

Applied Accelerated Artificial Intelligence
Prof. Satyadhyan Chickerur
Department of Computer Science and Engineering
Indian Institute of Technology, Palakkad

Fundamentals of Accelerating Deployments

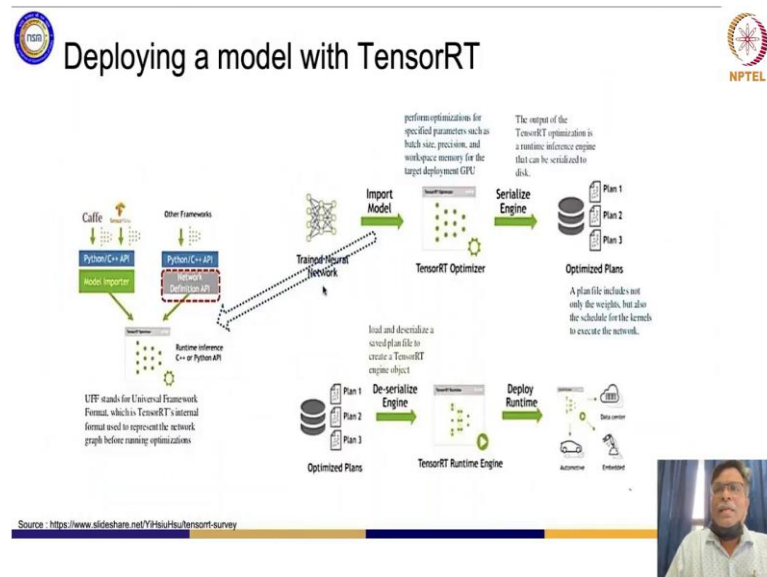
Lecture - 42

Accelerating neural network inference in PyTorch and TensorFlow Part 1

Good evening everyone. Today I will try to actually do some demos on Accelerating deployments ok. And as I showed you yesterday some videos maybe today we will show you the demonstration by again connecting a webcam there is some issue with the team viewer, but we will anyway show you some more videos today in case that does not do materialize right.

Because it is getting struck the team viewer video is getting stuck. So, no issues for that but, let us start with certain things wherein we will try to show you the maximum demo stuff ok.

(Refer Slide Time: 01:00)



So, yeah so, yesterday actually we ended up at this point right that how do you actually deploy a model with TensorRT ok. And we showed you this particular concept we showed you a demo and we actually discussed this yesterday that we have either a TensorFlow or a caffe or any other framework which basically is written using python ok.

And then you have a model importer which actually imports ok using a TensorRT ok right. So, the idea is you basically come up with a model which is a universal framework format or at present you can use ONNX as well right. So, I have got two options of actually converting your actual model which you have trained it on a.

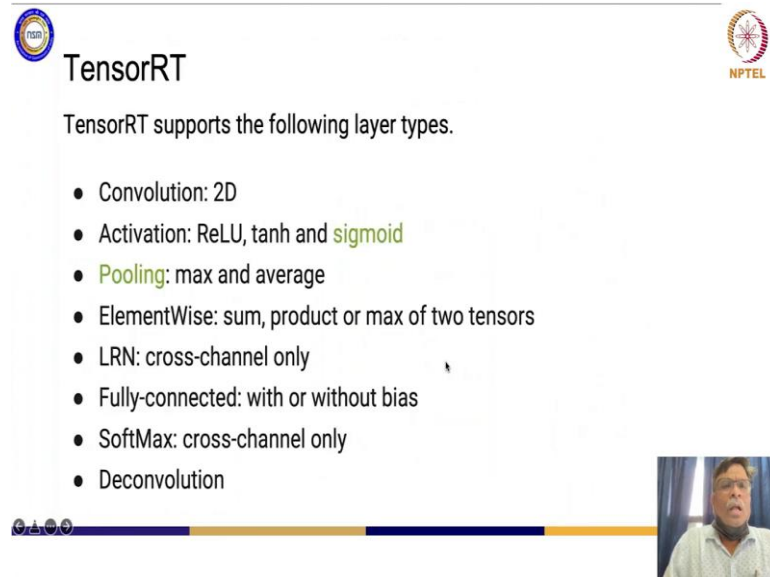
Students: Server.

Server right and then convert it into something which runs on a inferencing edge device right. So, you can use UFF or you can use ONNX. So, we can show we will show you ONNX also and we have done it using UFF also right. So, let us see how we can show you both of this right. So, the basic idea what we discussed yesterday was you have a trained neural network model you import that model and then try to use a TensorRT optimizer right for trying to optimize it. So, that it runs on your edge device right.

So, for that reason you actually have this process which is internally done right. You have to serialize it then you generate one or three optimized plans and then you get a plan file right with actually schedules ok for the kernel to execute whatever it is to be executed.

So, you take this optimized plan ok and then as shown you basically create a Tensor RT engine object and then you run it using a Tensor RT runtime engine and that is how you basically try to deploy on the edge device. So, basically you are trying to deploy the runtime. So, this is the overall model which we discussed yesterday right.

(Refer Slide Time: 03:51)



The slide features a blue circular logo on the top left and the NPTEL logo on the top right. The title 'TensorRT' is prominently displayed. Below the title, the text states 'TensorRT supports the following layer types.' followed by a bulleted list of supported operations. A small video inset of a speaker is visible in the bottom right corner of the slide area.

TensorRT

TensorRT supports the following layer types.

- Convolution: 2D
- Activation: ReLU, tanh and sigmoid
- Pooling: max and average
- ElementWise: sum, product or max of two tensors
- LRN: cross-channel only
- Fully-connected: with or without bias
- SoftMax: cross-channel only
- Deconvolution

So, now today we will try to understand that when we are trying to work with TensorRT, TensorRT supports the following layer types. What does it mean? It means if you are using convolution layers if you are using activation layers which involve ReLU, tanh, sigmoid and you have this pooling max and average stuff. Elementwise sum, product or max of two tensors ok and then LRN and then fully connected with or without bias SoftMax, deconvolution.

All of these layers could be converted from TensorFlow or PyTorch to a Tensor RT survey right this is what it means.

(Refer Slide Time: 04:40)

The slide is titled "Steps for TensorRT application." and features the logos of IIT Bombay and NPTEL in the top corners. It lists four steps for the process:

1. Convert the pretrained image segmentation PyTorch model into ONNX.
2. Import the ONNX model into TensorRT.
3. Apply optimizations and generate an engine.
4. Perform inference on the GPU.



Below the list is a flow diagram with three green boxes: "Import ONNX Model", "Build Engine", and "Perform Inference". An arrow labeled "ONNX Model" points to the first box, and arrows connect the boxes in sequence. A small video inset of a speaker is visible in the bottom right corner of the slide.

And now that effectively means that whenever you are trying to convert ok any model for inferencing the steps would be like something like this you actually convert for example, if you are trying to do a image segmentation application deployment. So, you convert the pretrained image that segmentation PyTorch model or for that matter the tensor flow model into ONNX or ok or UFF anything.

And then you import the ONNX model into Tensor RT model ok this conversion needs to be done and then you basically apply several more optimization generate an engine. So, technically speaking once your model gets converted into a Tensor RT runtime ok you can actually do the inferencing on that particular edge device right. So, we have written it on the GPU. So, you can do it on GPU you can do it on edge device you can do it on Jetson any of these right. So, this is how basically it is done ok.

So, there are various ways of doing it one of the ways what we are trying to tell you is this ok you can actually see in literature and various other places there are various ways of doing it one of the ways is you either convert your TensorFlow program into a TensorRT program through UFF or ok through ONNX. So, this is the basic idea ok.


