

Human Object Interaction Using Deep Learning

Group no.: 11

Aryan Menon (CEC21CS079)

Ken Jojo (CEC21CS057)

Sinan Noushad (CEC21CS088)

Midhun S Kumar (CEC21CS064)

Guided By : Judy Ann Joy

DEPARTMENT OF COMPUTER ENGINEERING

COLLEGE OF ENGINEERING CHERTHALA

January 8, 2025

INDEX

- 1 INTRODUCTION
- 2 RELEVANCE
- 3 PROBLEM STATEMENT
- 4 OBJECTIVES
- 5 LITERATURE SURVEY
- 6 CONCLUSION FROM LITERATURE SURVEY
- 7 PRODUCT FUNCTIONS
- 8 REQUIREMENTS
- 9 TOOLS AND LANGUAGE
- 10 SYSTEM DESIGN
- 11 SYSTEM ARCHITECTURE
- 12 USE CASE DIAGRAM
- 13 DATA FLOW DIAGRAM
- 14 GANTT CHART
- 15 COST ESTIMATION
- 16 IMPLEMENTATION
- 17 CONCLUSION
- 18 REFERENCES

INTRODUCTION

- Human-Object Interaction (HOI) using Deep Learning uses Convolutional Neural Networks (CNNs) to detect and understand interactions between humans and objects in images or videos.
- CNNs are particularly well-suited for this task due to their ability to effectively extract features from visual inputs.
- This project represents a fusion of cutting-edge deep learning techniques with the fundamental challenge of understanding how humans interact with objects in visual scenes.

RELEVANCE

- Human-Object Interaction (HOI) using Convolutional Neural Networks (CNNs) is crucial in computer vision for enhancing scene understanding and improving human-centric applications.
- By recognizing human-object interactions, CNNs enhance object recognition, prediction, and smart responses in surveillance, human-robot interaction, and assistive technologies.
- In autonomous systems like self-driving cars and smart homes, HOI improves decision-making and user experience. Despite challenges, HOI with CNNs is vital for advancing AI technologies.

PROBLEM STATEMENT

- To build an application that processes image inputs and predicts the type of Human-Object Interaction (HOI) occurring in the image using deep learning models. The system will analyze visual cues to identify and label specific interactions between people and objects.

OBJECTIVES

- Develop human-object interaction detection model
- Ensure robustness across lighting, poses, objects
- Enhance interpretive accuracy
- Scalable for diverse real-world applications

LITERATURE REVIEW

Recognizing Human-Object Interactions in Still Images by Modeling the Mutual Context of Objects and Human Poses, IEEE, 2012 [1]

Authors : Bangpeng Yao, Li Fei-Fei

- The paper employs a conditional random field (CRF) model to jointly capture the mutual context between human poses and objects in human-object interactions.
- The CRF model in the paper utilizes co-occurrence context via a co-occurrence frequency algorithm.
- Using the Gray-Level Co-Occurrence Matrix (GLCM), the algorithm analyzes the spatial relationships between pixel intensities in an image to quantify texture by calculating how frequently pairs of pixel values occur in a specified direction and distance.

LITERATURE REVIEW

Advantages

- The mutual context model enhances the accuracy of both object detection and human pose estimation.
- The model excels in detecting small or partially visible objects and estimating poses.

Disadvantages

- Jointly modeling objects and human poses adds complexity.
- The effectiveness of the model may depend on the availability of high-quality, labeled datasets.

Future Scope

- Extend the method to incorporate additional contextual information and various environmental conditions for more robust and generalizable HOI detection.

LITERATURE REVIEW

Detecting and Recognizing Human-Object Interactions, IEEE, 2018

[2]

Authors: Georgia Gkioxari, Ross Girshick, Piotr Dollár, Kaiming He

- The paper utilizes InteractNet, a custom architecture that integrates multiple branches to simultaneously detect objects and recognize human actions.
- The interaction branch of InteractNet specifically targets the detection of interactions by evaluating the spatial relationships between detected humans and objects, utilizing an action-specific density estimator.
- The estimator generates a probabilistic density function (PDF) that typically predicts a single probable location for simpler actions. For more complex interactions, it can extend to multiple probable locations.

