

Human Action Recognition using Detectron2 and LSTM

(Learn OpenCV and https://learnopencv.com)



Bibin Sebastian (<https://learnopencv.com/author/bibin-sebastian/>)

JULY 26, 2021

[Application](https://learnopencv.com/category/application/) (<https://learnopencv.com/category/application/>) [Deep Learning](https://learnopencv.com/category/deep-learning/) (<https://learnopencv.com/category/deep-learning/>) [Pose](https://learnopencv.com/category/pose/) (<https://learnopencv.com/category/pose/>)

[PyTorch](https://learnopencv.com/category/pytorch/) (<https://learnopencv.com/category/pytorch/>)

This exciting post comes to you from Bibin Sebastian, who took our [Deep Learning with Pytorch](https://courses.opencv.org/) (<https://courses.opencv.org/>) course by OpenCV, and then applied all the learning to create a Human Action Recognition application using PyTorch.

Isn't that amazing?

Let's read on to find out, how he build that. And be inspired to build one more.



Action: JUMPING

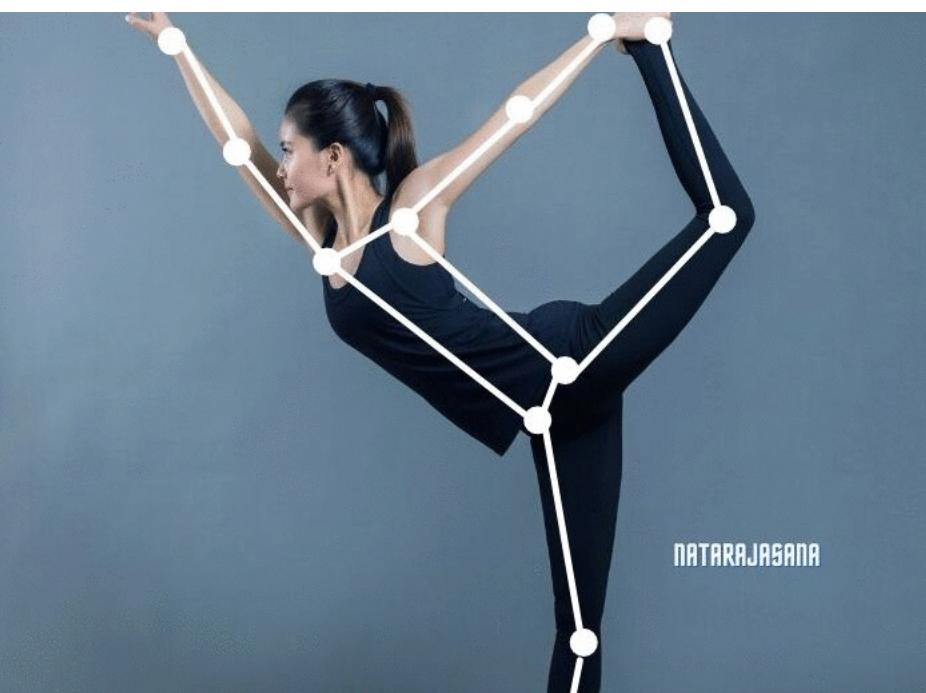
(<https://learnopencv.com/wp-content/uploads/2021/07/action-recognition.gif>)

Action recognition result based on key points detected on the human body

Human action recognition involves analyzing the video footage to predict or classify various actions performed by the person in that video. It is widely applied in diverse fields like surveillance, sports, fitness, and defense.

Let's say you want to build an application for teaching [Yoga](https://en.wikipedia.org/wiki/Yoga) (<https://en.wikipedia.org/wiki/Yoga>) online. It should offer a list of pre-recorded yoga session videos for users to watch. After watching a video on the app, users can upload the videos of their personal practice sessions. The app then evaluates their performance and gives feedback based on how well the user has performed the various yoga asanas (or poses). Wouldn't it be great to use action recognition to automate the evaluation of the video? Well, there's more you can do with it. Check out the video below.

The yoga application shown below uses human pose estimation to identify each yoga pose and identifies it as one of the following asanas – Natarajasana, Trikonasana, or Virabhadrasana.



