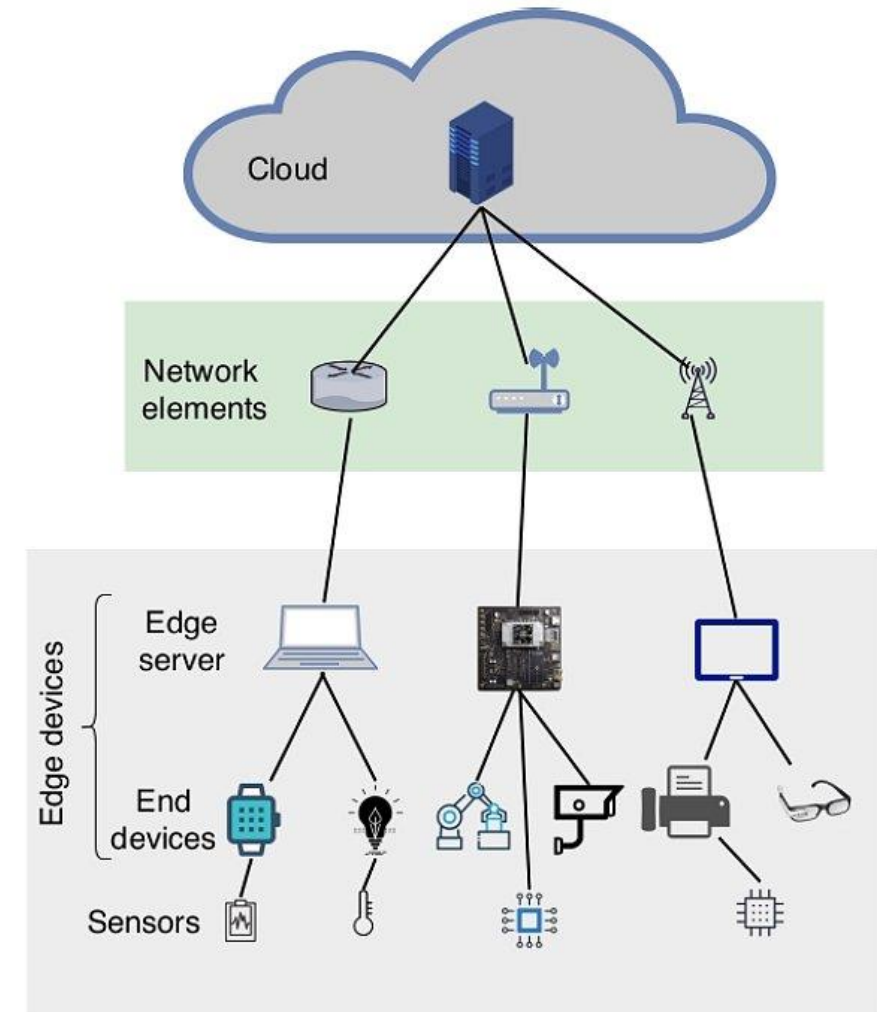


AI-on-the-edge-device

Sebastian Veigl

Was ist Edge AI?

- „Intelligenz“ wird in Geräte am Rand des Netzwerks verlagert -> nah an Datengenerierung
- Sensoren etc. führen intelligente Datenanalyse selbstständig durch
- Ergebnisse werden per Kommunikationsschnittstelle zur Verfügung gestellt
- Training wird zentral durchgeführt



<https://viso.ai/edge-ai/edge-ai-applications-and-trends/>

Vorteile / Nachteile von Edge AI

- **Vorteile:**

- Latenzzeiten
- Weniger Daten-Traffic
- Sicherheit
- Offline-fähig
- Kosten
- Energieverbrauch