LITERATURE REVIEW

Advantages

- Improved accuracy via human-centric cues.
- Unified detection of people, objects, and interactions.
- Proven effectiveness on benchmark datasets.

Disadvantages

- High computational cost.
- Limited in rare interactions.

Future Scope

- Future work could focus on training and evaluating the model across a wider range of environments and conditions.

LITERATURE REVIEW

Human-object Interaction Recognition Using Multitask Neural Network, IEEE, 2019 [3]

Authors: Weihao Yan, Yue Gao, Qiming Liu

- The paper employs the YOLO (You Only Look Once) algorithm for real-time object detection, utilizing a single convolutional network to predict multiple bounding boxes and class probabilities simultaneously, thereby streamlining the process of identifying objects within an image.
- An attention mechanism is integrated to on relevant regions in the input image, dynamically weighing the importance of features based on the context of detected human actions.
- The attention mechanism in images computes importance scores for different regions of the feature map, enabling the model to focus on relevant areas while processing the image.

LITERATURE REVIEW

Advantages

- Improved accuracy of HOI detection by effectively utilizing the relationships between body parts and objects.
- Integrated motion and object recognition.

Disadvantages

- Requires specialized equipment (RGBD camera and digital gloves).
- Complexity in setup and training.

Future Scope

- Extend the dataset and model to recognize a broader range of human-object interactions, enhancing its applicability to more diverse and complex scenarios.

LITERATURE REVIEW

A Human-Object Interaction Detection Method Inspired by Human Body Part Information, IEEE, 2020 [4]

Authors : Hailan Kuang, Zhikun Zheng, Xinhua Liu, Xiaolin Ma

- The paper employs a hierarchical model that integrates human body part information and object features to detect human-object interactions.
- The detection algorithm leverages a human body part-based feature extraction approach, where both local and global features are analyzed.
- The algorithm functions by first segmenting the human body into various parts, then extracting key features from these parts. It then analyzes the spatial configuration of the body parts in relation to objects to infer interactions, effectively mimicking human cognitive processes in recognizing actions.

LITERATURE REVIEW

Advantages

- Improved the accuracy of HOI detection by effectively utilizing the relationships between body parts and objects.
- The method makes better use of human pose information.

Disadvantages

- Incorporating detailed human body part information may increase the computational complexity.
- The effectiveness of the method relies heavily on accurate human pose estimation.

Future Scope

- Future work could include integrating advanced techniques to improve generalization across diverse scenarios.

LITERATURE REVIEW

HOTR: End-to-End Human-Object Interaction Detection with Transformers, IEEE, 2021 [5]

Authors: Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, Hyunwoo J. Kim

- The paper presents HOTR, a transformer-based model for end-to-end human-object interaction (HOI) detection that predicts HOI triplets directly, bypassing post-processing.
- The HOTR algorithm uses a transformer encoder-decoder architecture with parallel decoders to simultaneously predict human, object, and interaction, utilizing self-attention to capture relationships between them.
- The model works by first extracting global context via a shared encoder, then using HO Pointers to associate predicted interactions with object instances, enabling efficient localization and classification of interactions.

LITERATURE REVIEW

Advantages

- Achieves high performance with minimal processing time.
- Directly predicts $\langle \text{human, object, interaction} \rangle$ triplets, enhancing detection precision.

Disadvantages

- Relies on Object Detection Quality, accuracy depends on the effectiveness of object detection.
- Transformer models can be resource-intensive.

Future Scope

- Adapt the model for real-time applications and dynamic environments to enhance practical usability.

Spatial-Temporal Human-Object Interaction Detection, IEEE, 2021 [6]

Authors : Xu Sun, Yunqing He, Tongwei Ren, Gangshan Wu

- This paper uses the ST-HOID model to detect fine-grained human-object interactions and track trajectories of both humans and objects in videos.
- The algorithm employs visual tracking and multi-modal feature fusion, focusing on object trajectory detection and interaction reasoning.
- The reasoning module pairs human-object trajectories and classifies interactions by analyzing body motion, relative motion, and semantic context.

