# Model conversion

# 1. Jetson Orin YOLO11 (benchmark)

YOLO11 benchmark data comes from the Ultralytics team, which tests models in multiple formats (data is for reference only)

# 2. Enable optimal performance of the motherboard

## Enable Jetson clocks

Enabling Jetson Clocks will ensure that all CPU and GPU cores run at maximum frequency:

```
sudo jetson_clocks
```

# 3. Model conversion

According to the test parameters of different formats provided by the Ultralytics team, we can find that the inference performance is best when using TensorRT!
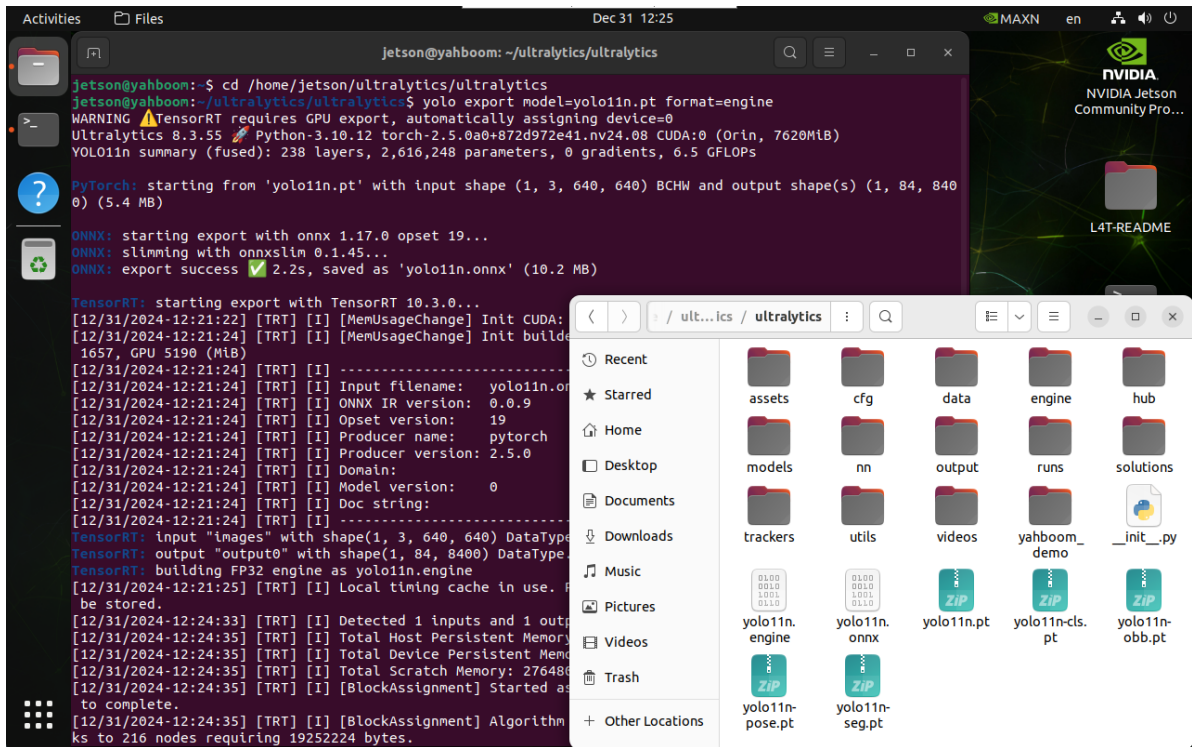
```
When using the export mode of YOLO11 for the first time, some dependencies will
be automatically installed. Just wait for it to be completed automatically!
```