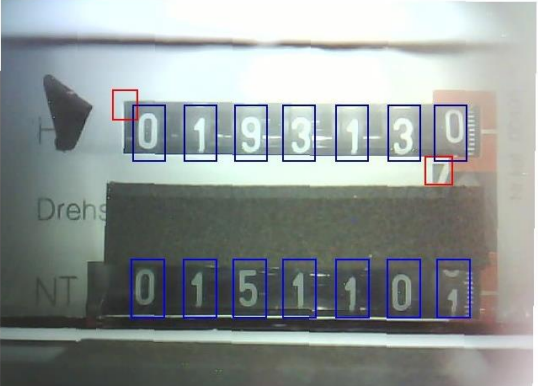
- **Nachteile:**

- Nur begrenzte Rechenleistung
- Irreversibler Datenverlust durch Vorverarbeitung

AI-on-the-edge-device

- Open-source Repository zur Anwendung von Edge AI
- Ursprünglich: Analyse von Zählerständen auf Strom-/Wasser-/Gaszählern
- Auswertung von digitalen und analogen Anzeigeelementen

<https://github.com/jomjol/AI-on-the-edge-device>

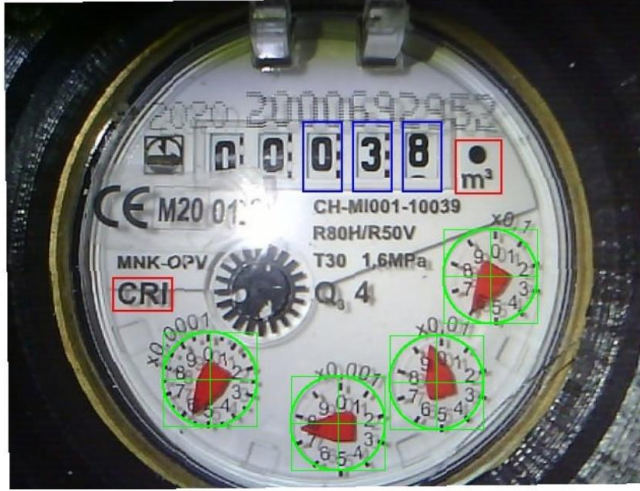


Value:	
NT	19313.0
HT	19013.0

Previous Value:	
NT	15110.1
HT	19013.0

Raw Value:	
NT	015110.N
HT	019013.0

Error:	
Neg. Rate - Read.	15110.0 - Raw.



