# DNN Partitioning

## APROPOS Summer School

### Delft

Axel Jantsch

June 28, 2023

## Resource limitations

|  | Embedded | Data center |
| --- | --- | --- |
| Computation [flop] | $30 - 1800 \cdot 10^{12}$ | $86 \cdot 10^{18}$ |
| Memory [bit] | $10^{10}$ | $10^{15}$ |
| Power [W] | 5-100 | $10^3 - 10^6$ |
| Energy [Wh] | 48-1000 | $200 \cdot 10^6$ |

**Computation Embedded** refers to an Nvidia Jetson Nano running 1 min and 1 hour, respectively.
**Computation server** refers to the computation needed for the 40 day experiment with AlphaGo Zero
**Energy embedded** refers to a mobile phone and to a car battery, respectively.
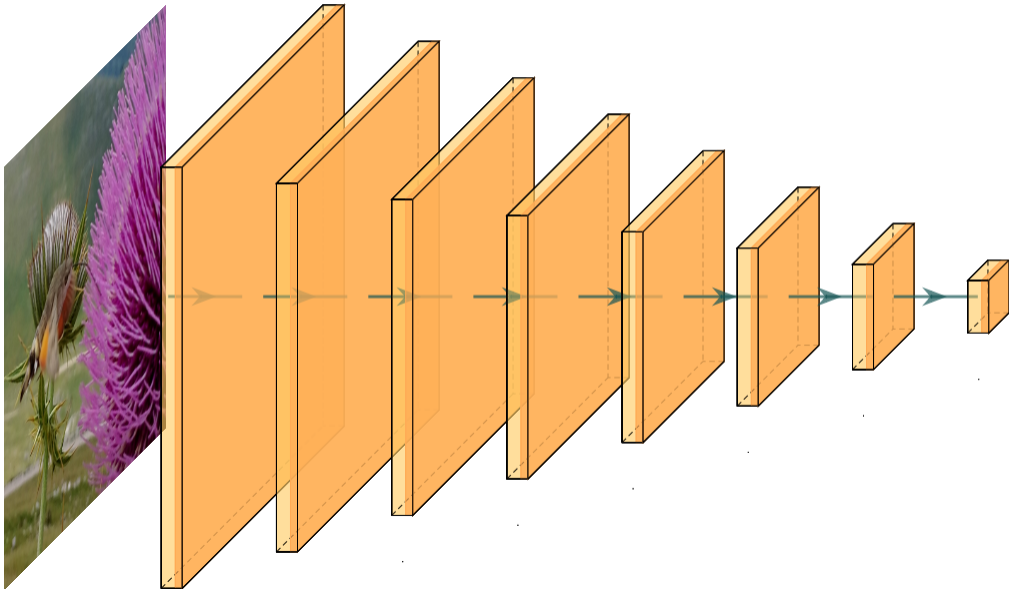**Energy server** refers to the 40 day experiment for AlphaGo Zero.

# Outline

1. Partitioning for Inference

2. Impact of image size and content

3. Waist Tightening

# Outline

# PARTITIONING FOR INFERENCE

- Energy depends on
    - computation platform
    - amount of computation done
    - communication protocol
    - amount of communication done
- Communication energy cost is very different for different protocols
- Communication energy (and latency tends) to dominate total energy and latency

Irida Shallari, Isaac Sánchez Leal, Silvia Krug, Axel Jantsch, and Mattias O'Nils. "Design space exploration on IoT node: Trade-offs in processing and communication". In: *IEEE Access* (2021)

# Energy and Latency Model

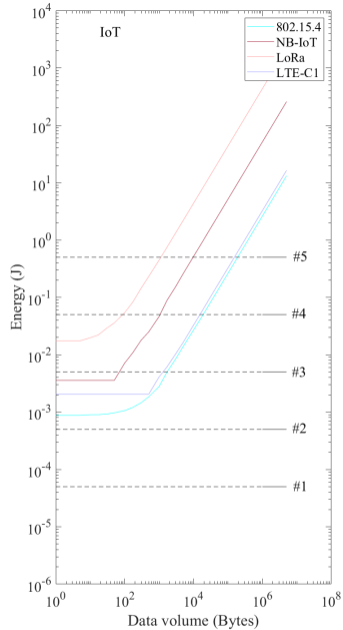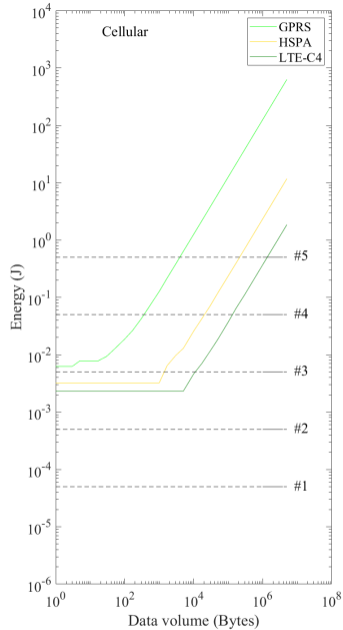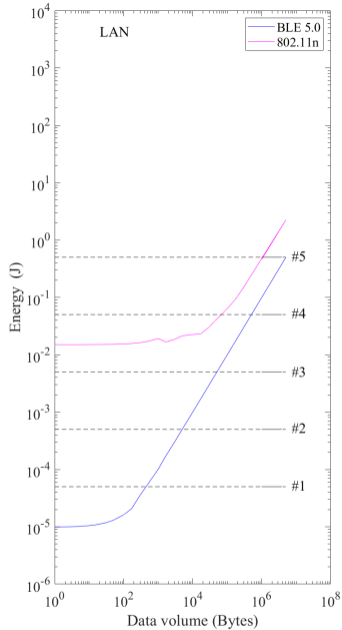$$E_{\text{Node}}(x) = E_S + E_P(T_{(0,x)}, P) + E_c(V_x, C)$$
$$L_{\text{Node}}(x) = L_s + L_P(T_{(0,x)}, P) + L_c(V_x, C)$$