## 3.1. CLI: pt → onnx → engine

Convert the PyTorch format model to TensorRT: The conversion process will automatically generate an ONNX model

```
cd /home/jetson/ultralytics/ultralytics
```

```
yolo export model=yolo11n.pt format=engine
# yolo export model=yolo11n-seg.pt format=engine
# yolo export model=yolo11n-pose.pt format=engine
# yolo export model=yolo11n-cls.pt format=engine
# yolo export model=yolo11n-obb.pt format=engine
```
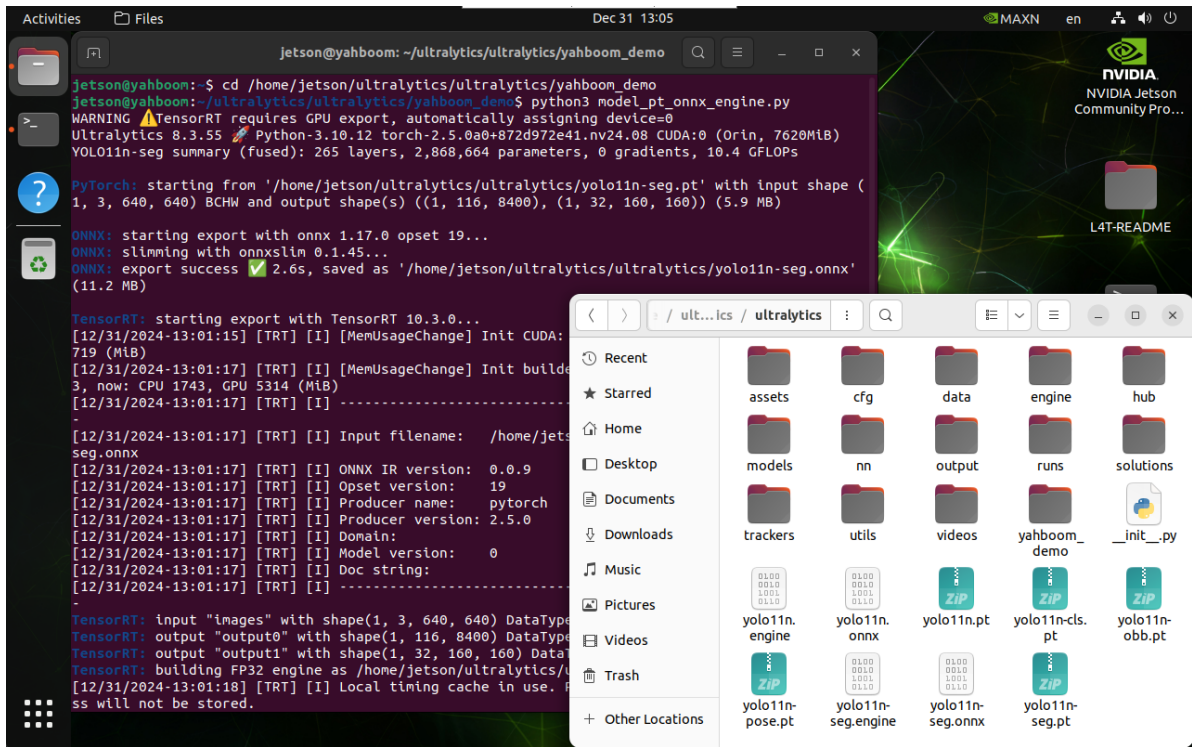


## 3.2、Python: pt → onnx → engine

Convert PyTorch format models to TensorRT: The conversion process will automatically generate ONNX models

```
cd /home/jetson/ultralytics/ultralytics/yahboom_demo
```

```
python3 model_pt_onnx_engine.py
```

```python
from ultralytics import YOLO

# Load a YOLO11n PyTorch model
# model = YOLO("/home/jetson/ultralytics/ultralytics/yolo11n.pt")
model = YOLO("/home/jetson/ultralytics/ultralytics/yolo11n-seg.pt")
# model = YOLO("/home/jetson/ultralytics/ultralytics/yolo11n-pose.pt")
# model = YOLO("/home/jetson/ultralytics/ultralytics/yolo11n-cls.pt")
# model = YOLO("/home/jetson/ultralytics/ultralytics/yolo11n-obb.pt")

# Export the model to TensorRT
model.export(format="engine")
```