Raw Value:	
038.5975	

Corrected Value:	
38.5975	

Checked Value:	
38.5975	

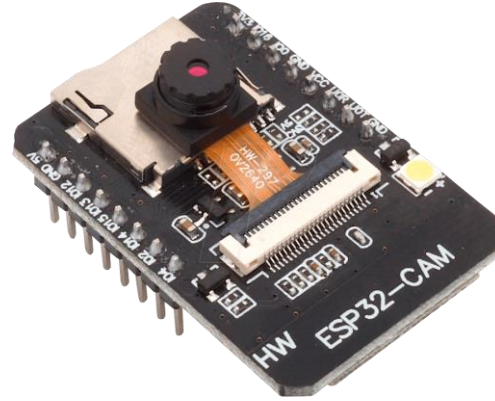
Start Time:	
20201118-075416	

Last Page Refresh: 06:57:39

<https://github.com/jomjol/AI-on-the-edge-device>

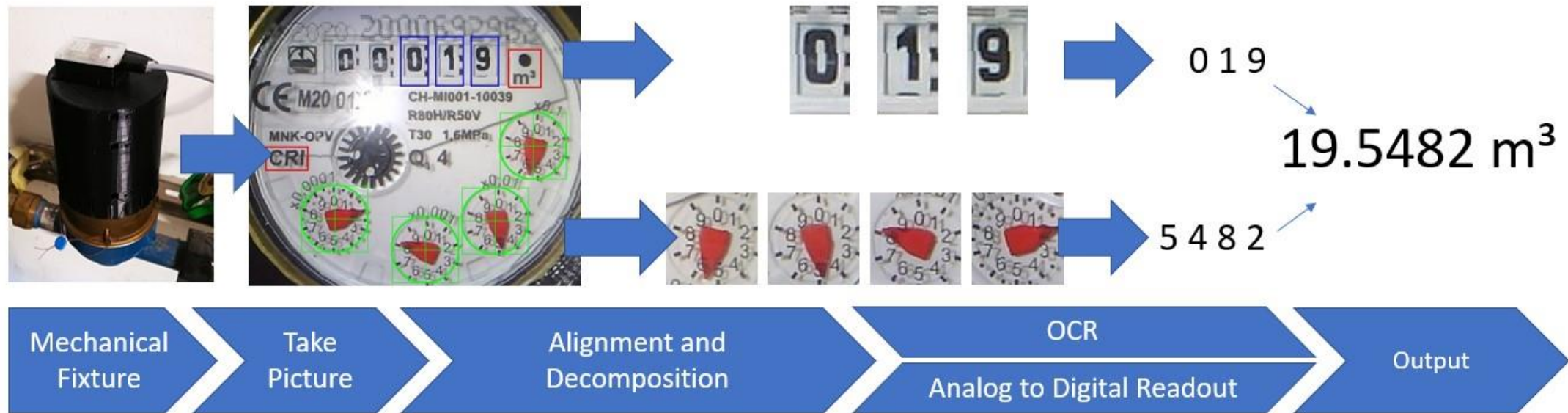
Hardware

- ESP32-Cam
 - ESP32
 - Kamera
 - SD-Kartenslot
 - Min. 4MB PSRAM
 - LED-Blitzlicht
- 5V Spannungsversorgung
- USB-TTL Adapter
 - zum Flashen der Firmware
- Halter



<https://github.com/jomjol/AI-on-the-edge-device>

Funktionsweise



<https://github.com/jomjol/AI-on-the-edge-device>

Zusätzliche Features

- Tensorflow Lite Integration
- Webserver
- REST API
- MQTT
- Fileserver
- OTA-updatefähig
- ...

Setup



Alignment Marker

► [CLICK HERE](#) for usage description. More infos in documentation: [Alignment](#)

Marker: Marker 1 Filename: /config/ref0.jpg

x: dx:

y: dy:

Selected Image Area:  Resulting Marker: 



Digit ROI

► [CLICK HERE](#) for usage description. More infos in documentation: [ROI Configuration](#)

☒ Digit ROI Processing

Number Sequence:

main

ROI:

dig2

Multiplier: x10
(only based on order)

Multiplier: x10
(order + decimal shift: 0)

x:

Δ x:

☒ Lock aspect ratio

y:

Δ y:

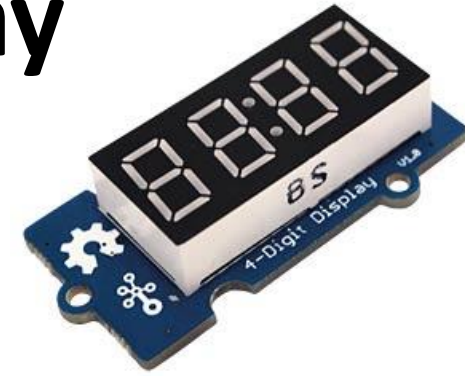
☒ Synchronize y, Δ x and Δ y between ROIs

