

Neural Network Optimization with OpenVINO™ NNCF

Alexander Suslov



Neural Network Applications

Self-Driving Car



This image is in the public domain

Robots



This image is in the public domain

Image processing



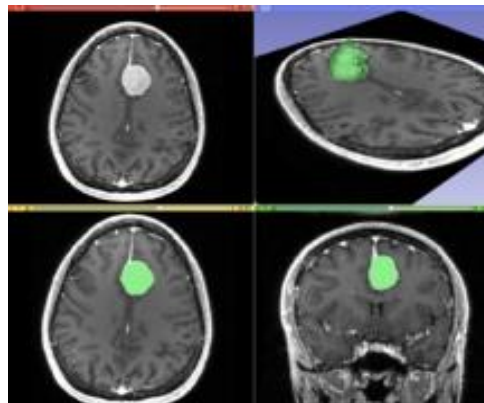
This image is in the public domain

Machine Translation



This image is in the public domain

Medical imaging



This image is in the public domain

3D scanning



This image is in the public domain

Where does Inference of Neural Networks Compute?

Standalone

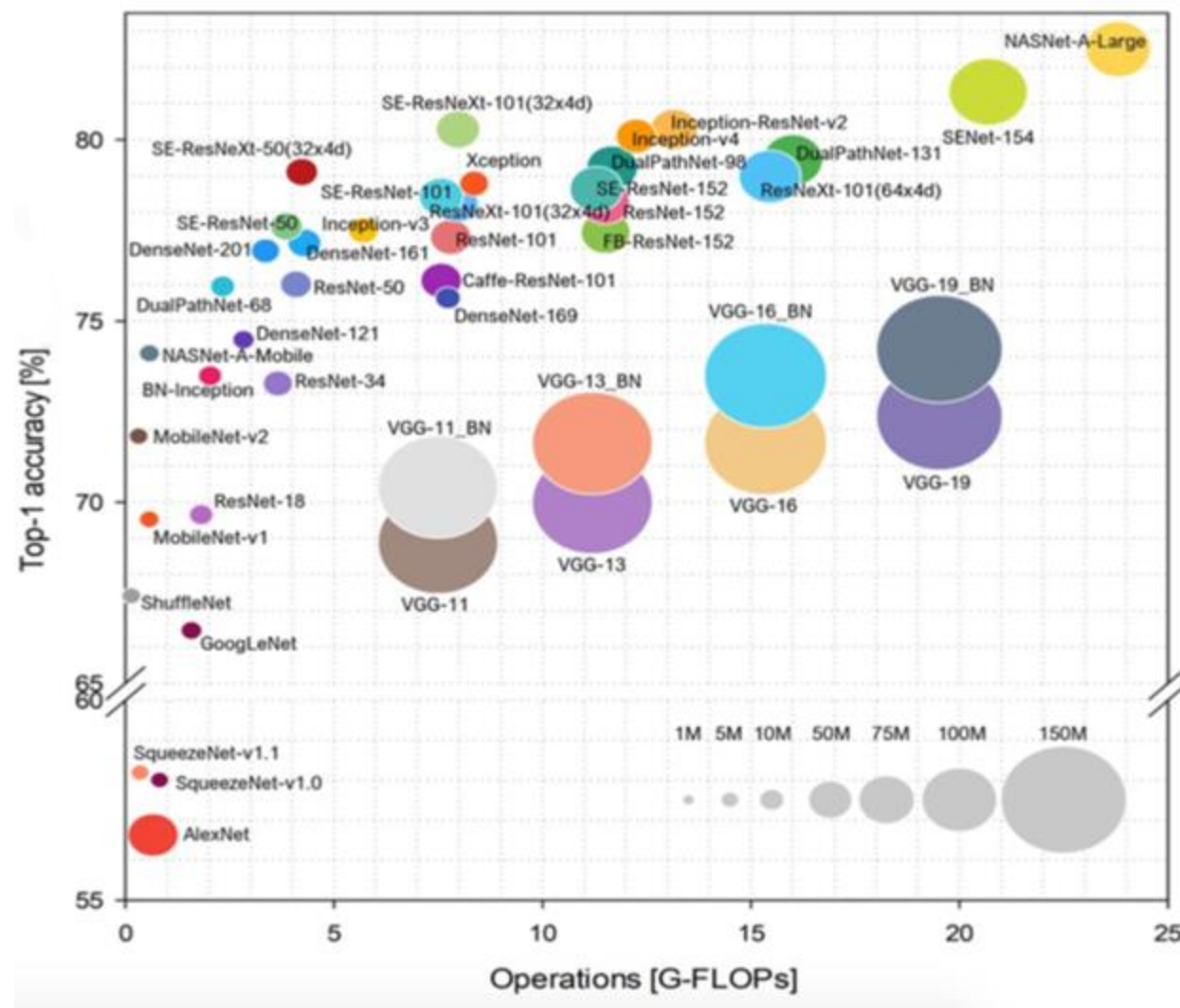


Client-Server



Models are Getting Larger

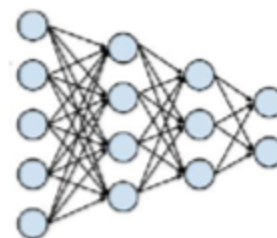
- While models are becoming more efficient, high accuracy still implies high complexity



From: [Benchmark Analysis of Representative Deep Neural Network Architectures](#), Simone Bianco et al,

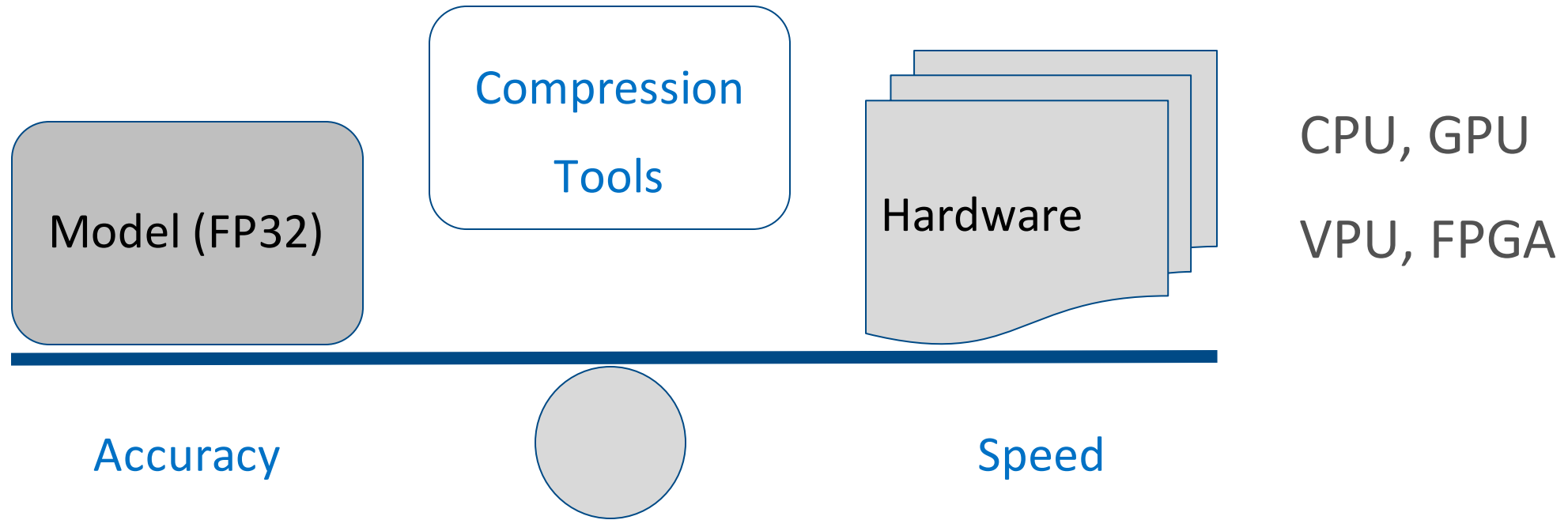
The First Challenge: Model Size

- Hard to distribute large models through over-the-air update
- The first run is slow due to loading weights.