LITERATURE SURVEY

Advantages

- The ST-HOID task allows for detailed detection of human-object interactions and object trajectories.
- The proposed method, combined with the VidOR-HOID dataset, advances the field by providing new benchmarks.

Disadvantages

- The effectiveness of the method may be limited to the VidOR-HOID dataset.
- The method's reliance on both object trajectory detection and interaction reasoning can increase computational complexity.

Future Scope

- Future work could focus on expanding the method to diverse datasets and integrating multi-model data to enhance overall interaction detection and understanding.

LITERATURE REVIEW

Hierarchical Reasoning Network for Human-Object Interaction Detection, IEEE, 2021 [7]

Authors: Yiming Gao, Zhanghui Kuang, Guanbin Li, Wayne Zhang, Liang Lin

- This paper uses the Hierarchical Reasoning Network (HRNet) to detect human-object interactions by modeling relations between human parts and objects across multiple scales.
- The HRNet algorithm employs a graph neural network that builds a multi-level graph, with nodes representing human parts, objects, and human-object pairs, and edges capturing visual and spatial relationships.
- The algorithm refines human-object interaction representations by propagating information through intra-level (within-scale) and inter-level (cross-scale) reasoning to capture both coarse and fine details of human-object interactions.

LITERATURE REVIEW

Advantages

- The model improves accuracy by capturing detailed relationships among human parts and objects at multiple scales.
- It sets new records on major datasets, proving its effectiveness.

Disadvantages

- The hierarchical reasoning and multi-level graph structure add complexity to the model, which may affect training and deployment.
- The extensive graph-based reasoning process can be computationally intensive, requiring significant resources.