(Refer Slide Time: 06:27)



contd..

The application uses the following components in TensorRT:

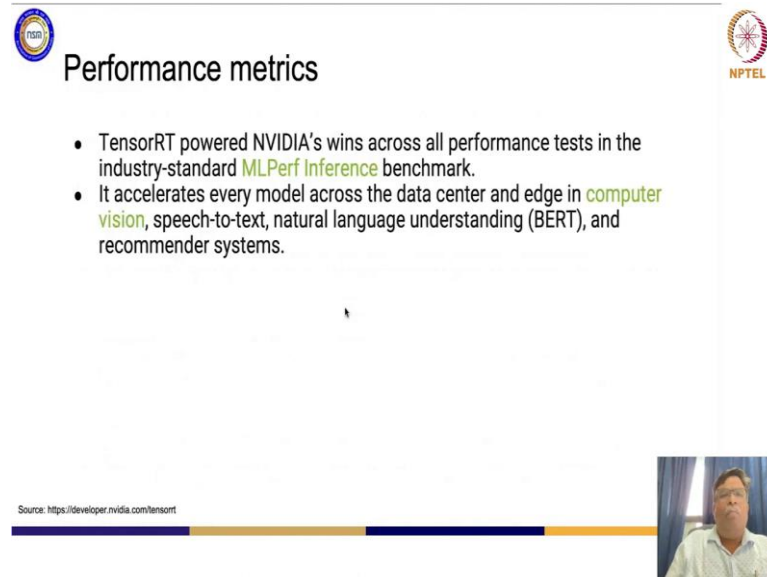
- **ONNX parser:** Takes a converted any trained model into the ONNX format as input and populates a network object in TensorRT.
- **Builder:** Takes a network in TensorRT and generates an engine that is optimized for the target platform.
- **Engine:** Takes input data, performs inferences, and emits inference output.
- **Logger:** Associated with the builder and engine to capture errors, warnings, and other information during the build and inference phases



So, you have got various components in TensorRT you have a ONNX parser which basically converts ok or you have a UFF parser you basically take any trained model into the ONNX format as input and you get a network object in tensor RT. This is possible ok and then you take a network in Tensor RT and you generate an engine that is optimized for the target platform.

So, this target platform can be anything. So, ultimately you basically take a network in TensorRT and generate an engine similarly your engine takes an input data perform inferences and then emits the inferencing output and logger basically tries to give you warnings capture errors and how you run and everything right.




(Refer Slide Time: 07:24)



Performance metrics

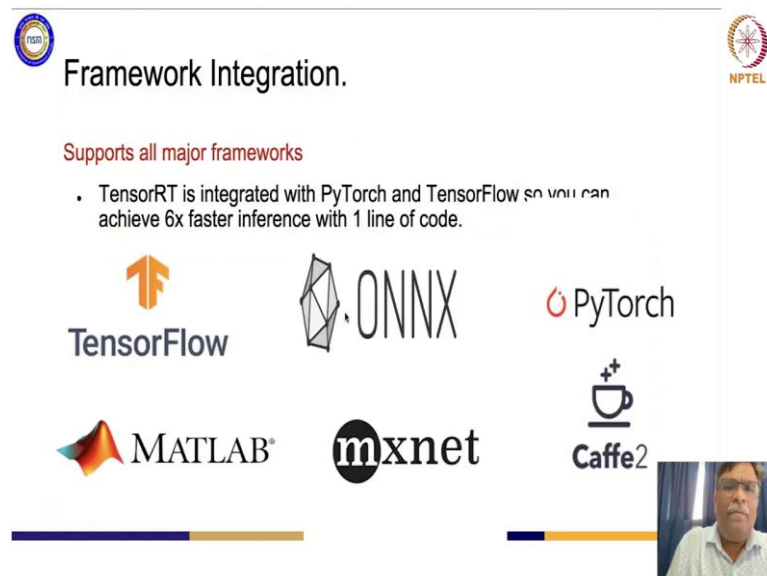
- TensorRT powered NVIDIA's wins across all performance tests in the industry-standard **MLPerf Inference** benchmark.
- It accelerates every model across the data center and edge in **computer vision**, speech-to-text, natural language understanding (BERT), and recommender systems.

Source: <https://developer.nvidia.com/tensorrt>



So, this basically is various component of a TensorRT. So, TensorRT powered NVIDIA edge devices actually if you see their benchmarking performances they are quite good they are very very good. So, you can accelerate every model across the data center and edge in computer vision applications speech to text NLP for example, BERT and various recommended systems right. So, all of this would actually be done ok on the edge devices using TensorRT optimizations right.










(Refer Slide Time: 08:04)



Framework Integration.

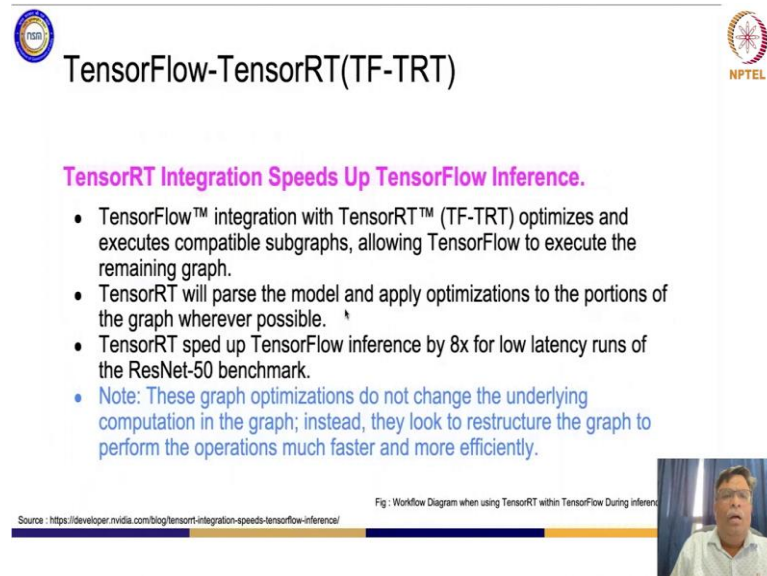
Supports all major frameworks

- TensorRT is integrated with PyTorch and TensorFlow so you can achieve 6x faster inference with 1 line of code.



So, it supports all major frameworks TensorRT is integrated with PyTorch TensorFlow and you can achieve 6x faster inferencing with 1 line of code. So, this is the basic idea you can actually integrate everything ok.

(Refer Slide Time: 08:26)



The slide features the NVIDIA logo on the top left and the NPTEL logo on the top right. The title is "TensorFlow-TensorRT(TF-TRT)". Below the title, there is a list of bullet points:

- TensorFlow™ integration with TensorRT™ (TF-TRT) optimizes and executes compatible subgraphs, allowing TensorFlow to execute the remaining graph.
- TensorRT will parse the model and apply optimizations to the portions of the graph wherever possible.
- TensorRT sped up TensorFlow inference by 8x for low latency runs of the ResNet-50 benchmark.
- Note: These graph optimizations do not change the underlying computation in the graph; instead, they look to restructure the graph to perform the operations much faster and more efficiently.

At the bottom left, there is a source link: <https://developer.nvidia.com/blog/tensorrt-integration-speeds-tensorflow-inference/>. At the bottom right, there is a small video inset showing a man speaking. A caption below the video reads: "Fig : Workflow Diagram when using TensorRT within TensorFlow During inference".