All images are in the public domain

The Second Challenge: Speed



Tradeoff between accuracy and performance

Solutions from OpenVINO™ Toolkit

■ Training Extensions

- https://github.com/openvinotoolkit/training_extensions

■ Compression Tools:

- Post-Training Optimization Toolkit (POT)
 - https://docs.openvinotoolkit.org/latest/pot_README.html
- Neural Network Compression Framework (NNCF)
 - <https://github.com/openvinotoolkit/nncf>

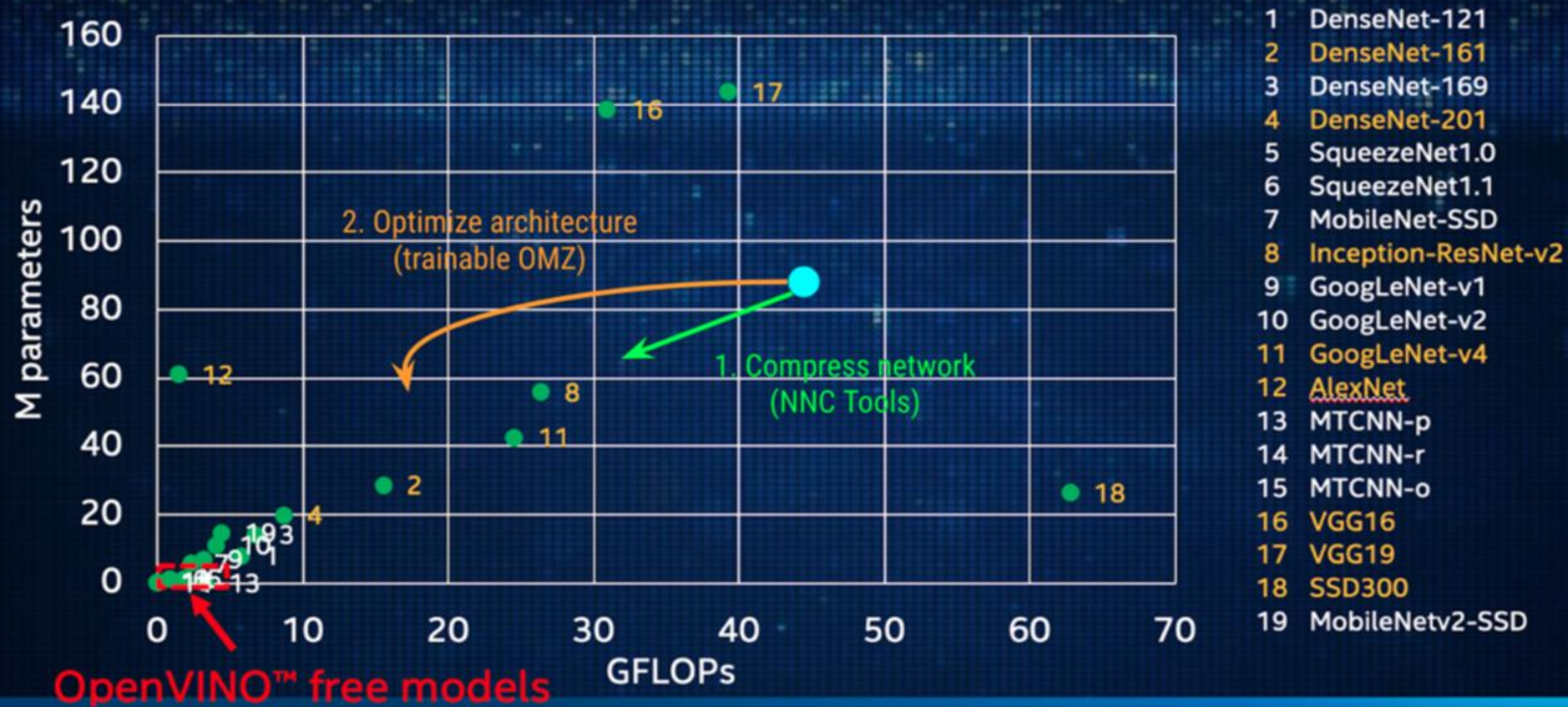
Solutions from OpenVINO™ Toolkit

■ Training Extensions

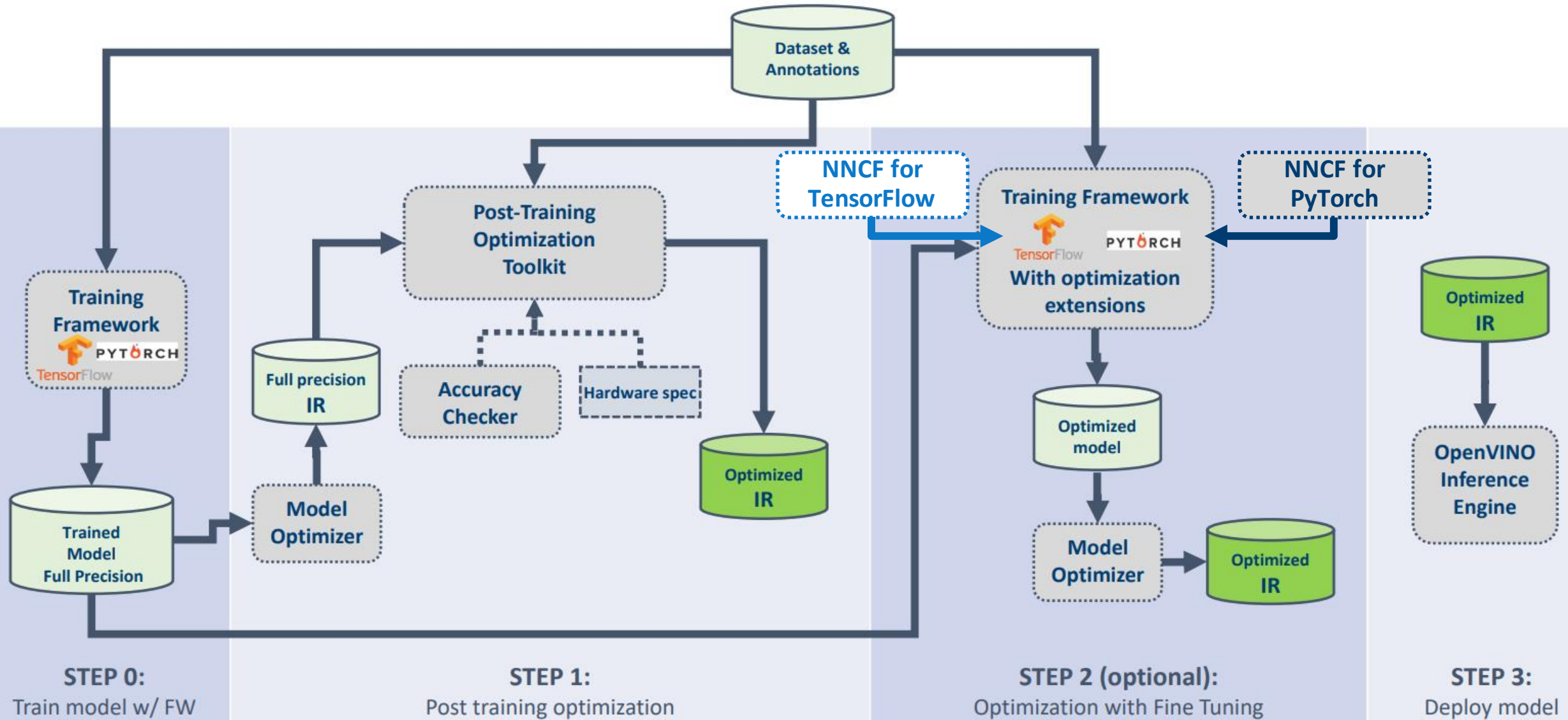
- https://github.com/openvinotoolkit/training_extensions

■ Compression Tools:

- Post-Training Optimization Toolkit (POT)
 - https://docs.openvinotoolkit.org/latest/pot_README.html
- Neural Network Compression Framework (NNCF)
 - <https://github.com/openvinotoolkit/nncf>



Optimization Customer Flow in OpenVINO™



Neural Network Compression Framework (NNCF)

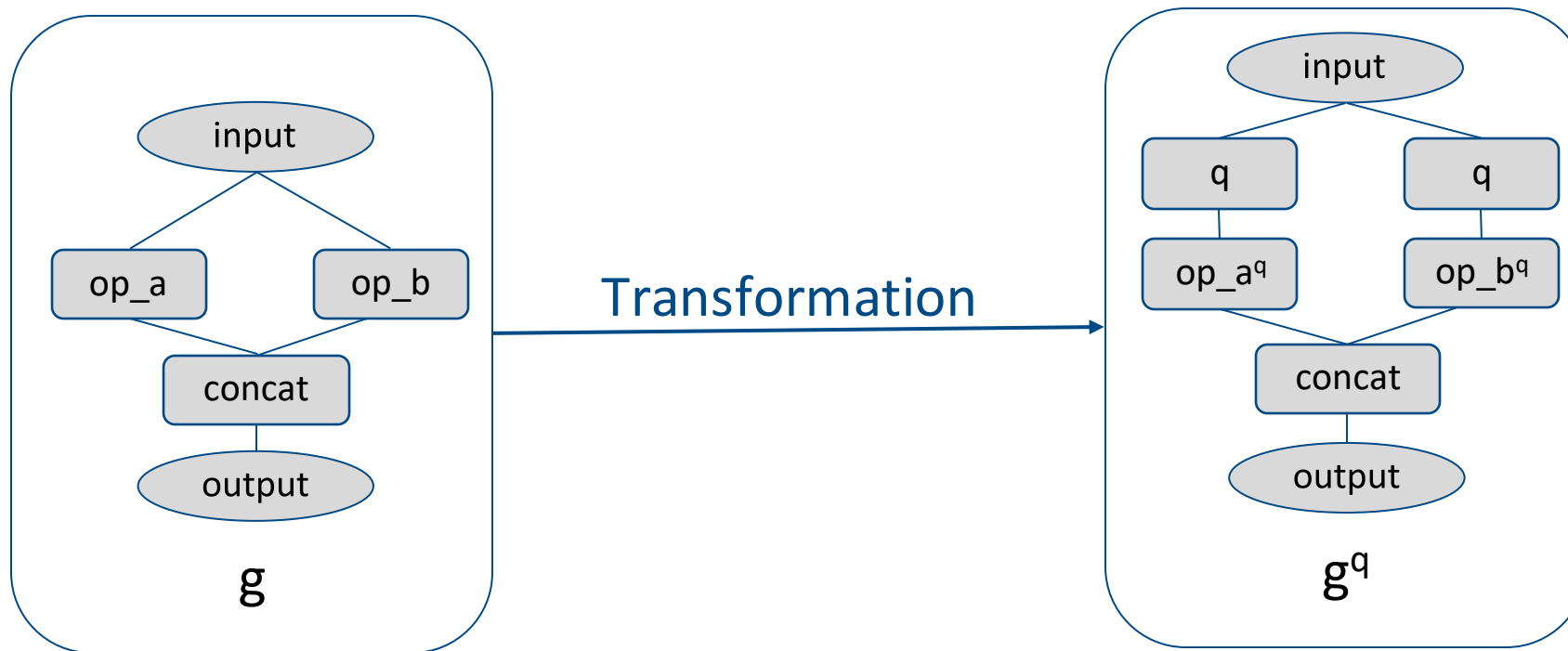
- Support of various compression algorithms, applied during a model fine-tuning process :
 - Quantization
 - Binarization
 - Sparsity
 - Filter pruning
- Automatic, configurable model graph transformation to obtain the compressed model.