Note: The model file generated by the conversion is located in the converted model file location
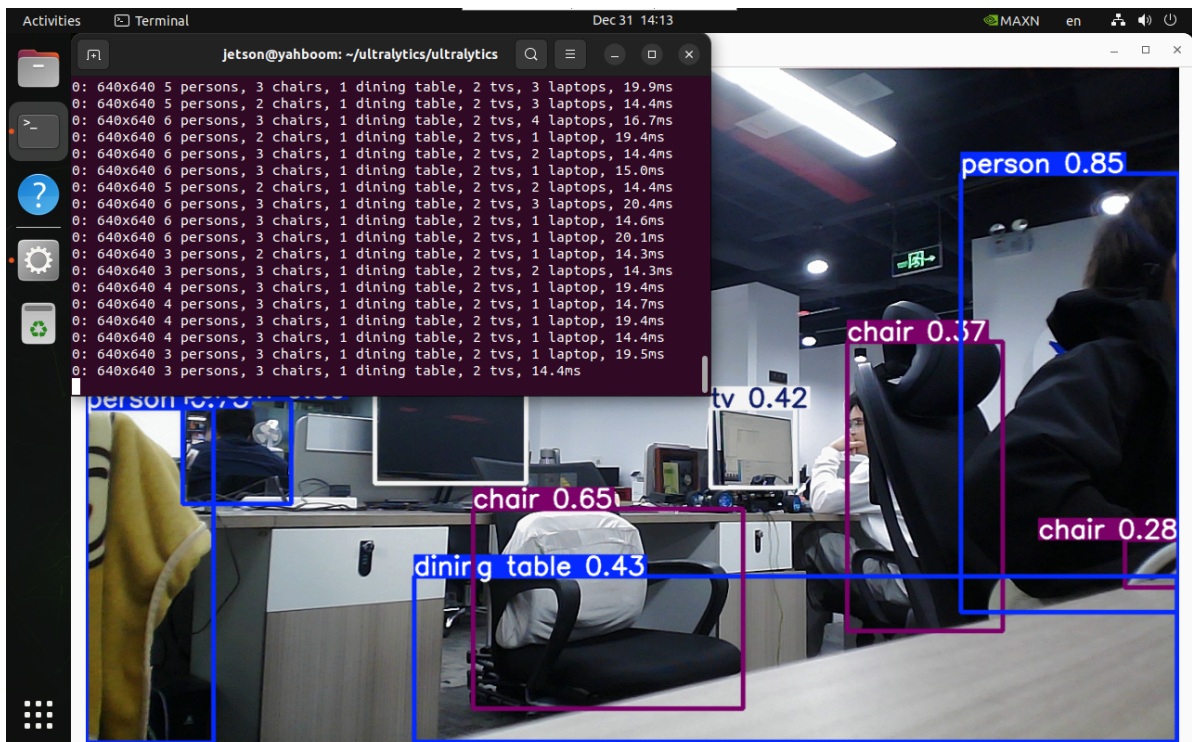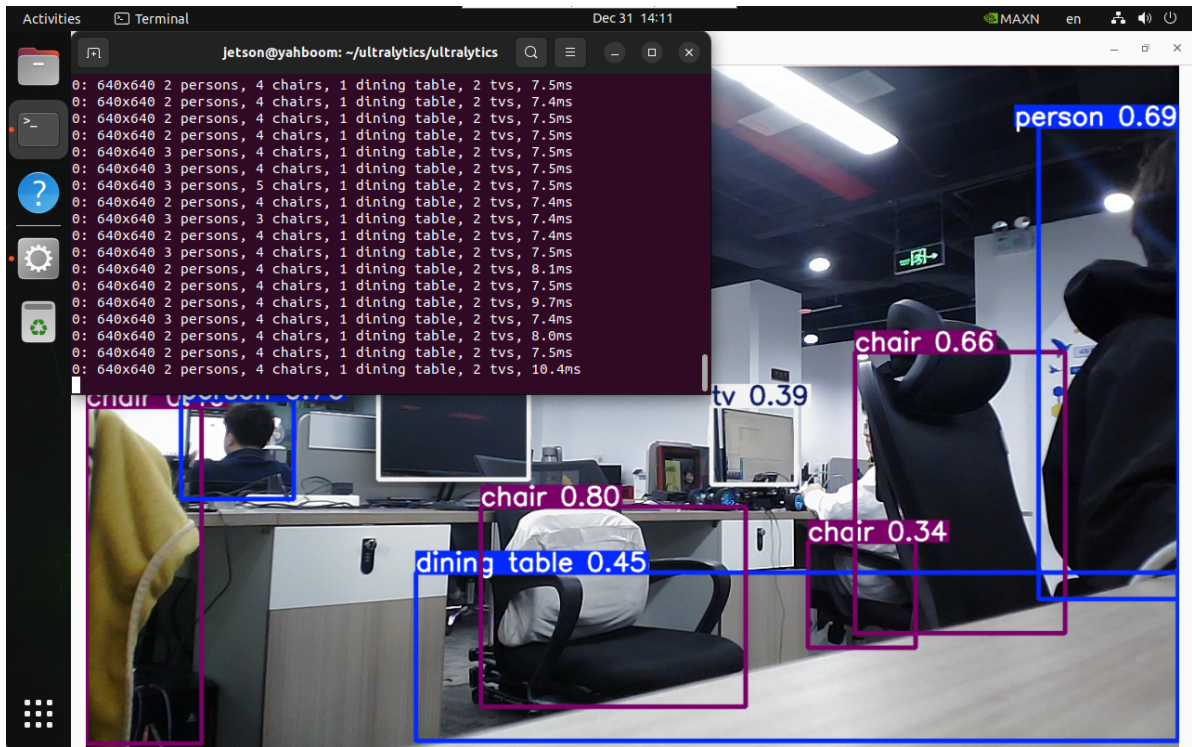
# 4. Model prediction

