Joint ICTP-IAEA School on Detector Signal Processing and Machine Learning for Scientific Instrumentation and Reconfigurable Computing

Model Optimization and Compression Techniques

smr4110 | Trieste, Italy - 2025

Romina Soledad Molina, Ph.D.

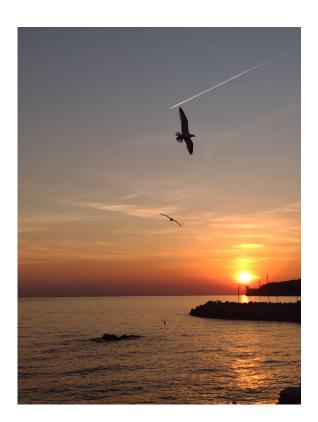






#### **Outline**

- ML and model compression techniques.
- Pruning.
- Quantization.
- Knowledge Distillation.
- How do we combine compression techniques?.
- Hyperparameters tuning and compression.
- Demo: MNIST-based binary classification QAP -





## Machine Learning and model compression techniques



• Question: What is compression?













• **Compression** is the process of reducing the size of data, often with the aim of preserving important information or, in some cases, enabling perfect reconstruction. The goal is to make data storage or transmission more efficient by using fewer resources (e.g., memory, storage, bandwidth).

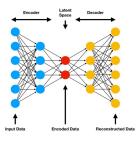


• Question: Why is important in Machine Learning?





 Compression, in the context of Machine Learning, involves techniques aiming to reduce the size of models or datasets while preserving performance.

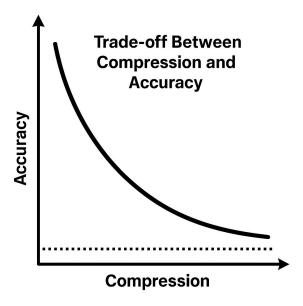




- **Compression**, in the context of Machine Learning, involves techniques aiming to reduce the size of models or datasets while preserving performance.
  - Model Compression: Reducing the size of machine learning models (such as neural networks) without significantly affecting their accuracy.
  - Data compression: Reducing the size of the data used for training, validation, and testing.



- Trade-off Between Compression and Accuracy
  - Balance between reducing the size of the model or data and maintaining its performance.





#### Applicability

 Not all compression techniques are suitable for every type of model or dataset.



#### Applicability

- Not all compression techniques are suitable for every type of model or dataset.
- The decision of which compression strategy to apply depends on factors such as the desired trade-off between model size and performance, the computational resources available, and the nature of the data being processed.



- Problem:
  - Very large models → high memory consumption and inference time.



The most accurate models (such as deep neural networks) are large and expensive in terms of memory and processing time.



The most accurate models (such as deep neural networks) are large and expensive in terms of memory and processing time.

Model	Parameters (millions)	Disk size (MB)
ResNet-50	~25.6M	~98 MB
BERT-base	~110M	~440 MB
BERT-large	~340M	~1.3 GB
VGG-16	~138M	~528 MB
VGG-19	~144M	~548 MB
YOLOv3	~62M	~236 MB
SpineNet-49S (small)	~11M	~45 MB
MobileNetV3-Large	~5.4M	~20 MB



**VGG-16** 

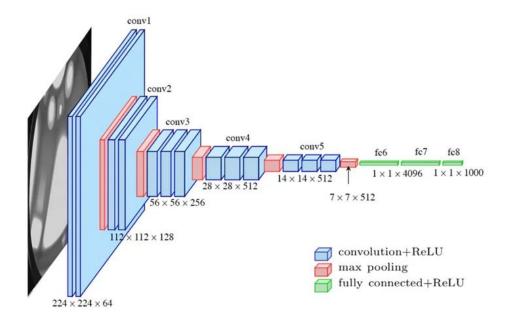
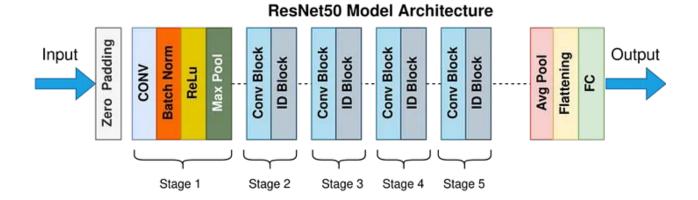


Image from https://lekhuyen.medium.com/an-overview-of-vgg16-and-nin-models-96e4bf398484



#### ResNet-50





SpineNet-49

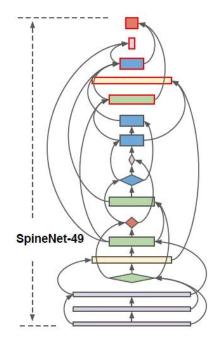


Image from Du, X., Lin, T. Y., Jin, P., Ghiasi, G., Tan, M., Cui, Y., ... & Song, X. (2020). Spinenet: Learning scale-permuted backbone for recognition and localization. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 11592-11601).

