# Image Classification Example



Input Image      Pre-trained Model      Prediction Output

**Golden Retriever**

# Image Classification Tensorflow Code

Pre-trained model

Pre-processing

Inference

```python
from tensorflow.keras.applications.inception_v3 import InceptionV3
from tensorflow.keras.preprocessing import image
from tensorflow.keras.applications.imagenet_utils import preprocess_input, decode_predictions

model = InceptionV3(weights='imagenet')
img = image.load_img(inp_file_name, target_size=(299, 299))
x = image.img_to_array(img)
x = np.expand_dims(x, axis=0)
x = preprocess_input(x, mode='tf')

preds = model.predict(x)
```
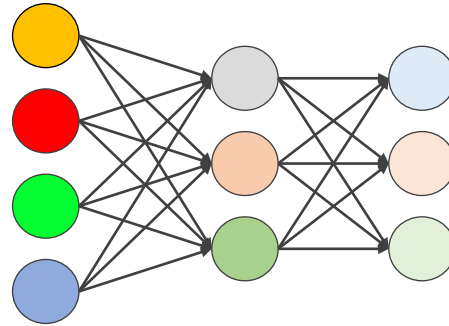
Pre-processing

[299,299]

filtering

**golden_retriever: 85.93 %**
Labrador_retriever: 5.18 %
kuvasz: 0.94 %

Input Image

InceptionV3
imagenet

Predictions
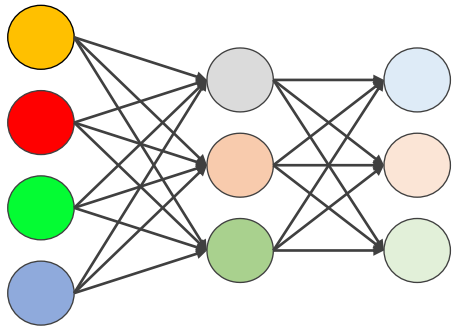
[dog.jpg]

[frozen.pb]

[Array]

# OpenVINO Model Optimizer

```python
# Setup model optimizer command ...
ir_name = "inceptionv3"
ir_data_type = "FP32"
ir_out_dir = f"{model_fname}/IR_models/{ir_data_type}"
ir_input_shape = "[1,299,299,3]"

mo_cmd = f"mo_tf.py \
        --saved_model_dir {model_fname} \
        --input_shape {ir_input_shape} \
        --data_type {ir_data_type} \
        --output_dir {ir_out_dir}  \
        --model_name {ir_name}"

#run the Optimizer command
output = subprocess.check_output(mo_cmd, shell=True)
print (output.decode('utf-8'))
```
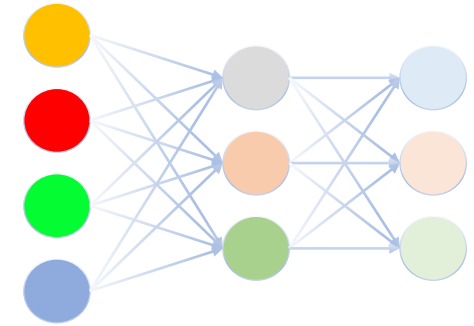
# OpenVINO
# Model Optimizer

[frozen.pb]

- Linear Operations Fusing

- ResNet optimization (stride optimization)

- Grouped Convolution Fusing
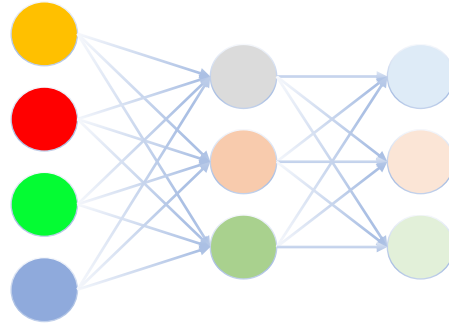
- Model Cutting

[model.bin, model.xml]

Pre-processing

[299,299]

filtering

**golden_retriever: 85.93 %**
Labrador_retriever: 5.18 %
kuvasz: 0.94 %

Input Image

OpenVINO
Optimized Model

Predictions

[dog.jpg]

[model.bin, model.xml]

[Array]