| | | | |
|---|---|---|---|
| $L_{\text{Node}}$ | Node latency per sample | $E_{\text{Node}}$ | Node energy per sample |
| $L_S$ | Sensing latency | $E_S$ | Sensing energy |
| $L_P$ | Processing latency | $E_P$ | Processing energy |
| $L_C$ | Communication latency | $E_C$ | Communication energy |

| | |
|---|---|
| $x$ | Partitioning point in $[0, \ldots, N]$ |
| $T_{(0,x)}$ | Computation tasks of stages $0 \ldots x$ |
| $P$ | Hardware platform |
| $V_x$ | Data volume at output of stage $x$ |
| $C$ | Communication protocol |

| Communication groups | | |
| --- | --- | --- |
| LAN | Cellular | IoT |
| BLE 5.0 | GPRS | 802.15.4 g |
| 802.11 | HSPA | NB-IoT |
| | LTE C. 4 | LoRa |
| | | LTE C. 1 |

Communication energy for different protocols.

| Tasks | Traditional systems | | CNN systems | |
| --- | --- | --- | --- | --- |
| | People Counting | Particle Detection | AlexNet | VGG16 |
| 0 | 307 200 | 307 200 | 307 200 | 307 200 |
| 1 | 8940 | 256 000 | 154 587 | 150 528 |
| 2 | 91 | 680 | 69 984 | 3 211 264 |
| 3 | 75 | 500 | 43 264 | 1 605 632 |
| 4 | 4 | 259 | 64 896 | 802 816 |
| 5 | | | 9216 | 401 408 |
| 6 | | | 4096 | 100 352 |
| 7 | | | 1000 | 25 088 |
| 8 | | | | 4096 |
| 9 | | | | 1000 |

(data volume after each processing stage in bytes)

People counting application

Particle detection application

AlexNet

Legend:
- ▼ Processing energy for each partition point.
- ● Processing and communication energy consumption for each partition point.
- – – Communication energy consumption.
- # 1-5 Energy constraints.

Series: BLE 5, 802.11n, GPRS, HSPA, LTE C. 4, 802.15.4, NB-IoT, LoRa, LTE C. 1

X-axis: Data points (Bytes)
Y-axis: Energy (Joules)

Constraint lines: #5, #4, #3, #2, #1

VGGNet 16

19

# Conclusions

- For high energy communication protocols the it is optimal to minimize transmitted data.

- For low energy communication protocols the sweet spot is not at the extremes.

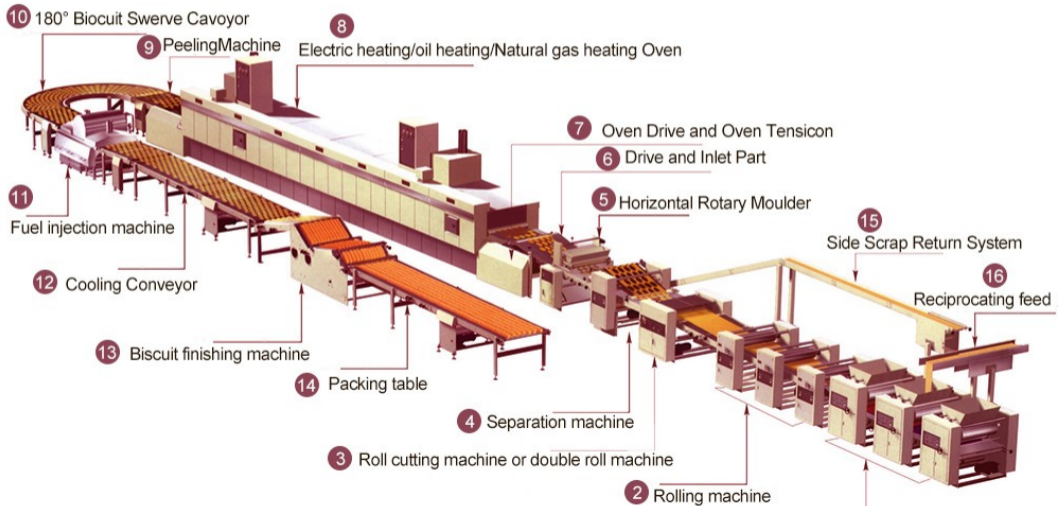- The optima depend on the application, the IoT platform and the communication protocol used.