- Problem:
  - Very large models → high memory consumption and inference time.

- Solution:
  - Compression



#### **Distilled versions**

Model	Parameters (millions)	Disk size (MB)	Accuracy
DistilBERT	66M (↓40%)	~250 MB	97% of BERT
SqueezeNet	1.2M (↓99%)	~4.8 MB (↓99%)	Similar to VGG-16
YOLOv8-Nano	3.2M (↓92%)	~8 MB (↓90%)	Good trade-off with YOLOv8-L
YOLO-Fastest	0.2M (↓99.5%)	~1 MB (↓99%)	Lower accuracy



# Pruning, Quantization, and Knowledge Distillation

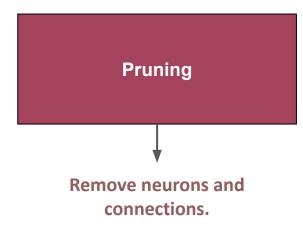


**Pruning** 

Quantization

**Knowledge distillation** 

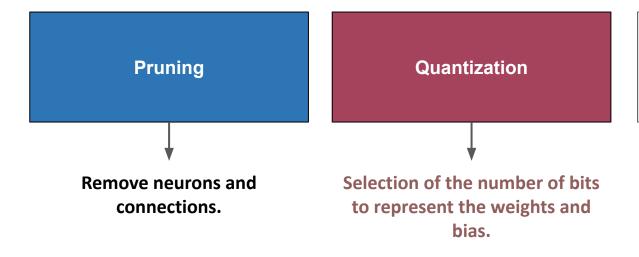




Quantization

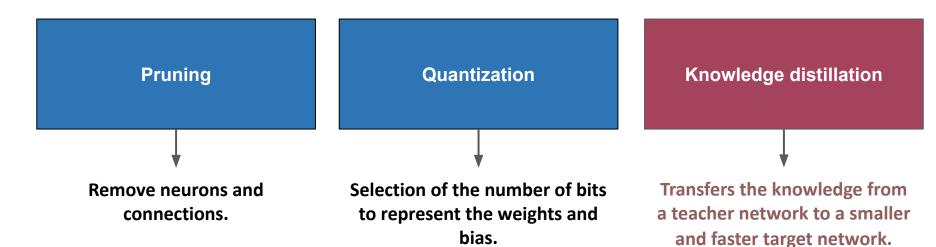
**Knowledge distillation** 



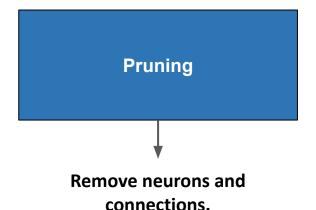


**Knowledge distillation** 









Selection of the number of bits to represent the weights and bias.

Quantization

Fully on-chip deployment

**Knowledge distillation** 

Transfers the knowledge from a teacher network to a smaller and faster target network.



## **Pruning**



**Pruning** 

This technique is used to reduce the size and complexity of a deep learning model by eliminating unnecessary weights or neurons.



**Pruning** 

- This technique is used to reduce the size and complexity of a deep learning model by eliminating unnecessary weights or neurons.
- The primary objective is to enhance the model's efficiency by decreasing memory usage and speeding up inference times, all while maintaining its performance.

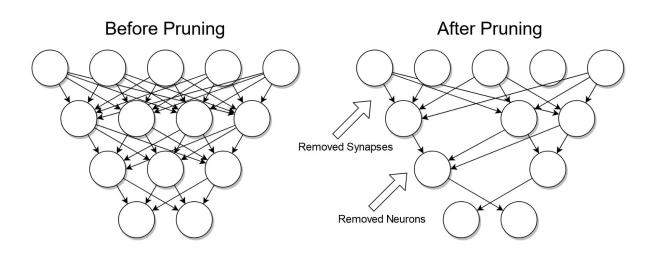


#### **Pruning**

- This technique is used to reduce the size and complexity of a deep learning model by eliminating unnecessary weights or neurons.
- The primary objective is to enhance the model's efficiency by decreasing memory usage and speeding up inference times, all while maintaining its performance.
- Prune weights: setting the weights of chosen individual parameters into zero.