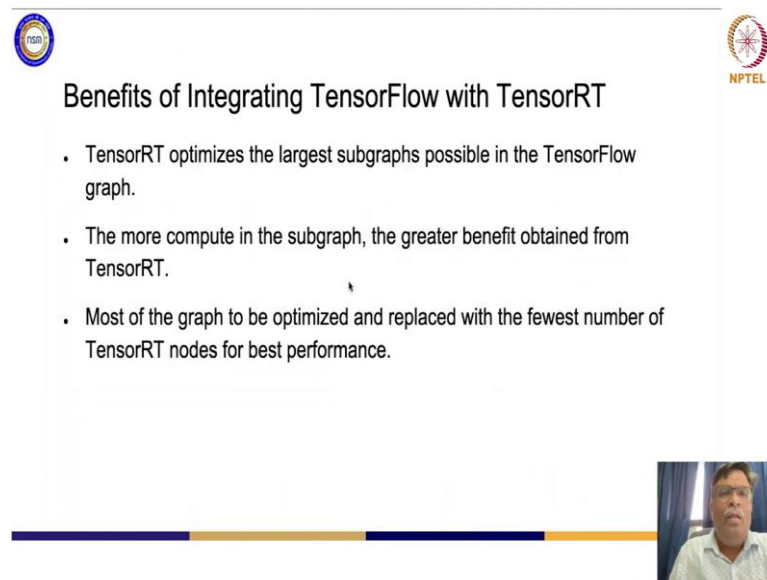
So, TensorRT then from TensorFlow: So, now, the idea is let us say we have a TensorFlow program and this TensorFlow program needs to be converted into a TensorRT program or a TRT program.

So, TF-TRT is what we generally call. So, TensorFlow integration with TensorRT ok; so, when you do this you basically are trying to work with sub graphs ok and then trying to understand how tensor flow ok graph can be converted into an optimized TensorRT graph. So, the effective idea here is your TensorRT is going to generate a model a parse model and then it applied optimizations to the portions of the graph wherever possible. So, we will show you various graphs right in the next few seconds or minutes right.

So, we will show you how basically a saved model a TensorFlow graph is converted into a optimized graph right that also you can visualize. So, technically speaking TensorFlow RT speeds up your inferencing by 8x for a ResNet 50 benchmarks right. So, there are certain things wherein you can see about graph optimizations ok do not change the underlying computations in the graph. So, that is very very important see your computations will not change ok, but the graph will be optimized.

So, basically you are restructuring the graph to perform the operations much faster and more efficient right efficiently right. So, this is the basic test software of it of why do you do actually this and how is that it is going to speed up ok that is also what you are going to actually see, ok.

(Refer Slide Time: 10:31)



The slide features a title "Benefits of Integrating TensorFlow with TensorRT" and three bullet points. In the top left corner is the TensorFlow logo, and in the top right corner is the NPTEL logo. A small video inset in the bottom right corner shows a man speaking. A decorative horizontal bar with blue and yellow segments is located below the slide content.

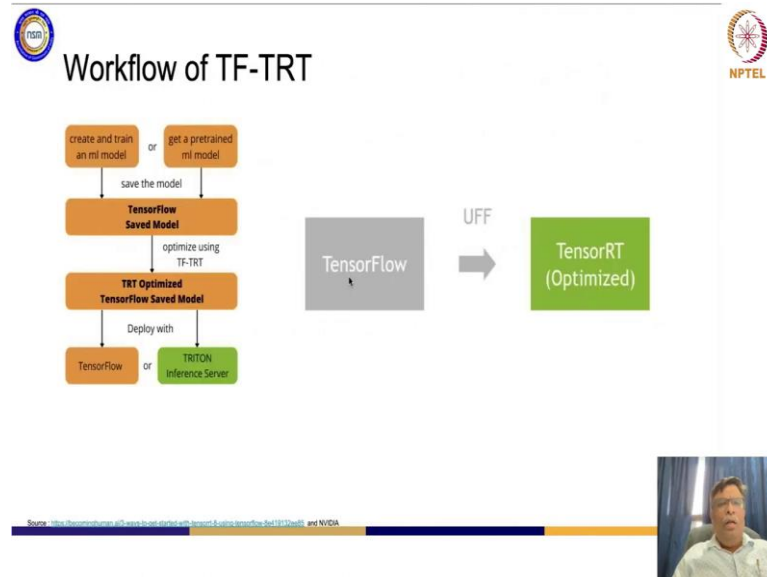
Benefits of Integrating TensorFlow with TensorRT

- TensorRT optimizes the largest subgraphs possible in the TensorFlow graph.
- The more compute in the subgraph, the greater benefit obtained from TensorRT.
- Most of the graph to be optimized and replaced with the fewest number of TensorRT nodes for best performance.

So, now when you integrate the TensorFlow with Tensor RT ok. So, Tensor RT basically optimizes the largest sub graph possible in the TensorFlow graph and if there is more compute in the sub graph; obviously, the TensorFlow optimizes it more ok. So, more computations are there in your model right more optimization you can expect out of TensorRT.

So, most of the graphs which are to be optimized and replaced ok with the fewest number of tensor node RT; so, you generate a TensorRT node actually. So, the idea is you can optimize and replace ok with certain specific less number of nodes as compared to your actual nodes in the graph right. So, that is what actually you are going to be doing right.

(Refer Slide Time: 11:36)

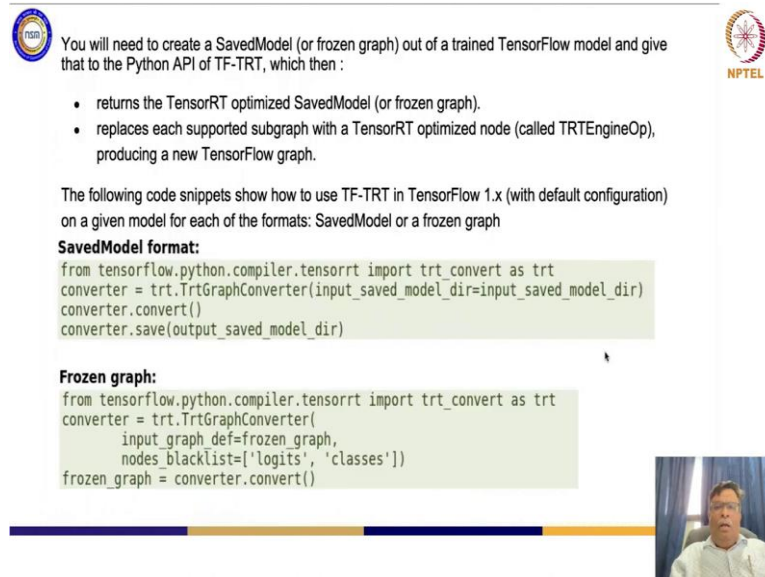


Now, this is how it actually happens you create and train an ML model or you get a pre trained ML model both are possible right you save that node; obviously, you are training your model and then you are saving your model once you save that model you get a TensorFlow save the model there is something which is called as TensorFlow frozen model ok. Now this TensorFlow frozen model cannot be edited or it cannot be changed that is what it means in model trainer right you can actually train a save model again and again.

But once you freeze that model right it has got all the weights parameters everything frozen right. So, you cannot actually use it for training again that is what it means and that is what you are going to optimize ok. So, you are going to optimize that ok saved model from the frozen model and then you deploy with either a train triton server or you do it with TensorFlow.

So, now, this is what I told you that your actual thing of understanding at the first level of complexity is this that you convert your TensorFlow program into a TensorRT optimized programs right by converting it intermediately to a UFF or ONNX right. So, open neural network exchange format ok. So, something like that. So, this is what it means that you are converting your TensorFlow program to a TensorFlow RT, TensorRT program right ok.