NATARAJASANA

(<https://learnopencv.com/wp-content/uploads/2021/07/yoga-poses.gif>)

Yoga poses (Natarajasana, Trikonasana or Virabhadrasana) identified based on keypoints detected on the human body

In this post, we will explain how to create such an application for human-action recognition (or classification), using **pose estimation** and **LSTM (Long Short-Term Memory)**. We will create a web application that takes in a video and produces an output video annotated with identified action classes. We will be using the Flask framework for the web application and [PyTorch lightning](https://github.com/PyTorchLightning/pytorch-lightning) (<https://github.com/PyTorchLightning/pytorch-lightning>) for model

Table of Contents

1. Detectron2
2. LSTM
3. Dataset
4. Flask
5. High-level Approach to the Solution
6. Training ML Models
7. Inferencing
8. Code Walkthrough
 - 8.1. Detectron2 Pose-Estimation Model
 - 8.2. LSTM Model Definition
 - 8.3. Web application
 - 8.4. Video Analysis
9. See the application in action
10. Conclusion

Beyond Flask, we will be deploying several other important toolsets like Detectron2, LSTM, Dataset etc. to reach our goal. Each is discussed in detail here.

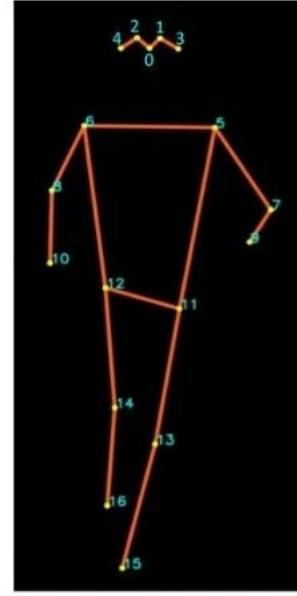
Detectron2

[Detectron2](https://github.com/facebookresearch/detectron2/) (<https://github.com/facebookresearch/detectron2/>) is Facebook AI Research's open source platform for object detection, dense pose, segmentation and other visual recognition tasks. The platform is now implemented in PyTorch, unlike its previous version, [Detectron](https://github.com/facebookresearch/Detectron/) (<https://github.com/facebookresearch/Detectron/>), which was implemented in Caffe2.

Here, we use a pre-trained 'R50-FPN' model from the Detectron2 model zoo for pose estimation. This model is already trained on the COCO dataset containing more than 200,000 images and 250,000 person instances, labelled with keypoints. The model outputs 17 keypoints for every human present in the input image frame, as shown in the image below.



Index	Key point
0	Nose
1	Left-eye
2	Right-eye
3	Left-ear
4	Right-ear
5	Left-shoulder
6	Right-shoulder
7	Left-elbow
8	Right-elbow
9	Left-wrist
10	Right-wrist
11	Left-hip
12	Right-hip
13	Left-knee
14	Right-knee
15	Left-ankle
16	Right-ankle



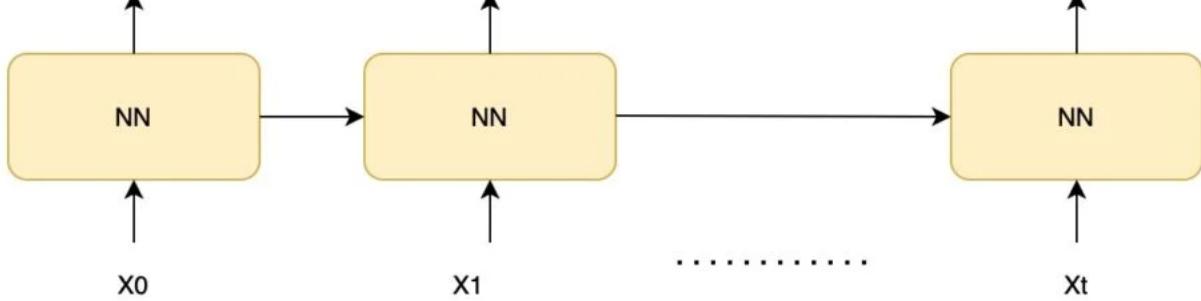
(<https://learnopencv.com/wp-content/uploads/2021/07/keypoints.jpg>).

17 keypoints on a human body. Left part of the image shows a person, the middle part shows a list of keypoints and the right part shows the location of keypoints on the person's body

Want to know more about Pose Estimation algorithm internals? checkout this [blog post](https://learnopencv.com/human-pose-estimation-using-keypoint-rcnn-in-pytorch/) (<https://learnopencv.com/human-pose-estimation-using-keypoint-rcnn-in-pytorch/>).

