 |
| momotus momota | 0.95555 |
| eudorcas thomsonii | 0.95028 |
| viverricula indica | 0.95 |
| aramides cajanea | 0.92666 |

## All labels

| | |
|---|---|
| Average precision ❓ | 0.676 |
| Precision ❓ | 79.3% |
| Recall ❓ | 50.4% |
| Created | Jun 29, 2021, 12:44:58 PM |
| Total images | 27,884 |
| Training images | 22,303 |
| Validation images | 2,804 |
| Test images | 2,777 |

To evaluate your model, set the **confidence threshold** to see how precision and recall are affected. The best confidence threshold depends on your use case. Read some example scenarios to learn how evaluation metrics can be used.

**Precision-recall curve** ❓

Confidence threshold: 0.05
Precision: 28.193%
Recall: 82.355%

**Precision-recall by threshold** ❓

— Recall    — Precision

## Confusion matrix

⬤ Item counts ⬇

This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray).

Show debug panel

wildlife_classifier_20210629_182330   ⊞ VIEW DATASET

**EVALUATE**   DEPLOY & TEST   BATCH PREDICTIONS   MODEL PROPERTIES

≡ Filter   Filter labels

Confidence threshold ?   ———●———   0.5

| All labels | 0 |
| ortalis guttata | 1 |
| mazama pandora | 1 |
| sapajus apella | 1 |
| penelope purpurascens | 1 |
| cryprturellus cinereus | 1 |
| unknown frog | 1 |
| tinamus sp | 1 |
| galictis vittata | 1 |
| mitu tomentosa | 0.99242 |
| grallaria andicolus | 0.98888 |
| ceratotherium simum | 0.97048 |
| xerus erythropus | 0.96231 |
| conepatus semistriatus | 0.95987 |
| momotus momota | 0.95555 |
| eudorcas thomsonii | 0.95028 |
| viverricula indica | 0.95 |
| aramides cajanea | 0.92666 |
| tapirus indicus | 0.92572 |
| hybomys univittatus | 0.91597 |
| turtur tympanistria | 0.91507 |

Item counts ⬇

## Confusion matrix

This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray).

Predicted label →

| True label | nesomys sp | rattus rattus | columba larvata | rhynchocyon cirnei | paradoxurus hermaphrodit | atherurus macrourus | equus grevyi | equus quagga | tinamus major | leptotila plumbeiceps |
|---|---|---|---|---|---|---|---|---|---|---|
| nesomys sp | 17% | 58% | — | — | — | — | — | — | — | — |
| rattus rattus | 7% | 79% | — | — | — | — | — | — | — | — |
| columba larvata | — | — | 25% | 63% | — | — | — | — | — | — |
| rhynchocyon cirnei | — | — | 17% | 67% | — | — | — | — | — | — |
| paradoxurus hermaphroditus | — | — | — | — | 30% | 50% | — | — | — | — |
| atherurus macrourus | — | — | — | — | 17% | 67% | — | — | — | — |
| equus grevyi | — | — | — | — | — | — | 27% | 45% | — | — |
| equus quagga | — | — | — | — | — | — | 33% | 50% | — | — |
| tinamus major | — | — | — | — | — | — | — | — | 29% | 71% |
| leptotila plumbeiceps | — | — | — | — | — | — | — | — | — | 100% |

Show debug panel

# Getting predictions

**High Accuracy**

category:
dicerorhinus sumatrensis

file: 'animals/0325/1529.jpg'

prediction:
dicerorhinus sumatrensis:
92.79% confidence

# Getting predictions

Use 'family' instead of species

category:
leopardus wiedii

file: 'animals/0000/1705.jpg'

prediction:
leopardus pardalis:
55.56% confidence
**leopardus wiedii:**
**33.45% confidence**

# Getting predictions

AutoML Vision



category:
dasypus novemcinctus

file: 'animals/0000/0425.jpg'

prediction:
procyon cancrivorus:
19.38% confidence

dasypus novemcinctus:
16.65% confidence

columba larvata:
10.84% confidence

Train more with far shots