# Outline

# Impact of image size and content

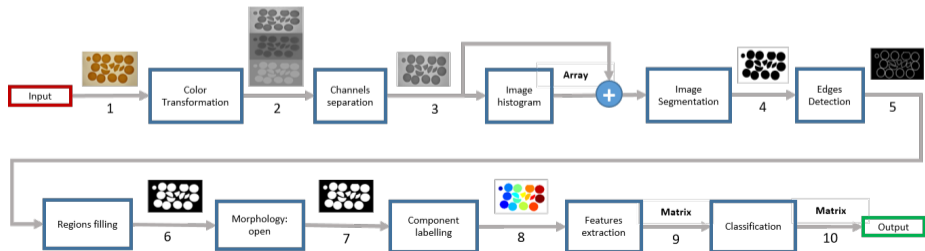Case study: Conventional image processing pipeline - Biscuit inspection system

Isaac Sánchez Leal, Irida Shallari, Silvia Krug, Axel Jantsch, and Mattias O'Nils. "Impact of Input Data on Intelligence Partitioning Decisions for IoT Smart Camera Nodes". In: *Electronics* 10.16 (2021)
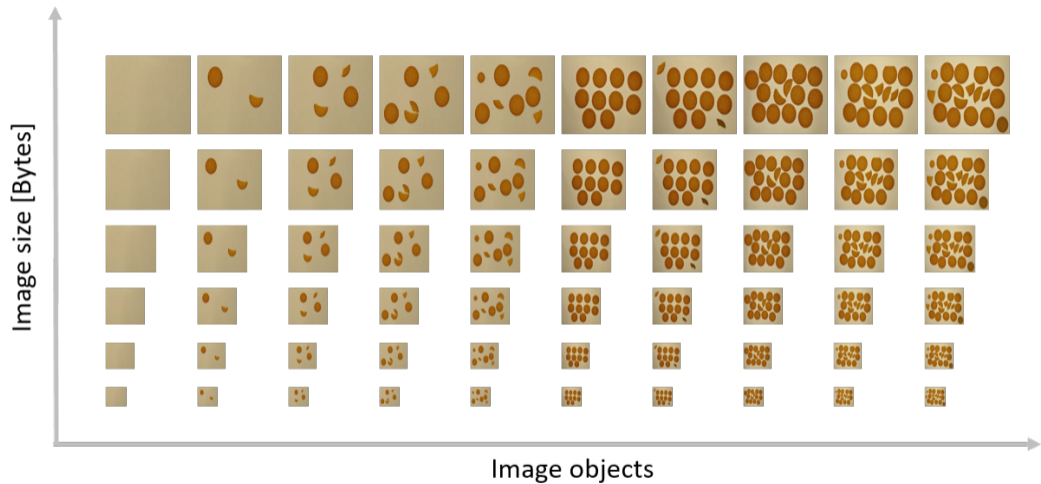
# Biscuit Production



10 180° Biocuit Swerve Cavoyor

9 PeelingMachine

8 Electric heating/oil heating/Natural gas heating Oven

7 Oven Drive and Oven Tensicon

6 Drive and Inlet Part

5 Horizontal Rotary Moulder

11 Fuel injection machine

12 Cooling Conveyor

13 Biscuit finishing machine

14 Packing table

15 Side Scrap Return System

16 Reciprocating feed

4 Separation machine

3 Roll cutting machine or double roll machine

2 Rolling machine

| # | Processing Task | Output type | Changes in data out due to: $\Delta Img.Size$ | $\Delta Img.Objects$ | Processing time behavior |
|----|----|----|----|----|----|
| t1 | Color transformation. | Color space YCbCr. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t2 | Channels separation. | Y channel. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t3 | Image Histogram. | Array with 256 elements. | Constant | Constant | $\Delta Linear_{Size}$ |
| t4 | Segmentation. | BW image without background. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t5 | Edges detection. | BW image with detected regions. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t6 | Regions filling. | BW image with filled regions. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t7 | Morphology: open. | BW image without particles. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t8 | Component Labelling. | Non-binary labels image. | $\Delta Linear$ | Constant | $\Delta Linear_{Size}$ |
| t9 | Features extraction. | 2D features matrix. | Constant | $\Delta Linear$ | $\Delta Linear_{Size, Objects}$ |
| t10 | Classification. | 2D coordinates matrix. | Constant | $\Delta Linear$ | $\Delta Linear_{Objects}$ |

Image size [Bytes]

Image objects
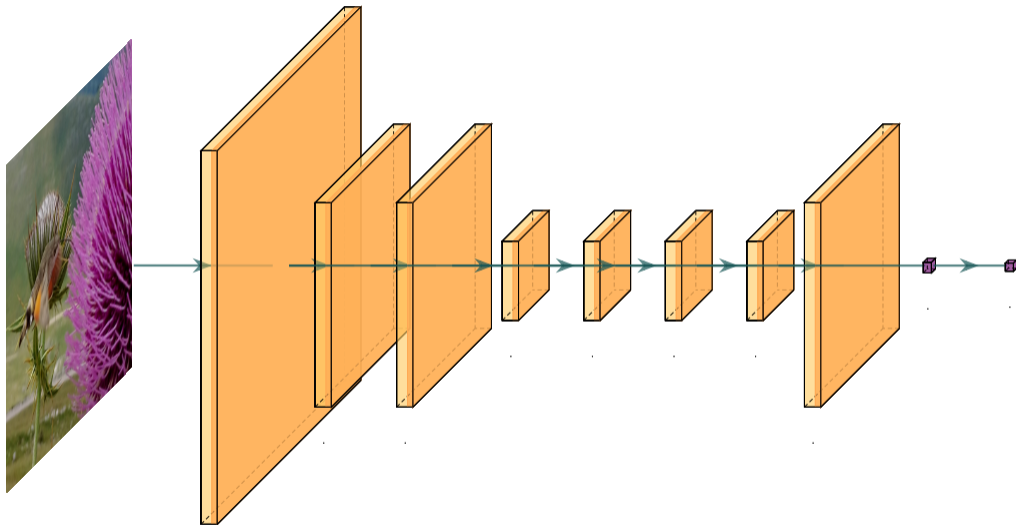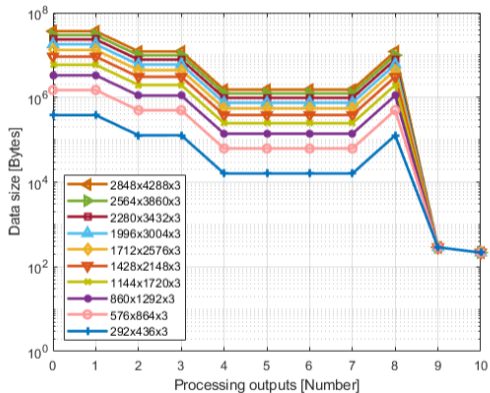
- Most tasks' processing time depends on the size of image: t1,t2,t4-t8
- Some tasks depend on number of objects: t9, t10, which are later in the pipeline
- Data Volume depends heavily on image size, only late in the pipeline (t9,t10) also on number of objects