(Refer Slide Time: 13:22)



The slide features a blue circular logo on the top left and the NPTEL logo on the top right. The main text explains the process of converting a TensorFlow model to a SavedModel or frozen graph using TF-TRT. It includes a bulleted list of actions and two code snippets for SavedModel and Frozen graph formats. A small video inset of a speaker is visible in the bottom right corner of the slide area.

You will need to create a SavedModel (or frozen graph) out of a trained TensorFlow model and give that to the Python API of TF-TRT, which then :

- returns the TensorRT optimized SavedModel (or frozen graph).
- replaces each supported subgraph with a TensorRT optimized node (called TRTEngineOp), producing a new TensorFlow graph.

The following code snippets show how to use TF-TRT in TensorFlow 1.x (with default configuration) on a given model for each of the formats: SavedModel or a frozen graph

SavedModel format:

```
from tensorflow.python.compiler.tensorrt import trt_convert as trt
converter = trt.TrtGraphConverter(input_saved_model_dir=input_saved_model_dir)
converter.convert()
converter.save(output_saved_model_dir)
```

Frozen graph:


```
from tensorflow.python.compiler.tensorrt import trt_convert as trt
converter = trt.TrtGraphConverter(
    input_graph_def=frozen_graph,
    nodes_blacklist=['logits', 'classes'])
frozen_graph = converter.convert()
```

So, the idea is that you create a SavedModel or a frozen graph. Out of a trained TensorFlow model and give that to the Python API of TF-TRT which is going to return TensorRT optimized model; replaces each supported sub graph with a TensorRT optimized node or a TRTEngine operation producing a new TensorFlow graph.


So, this is how you actually convert ok I use the SavedModel format and this is how you get a frozen graph format. So, if you see here ok. This basic a thing TensorRT import trt converter as trt and then you are using a converter you are just trying to actually save a saved model ok.

And then you actually get a frozen graph by using a converter convert ok and then you get a frozen graph ok. So, this is how basically you would be using it in your program ok this is just the idea of how you are basically using various converters which are available ok. For converting your various models into various formats like a SavedModel format or frozen graph format ok.


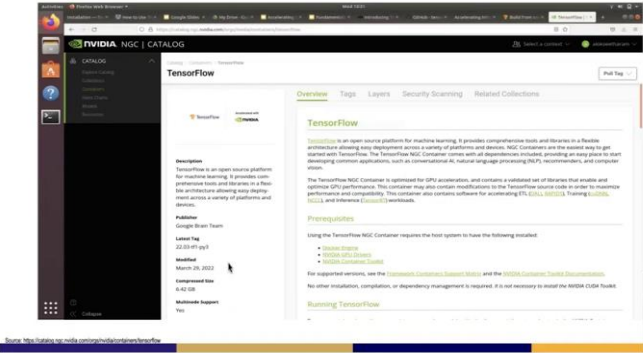
(Refer Slide Time: 14:59)



Installing TF-TRT




NVIDIA containers of TensorFlow are built with enabling TensorRT, which means TF-TRT is part of the TensorFlow binary in the container and can be used out of the box. The container has all the software dependencies required to run TF-TRT




So, one of the things when you want to work with TF-TRT is that you need to install TF-TRT you have NVIDIA containers as we told you in our previous sessions as to how do you pull NGC docker container. So, you pull a NGC docker container and you actually use ok for TensorRT enabled TensorFlow thing right. So, TF-TRT container ok and this basically is supposed to be pulled or installed ok from a TF-TRT container ok from the NGC cloud ok.

(Refer Slide Time: 15:45)




Conversion Parameters in TF-TRT



Parameters that can be passed to `saved_model_cli` and `TrtGraphConverter`.

- `precision_mode`: The precision mode to use (FP32, FP16, or INT8)
- `minimum_segment_size`: The minimum number of TensorFlow nodes required for a TensorRT subgraph to be valid.
- `is_dynamic_op`: TensorRT engines are converted and built at model run time instead of during the `converter.convert()` call. This is required if there are tensors with unknown or dynamic shapes.
- `use_calibration`: Only used if `precision_mode=INT8`. If True, a calibration graph will be created, and `converter.calibrate()` should be called. This is the recommended option. If False, all tensors that will not be fused must have quantization nodes.
- `max_batch_size`: Used when `is_dynamic_op=False`. This is the maximum batch size for TensorRT engines. At run time, smaller batch sizes can be used, but a larger batch size will result in an error.
- `maximum_cached_engines`: Used when `is_dynamic_op=True`. This limits the number of TensorRT engines that are cached, per `TRTEngineOp`.



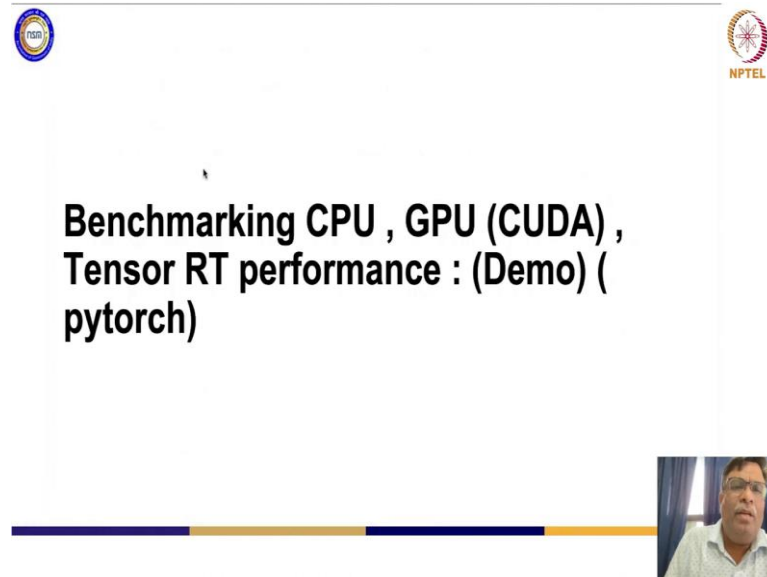
So, what are the various parameters when you try to actually use this ok? So, the parameters which can be passed ok so, this is a command line interface option saved model CLI to TRT graph converter you can work with any precision mode ok. So, the precision mode which you can use is FP32, FP16 or INT8 ok.

What is the minimum segment size? So, that basically means the minimum number of TensorFlow nodes required for a Tensor RT sub graph to be valid ok. This when you go on practicing and running your programs you will try to understand ok all of these particular parameters.

But just to have the completeness of trying to link ok. How these parameters are actually connected we are just trying to give you a very very brief idea to this ok. And then you can use calibration because calibration is needed yesterday I told you that when you are converting from FP32 to FP16 ok, it is done. But when you convert something from FP32 to FP16 to something like INT8 right there is a calibration required right. So, if the precision mode is INT8 the calibration is needed. So, effectively a calibration graph will be created ok.

So, this is what it means ok and then you have the maximum batch size and the maximum cached engine thing right. So, the limit of the number of RT engines which can be cached per TRT engine operation because it depends on the memory of your edge device right based on that how much you can cache is going to actually make your system slower and it can get stuck or something of that sort might happen. So, that is why you are going to use some of these parameters ok.