Neural Network Quantization

- For any neural network*
- What we do:
 - Store weights in n-bits?
 - Do calculations in n-bits?
- Why run quantized models?
 - Latency
 - Memory usage (i.e., 32bit float → 8bit fixed)
 - Power Consumption

Quantization: Overview

This is the process of transforming a neural network such that it can be represented and executed at a lower precision by discretizing the original neural network weights and activations.

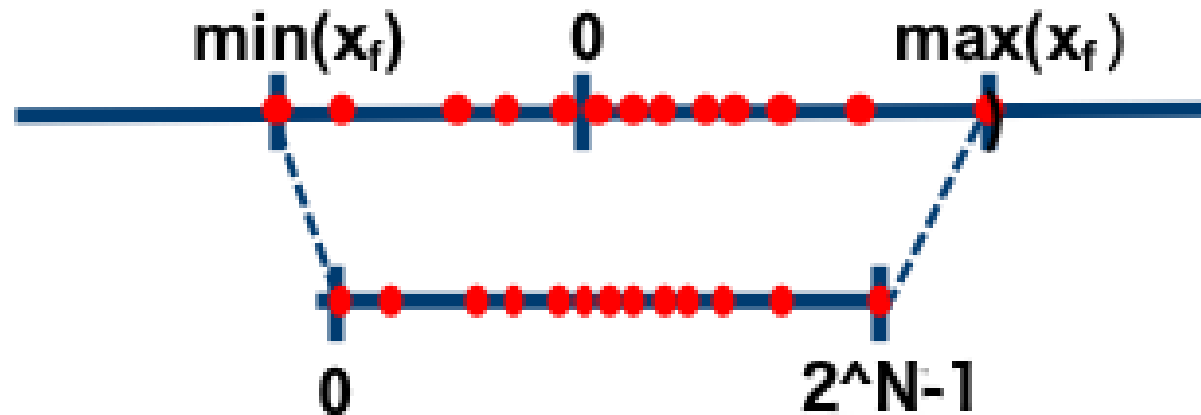


Quantization: Overview

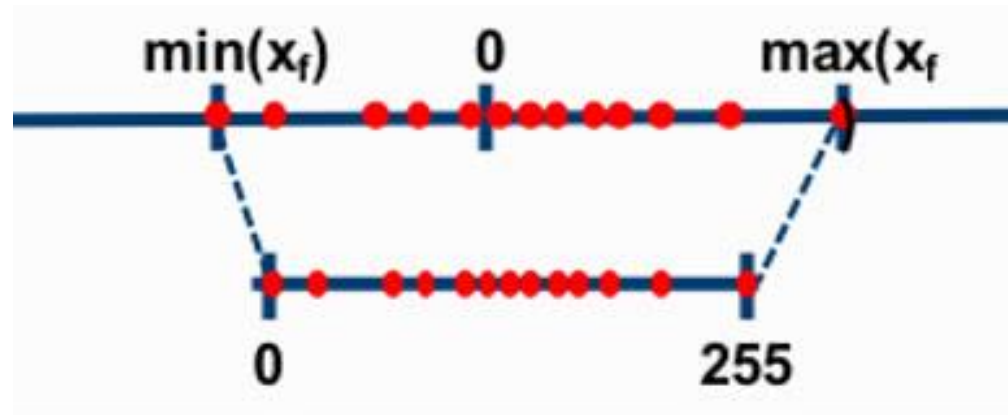
- **Quantization function** is a mapping of values from high to low precision
- **Neural network transformation** is the process of getting g' from g
- **Quantization algorithm** computes quantization parameters required by new neural network g^q and optimizes g^q via fine-tuning:
 - Post-training quantization without dataset
 - Post-training quantization with dataset
 - Quantization aware training

Quantization: Quantization function

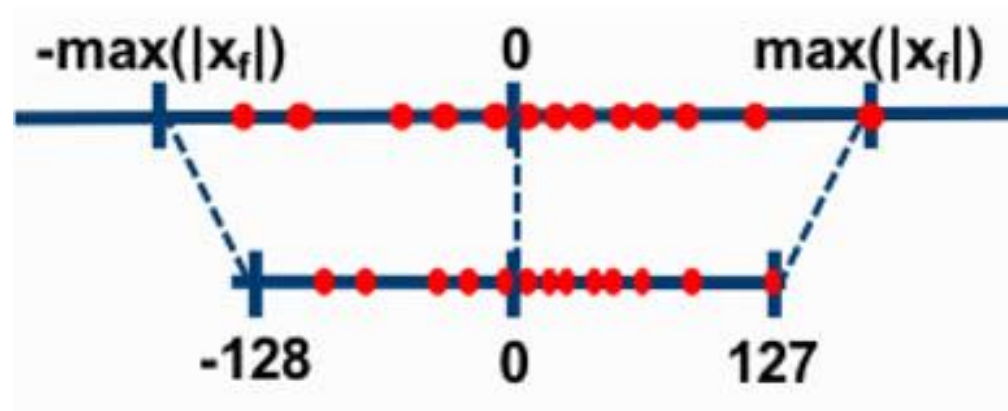
- Quantization refers to mapping values from fp32 to a lower precision format with specified parameters:
 - Precision
 - Quantization type
 - Granularity



Quantization: Quantization types



Asymmetric Mode



Symmetric Mode

Asymmetric Quantization

- Quantization:

$$x_{int} = \text{round}\left(\frac{x}{\Delta}\right) + z$$

$$x_Q = \text{clamp}(0, N_{levels} - 1, x_{int})$$

- Δ specifies the step size of the quantizer and floating point zero maps to zero-point.
- z - zero-point.
- $N_{levels} = 256$ or 8-bits of precision

- De-quantization:

$$x_{float} = (x_Q - z)\Delta$$

Asymmetric Quantization

- 2D convolution:

$$y(k, l, n) = \Delta_w \Delta_x \text{conv}(w_Q(k, l, m; n) - z_w, x_Q(k, l, m) - z_x)$$

$$y(k, l, n) = \text{conv}(w_Q(k, l, m; n), x_Q(k, l, m)) - z_w \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} x_Q(k, l, m) \\ - z_x \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} w_Q(k, l, m; n) + z_x z_w$$

Symmetric Quantization

■ Quantization, zero-point = 0

- Activations:

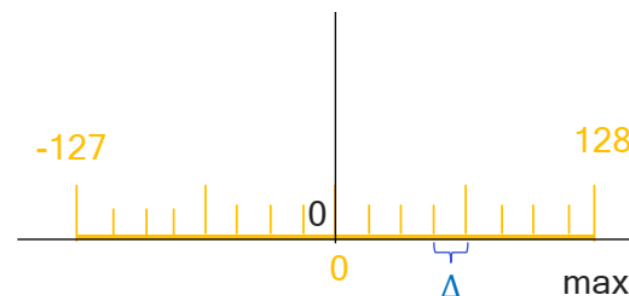
$$x_{int} = \text{round}\left(\frac{x}{\Delta}\right)$$

$$x_Q = \text{clamp}(-N_{levels}/2, N_{levels}/2 - 1, x_{int})$$

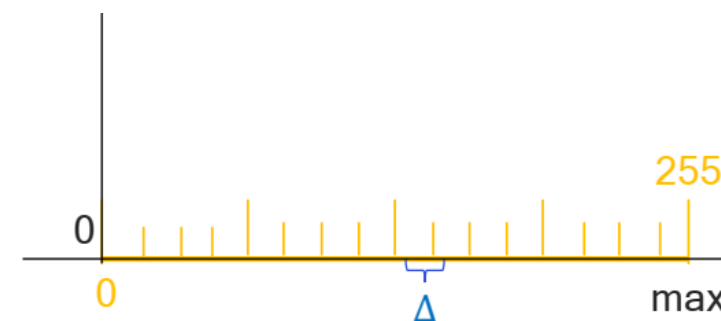
$$x_Q = \text{clamp}(0, N_{levels} - 1, x_{int})$$

if signed
if un-signed

Symmetric signed



Symmetric unsigned



- Weights:

$$x_Q = \text{clamp}(-(N_{levels}/2 - 1), N_{levels}/2 - 1, x_{int})$$

$$x_Q = \text{clamp}(0, N_{levels} - 2, x_{int})$$

if signed
if un-signed

Symmetric vs Asymmetric Quantization

Symmetric Quantization

Asymmetric Quantization

Same calculation

Unavoidable overhead

$$y(k, l, n) = \text{conv}(w_Q(k, l, m; n), x_Q(k, l, m))$$

$$y(k, l, n) = \text{conv}(w_Q(k, l, m; n), x_Q(k, l, m)) - z_w \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} x_Q(k, l, m)$$

$$- z_x \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} w_Q(k, l, m; n) + z_x z_w$$

Precompute

Symmetric vs Asymmetric Quantization

Symmetric Quantization

Asymmetric Quantization

Same calculation

Unavoidable overhead

$$y(k, l, n) = \text{conv}(w_Q(k, l, m; n), x_Q(k, l, m))$$

$$y(k, l, n) = \text{conv}(w_Q(k, l, m; n), x_Q(k, l, m)) - z_w \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} x_Q(k, l, m)$$

$$- z_x \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} w_Q(k, l, m; n) + z_x z_w$$

Precompute

Symmetric weights/Asymmetric activations are the best option

Quantization granularity

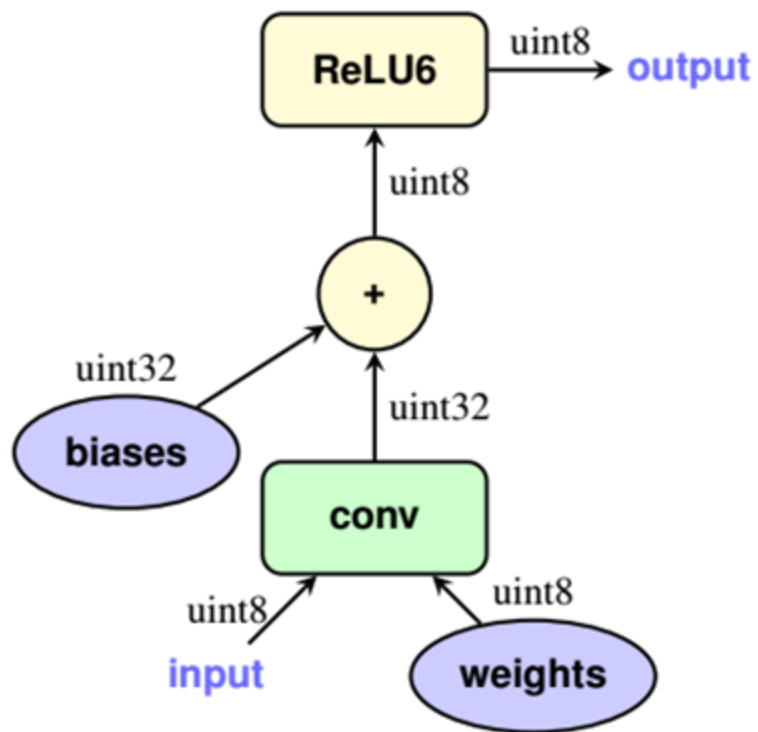
- Per-tensor quantization
 - Same quantization parameters for all elements in a tensor.
- Per-channel quantization:
 - Own quantization parameters for each output channel

Fake-quantization

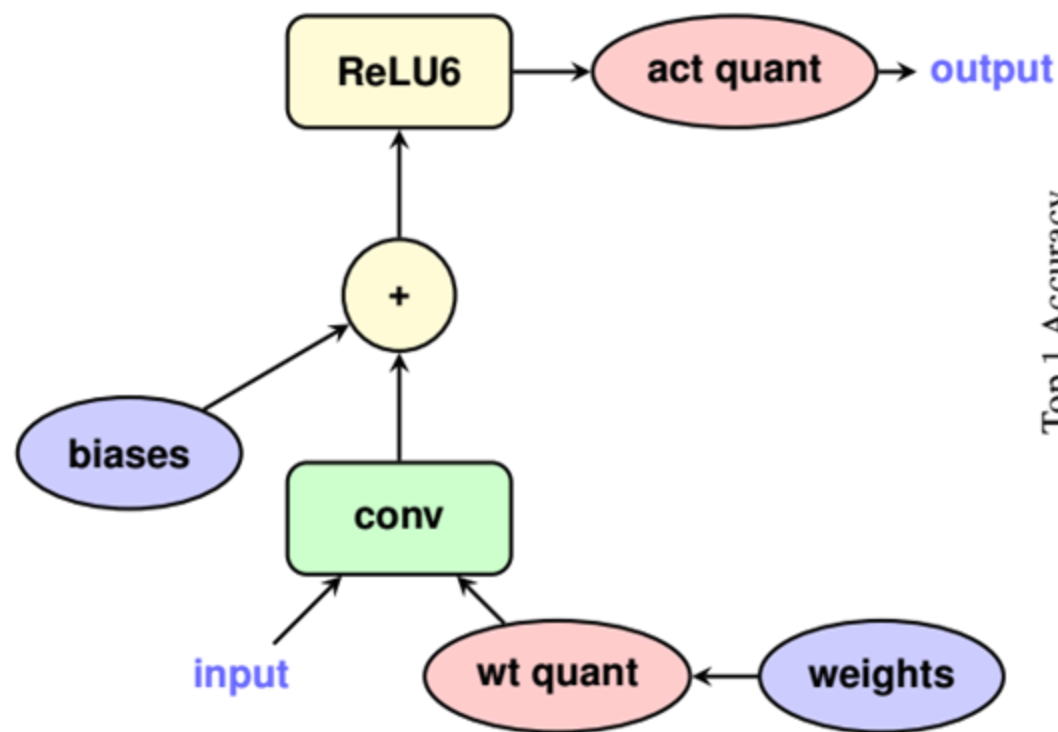
- Simulates quantization in floating point by quantizing and de-quantizing the input. The output is still in floating point, but with reduced precision.

$$\begin{aligned}\text{clamp}(r; a, b) &:= \min(\max(x, a), b) \\ s(a, b, n) &:= \frac{b - a}{n - 1} \\ q(r; a, b, n) &:= \left\lfloor \frac{\text{clamp}(r; a, b) - a}{s(a, b, n)} \right\rfloor s(a, b, n) + a, \\ &\quad (12)\end{aligned}$$