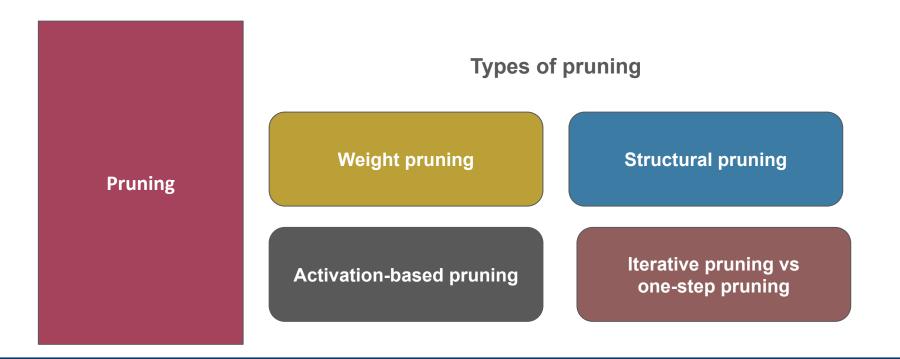


**Pruning** 



 $Image\ from\ https://en.wikipedia.org/wiki/Decision\_tree\_pruning\#: -: text = Pruning\%20 is \%20 adata\%20 compression, and \%20 redundant\%20 to \%20 classify\%20 instances.$ 







#### **Pruning**

#### Weight pruning

- Removes individual weights within a layer.
- Results in a sparse weight matrix with many zero values.
- How to take advantage of sparsity?

1	0	0	3	0	4
0	3	0	0	0	0
0	0	5	0	7	0
0	0	6	0	0	0
0	0	6	0	0	8

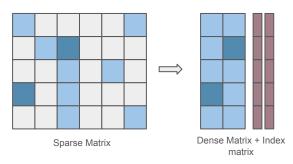
Sparse matrix



#### **Pruning**

#### Structural pruning

- Removes entire neurons, channels, or filters instead of individual weights.
- Produces a smaller, dense model, improving efficiency.
- Reduces model size.





## ML and model compression techniques for reconfigurable hardware accelerators

#### **Advantages**

- Model size reduction.
- Acceleration of the inference stage.
- Lower energy consumption.

#### **Pruning**

#### **Drawbacks**

- Loss of precision.
- Requires fine-tuning (to recover lost accuracy).
- Hardware compatibility.



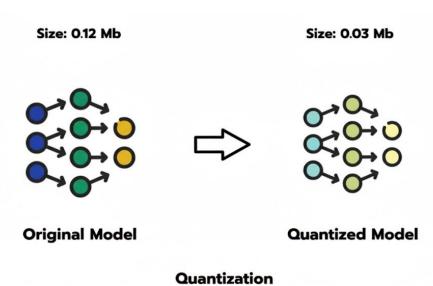
# Demo: Pruning for MLP and Fashion MNIST





- Quantization is a technique that reduces the numerical precision of a neural network's parameters by transforming floating-point values (e.g., 32-bit) into lower-precision representations, such as 16-bit or even 8-bit.
- The main goal is to reduce the model size and speed up inferences, especially on resource-constrained devices.







Quantization - PTQ -

- Post-Training Quantization (PTQ) is a quantization technique applied to the trained model.
- It converts the weights and activations from floating-point precision to a lower precision format.
- It is used when the goal is to reduce the model size without requiring retraining.



Quantization - QAT -

- Quantization-aware training (QAT) is a training technique in which the model learns to adapt to quantization before being deployed on hardware. Instead of training the model in full precision (32-bit floating point) and then quantizing it, quantization is introduced during training.
- It provides better results than PTQ but requires more computational effort.



Quantization - DQ -

- Dynamic Quantization (DQ) is a quantization method that focuses on quantizing only the model's weights, leaving the activations in floating-point format.
- This approach is particularly beneficial for models where fast inference takes priority over storage efficiency.