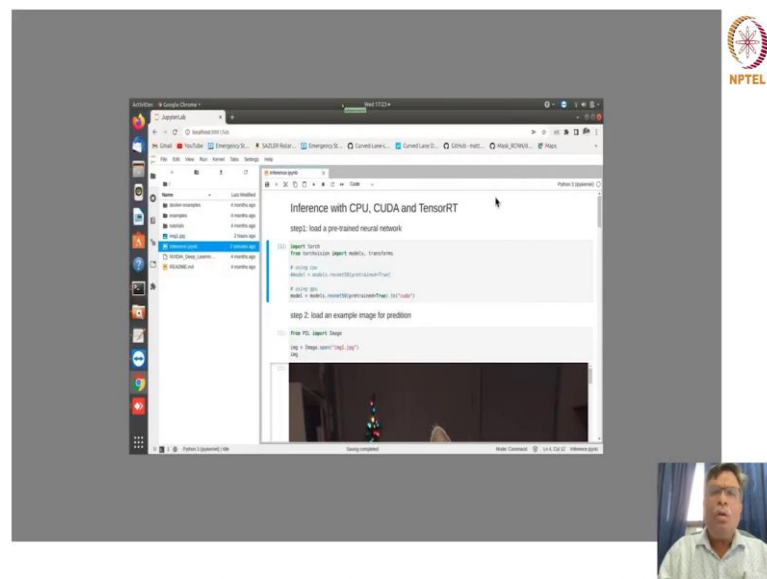
(Refer Slide Time: 17:56)



The slide features the IIT Bombay logo on the top left and the NPTEL logo on the top right. The main content is the title "Benchmarking CPU, GPU (CUDA), Tensor RT performance : (Demo) (pytorch)" centered on the slide. A small video inset of the presenter is located in the bottom right corner of the slide area.

So, let us try to do some benchmarking ok.

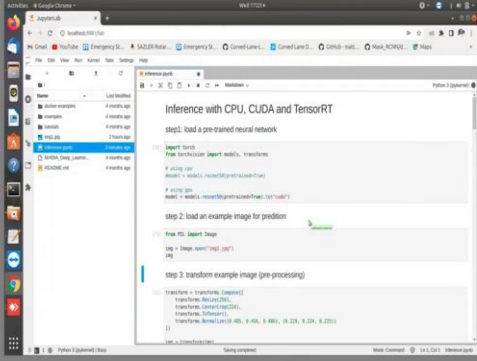
(Refer Slide Time: 17:59)



The screenshot shows a Jupyter Notebook interface. The notebook title is "Inference with CPU, CUDA and TensorRT". The code is divided into two sections: "step 1: load a pre-trained neural network" and "step 2: load an example image for prediction". The code in step 1 includes imports for torch, torchvision, and tensorrt, and the loading of a pre-trained model. The code in step 2 includes the loading of an image. A small video inset of the presenter is visible in the bottom right corner of the notebook screenshot.

So, let us try to understand the inferencing portion on just GPU, CUDA and TensorRT. So, this particular thing is a very very simple example of trying to import ok a pre trained neural network model. So, we are just trying to load an example image ok.

(Refer Slide Time: 18:26)




The screenshot shows a Jupyter Notebook titled "Inference with CPU, CUDA and TensorRT". The code is as follows:

```
step 1: load a pre-trained neural network
%% Import torch
from torch.nn import Linear, Softmax
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim
from torchvision import datasets, transforms

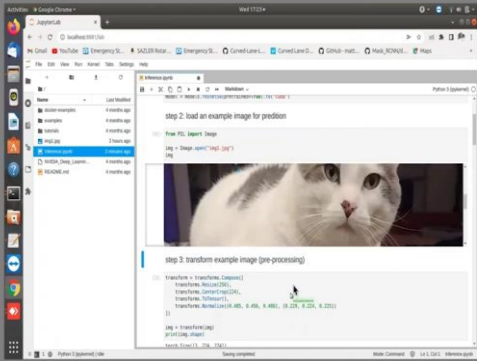
step 2: load an example image for prediction
from PIL import Image
img = Image.open('img.jpg')
img
```

The output shows the image's dimensions: `torch.Size([3, 224, 224])`. The NPTEL logo is visible in the top right corner.



For the prediction which we are able to see ok.


(Refer Slide Time: 18:32)



The screenshot shows the same Jupyter Notebook, but now with an image of a white cat. The code for step 3 is visible:

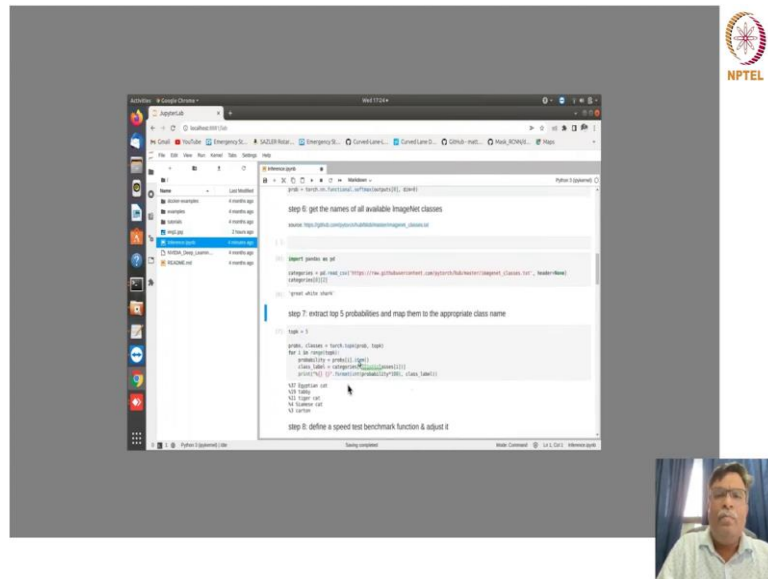
```
step 3: transform example image (pre-processing)
transforms = transforms.Compose([
    transforms.Resize(224),
    transforms.CenterCrop(224),
    transforms.ToTensor(),
    transforms.Normalize([0.4914, 0.4516, 0.5015], [0.2025, 0.2228, 0.2145])
])
img = transforms(img)
img
```

The output shows the transformed image: `torch.Size([3, 224, 224])`. The NPTEL logo is visible in the top right corner.



Now, what we are trying to show here is this is a image ok and then there is some transform which is happening. So, we are trying to pre process this image ok.

(Refer Slide Time: 19:26)



The screenshot shows a Jupyter Notebook with the following code:

```
step 6: get the names of all available Inception classes
source: https://github.com/tensorflow/tensorflow/blob/master/tensorflow/models/official/inception/inception.py

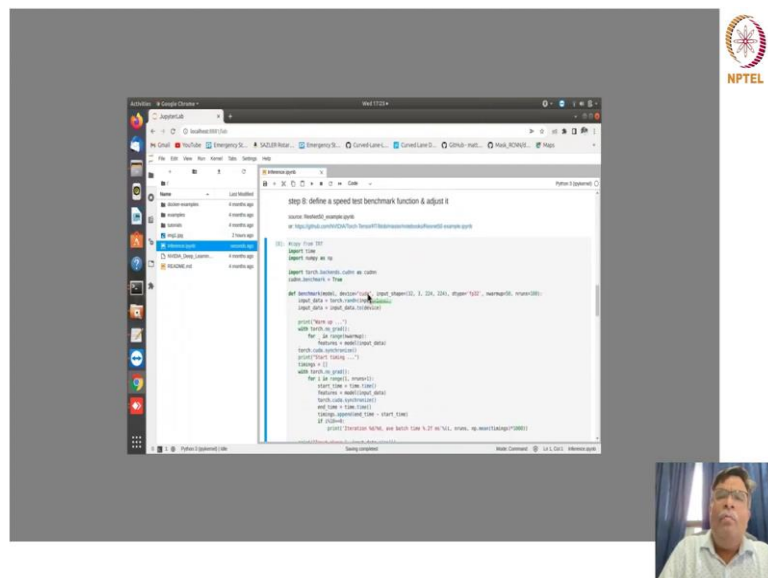
import pandas as pd
categories = all_model_urls["https://raw.githubusercontent.com/tensorflow/tensorflow/master/tensorflow/models/official/inception/inception.py"]
print(categories)

step 7: extract top 5 probabilities and map them to the appropriate class name
In [ ]:
def top_5_probs(image, model):
    probs, classes = sess.run([inception_top5_op, inception_top5_class_names],
                              feed_dict={inception_top5_input: image})
    return probs, classes

img = image_loader.load_image('cat.jpg')
img = image_loader.resize_image(img)
img = image_loader.preprocess_image(img)
probs, classes = top_5_probs(img, model)
print(classes)
```

This is for that particular thing. So, it has detected that it is basically a Egyptian cat its tabby, tiger cat semi stat and carton right; so, with these many prediction percentages right.