Quantization: Neural Network Transformation



(a) Integer-arithmetic-only inference



(b) Training with simulated quantization

Top 1 Accuracy

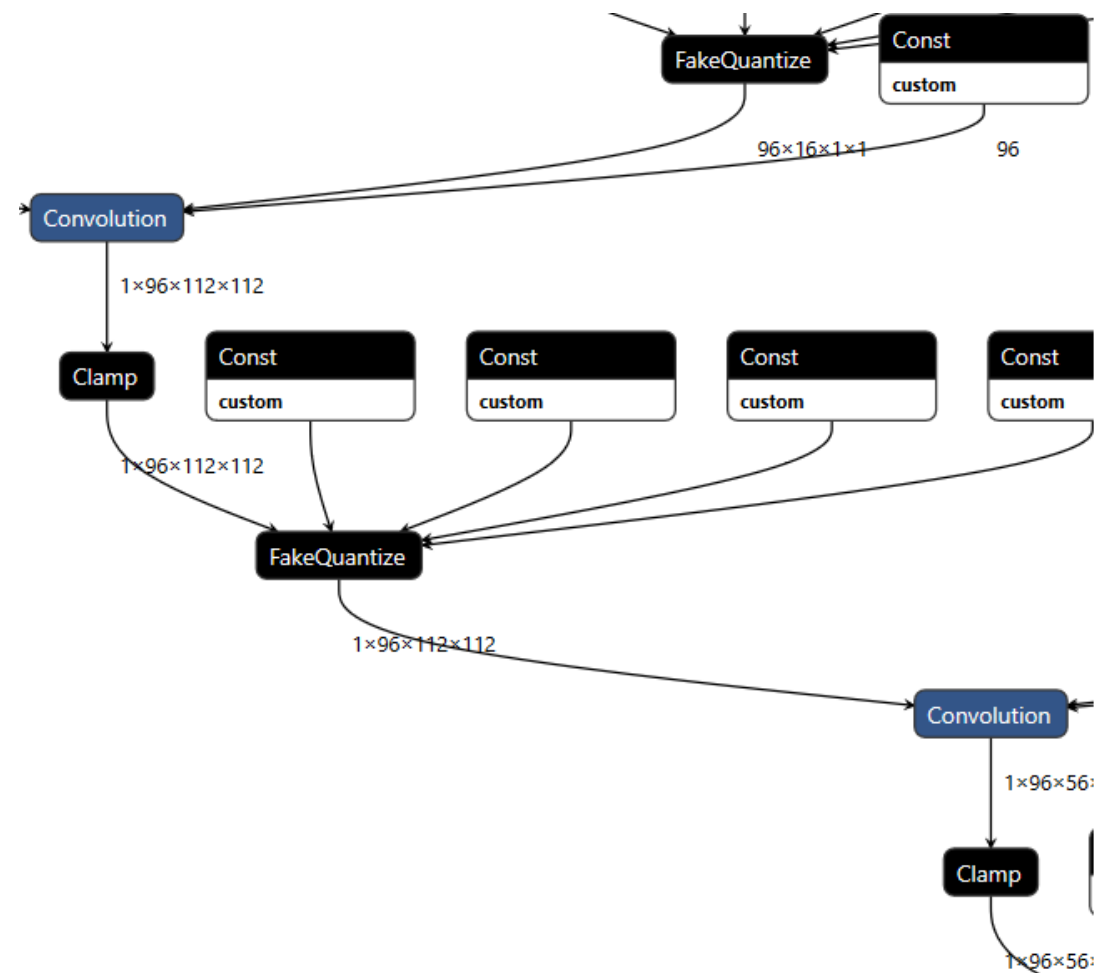
Quantization aware training

Step1: Create a training graph of the floating-point model.

Step2: Insert Fake Quantization Layers

Step3: Train in simulated quantized mode until convergence.

Step4: Export quantized model in the format that is supported by OpenVINO™ Toolkit

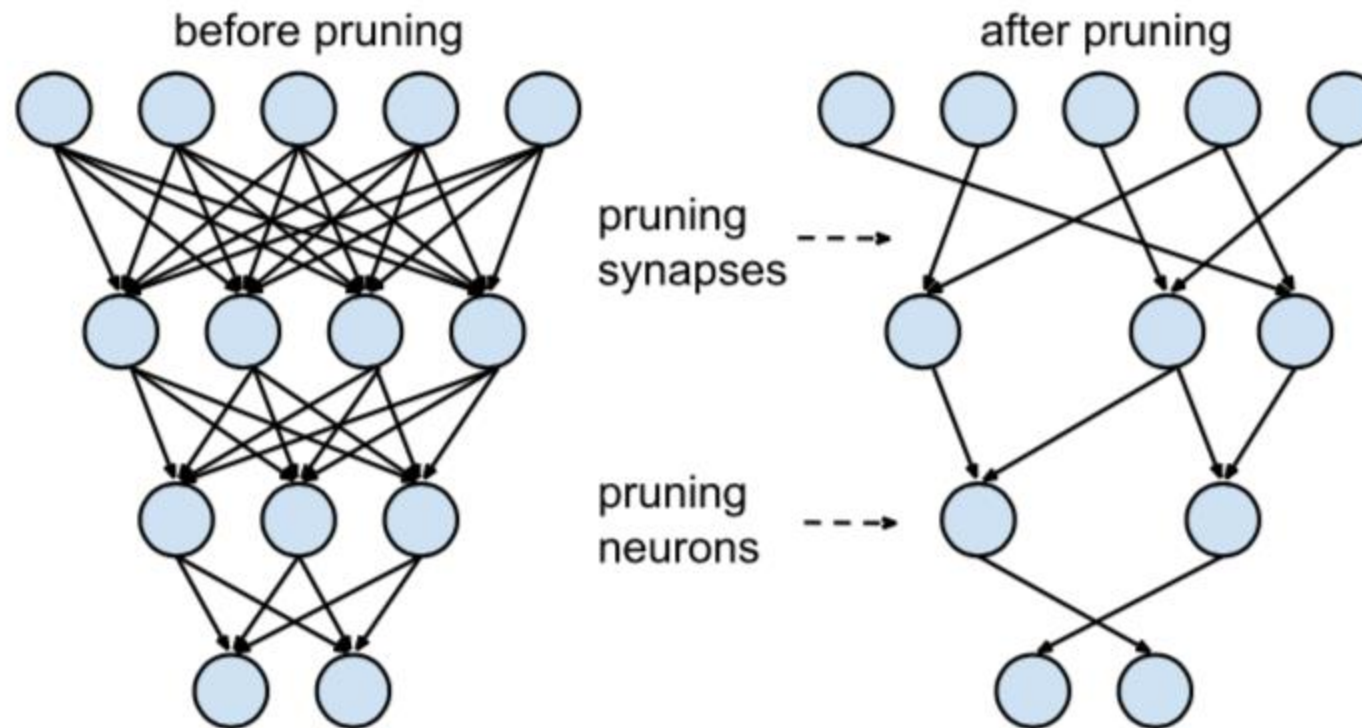


Quantization Results

Model	Compression algorithm	Dataset	PyTorch FP32 baseline	PyTorch compressed accuracy
ResNet-50	INT8	ImageNet	76.13	76.08
Inception V3	INT8	ImageNet	77.32	76.90
MobileNet V2	INT8	ImageNet	71.81	71.29
SqueezeNet V1.1	INT8	ImageNet	58.18	58.07
SSD300-BN	INT8	VOC12+07	78.28	78.08
SSD512-BN	INT8	VOC12+07	80.26	80.11
UNet	INT8	CamVid	71.95	71.66
ICNet	INT8	CamVid	67.89	67.85
UNet	INT8	Mapillary	56.23	56.1
BERT-base-chinese	INT8	XNLI	77.68	77.22
BERT-large (Whole Word Masking)	INT8	SQuAD v1.1	93.21 (F1)	92.68 (F1)
RoBERTa-large	INT8	MNLI	90.6 (matched)	89.25 (matched)
DistilBERT-base	INT8	SST-2	91.1	90.3
MobileBERT	INT8	SQuAD v1.1	89.98 (F1)	89.4 (F1)
GPT-2	INT8	WikiText-2 (raw)	19.73 (perplexity)	20.9 (perplexity)
RetinaNet-ResNet50-FPN	INT8	COCO2017	35.6 (avg bbox mAP)	35.3 (avg bbox mAP)
RetinaNet-ResNeXt101-64x4d-FPN	INT8	COCO2017	39.6 (avg bbox mAP)	39.1 (avg bbox mAP)
Mask-RCNN-ResNet50-FPN	INT8	COCO2017	40.8 (avg bbox mAP), 37.0 (avg segm mAP)	40.6 (avg bbox mAP), 36.5 (avg segm mAP)

