

---

---

# Optimizer models

---

---

## 1. Convert pytorch 32bit to 16bit

```
1 model = torch.load("models/torch/resnet18.pth").to('cpu')
2 model.float()
3 model.half()
4 torch.save(model, 'models/torch/resnet18_float16.pth')
```

```
1 size = get_size_file('models/torch/resnet18.pth')
2 print(f"Size of model resnet18: {size:.1f} MB")
3 size = get_size_file('models/torch/resnet18_float16.pth')
4 print(f"Size of model resnet18_float16: {size:.1f} MB")
```

[50] ✓ 0.0s

```
... Size of model resnet18: 42.7 MB
     Size of model resnet18_float16: 21.4 MB
```

### Điểm Mạnh:

- Dễ chuyển đổi
- Dung lượng mô hình **giảm đi 50%**
- Tốc độ dự đoán **tăng 3 lần**

```
1 print(f"Time infer per image of resnet18_float32: {np.mean(time_infer_32bit): .4f} s")
2 print(f"Time infer per image of resnet18_float16: {np.mean(time_infer_16bit): .4f} s")
```

[51] ✓ 0.0s

```
... Time infer per image of resnet18_float32: 0.0121 s
     Time infer per image of resnet18_float16: 0.0045 s
```

## 1. Convert pytorch 32bit to 16bit

Đây là đoạn code kiểm tra sự sai khác giữa mô hình float32 và float16.

Model float16: Bị NaN ở các batch\_size 8, 64, 128

=> **Nó không stable**

### Solution:

- Dùng **Mixed Precision Training**
- Dùng **ONNX 16bit**

```
1 import torch
2 import numpy as np
3
4 batch_size = [2, 4, 8, 16, 32, 64, 128]
5
6 for batch in batch_size:
7     with torch.cuda.stream(torch.cuda.current_stream()):
8         input = torch.randn(batch, 3, 224, 224)
9         input = input.cuda() # Chuyển input lên GPU bên trong stream
10        with torch.no_grad():
11            output_32bit = float_32(input)
12            output_16bit = float_16(input.half())
13        error = np.mean(output_32bit.cpu().numpy() - output_16bit.cpu().numpy())
14        print(f"Batch {batch}: {error:.5f}")
15
```

[53] ✓ 2.5s

```
... Batch 2: 0.00063
Batch 4: 0.00045
Batch 8: nan
Batch 16: 0.00043
Batch 32: 0.00065
Batch 64: nan
Batch 128: nan
```

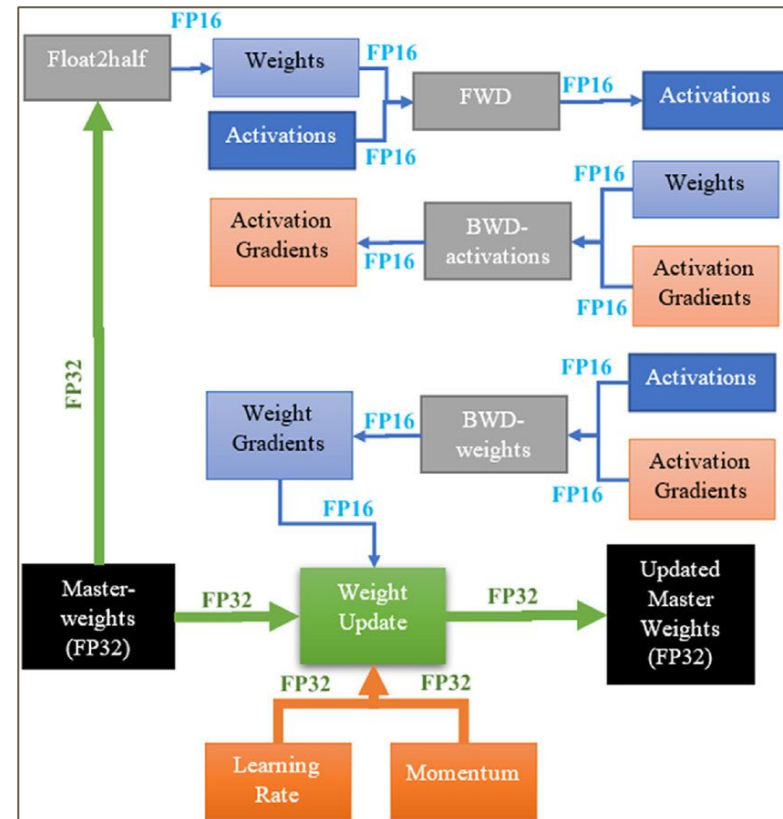
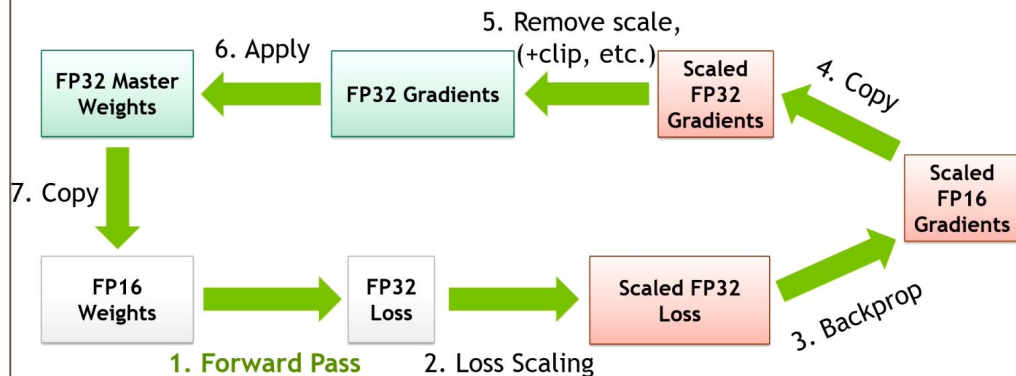
```
1 output_16bit
```