18 objects

Image size: $292 \times 436 \times 3$.

18 objects
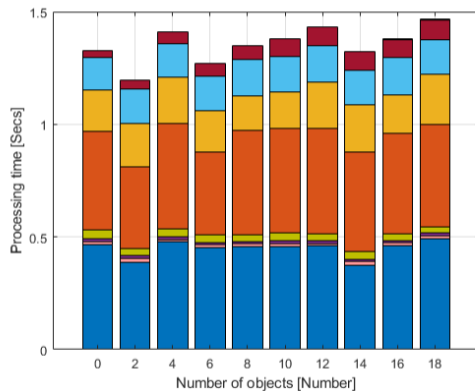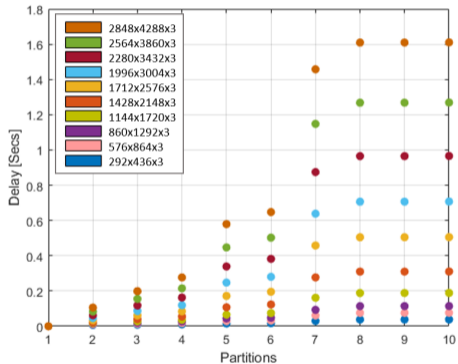
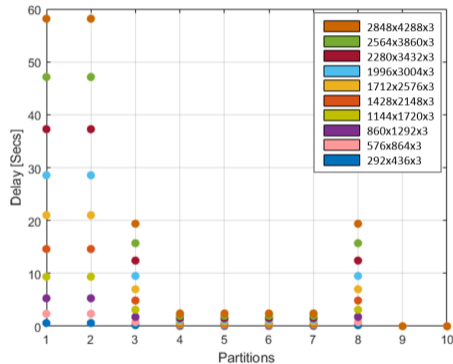Image size: $2848 \times 4288 \times 3$

Proccesing latency with RasberryPi.



Communication latency with LTE Cat.1.



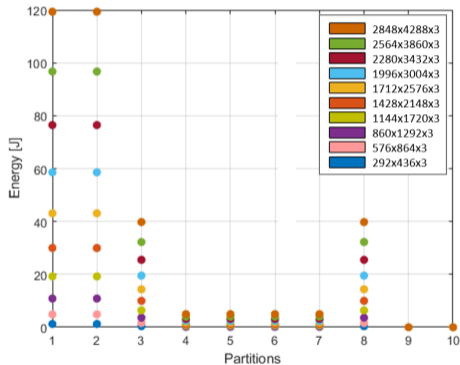18 objects

Proccesing energy on RasberryPi.

Communication energy with LTE Cat.1.



18 objects

# Optimization Problems

1. Miminizing latency:
$$L_{\text{Node}} \rightarrow Min$$

2. Minimizing energy:
$$E_{\text{Node}} \rightarrow Min$$

3. Minimizing energy under a latency constraint:
$$L_{\text{Node}} \leq MaxDelayConstraint$$
$$E_{\text{Node}} \rightarrow Min$$

# Outline

# Waist Tightening

## Assumptions

- Energy and delay model
- Optimization goal
- Constraints on energy, latency and accuracy

## Selecting the partitioning point

- Candidate points have low data volume
- For each candidate point
  - Quantization,Pruning and compression
  - Retraining
- Best point according to constraints and optimization criteria is selected.

# Energy and Latency Model

$$E_{\text{Node}}(x) = E_S + E_P(T_{(0,x)}, P) + E_c(V_x, C)$$
$$L_{\text{Node}}(x) = L_s + L_P(T_{(0,x)}, P) + L_c(V_x, C)$$

| | | | |
|---|---|---|---|
| $L_{\text{Node}}$ | Node latency per sample | $E_{\text{Node}}$ | Node energy per sample |
| $L_S$ | Sensing latency | $E_S$ | Sensing energy |
| $L_P$ | Processing latency | $E_P$ | Processing energy |
| $L_C$ | Communication latency | $E_C$ | Communication energy |

| | |
|---|---|
| $x$ | Partitioning point in $[0, \ldots, N]$ |
| $T_{(0,x)}$ | Computation tasks of stages $0 \ldots x$ |
| $P$ | Hardware platform |
| $V_x$ | Data volume at output of stage $x$ |
| $C$ | Communication protocol |

