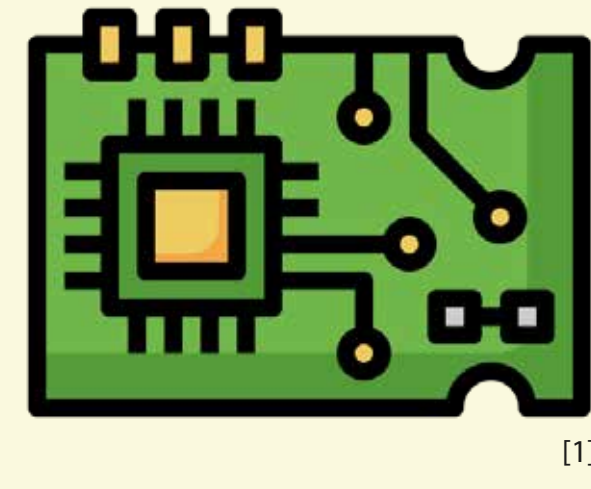


Inference of an average DNN releases 2.840 t of CO₂. [5]



How to optimize DNNs to be more efficient?



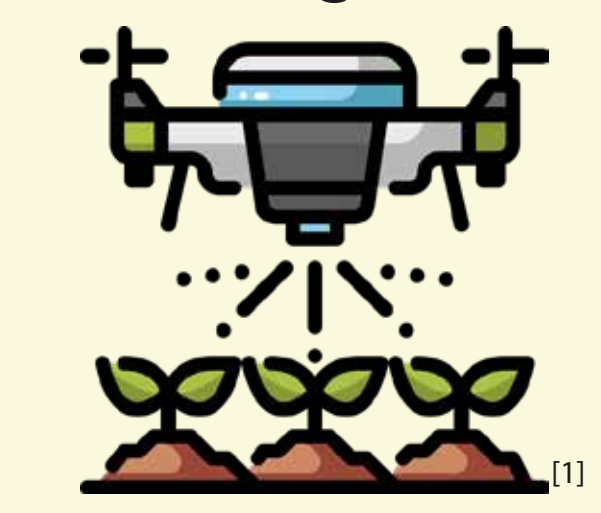
Scenarios

High speed cloud computing

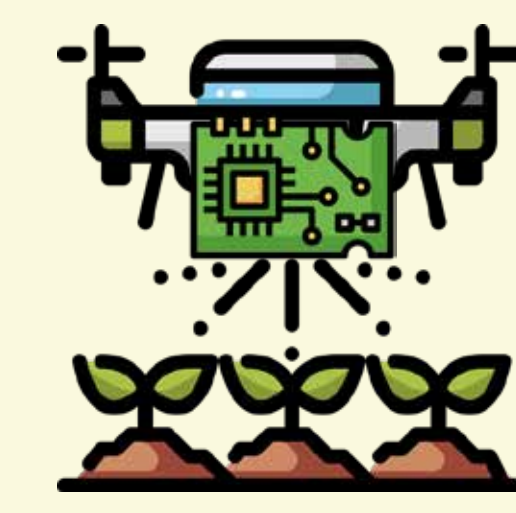


Sending images

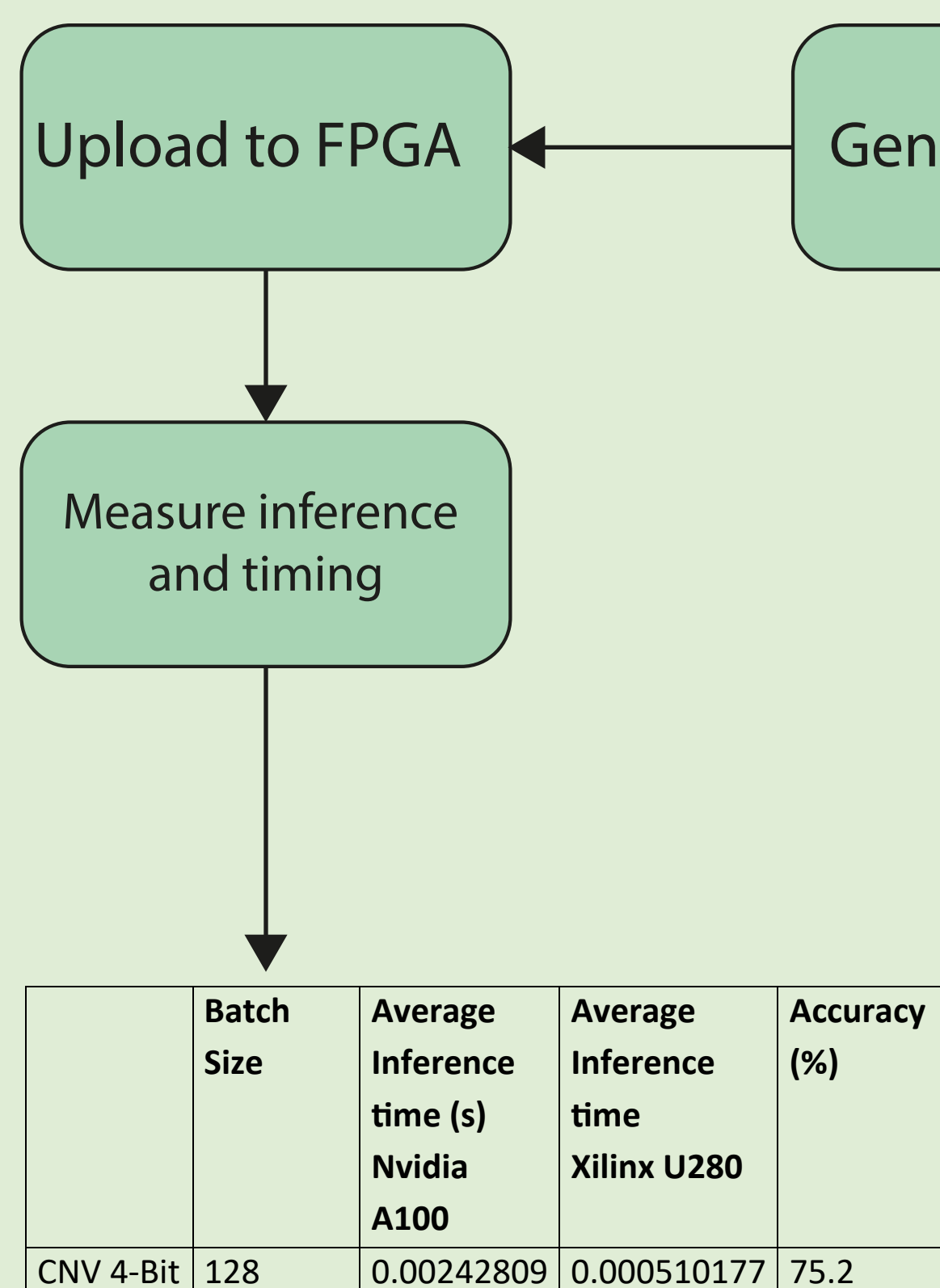
Receiving results



Efficient edge computing with FPGA as payload



Four images showing sugar beet plants from a bird's eye view.[2]



Model	Throughput (images/s)	Clock (MHz)
CNV 4-Bit	1740.74	100

Measure inference and timing

Move to GPU

GPU

FPGA

Hardware build

Network preparation

Generate bitfile

Upload to FPGA



Data preprocessing

Test set

Training set

Pretrained model

(QAT)

Model training

Training parameter

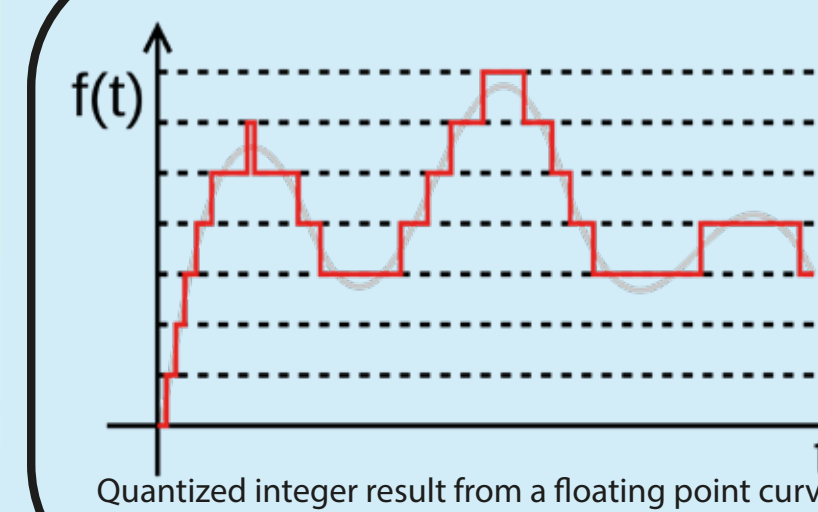
Model validation

Model optimizing (PTQ/ Pruning)

ONNX file

Optimizing

Results



Post-training quantization (PTQ) is a technique applied after a machine learning model has been trained, by converting the model's weights and activations from floating-point to integer numbers.

Quantization aware training (QAT) incorporates quantization directly into the training process of a model by simulating the effects of quantization while the model is still being trained.

Network	Accuracy	Inference (ms)	Size on disk (KB)
ResNet18 base	50.1	15,758	46835
ResNet18 static quantization	88.7	9,725	11832
ResNet50 base	94.5	53,329	102542
ResNet50 static quantization	94.3	18,8409	26576

References:

- [1] <https://www.flaticon.com>
 [2] R. Sharma L. Fink. V2 plant seedling dataset.[online]. Available: <https://www.kaggle.com/datasets/vbookshelf/v2-plant-seedlings-dataset>, 26.09.2023.
 [3] C. Rice. Gain, range, and quantization. [online]. Available: <https://community.sw.siemens.com/s/article/gain-range-and-quantization>, 26.09.2023.
 [4] <https://www.eki-project.tech/project>

Supported by:



Federal Ministry
for the Environment, Nature Conservation,
Nuclear Safety and Consumer Protection

based on a decision of
the German Bundestag