Future Scope

- Focus on optimizing the model for faster processing and reduced computational resource usage, making it more scalable for real-time applications.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
Recognizing Human-Object Interactions in Still Images by Modeling the Mutual Context of Objects and Human Poses	Bangpeng Yao, Li Fei-Fei, IEEE, 2012.	<ul style="list-style-type: none"> Integrates human body part information to improve Human-Object Interaction detection. Shows significant performance improvements on the V-COCO dataset. 	<ul style="list-style-type: none"> Enhances the accuracy of both object detection and human pose estimation. Excels in detecting small or partially visible objects and estimating poses. 	<ul style="list-style-type: none"> Jointly modeling objects and human poses adds complexity. Depend on the availability of high-quality, labeled datasets. 	<ul style="list-style-type: none"> Extend the method to include diverse contextual information for improved robustness.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
Detecting and Recognizing Human-Object Interactions	Georgia Gkioxari, Ross Girshick, Piotr Dollár, Kaiming He, IEEE, 2018	<ul style="list-style-type: none"> The proposed InteractNet model leverages the appearance of detected people—such as their pose, clothing, and actions. InteractNet jointly learns to detect people and objects, fusing these predictions to efficiently infer interaction triplets in an end-to-end system. 	<ul style="list-style-type: none"> Improved accuracy via human-centric cues. Unified detection of people, objects, and interactions. 	<ul style="list-style-type: none"> High computational cost. Limited in rare interactions. 	<ul style="list-style-type: none"> Future work could focus on training and evaluating the model across a wider range of environments and conditions.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
Human-object Interaction Recognition Using Multi-task Neural Network	Weihaoyan, Yue Gao, Qiming Liu, IEEE, 2019.	<ul style="list-style-type: none"> The proposed method uses a multitask 2D convolutional neural network that integrates human body motion, hand motion, and object recognition. A YOLOv3-based object recognition network is introduced to boost the accuracy of interaction. 	<ul style="list-style-type: none"> Improved accuracy of HOI detection by effectively utilizing the relationships between body parts and objects. Integrated motion and object recognition. 	<ul style="list-style-type: none"> Requires specialized equipment. Complexity in setup and training. 	<ul style="list-style-type: none"> Extend the dataset and model to recognize a broader range of human-object interactions, enhancing its applicability to more diverse and complex scenarios.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
A Human-Object Interaction Detection Method Inspired by Human Body Part Information	Hailan Kuang, Zhikun Zheng, Xinhua Liu, Xiaolin Ma, IEEE, 2020.	<ul style="list-style-type: none">The proposed HOI detection method incorporates human body part (HBP) information to leverage the relationship between body parts and objects.	<ul style="list-style-type: none">Improved accuracy via human-centric cues.Unified detection of people, objects, and interactions.	<ul style="list-style-type: none">High computational cost.Effectiveness of the method relies heavily on accurate human pose estimation.	<ul style="list-style-type: none">Evaluate the model in diverse environments and conditions.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
HOTR: End-to-End Human-Object Interaction Detection with Transformers	Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, Hyunwoo J. Kim, IEEE, 2021.	<ul style="list-style-type: none"> • Directly predicts ⟨human, object, interaction⟩ triplets with a transformer model. • Sets new benchmarks with ≤ 1 ms inference time post-detection. 	<ul style="list-style-type: none"> • Delivers top results with minimal processing time. • Directly predicts ⟨human, object, interaction⟩ triplets for better accuracy. 	<ul style="list-style-type: none"> • Accuracy relies on the quality of object detection. • Transformers can demand significant resources. 	<ul style="list-style-type: none"> • Optimize for real-time and dynamic environments.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
Spatial-Temporal Human-Object Interaction Detection	Xu Sun, Yunqing He, Tongwei Ren, Gangshan Wu, Xu Sun, Yunqing He, Tongwei Ren, Gangshan Wu, IEEE, 2021.	<ul style="list-style-type: none">• Detects fine-grained human-object interactions and trajectories in videos.• Provides 10,831 instances and shows better performance than existing methods.	<ul style="list-style-type: none">• Enables detailed detection of interactions and object trajectories.• Advances the field with new benchmarks using VidOR-HOID.	<ul style="list-style-type: none">• Effectiveness may be limited to VidOR-HOID.• Relies on object trajectory and interaction reasoning, raising computational demands.	<ul style="list-style-type: none">• Expand to diverse datasets and integrate multi-modal data for improved interaction detection.

LITERATURE SURVEY

Name	Author	Methodology	Merits	Limitations	Future Scope
Hierarchical Reasoning Network for Human-Object Interaction Detection	Yiming Gao, Zhanghui Kuang, Guanbin Li, Wayne Zhang, Liang Lin, IEEE, 2021.	<ul style="list-style-type: none">• Models human parts and objects with a multi-level graph and hierarchical reasoning.• Delivers leading results on HICO-DET, V-COCO, and HOI-A benchmarks.	<ul style="list-style-type: none">• Improves accuracy by capturing detailed relationships among human parts and objects.• It sets new records on major datasets, proving its effectiveness.	<ul style="list-style-type: none">• Adds challenges to training and deployment.• Requires significant computational resources.	<ul style="list-style-type: none">• Improve processing speed and reduce resource usage for real-time scalability.

CONCLUSION FROM LITERATURE SURVEY

- The field has seen significant advancements through the use of convolutional neural networks (CNNs) and transformer-based models for detecting and understanding HOI in both static images and videos.
- While models like InteractNet and HOTR have improved accuracy by leveraging human-centric cues and predicting HOI triplets, challenges remain in computational complexity, the need for large annotated datasets, and scalability across diverse environments.
- Future work will likely focus on improving generalization, efficiency, and real-time processing capabilities to make HOI detection more practical and robust across various applications.

PRODUCT FUNCTIONS

- **Data Collection and Preprocessing:** Collect images or videos of human-object interactions and preprocess them (resizing, normalization) for consistency.
- **Object Detection:** Use a CNN (e.g., YOLO) to detect and localize objects in the images.
- **Action Recognition:** Apply action recognition models to map key human joints for interaction analysis.
- **Interaction Recognition:** Combine object and keypoint data to recognize actions like picking up or holding objects using temporal models.
- **Interaction Classification:** Classify interactions with a machine learning model (e.g., SVM) to label types of human-object engagements.