Quantization - FIQ -

- Full Integer Quantization (FIQ) is a quantization technique where both the weights and activations are fully converted to integer format (INT8).
- This method is commonly used in specialized hardware, such as TPUs, and NPUs, to maximize efficiency and performance.



Quantization-aware pruning - QAP -

 Quantization-aware pruning (QAP) combines pruning with quantization-aware training. The goal is to reduce the model size after quantization, resulting in a more efficient network without sacrificing accuracy.



#### **Advantages**

- Model size reduction.
- Acceleration of the inference stage.
- Lower energy consumption.

#### Quantization

#### **Drawbacks**

- Loss of precision.
- Requires fine-tuning (to recover lost accuracy).
- Hardware compatibility.



Quantization

#### QKeras

- Extension of Keras designed to quantize neural network models
- Useful when training models with lower precision.
- Auto QKeras.
- https://github.com/google/qkeras



Quantization

Automatic heterogeneous quantization of deep neural networks for low-latency inference on the edge for particle detectors

Claudionor N. Coelho Jr.
Palo Alto Networks (California, USA)

Aki Kuusela, Shan Li, and Hao Zhuang Google LLC (California, USA)

Thea Aarrestad,\* Vladimir Loncar,† Maurizio Pierini, Adrian Alan Pol, and Sioni Summers

European Organization for Nuclear Research (CERN) (Geneva, Switzerland)

Jennifer Ngadiuba California Institute of Technology (Caltech) (California, USA) (Dated: June 22, 2021)

Although the quest for more accurate solutions is pushing deep learning research towards larger and more complex algorithms, edge devices demand efficient inference and therefore reduction in model size, latency and energy consumption. One technique to limit model size is quantization, which implies using fewer bits to represent weights and biases. Such an approach usually results in a decline in performance. Here, we introduce a method for designing optimally heterogeneously quantized versions of deep neural network models for minimum-energy, high-accuracy, nanosecond inference and fully automated deployment on chip. With a per-layer, per-parameter type automatic quantization procedure, sampling from a wide range of quantizers, model energy consumption and size are minimized while high accuracy is maintained. This is crucial for the event selection procedure in proton–proton collisions at the CERN Large Hadron Collider, where resources are strictly limited and a latency of  $\mathcal{O}(1)$   $\mu s$  is required. Nanosecond inference and a resource consumption reduced by a factor of 50 when implemented on field-programmable gate array hardware are achieved.

```
# MLP architecture
# Create the student OKERAS
studentQ MLP = keras.Sequential(
        Input(shape=(30,)),
        QDense(20, name='fc1',
                 kernel quantizer=quantized bits(9,1,alpha=1), bias quantizer=quantized bits(23,15,alpha=1)),
        QActivation(activation=quantized relu(16,15), name='relu1'),
        ODense(10, name='fc2',
                 kernel quantizer=quantized bits(9,1,alpha=1), bias quantizer=quantized bits(23,15,alpha=1)),
        OActivation(activation=quantized relu(16.15), name='relu2'),
        QDense(10, name='fc6',
                 kernel quantizer=quantized bits(9,1,alpha=1), bias quantizer=quantized bits(23,15,alpha=1)),
        OActivation(activation=quantized relu(16,15), name='relu3'),
        QDense(4, name='output',
                 kernel quantizer=quantized bits(32,15,alpha=1), bias quantizer=quantized bits(32,15,alpha=1)),
        Activation(activation='softmax', name='softmax')
    name="student".
print qstats(studentQ MLP)
```