LSTM

A type of Recurrent Neural Network ([RNN](https://en.wikipedia.org/wiki/Recurrent_neural_network) (https://en.wikipedia.org/wiki/Recurrent_neural_network)), [LSTM](https://en.wikipedia.org/wiki/Long_short-term_memory) (https://en.wikipedia.org/wiki/Long_short-term_memory) networks are capable of learning-order dependence in sequence-prediction problems. An RNN, as you can see below, has a chain of repeating neural-network modules.



(<https://learnopencv.com/wp-content/uploads/2021/07/recurrent-neural-network.jpg>)

A Recurrent Neural Network with its chain of repeated neural-network modules

In the NN (neural network):

- X_0, X_1, \dots, X_t are the inputs, and h_0, h_1, \dots, h_t are the predictions.
- Every prediction at time t (h_t) depends on the previous prediction and current input X_t .

RNN remembers the previous information and uses it optimally to process the current input. But the shortcoming of RNN is that it cannot remember long-term dependencies.

LSTM also has a similar chain structure, but its neural-network module can easily handle long-term dependencies.

We use LSTM to do action classification on a sequence of keypoint detections from a video.



(<http://opencv.org/courses/>)

Official OpenCV Courses

Start your exciting journey from an absolute Beginner to Mastery in AI, Computer Vision & Deep Learning!

[Learn More](#)
(<https://opencv.org/courses/>)

Dataset

To train the LSTM model we use [this dataset](#). (<https://github.com/stuartteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input#dataset-overview>)

What's so special about this dataset? It consists of keypoint detections, made using [OpenPose](#) (<https://github.com/CMU-Perceptual-Computing-Lab/openpose>) deep-learning model, on a subset of the [Berkeley Multimodal Human Action Database \(MHAD\)](#) (http://tele-immersion.citris-uc.org/berkeley_mhad) dataset.

OpenPose is the first, real-time, multi-person system to jointly detect human body, hand, facial, and foot key-points (in total 135 key-points) on single images. Keypoint detections are made over videos of 12 subjects (filmed from 4 angles), doing the following 6 actions for 5 repetitions:



(<https://learnopencv.com/wp-content/uploads/2021/07/human-actions.jpg>).

A series of six human actions (JUMPING, JUMPING_JACKS, BOXING, WAVING_2HANDS, WAVING_1HAND, CLAPPING_HANDS) from a subset of the MHAD dataset.

- JUMPING
- JUMPING_JACKS
- BOXING
- WAVING_2HANDS
- WAVING_1HAND
- CLAPPING_HANDS

Flask

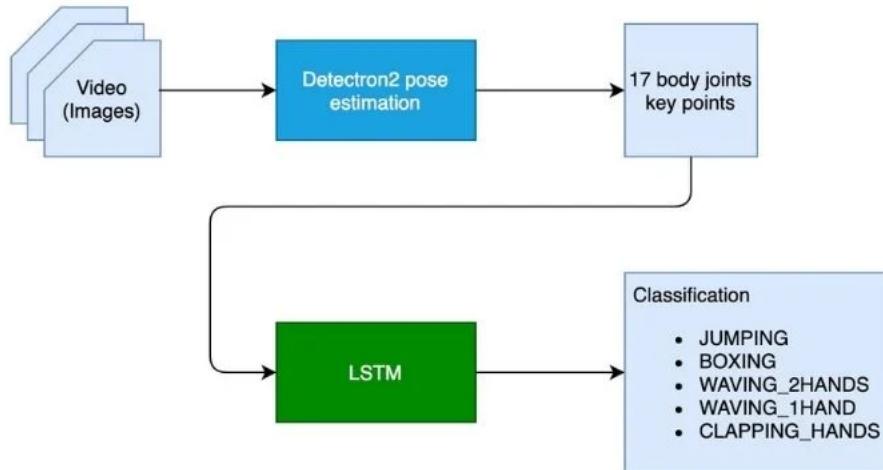
Flask is a popular Python web framework used for developing several web applications. This application internally uses Detectron2 and LSTM model to identify the actions.

High-level Approach to Activity Recognition

To classify an action, we first need locate various body parts in every frame, and then analyze the movement of the body parts over time.

The first step is achieved using Detectron2 which outputs the body posture (17 key points) after observing a single frame in a video.