*Input image* → Layers 1 to *j*-1 → Layer *j* → $f_j$ → Quantization → $d_j^b$ → Compression → $c_j$ → De-compression → $d_j^b$ → Conversion bin2float → $f_j'$ → Layers *j*+1 to *N* → *System output*

(Node partition: Layers 1 to *j*-1, Layer *j*, Quantization, Compression; Edge partition: De-compression, Conversion bin2float, Layers *j*+1 to *N*)

DNN Partitioning:

$$M_p = M_{1..j} \cup I \cup M_{j+1..N}$$

Optimization problem
Minimizing energy under accuracy constraint:

$$
\begin{aligned}
E_{\text{Node}} &\rightarrow Min \\
mAP(M_p) &\geq mAP(M) - L_{\text{th}}.
\end{aligned}
$$

| | |
|---|---|
| $M$ | Original DNN |
| $M_p$ | DNN partitioned at point *j* |
| $I$ | Interface |
| $mAP$ | mean average precision |
| $L_{\text{th}}$ | Threshold for acceptable loss of precision |

Eiraj Saqib, Isaac Sánchez Leal, Irida Shallari, Axel Jantsch, Silvia Krug, and Mattias O'Nils. "Optimizing the IoT Performance: A Case Study on Pruning a Distributed CNN". In: *Proceedings of the IEEE Sensors Applications Symposium (SAS)*. 2023

# Quantization

Two quantization methods:

$$Q_1 : \quad \begin{aligned} d_j^b \quad &= \left\lfloor \frac{f_j - S}{M - S} \times 2^{b-1} \right\rfloor, \quad \text{with } b \in [1 - 8]. \\ f_j' \quad &= d_j^b \times \left( \frac{M - S}{2^{b-1}} \right) + S, \quad \text{where } b \in [1 - 8]. \end{aligned}$$

$$Q_2 : \quad \begin{aligned} d_j^b \quad &= \lfloor (f_j - \mu) / M \rfloor \quad \text{with } b = 1. \\ f_j' \quad &= \left( d_j^b \times M \right) + \mu \quad \text{with } b = 1. \end{aligned}$$

| | | | |
|---|---|---|---|
| $f_j$ | Floating point number in feature map at output of stage $j$ | $M$ | Maximum value of the dynamic range |
| | | $m$ | Minimum value of the dynamic range |
| $d_j^b$ | Quantized number to $b$ bit | $S$ | $(M - m)/2$ |
| $f_j'$ | De-quantized value | $\mu$ | Mean value |

# Pruning

Two step pruning:

1. Pruning of filters in Layer $j$, subject to

$$mAP(M_p') \geq mAP(M) - L_{\text{th}}$$

2. Pruning of layers $i = 1, \ldots, j - 1$, subject to

$$mAP(M_p'') \geq mAP(M) - L_{\text{th}}$$

$M'$     DNN with layer $j$ pruned
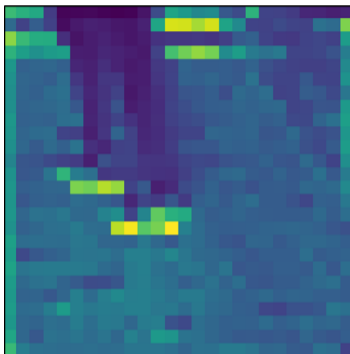$M''$    DNN with layers $i = 1, \ldots, j$ pruned

Compression to minimize the data transmitted over the channel
We use the zip and JPEG compression algorithms.

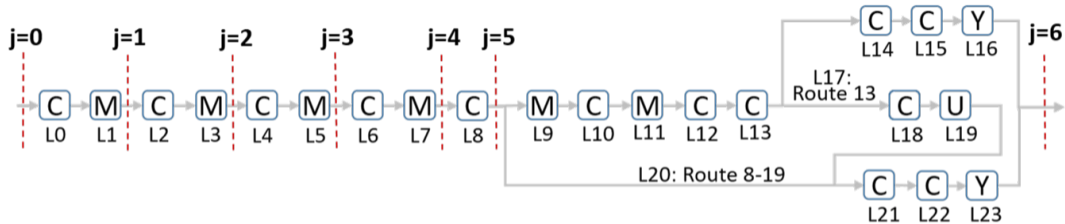Input image          32-bits          1-bit