REQUIREMENTS

- **Hardware Requirements :**

- Computer with
 - RAM: 8 GB
 - SSD: 512 GB

- **Software Requirements:**

- OS: WINDOWS 10

- **Dataset:**

- Verbs in COCO (V-COCO) is a dataset that builds off COCO for human-object interaction detection. V-COCO provides 10,346 images (2,533 for training, 2,867 for validating and 4,946 for testing) and 16,199 person instances.

TOOLS AND LANGUAGE

- **Language** : Python
- **Various libraries** : NumPy, PIL etc.
- **Machine learning framework** : TensorFlow, PyTorch etc.
- **FrontEnd** : Python Tkinter
- **Development Environment** : Google Colab, VS Code

SYSTEM ARCHITECTURE

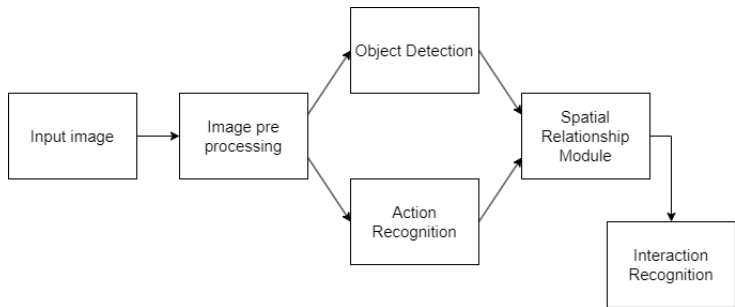


Figure 1: System Architecture

USE CASE DIAGRAM

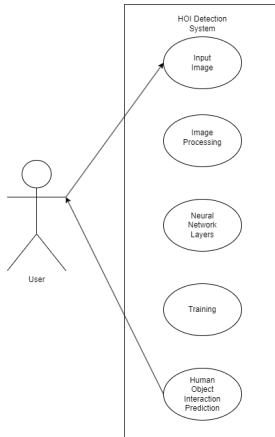


Figure 2: Use Case Diagram

DATA FLOW DIAGRAM LEVEL0

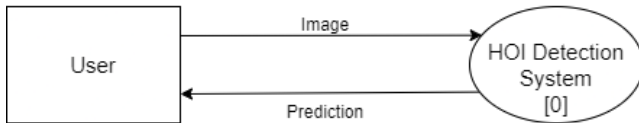


Figure 3: Level 0

DATA FLOW DIAGRAM LEVEL1

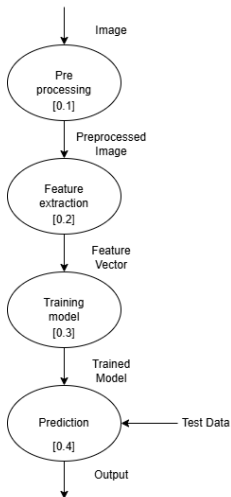


Figure 4: Level 1

GANTT CHART

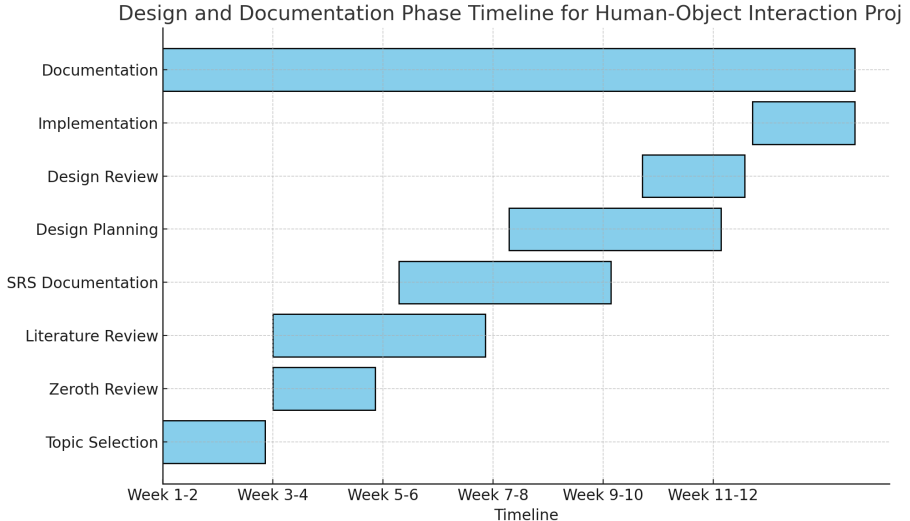


Figure 5: Gantt Chart

COST ESTIMATION

Lines of code = 211

KLOC = $211/1000 = 0.211$

Effort = $a(KLOC)^b$

$= 2.4(0.211)^{1.05}$

$= 0.468 \text{ PM}$

Duration = $c(\text{Effort})^d$

$= 2.5(0.468)^{0.38}$

$= 1.873 \cong 2 \text{ months}$

Staffing = $\text{Effort}/\text{Duration}$

$= 0.468/1.873$

$= 0.249 \cong 1 \text{ person}$

Rate per Person = 15,000 Rs

Total Rate for 1 staffs for 1 months = $15,000 * 1 * 1 = 15,000 \text{ Rs}$

IMPLEMENTATION

```
import os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import plotly.graph_objects as go
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix
import tensorflow as tf
from tensorflow.keras.preprocessing import image_dataset_from_directory