☒ Keep equidistance of between all ROIs

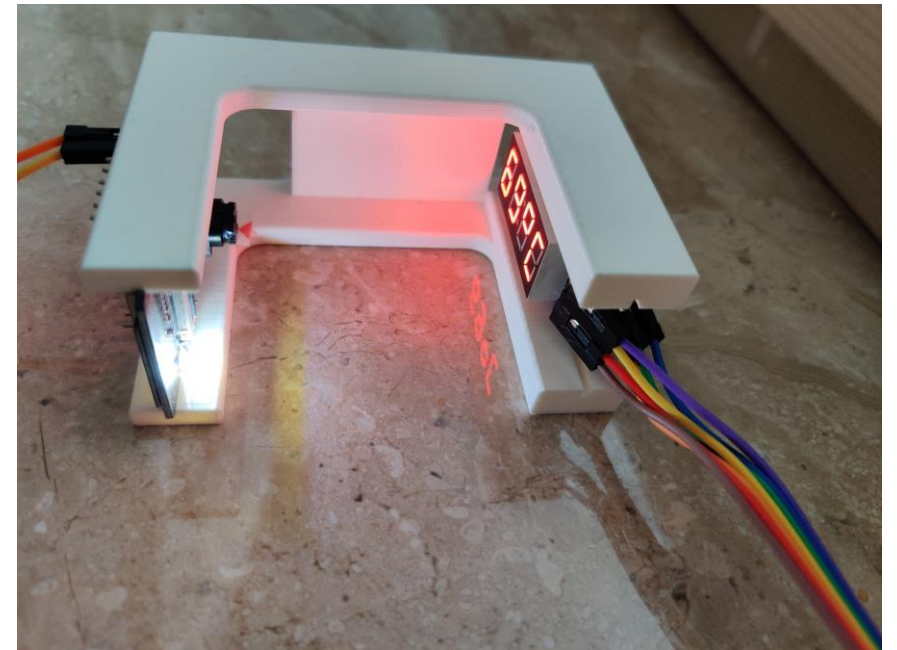
Reference Image:



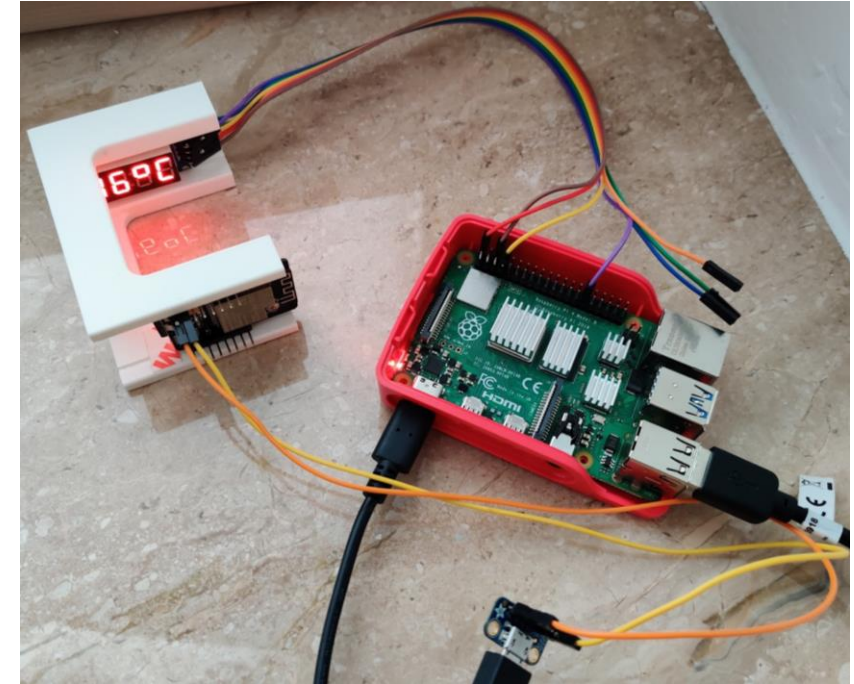
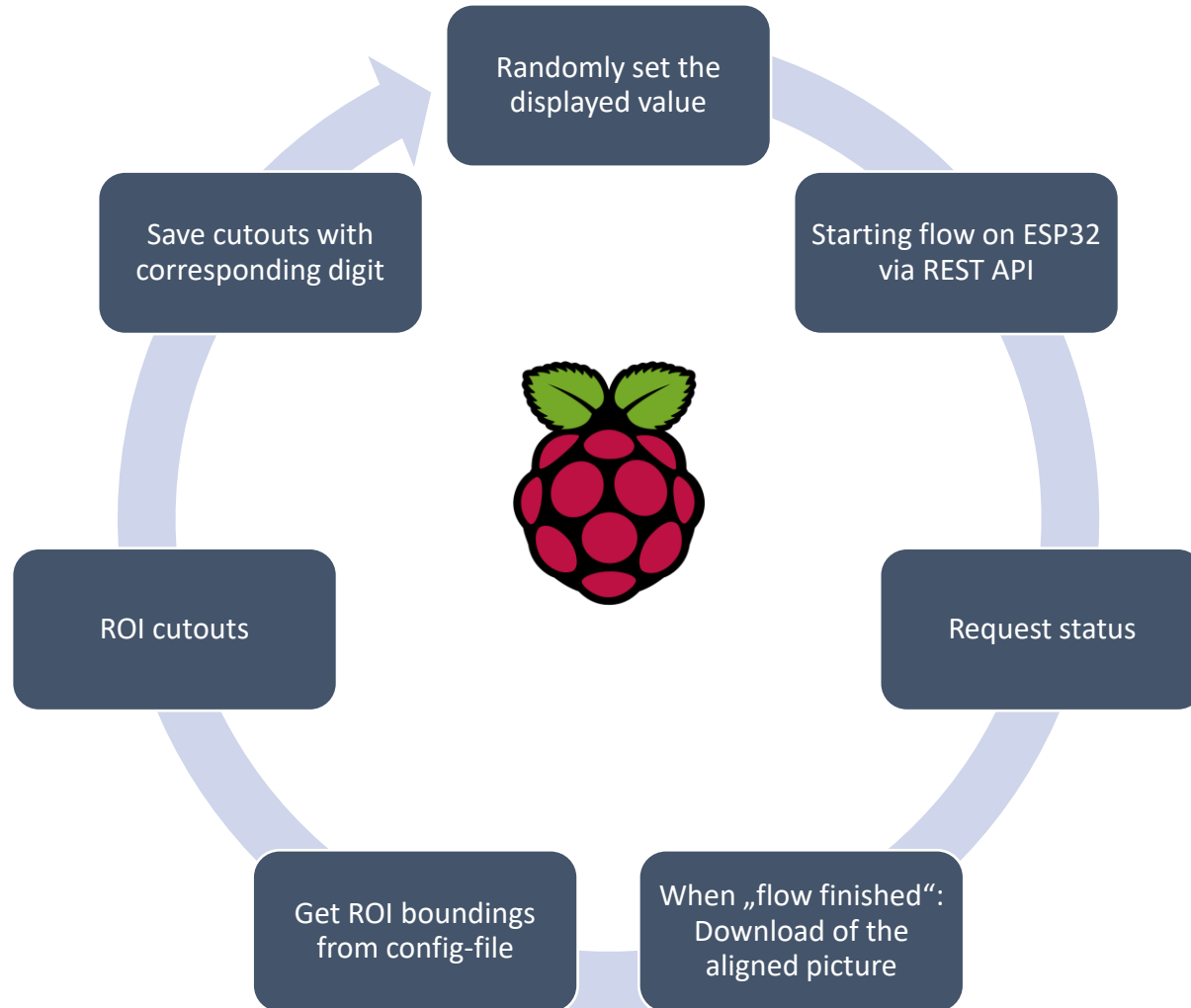
Versuche mit 7-Segment Display