✓ 0.0s

```
tensor([[nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan],
        [nan, nan, nan, nan, nan, nan]], device='cuda:0', dtype=torch.float16)
```

## Mixed Precision Training

### MIXED PRECISION TRAINING



## Mixed Precision Training

[https://pytorch.org/docs/stable/notes/amp\\_examples.html](https://pytorch.org/docs/stable/notes/amp_examples.html)

### Ưu điểm:

- **Giảm dung lượng GPU yêu cầu cho các phép tính:** Do các phép tính diễn ra trên fp16 nên dung lượng tính toán giảm nhiều.
- **Hiệu quả trên một số phần cứng hiện đại :** Một số phần cứng GPU Nvidia, được tối ưu hóa cho các phép tính fp16, cho kết quả infer nhanh hơn.
- **Tăng tốc độ huấn luyện:** Do phần forward tính toán trên fp16 nên việc tính toán nhanh.

### Nhược điểm:

- **Cần nhiều dung lượng GPU:** do việc nó cần tạo ra hai bản sao cho fp32 và fp16.
- **Không phải GPU Nvidia nào hỗ trợ tốt fp16:** một số phần cứng chạy trên fp32 nhanh hơn fp16.

**Nhận xét:** Nên dùng MPT khi muốn training nhanh mô hình, trong điều kiện tài nguyên thoải mái.

1	N/A	N/A	1401354	C	...rtualenvs/pose_estimator/bin/python	1308MiB
1	N/A	N/A	1405967	C	...rtualenvs/pose_estimator/bin/python	996MiB

## Mixed Precision Training

```
1 import torch
2 from torch.cuda.amp import GradScaler, autocast
3
4 dataloader = ... # Tạo dataloader
5 epochs = ... # Số epoch cần huấn luyện
6 criterion = ... # Khởi tạo hàm loss
7 model = ... # Khởi tạo mô hình
8 optimizer = ... # Khởi tạo trình tối ưu hóa
9 scaler = GradScaler() # Khởi tạo scaler cho mixed precision
10
11 for epoch in range(epochs):
12     for batch in dataloader:
13         inputs, targets = batch
14         inputs, targets = inputs.cuda(), targets.cuda()
15
16         optimizer.zero_grad()
17
18         with autocast():
19             outputs = model(inputs)
20             loss = criterion(outputs, targets)
21
22         scaler.scale(loss).backward()
23         scaler.step(optimizer)
24         scaler.update()
```

```
1 import torch
2
3 dataloader = ... # Tạo dataloader
4 epochs = ... # Số epoch cần huấn luyện
5 criterion = ... # Khởi tạo hàm loss
6 model = ... # Khởi tạo mô hình
7 optimizer = ... # Khởi tạo trình tối ưu hóa
8 scheduler = ... # Khởi tạo scheduler
9
10 for epoch in range(epochs):
11     for batch in dataloader:
12         inputs, targets = batch
13         inputs, targets = inputs.cuda(), targets.cuda()
14
15         optimizer.zero_grad()
16         outputs = model(inputs)
17         loss = criterion(outputs, targets)
18         loss.backward()
19         optimizer.step()
20
21     scheduler.step() # Cập nhật learning rate sau mỗi epoch
22
```

Do tính toán fp16 dễ bị tràn số ( overflow computing ), dùng một cái scaler để chỉnh gradient trước khi cập nhật trọng số , giúp quá trình cập nhật ổn định hơn.