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense, Dropout
from tensorflow.keras.callbacks import EarlyStopping
```

```
[29] # Data Preprocessing
train_df['label'] = train_df['label'].astype('category')
train_df['label'] = train_df['label'].cat.codes
train_df['filepath'] = train_df['filename'].apply(lambda x: os.path.join(train_path, x))

[30] # Split training and validation set
train_set, val_set = train_test_split(train_df, test_size=0.2, stratify=train_df['label'], random_state=42)

[31] def load_image(filepath, label):
    image = tf.io.read_file(filepath)
    image = tf.image.decode_jpeg(image, channels=3)
    image = tf.image.resize(image, [128, 128])
    image = image / 255.0
    return image, label

# Create TensorFlow datasets
train_dataset = tf.data.Dataset.from_tensor_slices((train_set['filepath'].values, train_set['label'].values))
train_dataset = train_dataset.map(load_image).batch(32).shuffle(buffer_size=len(train_set))

val_dataset = tf.data.Dataset.from_tensor_slices((val_set['filepath'].values, val_set['label'].values))
val_dataset = val_dataset.map(load_image).batch(32)

[32] from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense, Dropout, Input

# Model Building
model = Sequential([
    Input(shape=(128, 128, 3)), # Use Input layer to specify input shape
    Conv2D(32, (3, 3), activation='relu'),
    MaxPooling2D((2, 2)),
    Conv2D(64, (3, 3), activation='relu'),
    MaxPooling2D((2, 2)),
    Conv2D(128, (3, 3), activation='relu'),
    MaxPooling2D((2, 2)),
    Flatten(),
    Dense(128, activation='relu'),
    Dropout(0.5),
    Dense(15, activation='softmax')
])

model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])
```

```

# Model Training
history = model.fit(train_dataset, validation_data=val_dataset, epochs=12)