The second step of analyzing the motion of the body over time and making a prediction is done using the LSTM network. So, keypoints from a sequence of frames are sent to LSTM for action classification, as shown below.



(<https://learnopencv.com/wp-content/uploads/2021/07/data-flow-diagram.jpg>).

End to end action recognition workflow using Detectron2 and LSTM

Training ML Models for Action Recognition

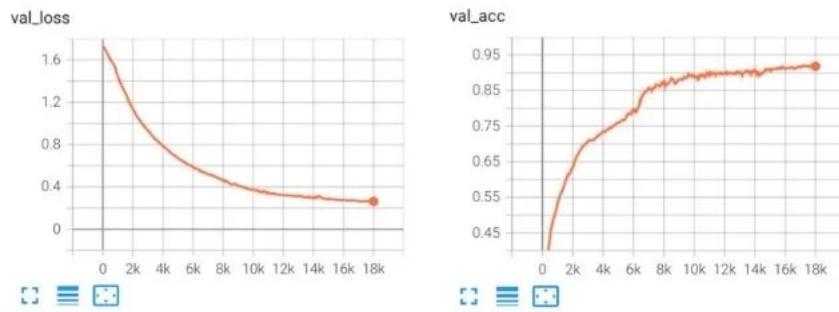
1. As we had mentioned before, for keypoint detection, we are using the pre-trained 'R50-FPN' model from Detectron2 [model zoo](#) (https://github.com/facebookresearch/detectron2/blob/master/MODEL_ZOO.md#coco-person-keypoint-detection-baselines-with-keypoint-r-cnn). So no further training is required.
2. The LSTM model which is used for action classification based on keypoints is trained with [pytorch lightning](#) (<https://github.com/PyTorchLightning/pytorch-lightning>).

$$\begin{bmatrix} [x_0, y_0, x_1, y_1, \dots, x_{10}, y_{10}] \\ [x_0, y_0, x_1, y_1, \dots, x_{16}, y_{16}] \\ \vdots \\ [x_0, y_0, x_1, y_1, \dots, x_{16}, y_{16}] \end{bmatrix}$$

Every row contains 17 keypoint values. Every keypoint is represented as (x,y) values, hence a total of 34 values per row.

Note: Unlike the 18 keypoints of a human body detected by the OpenPose model in the original dataset, our application has just 17 keypoints detected by Detectron2. So, we convert to a 17 keypoints format before training our LSTM model.

We trained our model for **400 epochs** and got a validation accuracy of **0.913**. Validation accuracy and loss curves are shown below. The trained model is checked into the codebase, and the same is used during inferencing.



(<https://learnopencv.com/wp-content/uploads/2021/07/validation-loss-and-accuracy.jpg>).

Validation Loss and Validation Accuracy

You can obtain the LSTM model training by clicking on the button below.

Download Code To easily follow along this tutorial, please download code by clicking on the button below. It's FREE!

[Download Code](#)

Inferencing

The inference pipeline consists of both the Detectron2 model and a custom LSTM model.

- Our application accepts a video input, iterates through the frames and then uses Detectron2 to do keypoint detection on every frame.
- Keypoint results are next appended to a buffer of size 32, which operates in a sliding window fashion.
- Contents of the buffer are finally sent to our trained LSTM model for action identification.
- Also, the Flask-based web application has a UI to accept video input from the user.
- Actions detected by our inference pipeline are annotated on the video and displayed as the result.

Testing on 'Tesla T4' GPU shows that Detectron2 takes about **0.14 seconds**, and LSTM **0.002 seconds** respectively for inferencing. Hence, the combined **FPS** (frames per second) of our inference pipeline execution comes to about **6 frames per second**, i.e., if we process every frame in the video.

The above FPS rate might work for an application doing video analysis offline. But what if you are inferencing on a real-time video stream? In general, real-time video streams have a frame rate of 30 FPS or more (depending on the camera). In such cases, the FPS of the inference pipeline has to be higher than or at least equal to the video stream FPS to process the frames without any lag. Though few, you do have options to improve the FPS of the inference pipeline.