Cristian Vilar Giménez, Silvia Krug, Faisal Z. Qureshi, and Mattias O'Nils. "Evaluation of 2D-/3D-Feet-Detection Methods for Semi-Autonomous Powered Wheelchair Navigation". In: *Journal of Imaging* 7.12 (2021)
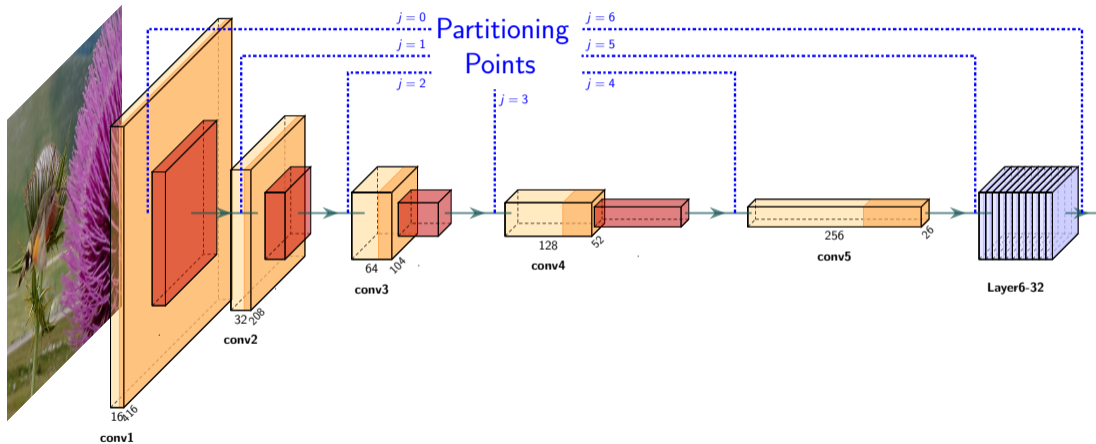
Isaac Sanchez Leal, Eiraj Saqib, Irida Shallari, Axel Jantsch, Silvia Krug, and Mattias O'Nils. "Waist Tightening of CNNs: A Case study on Tiny YOLOv3 for Distributed IoT Implementations". In: *Proceedings of the Real-time And intelliGent Edge computing workshop (RAGE)*. San Antonio, Texas, May 2023
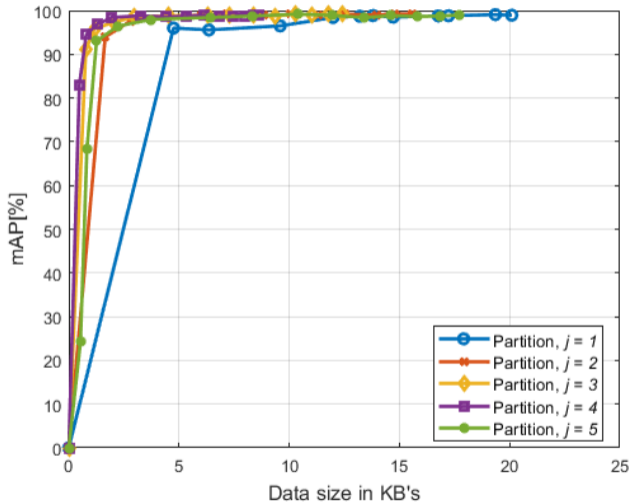
j=0   j=1   j=2   j=3   j=4  j=5

C   M   C   M   C   M   C   M   C      M   C   M   C   C
L0  L1  L2  L3  L4  L5  L6  L7  L8     L9  L10 L11 L12 L13

L17:
Route 13

C   C   Y
L14 L15 L16

j=6

C   U
L18 L19

L20: Route 8-19

C   C   Y
L21 L22 L23

Li: Layer "i"      C: Convolutional      U: Up-sampling
j: Partition       M: Maxpool            Y: YOLO region