## CLI usage

CLI currently only supports calling USB cameras. Nuwa camera users can directly modify the previous python code to call onnx and engine models!

```
cd /home/jetson/ultralytics/ultralytics
```

```
yolo predict model=yolo11n.onnx source=0 save=False show
```
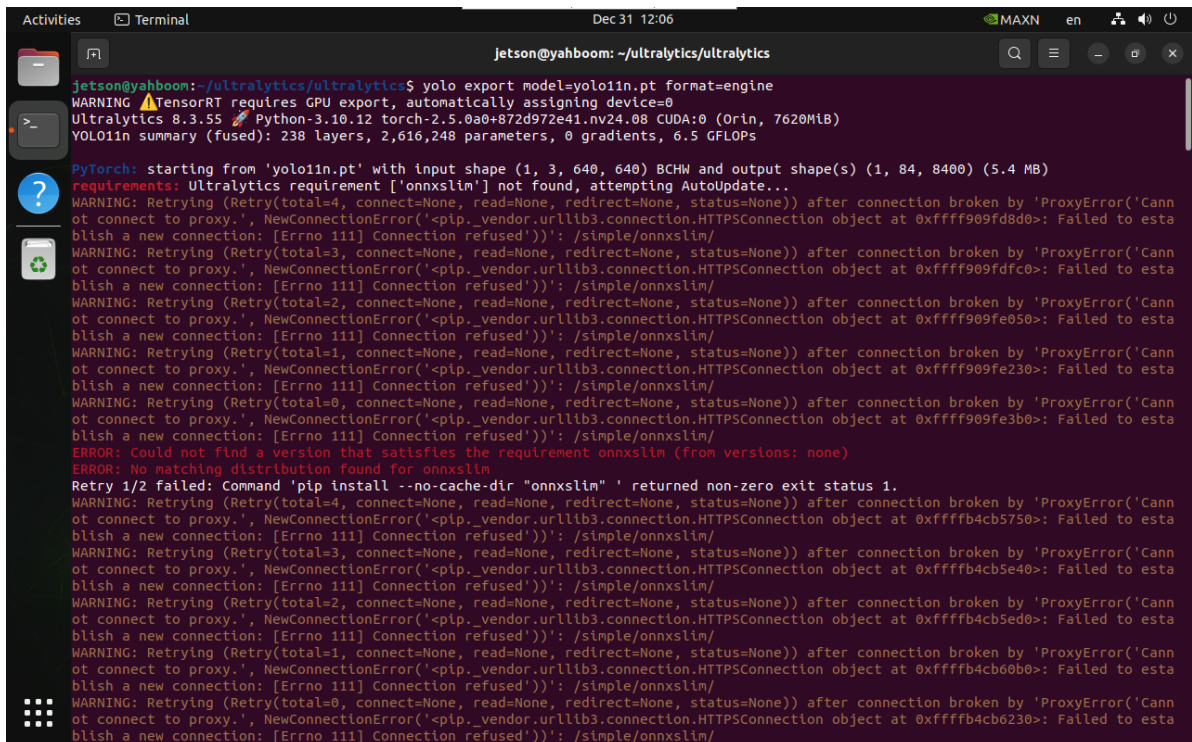


```
yolo predict model=yolo11n.engine source=0 save=False show
```

# Frequently asked questions

## ERROR: onnxslim



Solution: Enter the onnxslim installation command in the terminal

```
sudo pip3 install onnxslim
```

# References

https://docs.ultralytics.com/guides/nvidia-jetson/

https://docs.ultralytics.com/integrations/tensorrt/