# Demo: Quantization for MLP and Fashion MNIST



5

9 Mar 201

2531v1

**Knowledge** distillation

#### Distilling the Knowledge in a Neural Network

Geoffrey Hinton\*†
Google Inc.
Mountain View
qeoffhinton@qoogle.com

Oriol Vinyals†
Google Inc.
Mountain View
vinyals@qoogle.com

Jeff Dean Google Inc. Mountain View jeff@google.com

#### Abstract

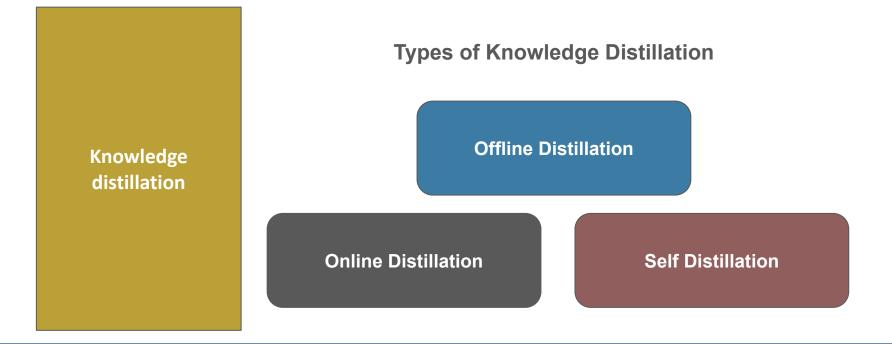
A very simple way to improve the performance of almost any machine learning algorithm is to train many different models on the same data and then to average their predictions [3]. Unfortunately, making predictions using a whole ensemble of models is cumbersome and may be too computationally expensive to allow deployment to a large number of users, especially if the individual models are large neural nets. Caruana and his collaborators [1] have shown that it is possible to compress the knowledge in an ensemble into a single model which is much easier to deploy and we develop this approach further using a different compression technique. We achieve some surprising results on MNIST and we show that we can significantly improve the acoustic model of a heavily used commercial system by distilling the knowledge in an ensemble of models into a single model. We also introduce a new type of ensemble composed of one or more full models and many specialist models which learn to distinguish fine-grained classes that the full models confuse. Unlike a mixture of experts, these specialist models can be trained rapidly and in parallel.



**Knowledge** distillation

 Knowledge Distillation (KD) is devoted to transferring knowledge from a teacher network to a smaller and faster target network (named a distilled or student network) that can reproduce the teacher's behavior while being computationally less expensive.







## **Knowledge** distillation

#### Offline Distillation

- Use of a pre-trained teacher model to guide the student model
- Knowledge
  - The output probabilities not only indicate the correct class but also contain the relative probabilities of incorrect classes.



Image from https://medium.com/@aryamaanthakur/knowledge-distillation-make-your-neural-networks-smaller-398485f811c6



## **Knowledge** distillation

#### Offline Distillation

 Instead of learning directly from the books (the original data), the child watches how the teacher answers the questions.

Question	Teacher says	Student learns
Is this a cat?	90% yes, 10% no	Learns that it really looks like a cat
Is this a dog?	5% yes, 95% no	Learns that it's not a dog
Is this a rabbit?	50% yes, 50% no	Learns that it's hard to tell the difference



**Knowledge** distillation

#### Offline Distillation

• In this way, the child learns not only the answer, but also how confident the teacher was. That helps the child understand the boundaries better and learn faster with less effort.



## **Knowledge** distillation

#### Offline Distillation

#### Temperature (T)

If **T is low** (for example, 1): the teacher gives very confident and sharp answers:

"This is 100% a cat! Nothing else". The child only sees that it's a cat and doesn't learn any nuance.

If **T is high** (for example, 5): the teacher gives **softer answers**, showing some uncertainty:

"It looks 70% like a cat, 20% like a rabbit, 10% like a dog". The child learns that some animals look similar and understands subtle differences better.

This helps the smaller model (the student) **learn the structure of the knowledge**, not just simple yes/no answers.



## **Knowledge** distillation

#### Offline Distillation

#### Alpha (α)

Now the child has two teachers:

The big teacher (the teacher model) and the book (the real data with true labels).

The  $\alpha$  (alpha) value tells how much attention the child pays to each one:

Total loss= $\alpha \times Teacher$  loss+ $(1-\alpha) \times Real$  data loss

If  $\alpha$  = 0.9, the child listens more to the teacher (learning from the teacher's soft answers). If  $\alpha$  = 0.1, the child relies more on the real data (the hard labels, 0 or 1).



## **Knowledge** distillation

#### **General Steps**

- Train the teacher model
  - soft labels (probability distributions) for each sample.