(Refer Slide Time: 19:47)



The screenshot shows a Jupyter Notebook with the following code:

```
step 8: define a speed test benchmark function & adjust it
In [ ]:
def cat_benchmark(image, model):
    """Benchmark the inference time of the ResNet50 model on a cat image.
    Args:
        image: A numpy array representing the image.
        model: A ResNet50 model object.
    Returns:
        A tuple containing the inference time in seconds and the predicted class name.
    """
    # Preprocess the image
    img = image_loader.preprocess_image(image)

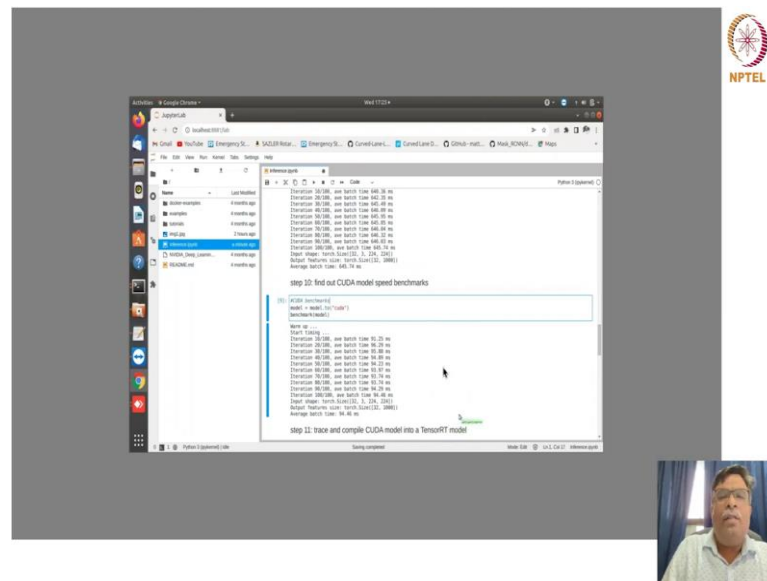
    # Inference
    probs, classes = top_5_probs(img, model)

    # Print the results
    print("Predicted class: {}".format(classes[0]))
    print("Inference time: {}".format(time.time() - start_time))

# Run the benchmark
cat_benchmark(image_loader.load_image('cat.jpg'), model)
```

So, now we will try to understand the speed test with the ResNet 50 ok example. So, what we are trying to do is this is basically a TRT program ok and then we are trying to we are trying to just show you ok that how basically for each of this ok.

(Refer Slide Time: 20:30)



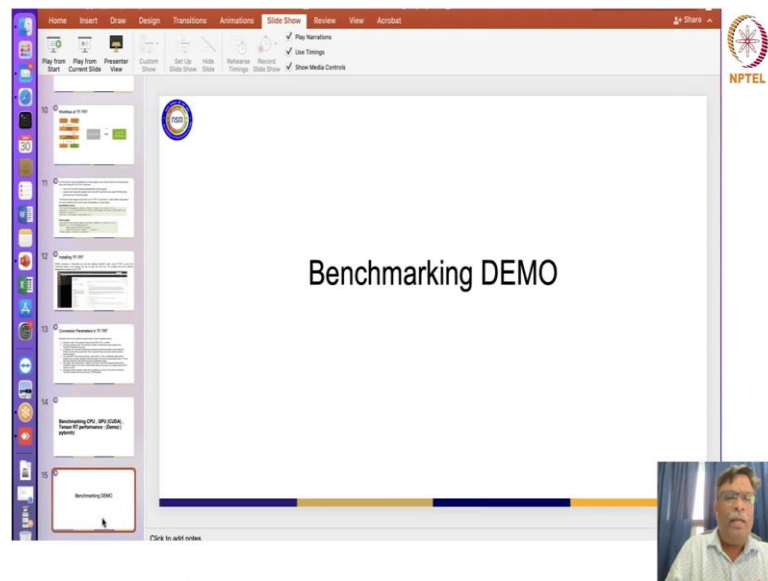
So, this is just a very very preliminary benchmarking thing which will help you to appreciate the optimize stuff right. So, that is what we are trying to do now. So, here so, if you run this ok which we have just now run it basically gives you 645.74 milliseconds right the average batch time and we are running it for CUDA here, which is executing which is telling you that it is giving you 91.85 milliseconds.

So, the idea here is you are trying to run the same model on a CPU which gives you 645.45 milliseconds 7 sorry 654 644.74 milliseconds. Whereas on CUDA ok the same model gives you 91.85 milliseconds for actually training right now TensorRT model if you see take a little bit of time yeah.

So, if you see here its 37 19 11 4 and 3. The same thing we get because we yesterday we are talking of even if it is reduced prediction sorry reduced precision right to get almost the same thing right there is nothing much of a difference ok right.

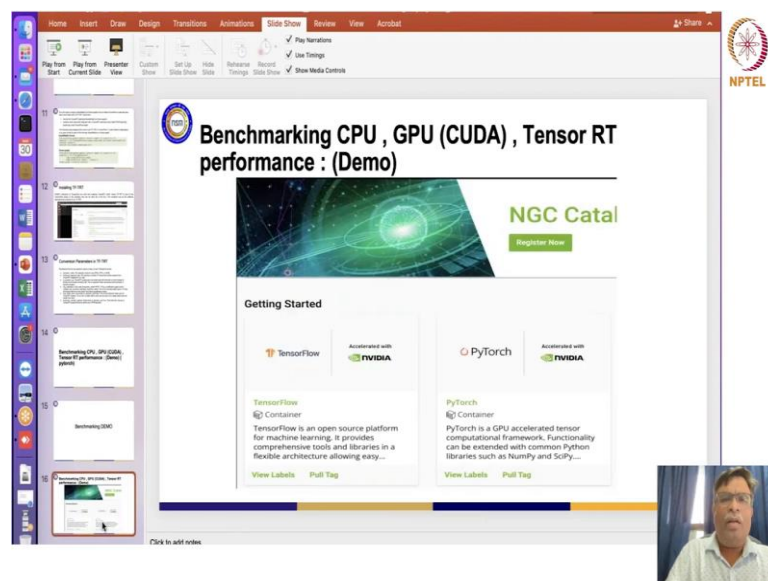
So, this is what we thought we will show you first. And then now so, now, once we have done this.