https://github.com/openvinotoolkit/nncf_pytorch

Pruning Neural Networks

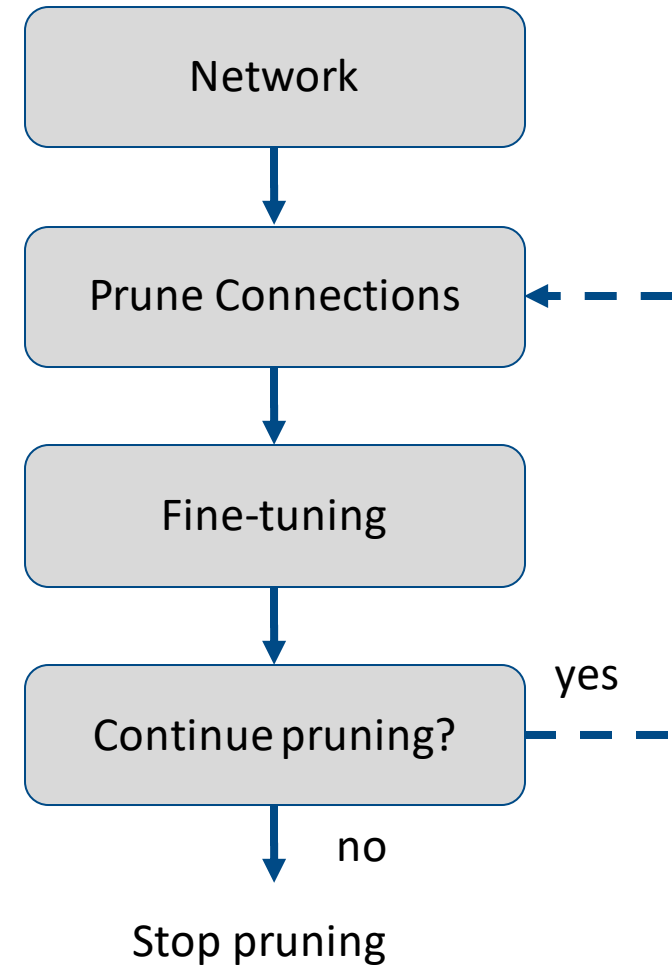


Pruning Neural Networks

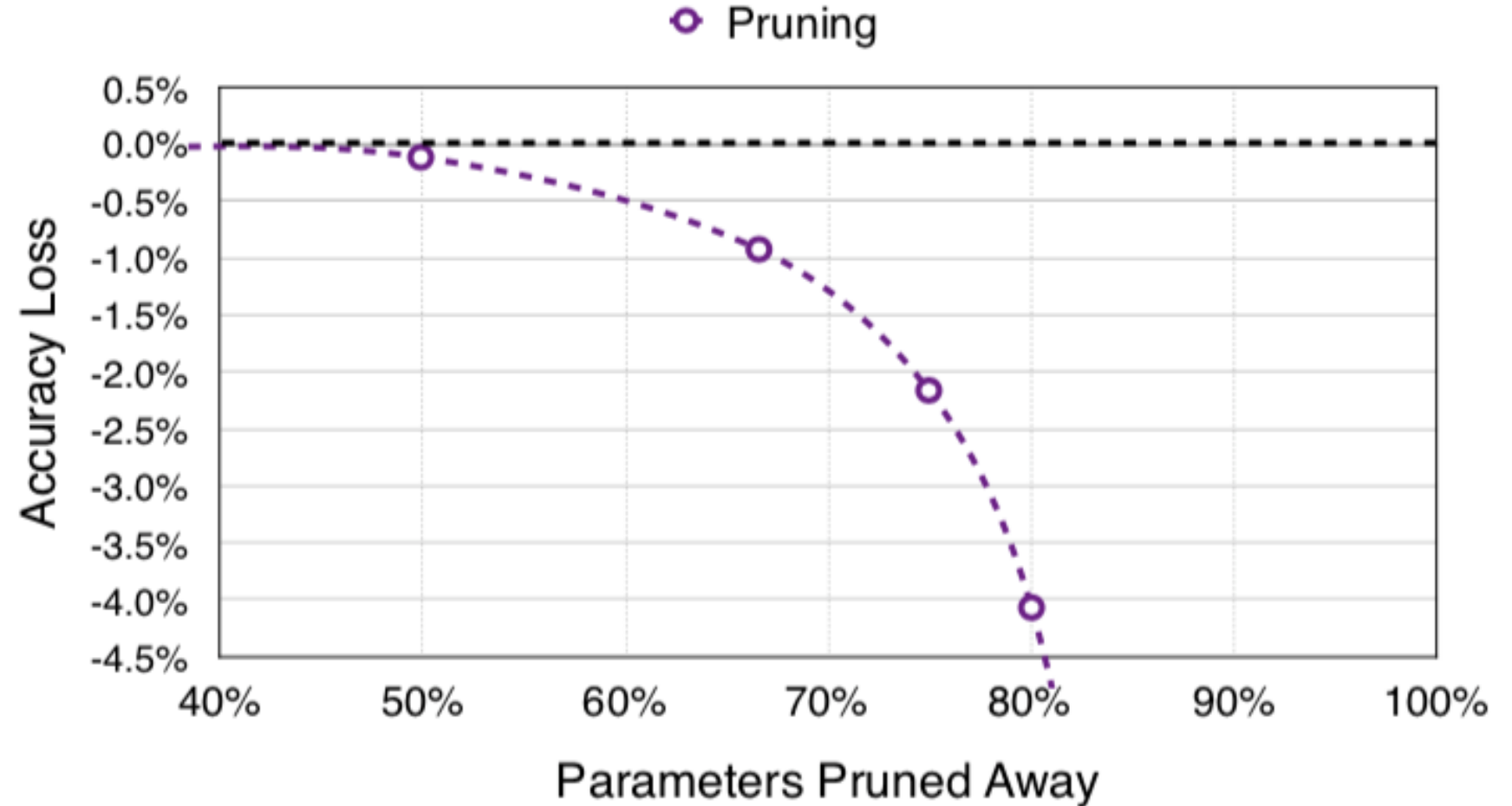
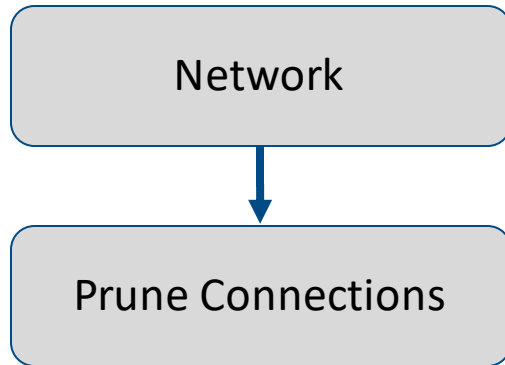
- Neural network pruning benefits
 - Reducing the binary model size
 - Smaller models reduce memory bandwidth bottlenecks
 - Faster kernels (depends on hardware support)

Pruning Neural Networks

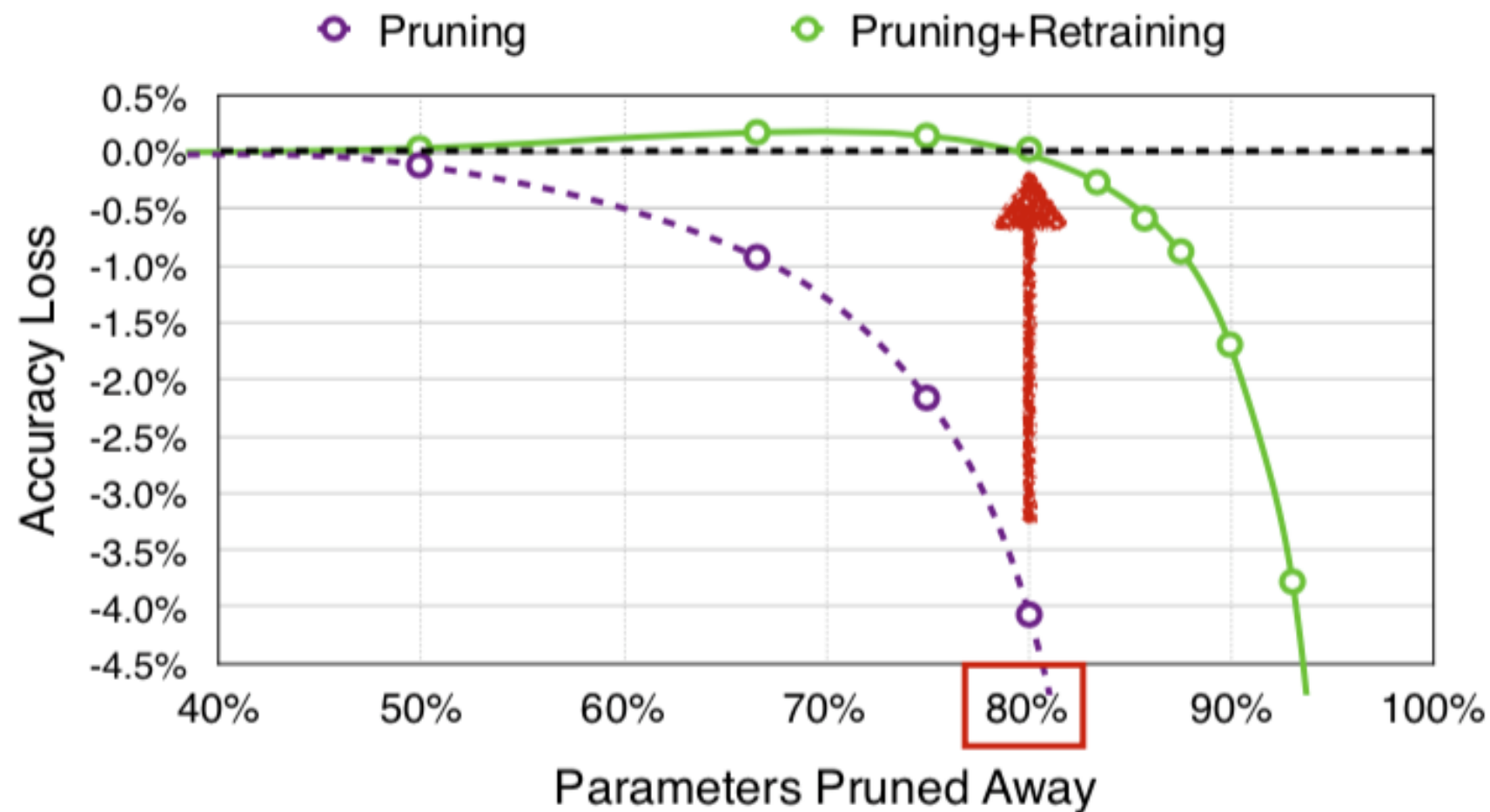
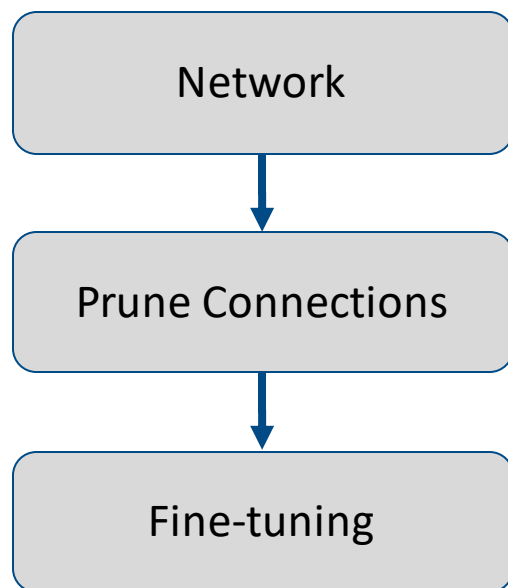
- Criteria for pruning
 - Connections with low weights
 - Neurons or filters with low impact
- External limitations
 - Spatial structure of pruned weights