## JIT ( Just in Time )

JIT trong **TorchScript** là quá trình biến đổi một mô hình Pytorch từ dạng **đồ thị tính toán động** (Dynamic Computation Graph) sang dạng **đồ thị tính toán tĩnh** (Static Computation Graph).

- Dynamic Computation Graph “define-by-run”: chỉ có pytorch dùng
- Static Computation Graph “define-and-run” : Tensorflow, TensorRT, TorchScript, ONNX

## Static vs Dynamic Graphs

**TensorFlow:** Build graph once, then run many times (**static**)

```
N, D, H = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=(N, D))
y = tf.placeholder(tf.float32, shape=(N, D))
w1 = tf.Variable(tf.random_normal((D, H)))
w2 = tf.Variable(tf.random_normal((H, D)))

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_sum(diff ** 2, axis=1))
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

learning_rate = 1e-5
new_w1 = w1.assign(w1 - learning_rate * grad_w1)
new_w2 = w2.assign(w2 - learning_rate * grad_w2)
updates = tf.group(new_w1, new_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = {x: np.random.randn(N, D),
              y: np.random.randn(N, D),}
    losses = []
    for t in range(50):
        loss_val, _ = sess.run([loss, updates],
                               feed_dict=values)
```

Build  
graph

Run each  
iteration

**PyTorch:** Each forward pass defines a new graph (**dynamic**)

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

New graph each iteration



### Torchscript

```
1 import torch
2 # Load models pytorch
3 model_pytorch_cuda = torch.load("models/torch/resnet18.pth").to('cuda')
4 # convert model to static graph
5 sample_input_cuda = torch.randn(1,3,224,224).cuda()
6 traced_cuda = torch.jit.trace(model_pytorch_cuda, sample_input_cuda)
7 torch.jit.save(traced_cuda, "cuda.pt")
8
9
10 # Load jit model
11 model_jit = torch.jit.load("cuda.pt")
12
13 # Inference model
14 model_jit(sample_input_cuda)
```

[41] ✓ 0.4s

```
... tensor([[ 0.3873,  2.0765, -0.4736, -6.3287,  2.5076, -3.1566]],
          device='cuda:0', grad_fn=<AddmmBackward0>)
```



# Torchscript

```
1 import torch
2 # Load models pytorch
3 model_pytorch_cuda = torch.load("models/torch/resnet18.pth").to('cuda')
4 # convert model to static graph
5 sample_input_cuda = torch.randn(1,3,224,224).cuda()
6 traced_cuda = torch.jit.trace(model_pytorch_cuda, sample_input_cuda)
7 torch.jit.save(traced_cuda, "cuda.pt")
8
9
10 # Load jit model
11 model_jit = torch.jit.load("cuda.pt")
12
13 # Inference model
14 model_jit(sample_input_cuda)
```

[43] ✓ 0.4s

```
... tensor([[ -0.3191,  2.0507, -0.0977, -6.1900,  3.0482, -3.3211]],
        device='cuda:0', grad_fn=<AddmmBackward0>)
```

```
1 # Một cách khác, nhưng cách này giúp tối ưu, nhưng không lưu được mô hình
2 frozen_mod = torch.jit.optimize_for_inference(torch.jit.script(model_pytorch_cuda.eval()))
3 frozen_mod(sample_input_cuda)
```

[44] ✓ 0.2s

```
... tensor([[ -0.3191,  2.0507, -0.0977, -6.1900,  3.0482, -3.3211]],
        device='cuda:0')
```