- [Pruning and Quantization of the model](https://arxiv.org/abs/2101.09671) (<https://arxiv.org/abs/2101.09671>) can accelerate the speed of execution.
- Skipping frames and inferencing at intervals: You can even skip a couple of frames and opt to infer only at intervals. For example, our testing shows that when inferencing only on every 5th frame in the sequence, FPS increases to **27 frames per second**. But we saw that the accuracy dropped. So choose an interval accordingly.
- Multi threading: Have separate threads for receiving video and inferencing.
 - Receiver thread can focus on reading the video frames from the stream, and adding them to a queue.
 - A separate child thread can read frames from the queue to do inferencing. A child thread may lag behind while processing the frames, but it won't block the receiver thread from reading the video streams.

Coding Detectron2, LSTM Models For Video Analysis on Web Application

Let's now understand how the important components of the application are coded.

```

1 # obtain detectron2's default config
2 cfg = get_cfg()
3 # load the pre trained model from Detectron2 model zoo
4 cfg.merge_from_file(model_zoo.get_config_file("COCO-Keypoints/keypoint_rcnn_R_50_FPN_3x.yaml"))
5 # set confidence threshold for this model
6 cfg.MODEL.ROI_HEADS.SCORE_THRESH_TEST = 0.5
7 # load model weights
8 cfg.MODEL.WEIGHTS = model_zoo.get_checkpoint_url("COCO-Keypoints/keypoint_rcnn_R_50_FPN_3x.yaml")
9 # create the predictor for pose estimation using the config
10 pose_detector = DefaultPredictor(cfg)

```

2. LSTM model definition

Our LSTM model is initialised with a hidden dimension (hidden_dim) of 50 and is trained with PyTorch Lightning. We have used the [Adam](https://pytorch.org/docs/stable/generated/torch.optim.Adam.html#torch.optim.Adam) (<https://pytorch.org/docs/stable/generated/torch.optim.Adam.html#torch.optim.Adam>) optimiser and also configured the [ReduceLROnPlateau](https://pytorch.org/docs/stable/generated/torch.optim.lr_scheduler.ReduceLROnPlateau.html) (https://pytorch.org/docs/stable/generated/torch.optim.lr_scheduler.ReduceLROnPlateau.html). scheduler to reduce the learning rate, based on the value of val_loss.

```

1 # We have 6 output action classes.
2 TOT_ACTION_CLASSES = 6
3
4 #lstm classifier definition
5 class ActionClassificationLSTM(pl.LightningModule):
6
7     # initialise method
8     def __init__(self, input_features, hidden_dim, learning_rate=0.001):
9         super().__init__()
10        # save hyperparameters
11        self.save_hyperparameters()
12        # The LSTM takes word embeddings as inputs, and outputs hidden states
13        # with dimensionality hidden_dim.
14        self.lstm = nn.LSTM(input_features, hidden_dim, batch_first=True)
15        # The linear layer that maps from hidden state space to classes
16        self.linear = nn.Linear(hidden_dim, TOT_ACTION_CLASSES)
17
18    def forward(self, x):
19        # invoke lstm layer
20        lstm_out, (ht, ct) = self.lstm(x)
21        # invoke linear layer
22        return self.linear(ht[-1])
23
24    def training_step(self, batch, batch_idx):
25        # get data and labels from batch
26        x, y = batch
27        # reduce dimension
28        y = torch.squeeze(y)
29        # convert to long
30        y = y.long()
31        # get prediction
32        y_pred = self(x)
33        # calculate loss
34        loss = F.cross_entropy(y_pred, y)
35        # get probability score using softmax
36        prob = F.softmax(y_pred, dim=1)
37        # get the index of the max probability
38        pred = prob.data.max(dim=1)[1]
39        # calculate accuracy
40        acc = torchmetrics.functional.accuracy(pred, y)
41        dic = {
42            'batch_train_loss': loss,
43            'batch_train_acc': acc
44        }
45        # log the metrics for pytorch lightning progress bar or any other operations
46        self.log('batch_train_loss', loss, prog_bar=True)
47        self.log('batch_train_acc', acc, prog_bar=True)
48        #return loss and dict
49        return {'loss': loss, 'result': dic}
50
51    def training_epoch_end(self, training_step_outputs):
52        # calculate average training loss end of the epoch
53        avg_train_loss = torch.tensor([x['result']['batch_train_loss'] for x in training_step_outputs]).mean()
54        # calculate average training accuracy end of the epoch
55        avg_train_acc = torch.tensor([x['result']['batch_train_acc'] for x in training_step_outputs]).mean()
56        # log the metrics for pytorch lightning progress bar and any further processing
57        self.log('train_loss', avg_train_loss, prog_bar=True)
58        self.log('train_acc', avg_train_acc, prog_bar=True)
59
60    def validation_step(self, batch, batch_idx):
61        # get data and labels from batch
62        x, y = batch
63        # reduce dimension
64        y = torch.squeeze(y)
65        # convert to long
66        y = y.long()
67        # get prediction
68        y_pred = self(x)
69        # calculate loss
70        loss = F.cross_entropy(y_pred, y)
71        # get probability score using softmax
72        prob = F.softmax(y_pred, dim=1)
73        # get the index of the max probability
74        pred = prob.data.max(dim=1)[1]
75        # calculate accuracy
76        acc = torchmetrics.functional.accuracy(pred, y)
77        dic = {
78            'batch_val_loss': loss,
79            'batch_val_acc': acc
80        }
81        # log the metrics for pytorch lightning progress bar and any further processing
82        self.log('batch_val_loss', loss, prog_bar=True)
83        self.log('batch_val_acc', acc, prog_bar=True)
84        #return dict
85        return dic
86
87    def validation_epoch_end(self, validation_step_outputs):
88        # calculate average validation loss end of the epoch