```

Epoch	315/315	Time	Accuracy	Loss	Val Accuracy	Val Loss
Epoch 1/12	361s	1s/step	0.0903	2.7202	0.1921	2.5180
Epoch 2/12	371s	1s/step	0.1819	2.4830	0.2532	2.3218
Epoch 3/12	347s	1s/step	0.2531	2.2906	0.2742	2.2421
Epoch 4/12	383s	1s/step	0.2974	2.1367	0.2885	2.2041
Epoch 5/12	345s	1s/step	0.3432	1.9993	0.2968	2.2037
Epoch 6/12	401s	1s/step	0.3949	1.8284	0.2948	2.1925
Epoch 7/12	367s	1s/step	0.4491	1.6551	0.3087	2.2195
Epoch 8/12	347s	1s/step	0.5172	1.4565	0.3135	2.2824
Epoch 9/12	386s	1s/step	0.5720	1.2741	0.3135	2.4861
Epoch 10/12	376s	1s/step	0.6142	1.1259	0.3052	2.5643
Epoch 11/12	346s	1s/step	0.6550	1.0144	0.3048	2.7817
Epoch 12/12	385s	1s/step	0.6827	0.9141	0.3111	2.8968


```
# Show 10 images with actual and predicted labels from validation set
fig, axes = plt.subplots(2, 5, figsize=(20, 10))
axes = axes.flatten()
sampled_val_set = val_set.sample(10, random_state=42)
for i, (index, row) in enumerate(sampled_val_set.iterrows()):
    img = plt.imread(row['filepath'])
    actual_label = categories[row['label']]
    predicted_label = categories[val_predictions[val_set.index.get_loc(index)]]
    axes[i].imshow(img)
    axes[i].set_title(f"Actual: {actual_label}\nPredicted: {predicted_label}")
    axes[i].axis("off")
plt.tight_layout()
plt.show()
```

Actual: listening_to_music
Predicted: texting



Actual: cycling
Predicted: cycling



Actual: clapping
Predicted: fighting



Actual: laughing
Predicted: laughing



Actual: texting
Predicted: texting



Actual: clapping
Predicted: drinking



Actual: fighting
Predicted: fighting



Actual: drinking
Predicted: dancing



Actual: hugging
Predicted: hugging



Actual: calling
Predicted: calling



```

1 from google.colab import files
2 import numpy as np
3 import time
4 import cv2
5 from google.colab.patches import cv2_imshow
6
7 # Upload the image file
8 uploaded = files.upload()
9
10 # This will prompt you to upload the image file. Use the filename returned by the upload to read the image.
11 image_path = None
12 for file_name in uploaded.keys():
13     image_path = file_name
14     print(f"Uploaded file: {file_name}")
15
16 # Define the paths for the YOLO files
17 labelsPath = 'coco.names'
18 weightsPath = 'yolov3.weights'
19 configPath = 'yolov3.cfg'
20
21 # Load the COCO class labels our YOLO model was trained on
22 LABELS = open(labelsPath).read().strip().split("\n")
23
24 # Initialize a list of colors to represent each possible class label
25 np.random.seed(42)
26 COLORS = np.random.randint(0, 255, size=(len(LABELS), 3), dtype="uint8")
27
28 # Load our YOLO object detector trained on COCO dataset (80 classes)
29 print("[INFO] loading YOLO from disk...")
30 net = cv2.dnn.readNetFromDarknet(configPath, weightsPath)
31
32 # Verify the image path and read the image
33 image = cv2.imread(image_path)
34 if image is None:
35     print(f"Failed to load image at {image_path}")
36 else:
37     print(f"Successfully loaded image at {image_path}")
38
39 # Proceed with the rest of the script
40 (H, W) = image.shape[:2]
41
42 # Determine only the "output" layer names that we need from YOLO
43 ln = net.getLayerNames()
44 ln = [ln[i - 1] for i in net.getUnconnectedOutLayers()]