**Model:** Resnet18

Pytorch (GPU): **0.0240 ms**

Torchscript (GPU): 0.0313 ms

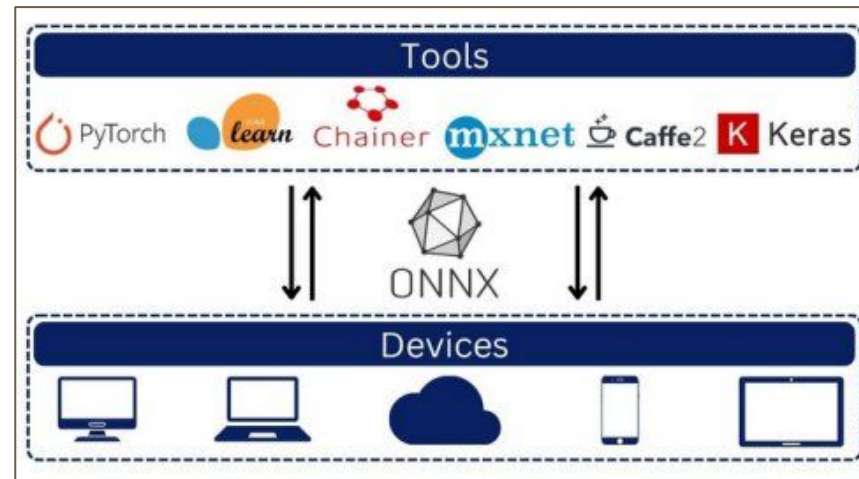
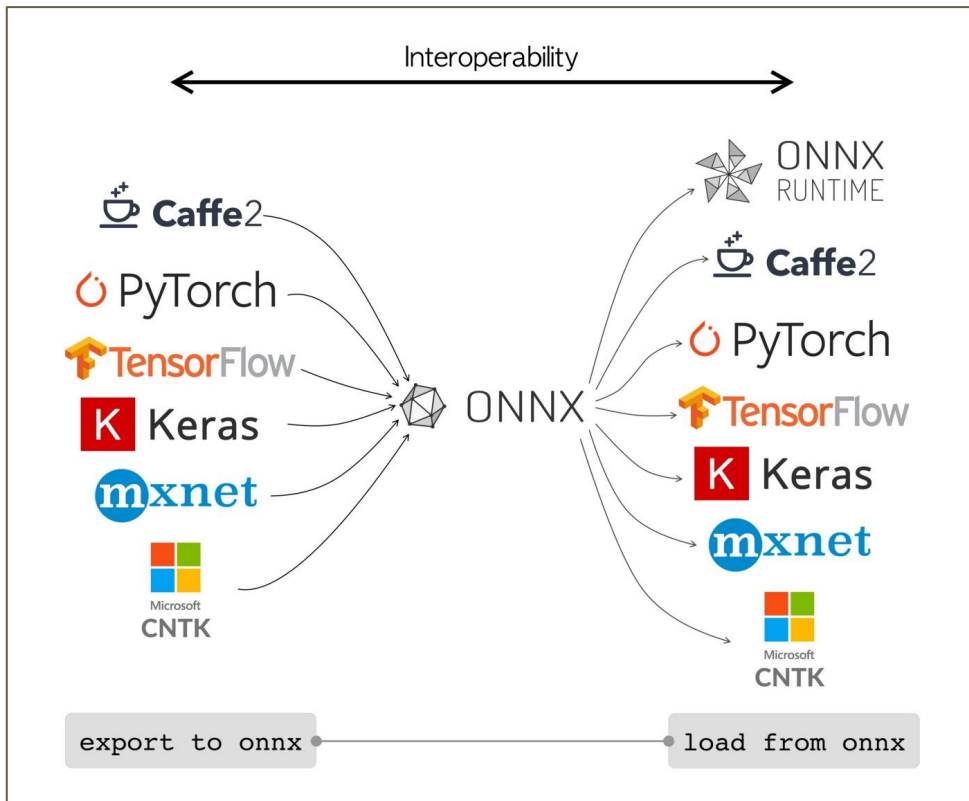
Torchscript\_Optimizer (GPU):  
0.0287 ms