```

```

95     self.log('val_loss', avg_val_loss, prog_bar=True)
96     self.log('val_acc', avg_val_acc, prog_bar=True)
97
98 def configure_optimizers(self):
99     # adam optimiser
100    optimizer = optim.Adam(self.parameters(), lr=self.hparams.learning_rate)
101    # learning rate reducer scheduler
102    scheduler = optim.lr_scheduler.ReduceLROnPlateau(optimizer, mode='min', factor=0.5, patience=10, min_lr=1e-15,
103    verbose=True)
104    # scheduler reduces learning rate based on the value of val_loss metric
105    return {"optimizer": optimizer,
106            "lr_scheduler": {"scheduler": scheduler, "interval": "epoch", "frequency": 1, "monitor": "val_loss"}}

```

3. Web Application

Our web application has several defined routes. The one below processes the input video. This route gets called when the user submits a video from the webpage for analysis.

```

1 # route definition for video upload for analysis
2 @app.route('/analyze/<filename>')
3 def analyze(filename):
4     # invokes method analyse_video
5     return Response(analyse_video(pose_detector, lstm_classifier, filename), mimetype='text/event-stream')

```

4. Video Analysis

Once our web application receives a video from the user,

- The below function parses through the video, and invokes: Detectron2 for keypoint detection, and LSTM for action classification. You can see that we have used the 'buffer_window' to store 32 consecutive keypoint results from frames, and the same is used for inferencing action classes.
- Next, we use OpenCV to read the input video, and also to create the output video with classification results and pose-estimation results.

If you want a higher FPS rate for video analysis, you can configure a higher value for SKIP_FRAME_COUNT.

```

1 # how many frames to skip while inferencing
2 # configuring a higher value will result in better FPS (frames per rate), but accuracy might get impacted
3 SKIP_FRAME_COUNT = 0
4
5 # analyse the video
6 def analyse_video(pose_detector, lstm_classifier, video_path):
7     # open the video
8     cap = cv2.VideoCapture(video_path)
9     # width of image frame
10    width = int(cap.get(cv2.CAP_PROP_FRAME_WIDTH))
11    # height of image frame
12    height = int(cap.get(cv2.CAP_PROP_FRAME_HEIGHT))
13    # frames per second of the input video
14    fps = int(cap.get(cv2.CAP_PROP_FPS))
15    # total number of frames in the video
16    tot_frames = int(cap.get(cv2.CAP_PROP_FRAME_COUNT))
17    # video output codec
18    fourcc = cv2.VideoWriter_fourcc(*'mp4v')
19    # extract the file name from video path
20    file_name = ntpath.basename(video_path)
21    # video writer
22    vid_writer = cv2.VideoWriter('res_{}.format(23
23        file_name), fourcc, 30, (width, height))
24    # counter
25    counter = 0
26    # buffer to keep the output of detectron2 pose estimation
27    buffer_window = []
28    # start time
29    start = time.time()
30    label = None
31    # iterate through the video
32    while True:
33        # read the frame
34        ret, frame = cap.read()
35        # return if end of the video
36        if ret == False:
37            break
38        # make a copy of the frame
39        img = frame.copy()
40        if(counter % (SKIP_FRAME_COUNT+1) == 0):
41            # predict pose estimation on the frame
42            outputs = pose_detector(frame)
43            # filter the outputs with a good confidence score
44            persons, pIndices = filter_persons(outputs)
45            if len(persons) >= 1:
46                # pick only pose estimation results of the first person.
47                # actually, we expect only one person to be present in the video.
48                p = persons[0]
49                # draw the body joints on the person body
50                draw_keypoints(p, img)
51                # input feature array for lstm
52                features = []
53                # add pose estimate results to the feature array
54                for i, row in enumerate(p):
55                    features.append(row[0])
56                    features.append(row[1])
57
58                # append the feature array into the buffer
59                # note that max buffer size is 32 and buffer_window operates in a sliding window fashion
60                if len(buffer_window) < WINDOW_SIZE:
61                    buffer_window.append(features)
62                else:
63                    # convert input to tensor
64                    model_input = torch.Tensor(np.array(buffer_window, dtype=np.float32))
65                    # add extra dimension
66                    model_input = torch.unsqueeze(model_input, dim=0)
67                    # predict the action class using lstm