```

```

45
46 # Construct a blob from the input image and then perform a forward pass
47 # of the YOLO object detector, giving us our bounding boxes and
48 # associated probabilities
49 blob = cv2.dnn.blobFromImage(image, 1 / 255.0, (416, 416), swapRB=True, crop=False)
50 net.setInput(blob)
51 start = time.time()
52 layerOutputs = net.forward(ln)
53 end = time.time()
54
55 # Show timing information on YOLO
56 print("[INFO] YOLO took {:.6f} seconds".format(end - start))
57
58 # Initialize our lists of detected bounding boxes, confidences, and
59 # class IDs, respectively
60 boxes = []
61 confidences = []
62 classIDs = []
63
64 # Loop over each of the layer outputs
65 for output in layerOutputs:
66     # Loop over each of the detections
67     for detection in output:
68         # Extract the class ID and confidence (i.e., probability) of the current object detection
69         scores = detection[5:]
70         classID = np.argmax(scores)
71         confidence = scores[classID]
72
73         # Filter out weak predictions by ensuring the detected probability is greater than the minimum probability
74         if confidence > 0.5:
75             # Scale the bounding box coordinates back relative to the size of the image
76             box = detection[0:4] * np.array([W, H, W, H])
77             (centerX, centerY, width, height) = box.astype("int")
78
79             # Use the center (x, y)-coordinates to derive the top and
80             # and left corner of the bounding box
81             x = int(centerX - (width / 2))
82             y = int(centerY - (height / 2))
83
84             # Update our list of bounding box coordinates, confidences, and class IDs
85             boxes.append([x, y, int(width), int(height)])
86             confidences.append(float(confidence))
87             classIDs.append(classID)
88
89 # Apply non-maxima suppression to suppress weak, overlapping bounding boxes
90 idxs = cv2.dnn.NMSBoxes(boxes, confidences, 0.5, 0.3)
91

```

```

92 # Ensure at least one detection exists
93 if len(idxs) > 0:
94     # Loop over the indexes we are keeping
95     for i in idxs.flatten():
96         # Extract the bounding box coordinates
97         (x, y) = (boxes[i][0], boxes[i][1])
98         (w, h) = (boxes[i][2], boxes[i][3])
99
100         # Draw a bounding box rectangle and label on the image
101         color = [int(c) for c in COLORS[classIDs[i]]]
102         cv2.rectangle(image, (x, y), (x + w, y + h), color, 2)
103         text = "{}: {:.4f}".format(LABELS[classIDs[i]], confidences[i])
104         cv2.putText(image, text, (x, y - 5), cv2.FONT_HERSHEY_SIMPLEX, 0.5, color, 2)
105
106 # Show the output image using cv2_imshow
107 cv2_imshow(image)
108

```

Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving hyundai-creta.jpeg to hyundai-creta.jpeg
Uploaded file: hyundai-creta.jpeg
[INFO] loading YOLO from disk...
Successfully loaded image at hyundai-creta.jpeg
[INFO] YOLO took 2.714727 seconds



CONCLUSION

- The proposed system leverages state-of-the-art technologies and innovative techniques to address the limitations of existing approaches.
- By enhancing accuracy, improving efficiency, providing better contextual understanding, and ensuring scalability and robustness, the proposed system offers a superior solution for detecting and understanding human-object interactions.
- Utilizing models like Faster R-CNN ensures quick and accurate localization of humans and objects.

REFERENCES

- 1 Bangpeng Yao, Li Fei-Fei "Recognizing Human-Object Interactions in Still Images by Modeling the Mutual Context of Objects and Human Poses", IEEE, 2012
- 2 Georgia Gkioxari, Ross Girshick, Piotr Dollár, Kaiming He "Detecting and Recognizing Human-Object Interactions", IEEE, 2018
- 3 Weihao Yan, Yue Gao, Qiming Liu "Human-object Interaction Recognition Using Multitask Neural Network", IEEE, 2019
- 4 Hailan Kuang, Zhikun Zheng, Xinhua Liu, Xiaolin Ma "A Human-Object Interaction Detection Method Inspired by Human Body Part Information", IEEE, 2020

- 5 Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, Hyunwoo J. Kim "HOTR: End-to-End Human-Object Interaction Detection with Transformers", IEEE, 2021
- 6 Xu Sun, Yunqing He, Tongwei Ren, Gangshan Wu "Spatial-Temporal Human-Object Interaction Detection", IEEE, 2021
- 7 Yiming Gao, Zhanghui Kuang, Guanbin Li, Wayne Zhang, Liang Lin "Hierarchical Reasoning Network for Human-Object Interaction Detection", IEEE, 2021
- Manli Zhu, Edmond S. L. Ho, and Hubert P. H. Shum "A Skeleton-Aware Graph Convolutional Network For Human-Object Interaction Detection", IEEE, 2022

THANK YOU