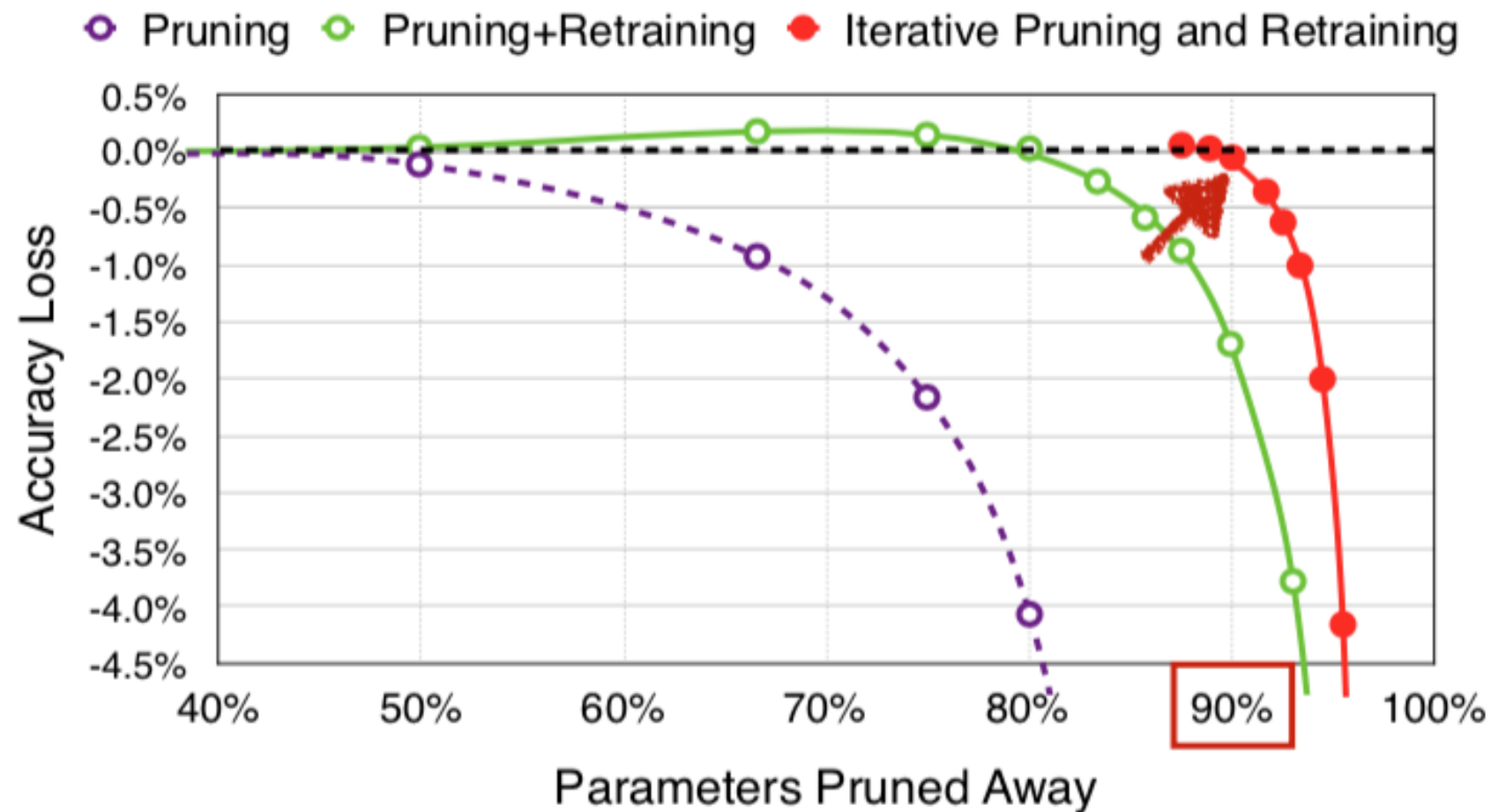
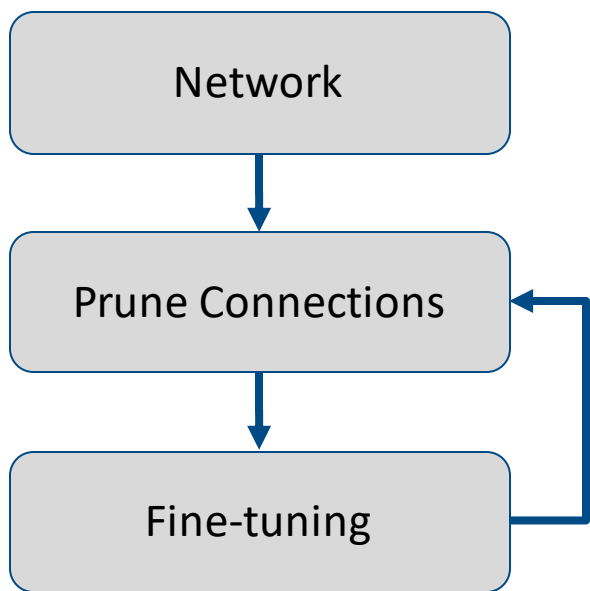
Pruning Neural Networks



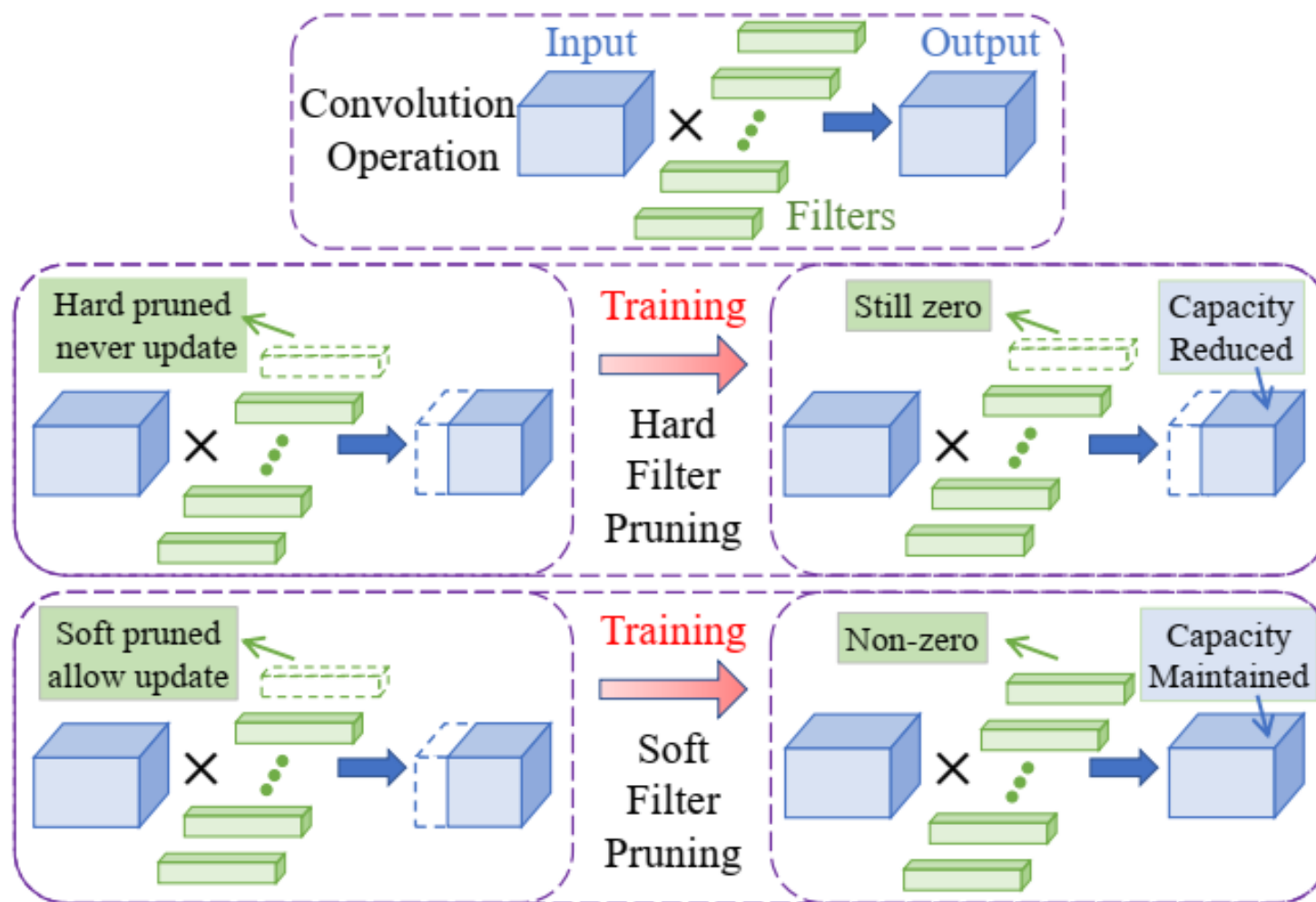
Retrain to Recover Accuracy



Iteratively Retrain to Recover Accuracy



Filter Pruning: Main Approaches



Filter Pruning: Filter Importance Criteria

- L1, L2

$$||F||_p = \sqrt[p]{\sum_{c,k_1,k_2=1}^{C,K,K} |F(c,k_1,k_2)|^p}$$

- “Geometric Median”