47

Partitioning Points

$j = 0$
$j = 1$
$j = 2$
$j = 3$
$j = 4$
$j = 5$
$j = 6$

conv1
16×16

conv2
32 208

conv3
64 104

conv4
128 52

conv5
256 26

Layer6-32

48
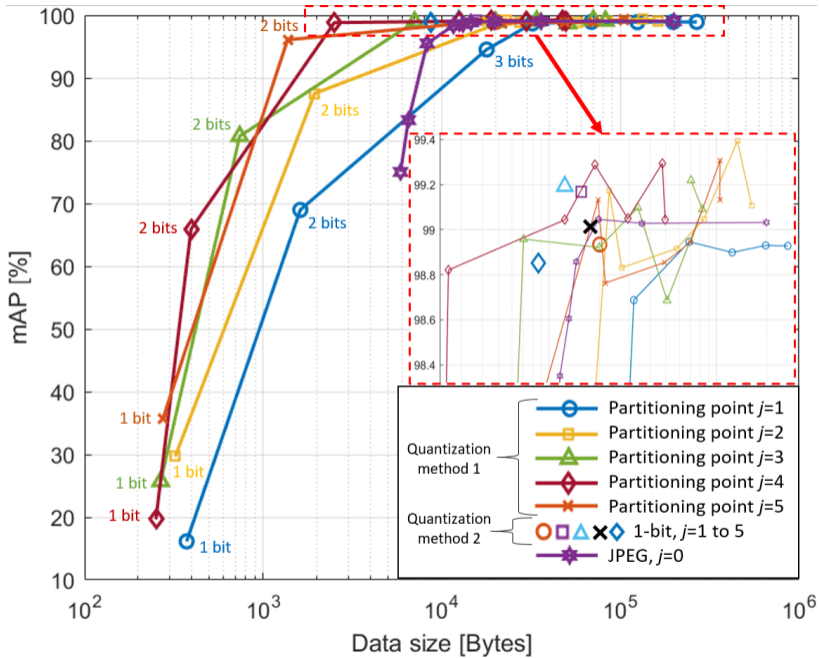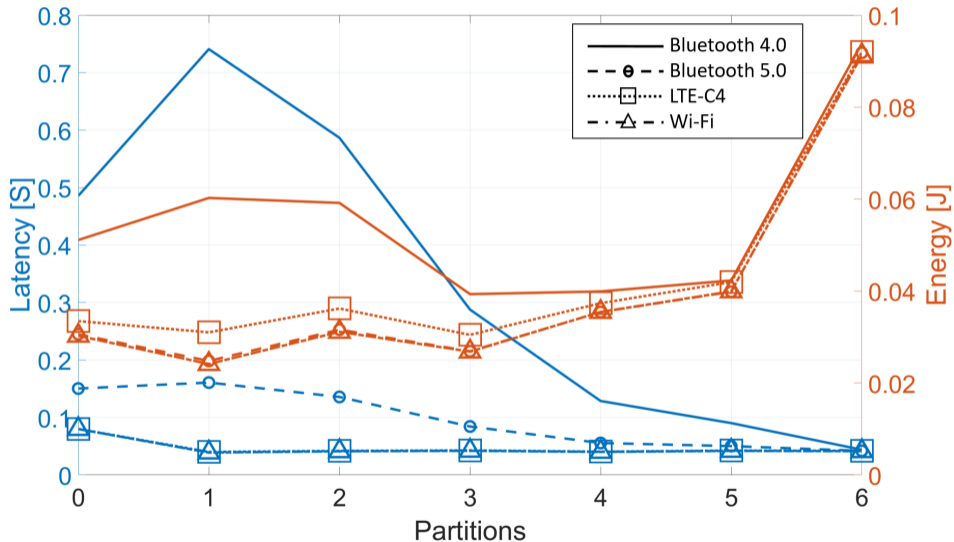
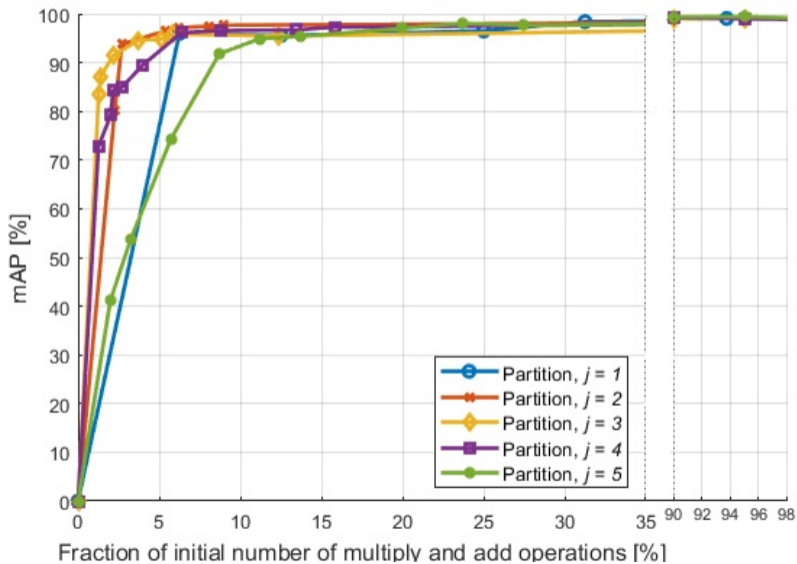Data volume at partition point after quantization, pruning and compression
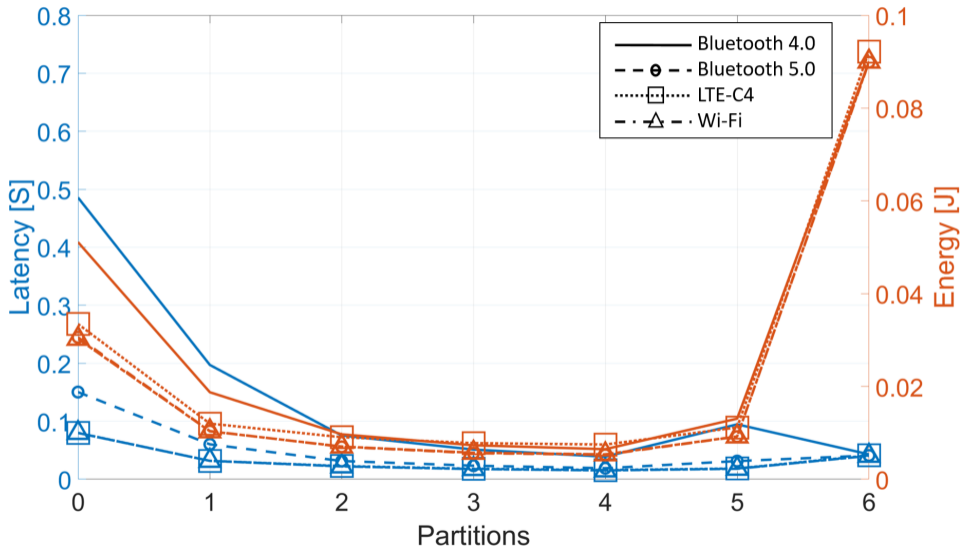
# Quantization $Q_2$ and Compression

# Quantization $Q_2$, Pruning and Compression

| | Quantization and Compression | | Quantization, Compression and Pruning | |
|---|---|---|---|---|
| | **All In-Edge** | **All In-Node** | **All In-Edge** | **All In-Node** |
| **Energy saving** | x1.26 | x3.8 | x5.74 | x17 |
| **System speed-up** | x2.05 | x1.05 | x5.24 | x2.65 |

Optimal partitioning versus all-in-edge and all-in-node reference solutions.