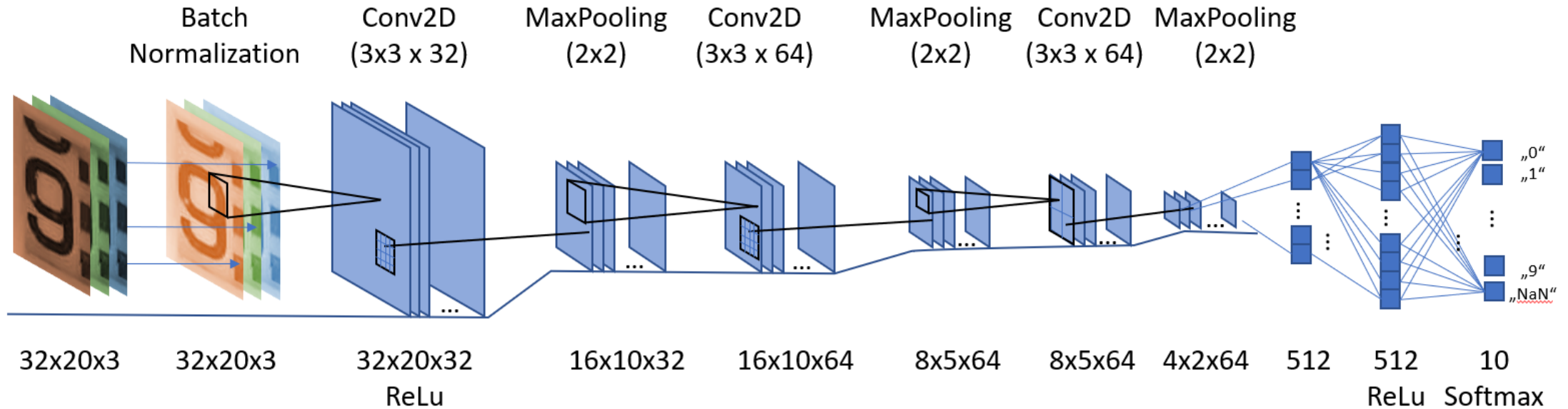
- Auswertung der angezeigten Temperatur auf einem 7-Segment Display mit I2C-Schnittstelle (TM1637)
 - Mitgelieferte Modelle zu ungenau (~31%)
- > Training eines eigenen Modells



Trainingsaufbau



Aufbau des Netzes



Starten der JupyterLab Umgebung

Zum Ausprobieren:

```
git clone https://github.com/SebastianVeigl/AI_on_the_edge_segment_train
cd AI_on_the_edge_segment_train/
pip install -r training/requirements.txt
```