## **Knowledge** distillation

#### **General Steps**

- Train the teacher model
  - o soft labels (probability distributions) for each sample.
- Generate Soft Labels (Logits) from the Teacher.
  - $\circ$  softmax with temperature T (a higher T makes the probabilities smoother)

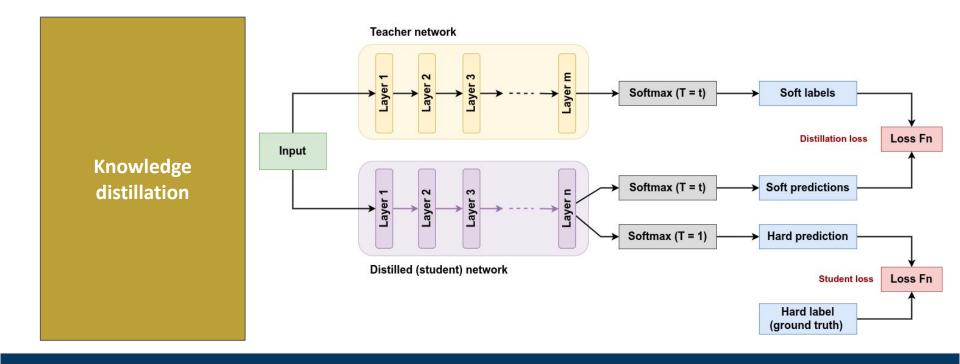


## **Knowledge** distillation

#### **General Steps**

- Train the teacher model
  - o soft labels (probability distributions) for each sample.
- Generate Soft Labels (Logits) from the Teacher.
  - $\circ$  softmax with temperature T (a higher T makes the probabilities smoother)
- Train the student with two losses:
  - Distillation loss: Matches Student's soft labels to the Teacher's soft labels.
  - Classification loss: Matches Student's **predictions** to ground **truth labels**.







**Knowledge** distillation

#### **Online Distillation**

- The teacher and student models are updated simultaneously in a single training process.
- When the teacher is not available.



# **Knowledge** distillation

#### **Online Distillation**

- The teacher and student models are updated simultaneously in a single training process.
- When the teacher is not available.

#### **Self Distillation**

- Similar to Online Distillation.
- The technique uses the same model as the teacher as well as the student.



# **Knowledge** distillation

#### **Advantages**

- Reduction of model size.
- Acceleration of the inference stage.
- Improves computational efficiency.
- Leverages the knowledge of the master model.

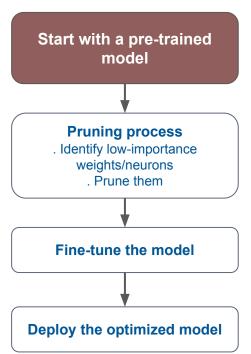
#### **Drawbacks**

- Complexity of proper implementation.
- Not always effective; it depends on the generalization ability of the master model.
- Expensive training.