# OpenVINO Inference



```python
from openvino.inference_engine import IECore

model_xml = f'{ir_out_dir}/{ir_name}.xml'
model_bin = f'{ir_out_dir}/{ir_name}.bin'

# Load network to the plugin
ie = IECore()
net = ie.read_network(model=model_xml, weights=model_bin)
exec_net = ie.load_network(network=net, device_name="CPU")
del net

input_layer = next(iter(exec_net.input_info))
output_layer = next(iter(exec_net.outputs))

# Run the Inference on the Input image...
res = exec_net.infer(inputs={input_layer: input_image})
res = res[output_layer]
```

CPU
GPU
FPGA
MYRIAD
MULTI

GPUs  CPUs  FPGAs  VPUs
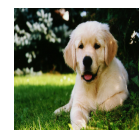
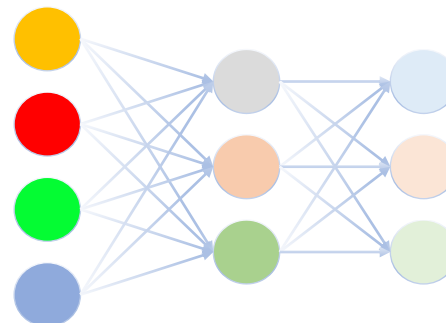golden_retriever: 89.87 %
Labrador_retriever: 1.84 %
tennis_ball: 0.79 %

Pre-processing

[299,299]

Input Image

OpenVINO Optimized
Pre-trained Model

Prediction Outcomes

[dog.jpg]

[model.bin, model.xml]

[Array]

# inception-v3-TF

## Throughput (higher is better)

## Latency (lower is better)

CPU INFERENCE ENGINES

| Engine | INT8 | FP32 | Latency (ms) |
|---|---|---|---|
| Intel® Atom™ x5-E3940 | 2.93 | 6.16 | 170.36 |
| Intel® Core™ i3-8100 | 33.92 | 54.99 | 18.93 |
| Intel® Core™ i5-8500 | 48.87 | 80.95 | 13.03 |
| Intel® Core™ i7-8700T | 56.04 | 102.56 | 10.85 |
| Intel® Core™ i9-10920X | 146.84 | 581.14 | 3.1 |
| Intel® Core™ i5-1145G7E CPU-only | 37.59 | 138.26 | 7.68 |
| Intel® Core™ i5-1145G7E GPU-only | 55.13 | 159.65 | 8.74 |
| Intel® Core™ i5-1145G7E GPU+CPU | 78.4 | 259.91 | 8.74 |
| Intel® Xeon® E-2124G | 39.34 | 62.88 | 16.53 |
| Intel® Xeon® Silver 4216R | 156.5 | 595.53 | 4.28 |
| Intel® Xeon® Gold 5218T | 166.96 | 632.13 | 4.17 |
| Intel® Xeon® Platinum 8270 | 482.8 | 1927.33 | 3.14 |

Legend: INT8 · FP32 · FP16 · Milliseconds

## Value (higher is better)

| Processor | FPS/$ |
|---|---|
| Intel® Atom™ x5-E3940 | 0.146 |
| Intel® Core™ i3-8100 | 0.469 |
| Intel® Core™ i5-8500 | 0.421 |
| Intel® Core™ i7-8700T | 0.338 |
| Intel® Core™ i9-10920X | 0.83 |
| Intel® Core™ i5-1145G7E CPU-only | 0.447 |
| Intel® Core™ i5-1145G7E GPU-only | 0.516 |
| Intel® Core™ i5-1145G7E GPU+CPU | 0.841 |
| Intel® Xeon® E-2124G | 0.295 |
| Intel® Xeon® Silver 4216R | 0.594 |
| Intel® Xeon® Gold 5218T | 0.468 |
| Intel® Xeon® Platinum 8270 | 0.26 |

## Efficiency (higher is better)

| Processor | FPS/Thermal Design Power |
|---|---|
| Intel® Atom™ x5-E3940 | 0.647 |
| Intel® Core™ i3-8100 | 0.845 |
| Intel® Core™ i5-8500 | 1.245 |
| Intel® Core™ i7-8700T | 2.93 |
| Intel® Core™ i9-10920X | 3.522 |
| Intel® Core™ i5-1145G7E CPU-only | 4.937 |
| Intel® Core™ i5-1145G7E GPU-only | 5.701 |
| Intel® Core™ i5-1145G7E GPU+CPU | 9.282 |
| Intel® Xeon® E-2124G | 0.885 |
| Intel® Xeon® Silver 4216R | 2.382 |
| Intel® Xeon® Gold 5218T | 3.01 |
| Intel® Xeon® Platinum 8270 | 4.7 |