```

```

73     buffer_window.pop(0)
74     buffer_window.append(features)
75     label = LABELS[pred_index.numpy()[0]]
76     #print("Label detected ", label)
77
78     # add predicted label into the frame
79     If label is not None:
80         cv2.putText(img, 'Action: {}'.format(label),
81             (int(width-400), height-50), cv2.FONT_HERSHEY_COMPLEX, 0.9, (102, 255, 255), 2)
82         # increment counter
83         counter += 1
84         # write the frame into the result video
85         vid_writer.write(img)
86         # compute the completion percentage
87         percentage = int(counter*100/tot_frames)
88         # return the completion percentage
89         yield "data:" + str(percentage) + "\n\n"
90     analyze_done = time.time()
91     print("Video processing finished in ", analyze_done - start)

```

See the Application in Action

You need a GPU to run these ML models. To make it easy for you to see the application in action, we have provided a [google colab](#) (<https://research.google.com/colaboratory/>) notebook in the codebase. Just upload the notebook in Google Colab, enable the GPU and run the notebook.



Video Analysis Progress bar

100%

Video Analysis Result

[Download Result File](#)



Action: WAVING_1HAND

(<https://learnopencv.com/wp-content/uploads/2021/07/application-ui.jpg>).

User interface (UI) of the application

Conclusion

Let's summarise what we have learned so far.

- You learned all that is involved in human-action recognition and also about their diverse applications
- Next, we discussed how to create an application for action recognition.
- You clearly understood why we chose Detectron2 and LSTM for our solution as we went over the capabilities of each.

inference pipeline that were annotated on the input video.

- Understood why the FPS must be optimised for real-time video stream inferencing. For this you explored various approaches like pruning and quantization of the model, skipping frames and inferencing at intervals, and also, multi threading.
- Next you examined in detail the given code for all the important components in the application.
- Finally, you went ahead and built a sample Flask-based web application to do inferencing on any video input.

Hope the contents were useful and you were able to learn something new from this blog post.

Subscribe & Download Code

If you liked this article and would like to download code (C++ and Python) and example images used in this post, please [click here](#). Alternately, sign up to receive a free [Computer Vision Resource Guide](#). In our newsletter, we share OpenCV tutorials and examples written in C++/Python, and Computer Vision and Machine Learning algorithms and news.

[Download Example Code](#)

Subscribe Now

Disclaimer

All views expressed on this site are my own and do not represent the opinions of OpenCV.org or any entity whatsoever with which I have been, am now, or will be affiliated.

 (<https://www.facebook.com/Learnopencv-277284839758/>) [tag me on my posts!](#) [click here](#)

Course

[OpenCV Courses](#)

[CV4Faces \(Old\)](#)

About LearnOpenCV

In 2007, right after finishing my Ph.D., I co-founded TAAZ Inc. with my advisor Dr. David Kriegman and Kevin Barnes. The scalability, and robustness of our computer vision and machine learning algorithms have been put to rigorous test by more than 100M users who have tried our products.

[Read More](#)

(<https://learnopencv.com/about/>)

Getting Started

[Installation](#)

[PyTorch](#)

[Getting Started with OpenCV](#)

[Keras & Tensorflow](#)

[Resource Guide](#)

Information

[Privacy Policy](#)

[Terms and Conditions](#)

Copyright © 2021 – BIG VISION LLC