(Refer Slide Time: 22:47)

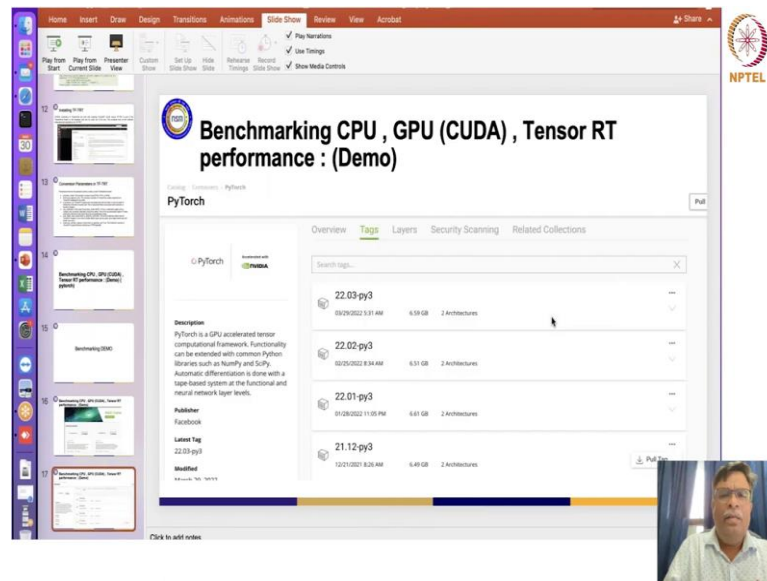


The benchmarking demo let us try to understand right step by step of how we have done it actually.

(Refer Slide Time: 22:49)

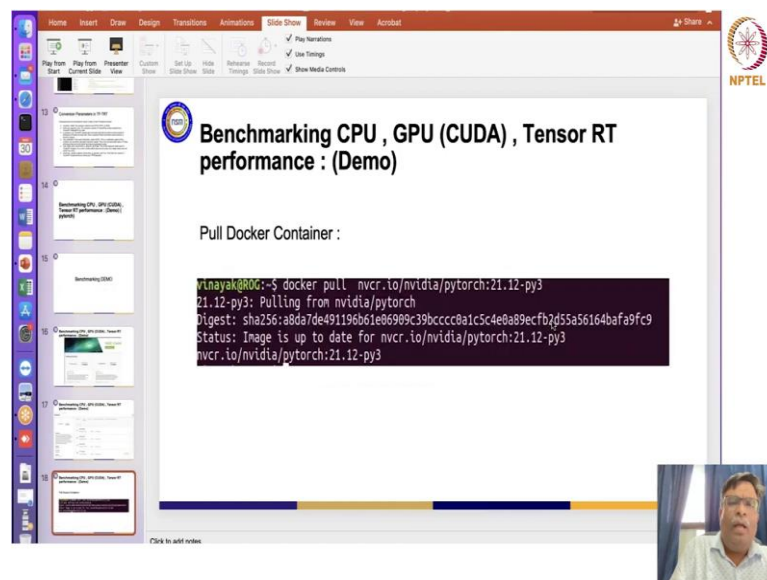


(Refer Slide Time: 22:52)



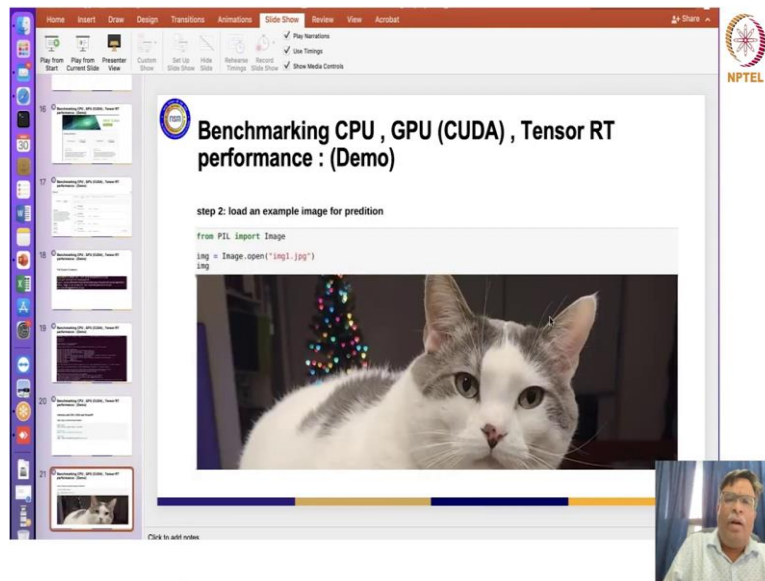
So, what we do is we basically as I told you we are supposed to be pulling the docker containers ok from the NGC site right. So, we are we have tried to download the PyTorch this thing ok.

(Refer Slide Time: 23:14)



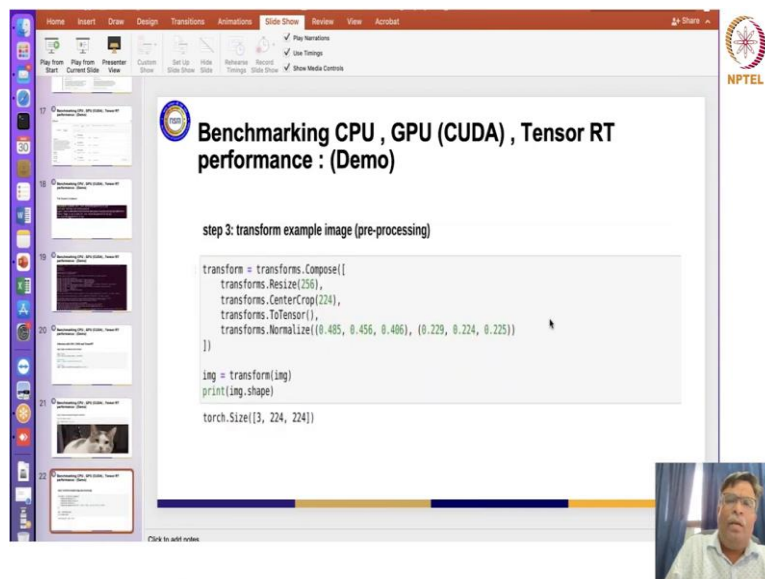
And then we basically pulled the docker container. So, we are actually telling you step by step of how we did it right all of this is executed. So, whatever we try to do till now we have actually tried to put it in the slide so, that tomorrow when you want to do it ok, you can actually see it automatically and do it step by step right.

(Refer Slide Time: 23:44)



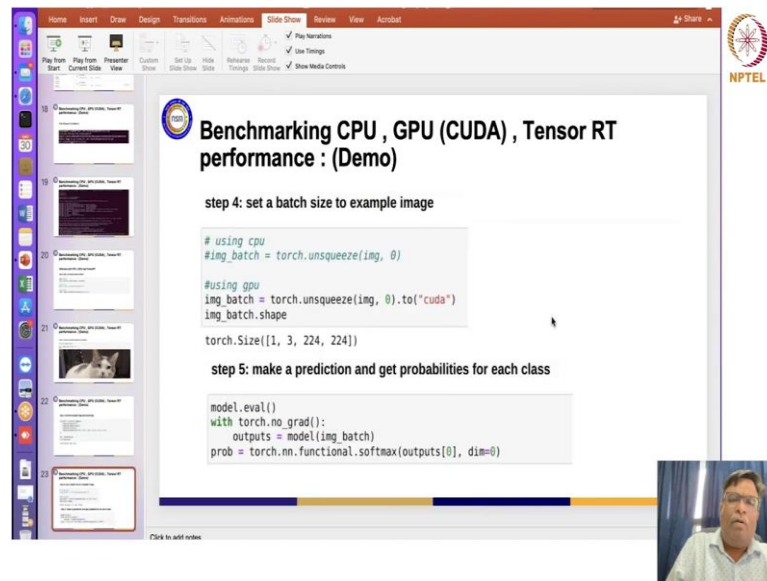
You get the image for production.