Pytorch (CPU): 0.4395 ms

Torchscript (CPU): 0.3081 ms

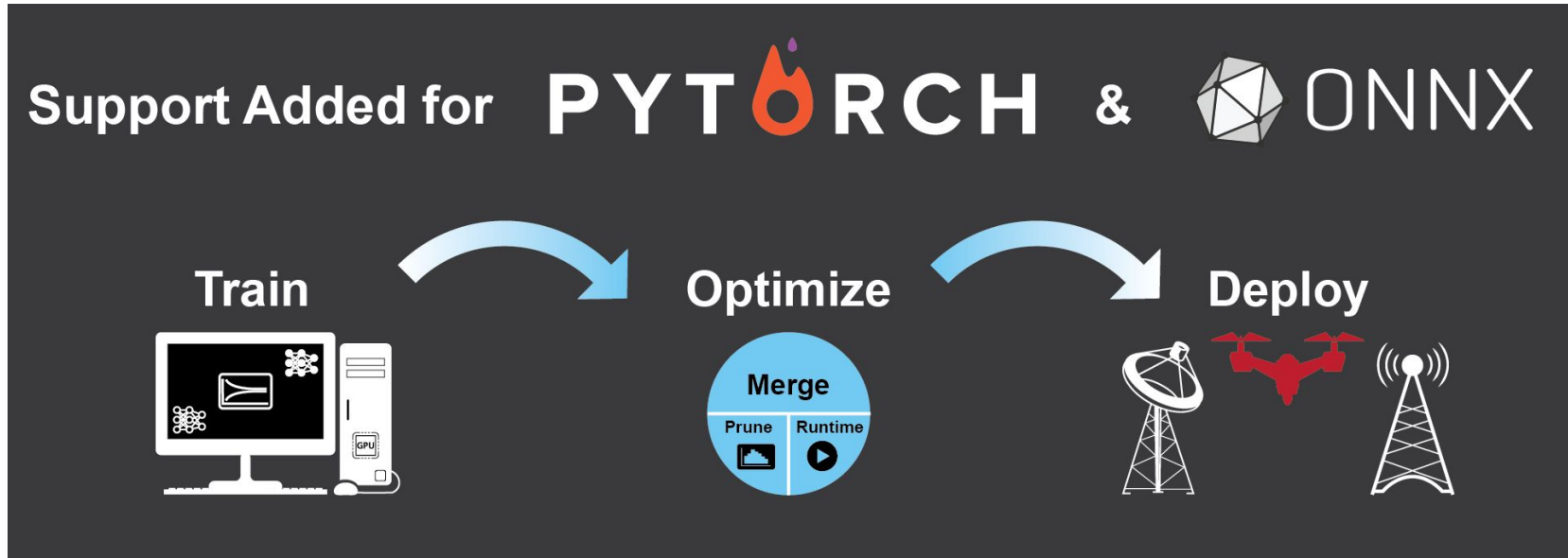
Torchscript\_Optimizer (CPU):  
**0.1935 ms**

## ONNX ( Open Neural Network Exchange )



- Là một framework
- Giúp chuyển đổi giữa các framework với nhau
- Tối ưu hóa mô hình và triển khai dễ dàng trên máy tính, di động, thiết bị nhúng.

## ONNX ( Open Neural Network Exchange )



## Pytorch to ONNX

```
1 def convert_pytorch2onnx(model, input_samples, path_onnx, mode = 'float32bit', device = 'cuda'):  
2     if mode == 'float16bit':  
3         print("float16bit")  
4         model.float()  
5         model.half() # Chuyển mô hình sang float16  
6         input_samples = input_samples.half()  
7     model.to(device)  
8     model.eval()  
9     input_samples = input_samples.to(device)  
10    torch.onnx.export(  
11        model, # Model load bằng pytorch  
12        input_samples, # input với kích thước mong muốn  
13        path_onnx, # path_onnx bạn lưu  
14        verbose=False, # Hiện thị thông báo trong quá trình chuyển đổi  
15        opset_version=12, # Đọc xem mỗi version sẽ có hỗ trợ onnx cho lớp nào  
16        do_constant_folding=True, # Dùng Constant-folding optimizer giúp cải thiện tốc độ và dung lượng  
17        input_names = ['images'], # the model's input names  
18        output_names = ['output'], # the model's output names  
19        dynamic_axes={  
20            'images' : {0 : 'batch_size'},  
21            'output' : {0 : 'batch_size'}  
22        }  
23    )
```