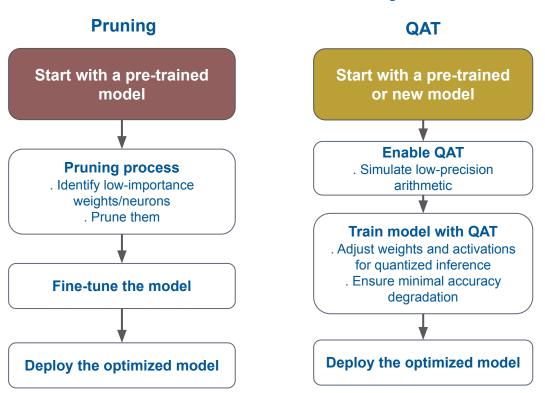




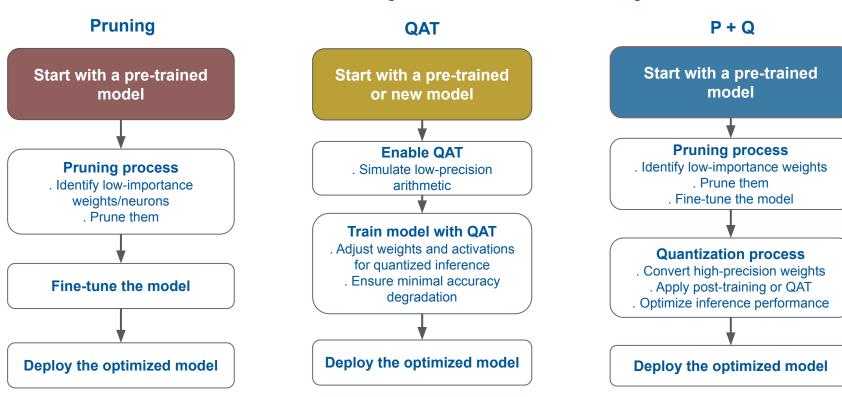
#### **Pruning**











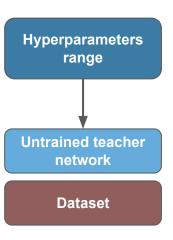


## Hyperparameters tuning and compression



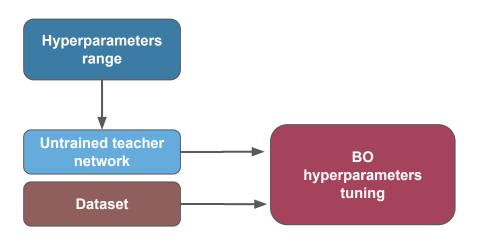
## **DNN** training and compression

#### **Stage 1 - Teacher training**





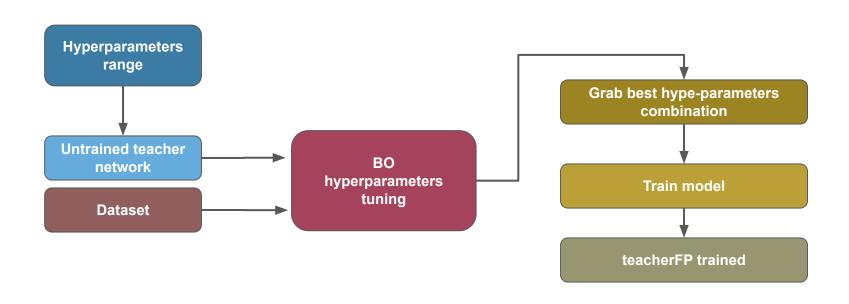
# **DNN training and compression**Stage 1 - Teacher training





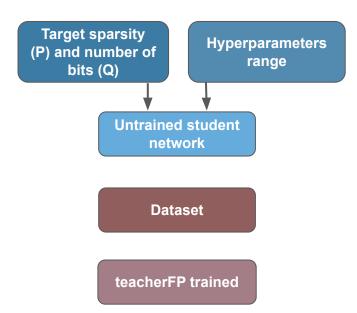
## **DNN** training and compression

#### **Stage 1 - Teacher training**



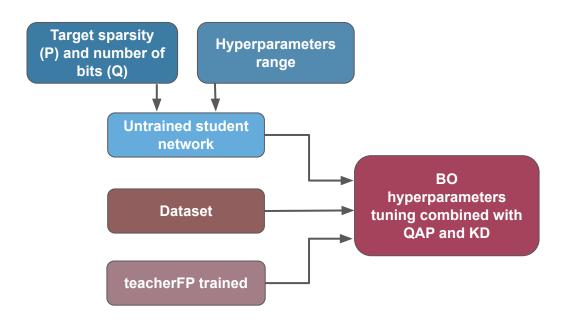


# **DNN training and compression Stage 2 - Student training**



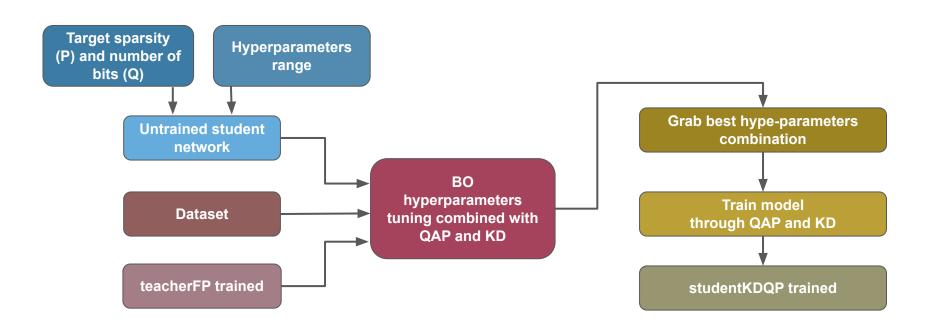


# **DNN training and compression Stage 2 - Student training**





# **DNN training and compression**Stage 2 - Student training





# Demo: MNIST-based binary classification -QAP -

Joint ICTP-IAEA School on Detector Signal Processing and Machine Learning for Scientific Instrumentation and Reconfigurable Computing

Model Optimization and Compression Techniques

smr4110 | Trieste, Italy - 2025

Romina Soledad Molina, Ph.D.