$$G(F_i) = \sum_{F_j \in \{F_1, \dots, F_m\}, j \neq i} ||F_i - F_j||_2$$

Filter Pruning: Results

Models	Compression algorithm	Dataset	Top-1 Accuracy FP32 model (%)	Top-1 Accuracy Pruned model (%)
ResNet-18	Filter pruning, 30%, magnitude criterion	ImageNet	69.76	68.69
ResNet-18	Filter pruning, 30%, geometric median criterion	ImageNet	69.76	68.97
ResNet-34	Filter pruning, 30%, magnitude criterion	ImageNet	73.31	72.54
ResNet-34	Filter pruning, 30%, geometric median criterion	ImageNet	73.31	72.60
ResNet-50	Filter pruning, 30%, magnitude criterion	ImageNet	76.13	75.7
ResNet-50	Filter pruning, 30%, geometric median criterion	ImageNet	76.13	75.7

https://github.com/openvinotoolkit/nncf_pytorch

Sparsification: Main Approaches

■ Magnitude Sparsity

- After each training epoch the method calculates a threshold based on the current sparsity ratio and uses it to zero weights which are lower than this threshold.

■ Regularization-Based (RB) Sparsity

- The sparsification algorithm based on probabilistic approach and loss regularization.

Ordinary convolution

$$output = conv(x, w)$$

Sparsifying weights we reparametrize weights as:

$$\bar{w} = w * z$$

Where:

w – weights, $w \in R$

z – binary mask, $z \in [0, 1]$

We train these masks using modified loss

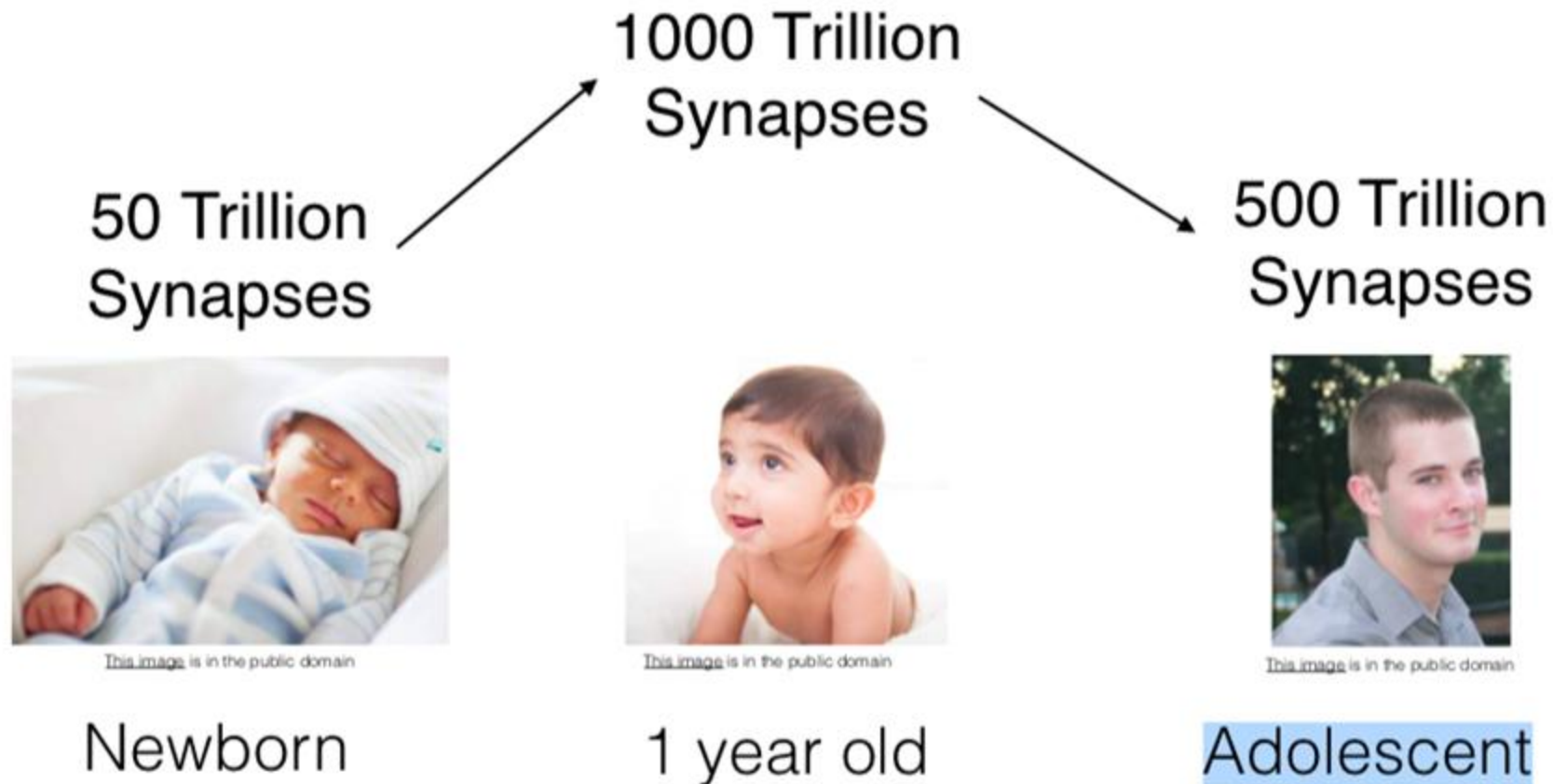
$$Loss = Loss_{task} + \alpha \left(\sum_l^{Layers} ||z|| - target \right)^2$$

Sparsification: Results

Models	Compression algorithm	Dataset	Top-1 Accuracy FP32 model (%)	Top-1 Accuracy Pruned model (%)
inception-v3	RB-sparsity, 50% sparsity rate	ImageNet	77.46	77.25
inception-v3	Magnitude sparsity, 50% sparsity rate	ImageNet	77.46	77.24
inception-v3	RB-sparsity, 92% sparsity rate	ImageNet	77.46	76.6
mobilenet-v2	RB-sparsity, 50% sparsity rate	ImageNet	71.8	71.2
mobilenet-v2	Magnitude sparsity, 50% sparsity rate	ImageNet	71.8	70.8
mobilenet-v2	RB-sparsity, 78% sparsity rate	ImageNet	71.8	69.98

https://github.com/openvinotoolkit/nncf_pytorch

Pruning Happens in Human Brain



Christopher A Walsh. Peter Huttenlocher (1931-2013). Nature, 502(7470):172–172, 2013.

Slide credits by Song Han

Usage NNCF: Injection

```
import torch
import nncf # Important - should be imported directly after torch
from nncf import create_compressed_model, NNCFConfig, register_default_init_args

# Instantiate your uncompressed model
from torchvision.models.resnet import resnet50
model = resnet50()

# Load a configuration file to specify compression
nncf_config = NNCFConfig.from_json("resnet50_int8.json")

# Provide data loaders for compression algorithm initialization, if necessary
nncf_config = register_default_init_args(nncf_config, train_loader, loss_criterion)

# Apply the specified compression algorithms to the model
comp_ctrl, compressed_model = create_compressed_model(model, nncf_config)

# Now use compressed_model as a usual torch.nn.Module to fine-tune compression parameters along with the model weights

# ... the rest of the usual PyTorch-powered training pipeline

# Export to ONNX or .pth when done fine-tuning
comp_ctrl.export_model("compressed_model.onnx")
torch.save(compressed_model.state_dict(), "compressed_model.pth")
```

Usage NNCF: Config file

Quantization

```
"compression": {  
  "algorithm": "quantization",  
  "weights": {  
    "mode": "symmetric",  
    "bits": 8,  
    "per_channel": true  
  },  
  "activations": {  
    "mode": "asymmetric",  
    "bits": 8,  
    "per_channel": false  
  }  
}
```

Filter Pruning

```
"compression": {  
  "algorithm": "filter_pruning",  
  "params": {  
    "schedule": "baseline",  
    "pruning_target": 0.3  
  }  
}
```

Sparsity

```
"compression": {  
  "algorithm": "magnitude_sparsity",  
  "params": {  
    "schedule": "multistep",  
    "multistep_steps": [5, 10, 20, 30, 40],  
    "multistep_sparsity_levels": [0.1, 0.2, 0.3,  
                                  0.4, 0.5, 0.6],  
    "sparsity_freeze_epoch": 50  
  }  
}
```

Q&A