## ONNX Runtime

ONNX Runtime: Là một trình thực thi mô hình Deep Learning

- Đa nền tảng: Chạy được trên Linux, Windows, MacOS
- Tối ưu hóa hiệu suất trên VPU ( Thiết bị lượng tử ), GPU, CPU
- API đa ngôn ngữ: C, Python, C++, Java
- Tích hợp với các công nghệ khác: TensorRT, openVINO

```
24
25 def load_onnx_model(path_onnx, providers=['CUDAExecutionProvider', 'CPUExecutionProvider']):
26     # Create an ONNX Runtime inference session for the ONNX model
27     ort_session = onnxruntime.InferenceSession(
28         path_onnx,
29         providers=providers
30     )
31     return ort_session
32
33 def onnx_infer(ort_session, input_data):
34     ort_inputs = {ort_session.get_inputs()[0].name: input_data}
35     ort_output = ort_session.run(None, ort_inputs)
36     return ort_output
```

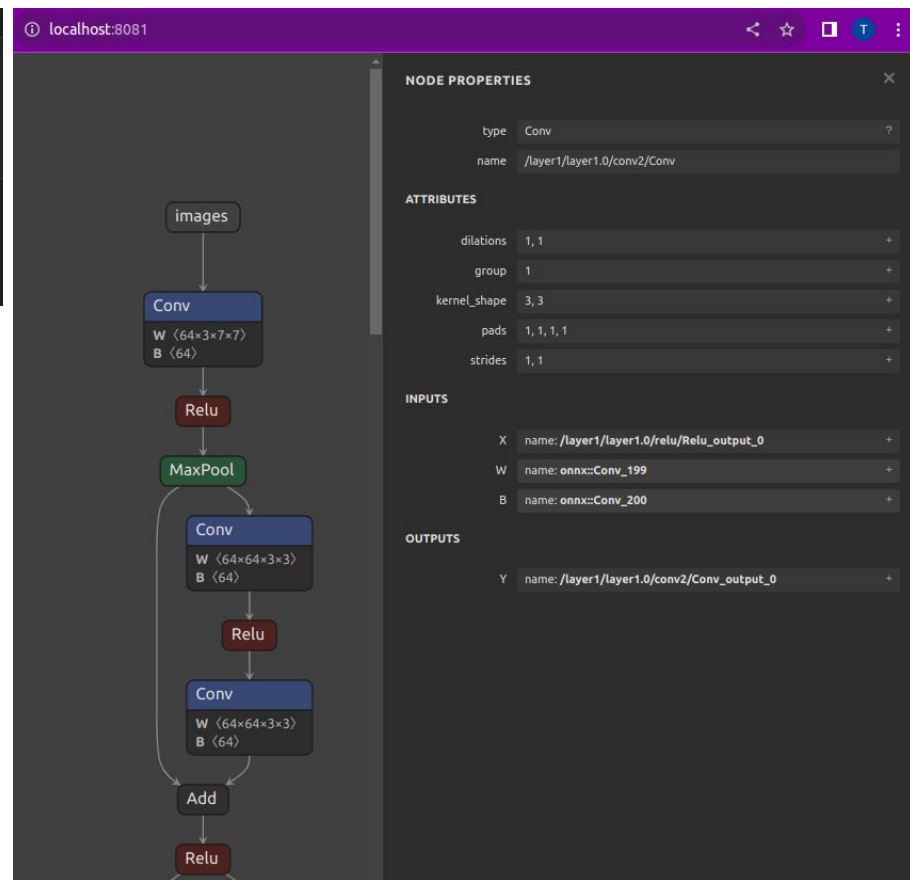
## ONNX Visualization

```
1 #!pip install netron
2 import netron
3 netron.start('models/onnx/vgg19_float32bit.onnx', 8081)

[38] ✓ 0.5s

... Serving 'models/onnx/vgg19_float32bit.onnx' at http://localhost:8081

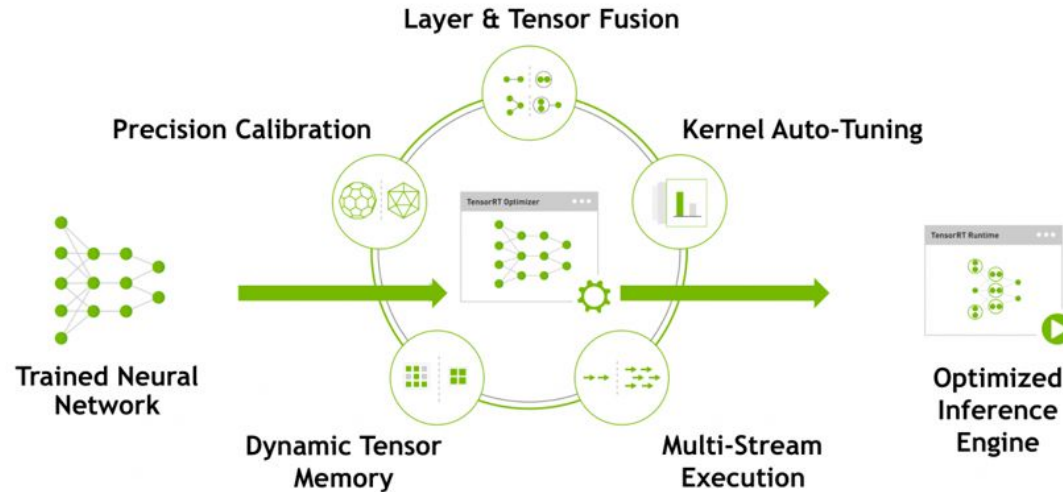
... ('localhost', 8081)
```