(Refer Slide Time: 23:46)



Then this is the transformed image.

(Refer Slide Time: 23:49)



Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

step 4: set a batch size to example image

```
# using cpu
#img_batch = torch.unsqueeze(img, 0)

#using gpu
img_batch = torch.unsqueeze(img, 0).to("cuda")
img_batch.shape

torch.Size([1, 3, 224, 224])
```

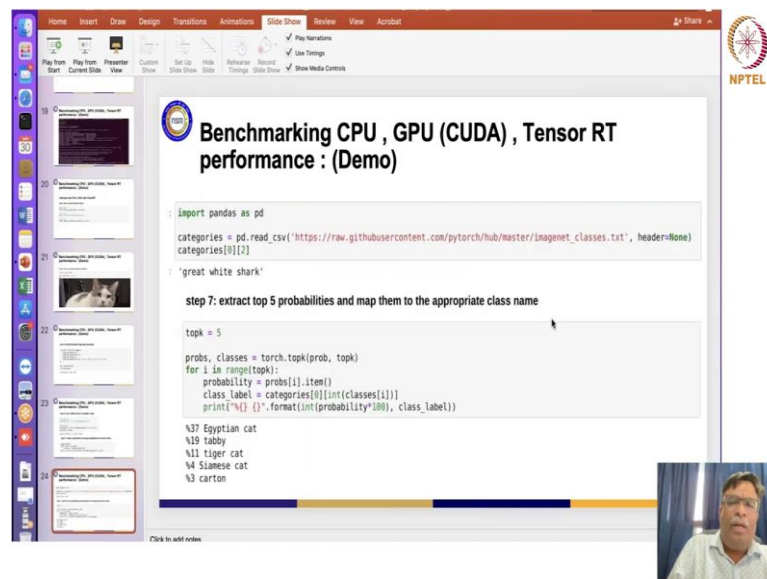
step 5: make a prediction and get probabilities for each class

```
model.eval()
with torch.no_grad():
    outputs = model(img_batch)
prob = torch.nn.functional.softmax(outputs[0], dim=0)
```

NPTEL

Then you see the batch size probability prediction for CPU ok.

(Refer Slide Time: 23:53)



Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

```
import pandas as pd
categories = pd.read_csv('https://raw.githubusercontent.com/pytorch/hub/master/imagenet_classes.txt', header=None)
categories[0][2]
'great white shark'
```

step 7: extract top 5 probabilities and map them to the appropriate class name

```
topk = 5
probs, classes = torch.topk(prob, topk)
for i in range(topk):
    probability = probs[i].item()
    class_label = categories[0][int(classes[i])]
    print("{} {} {}".format(int(probability*100), class_label))

%37 Egyptian cat
%19 tabby
%11 tiger cat
%4 Siamese cat
%3 carton
```

NPTEL

(Refer Slide Time: 23:54)

Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

step 8: define a speed test benchmark function & adjust it

source: <https://github.com/PyTorch/torchbenchmark/blob/master/examples.py>

```
#!/usr/bin/env python3
import time
import numpy as np

import torch.backends.cudnn as cudnn
cudnn.benchmark = True

def benchmark(model, device="cuda", input_shape=(1, 1, 224, 224), dtype=torch.float32, warmup=10, runs=1000):
    print("Warm up ...")
    with torch.no_grad():
        for _ in range(warmup):
            features = model(input_data)
            torch.cuda.synchronize()
            print("Start timing ...")
            timings = []
            with torch.no_grad():
                for i in range(runs):
                    start_time = time.time()
                    features = model(input_data)
                    torch.cuda.synchronize()
                    end_time = time.time()
                    timings.append(end_time - start_time)
            if i % 10 == 0:
                print(f"Iteration {i}/rd, ave batch time %.2f ms %d, runs, np.mean(timings)*1000)")

    print("Input shape:", input_shape)
    print("Output features size:", features.size())
    print("Average batch time: %.2f ms %d (np.mean(timings)*1000)"))
```

And then you do it for GPU.

(Refer Slide Time: 23:59)

Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

step 9: find out CPU model speed benchmarks

```
] : #CPU benchmarks
benchmark(model, device="cpu")

Warm up ...
Start timing ...
Iteration 10/100, ave batch time 640.36 ms
Iteration 20/100, ave batch time 642.35 ms
Iteration 30/100, ave batch time 645.49 ms
Iteration 40/100, ave batch time 646.09 ms
Iteration 50/100, ave batch time 645.95 ms
Iteration 60/100, ave batch time 645.85 ms
Iteration 70/100, ave batch time 646.04 ms
Iteration 80/100, ave batch time 646.32 ms
Iteration 90/100, ave batch time 646.03 ms
Iteration 100/100, ave batch time 645.74 ms
Input shape: torch.Size([1, 1, 224, 224])
Output features size: torch.Size([1, 1000])
Average batch time: 645.74 ms
```

The CPU benchmarking, the GPU benchmarking.

(Refer Slide Time: 24:01)

Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

step 10: find out CUDA model speed benchmarks

```
#CUDA benchmarks
model = model.to("cuda")
benchmark(model)

Warm up ...
Start timing ...
Iteration 10/100, ave batch time 92.37 ms
Iteration 20/100, ave batch time 92.42 ms
Iteration 30/100, ave batch time 92.32 ms
Iteration 40/100, ave batch time 94.92 ms
Iteration 50/100, ave batch time 97.94 ms
Iteration 60/100, ave batch time 97.01 ms
Iteration 70/100, ave batch time 99.83 ms
Iteration 80/100, ave batch time 101.69 ms
Iteration 90/100, ave batch time 104.36 ms
Iteration 100/100, ave batch time 106.73 ms
Input shape: torch.Size([32, 3, 224, 224])
Output features size: torch.Size([32, 1000])
Average batch time: 106.73 ms
```

(Refer Slide Time: 24:03)

Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

step 11: trace and compile CUDA model into a TensorRT model

```
traced_model = torch.jit.trace(model, [torch.randn((32, 3, 224, 224)).to("cuda")])

import torch_tensorrt

trt_model = torch_tensorrt.compile(
    traced_model,
    inputs=[torch_tensorrt.Input((32, 3, 224, 224), dtype=torch.float32)],
    enabled_precisions = {torch.float32}
)
```

(Refer Slide Time: 24:04)

Benchmarking CPU , GPU (CUDA) , Tensor RT performance : (Demo)

```
benchmark(trt_model)

Warm up ...
Start timing ...
Iteration 10/100, ave batch time 76.66 ms
Iteration 20/100, ave batch time 75.95 ms
Iteration 30/100, ave batch time 73.84 ms
Iteration 40/100, ave batch time 72.74 ms
Iteration 50/100, ave batch time 72.09 ms
Iteration 60/100, ave batch time 71.67 ms
Iteration 70/100, ave batch time 71.44 ms
Iteration 80/100, ave batch time 71.28 ms
Iteration 90/100, ave batch time 71.07 ms
Iteration 100/100, ave batch time 71.52 ms
Input shape: torch.Size([32, 3, 224, 224])
Output features size: torch.Size([32, 1000])
Average batch time: 71.52 ms
```

And then TRT model benchmarking ok.

(Refer Slide Time: 24:07)

Demo of Workflow for optimizing Tensorflow model to TensorRT (Tensorflow) – Reduced Precision graph generation

Something like this.