## Summary

- Effective DNN partitioning is feasible

## Summary

- Effective DNN partitioning is feasible

- DNN partitioning opens a considerable design space for DNN based IoT applications
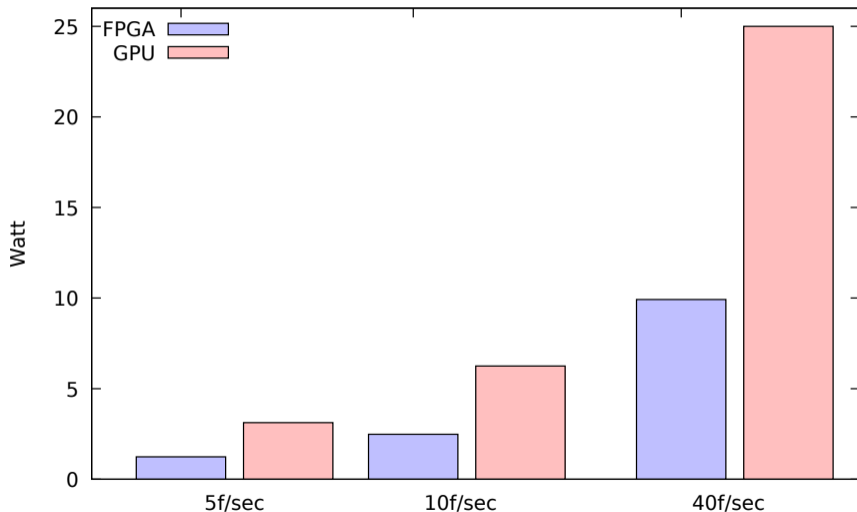
## Summary

- Effective DNN partitioning is feasible

- DNN partitioning opens a considerable design space for DNN based IoT applications

- Next steps is to explore more aggressive DNN adaptions for partitioning

# References I

[1] Irida Shallari, Isaac Sánchez Leal, Silvia Krug, Axel Jantsch, and Mattias O'Nils. "Design space exploration on IoT node: Trade-offs in processing and communication". In: *IEEE Access* (2021).

[2] Isaac Sánchez Leal, Irida Shallari, Silvia Krug, Axel Jantsch, and Mattias O'Nils. "Impact of Input Data on Intelligence Partitioning Decisions for IoT Smart Camera Nodes". In: *Electronics* 10.16 (2021).

[3] Isaac Sanchez Leal, Eiraj Saqib, Irida Shallari, Axel Jantsch, Silvia Krug, and Mattias O'Nils. "Waist Tightening of CNNs: A Case study on Tiny YOLOv3 for Distributed IoT Implementations". In: *Proceedings of the Real-time And intelliGent Edge computing workshop (RAGE)*. San Antonio, Texas, May 2023.

[4] Cristian Vilar Giménez, Silvia Krug, Faisal Z. Qureshi, and Mattias O'Nils. "Evaluation of 2D-/3D-Feet-Detection Methods for Semi-Autonomous Powered Wheelchair Navigation". In: *Journal of Imaging* 7.12 (2021).

[5] Eiraj Saqib, Isaac Sánchez Leal, Irida Shallari, Axel Jantsch, Silvia Krug, and Mattias O'Nils. "Optimizing the IoT Performance: A Case Study on Pruning a Distributed CNN". In: *Proceedings of the IEEE Sensors Applications Symposium (SAS)*. 2023.
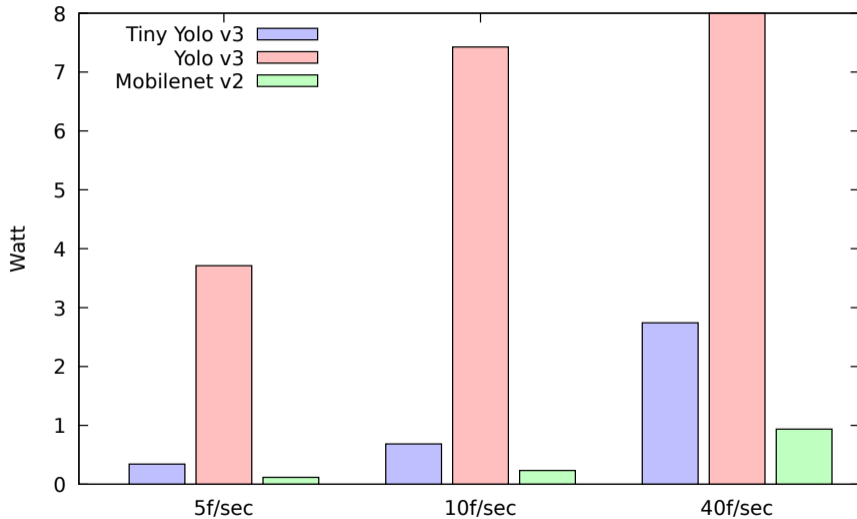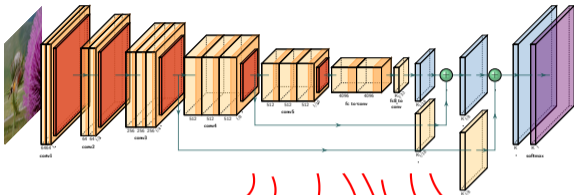
VGG16 applied to the ImageNet data set based on published papers.

A1

# Power Consumption in Inference



Object detection on the NCS2 platform; own measurements.

# Design Space

**DNN Choices**

Convolutional layers
Filter kernels
Number of filters
Pooling layers
Filter shape
Stride
Fully connected layer
Number of layers
Regularization
etc.

**Mapping Choices**

Neuron pruning
Data type selection
Approximation
Retraining
Connection pruning
Weight sparsifying
Regularization
etc.

**Platform Choices**

Platform Selection
Reconfiguration
Batch processing
Deep pipelining
Resource reuse
Hierarchical control
Processing unit selection
Memory allocation
Memory reuse
etc.

Intel® Vision Products with
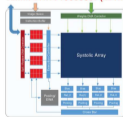Intel® Arria® 10 FPGA

Xilinx DNN Processor (xDNN)

Nvidia Turing

ARM NN