## TensorRT ( Tensor Run Time )

**TensorRT** là high-performance deep learning cho inference model. Đặc biệt chuyên tối ưu hóa hiệu suất của mô hình trên GPU của NVIDIA.

**TensorRT** giúp **Pruning** ( Loại bỏ những trọng số không cần thiết), **Quantization** ( Giảm thiểu số lượng bit biểu diễn trọng số ) và **Fusion** ( Ghép các phép tính lại với nhau ).

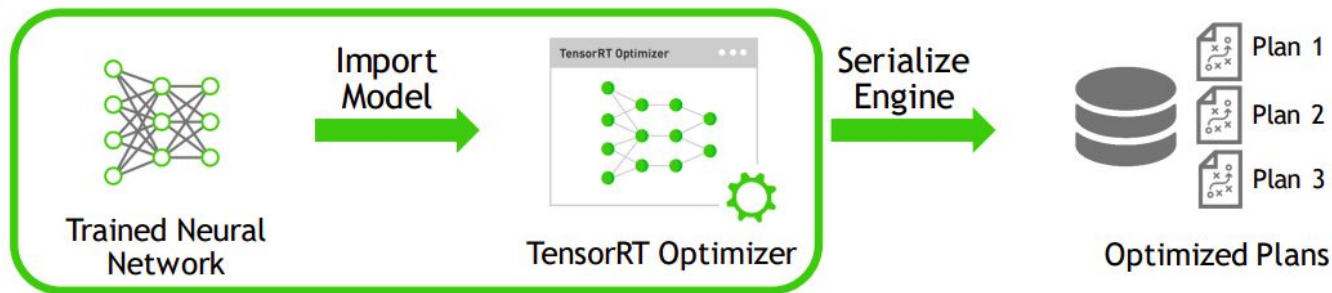




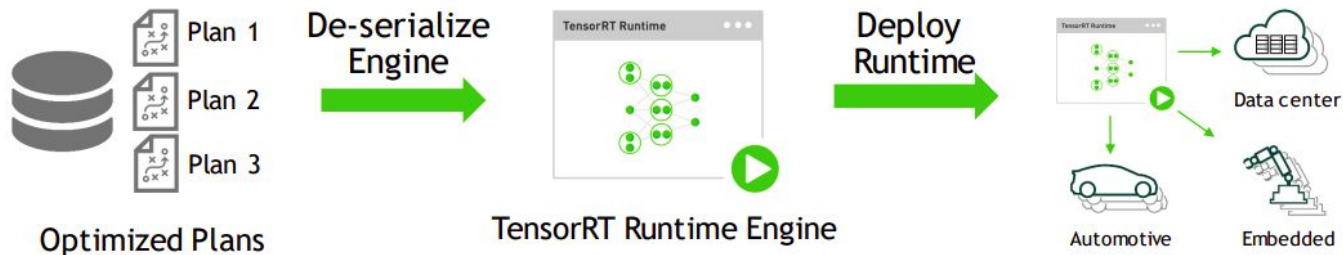
## TensorRT in Deployments

# TENSORRT DEPLOYMENT WORKFLOW

### Step 1: Optimize trained model



### Step 2: Deploy optimized plans with runtime



## TensorRT in Deployments

TensorRT supported layers:

- Convolution
- LSTM and GRU
- Activation: ReLU, Tanh, Sigmoid
- Pooling: max and average
- Scaling
- Element wise operations
- LRN
- Fully-connected
- Softmax
- Deconvolution

TensorRT cung cấp Custom Layer API.

- Dùng C++ để code cuda để định nghĩa layer