jupyter-lab training

Settings Help

requirements.txt x Image_Preparation.ipynb x +

Image Preparation

Image preparation

The original image size is 55x90 pixels with a color depth of 3 (RGB). The below code can be used to convert the image to grayscale. The below code can be used as well.

```
[1]: import glob
import os
from PIL import Image

input_dir = r'../digits'
output_dir = 'digits_resized'

if not os.path.exists(output_dir):
    os.mkdir(output_dir)

target_size_x = 20
target_size_y = 32

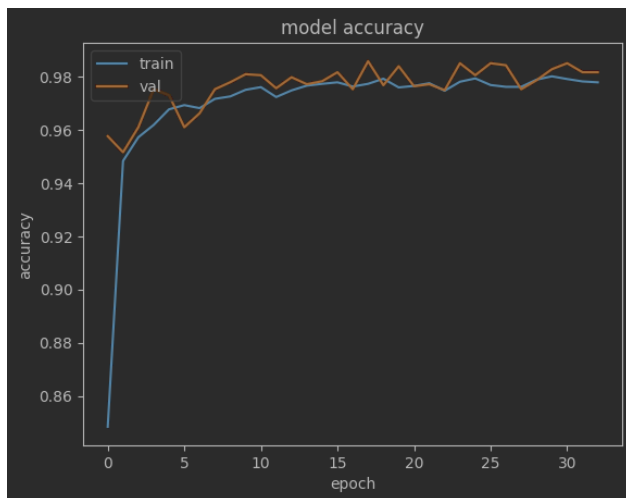
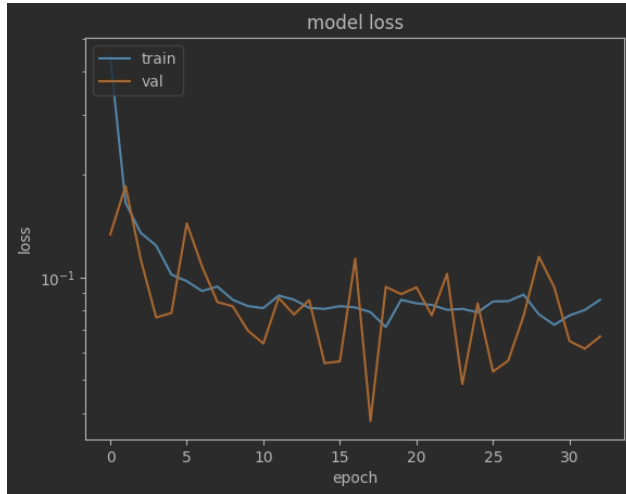
[2]: files = glob.glob(output_dir + '/*.jpg')
i = 0
for f in files:
    os.remove(f)
    i=i+1
print(str(i) + " files have been deleted.")

9574 files have been deleted.

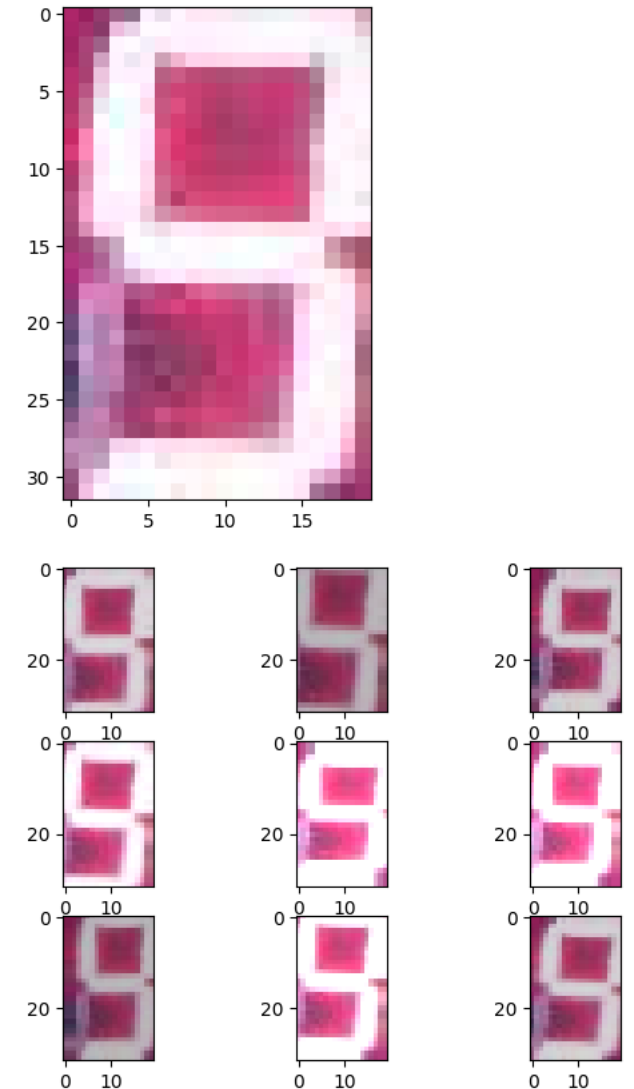
[3]: files = glob.glob(input_dir + '/*.jpeg', recursive=True)
files = files + glob.glob(input_dir + '/*.png')
files = files + glob.glob(input_dir + '/*.bmp')
count = 0
for aktfile in files:
    count = count + 1
    if not count % 250:
        print(str(count) + " ...")
    test_image = Image.open(aktfile)
    test_image = test_image.convert('RGB')
    test_image = test_image.resize((target_size_x, target_size_y), Image.NEAREST)
    base=os.path.basename(aktfile)
    base = os.path.splitext(base)[0] + ".jpg"
    save_name = output_dir + '/' + base
    # print("in: " + aktfile + " - out: " + save_name)
    test_image.save(save_name, "JPEG")
print(count)

250 ...
500
```

Training des neuronalen Netzes



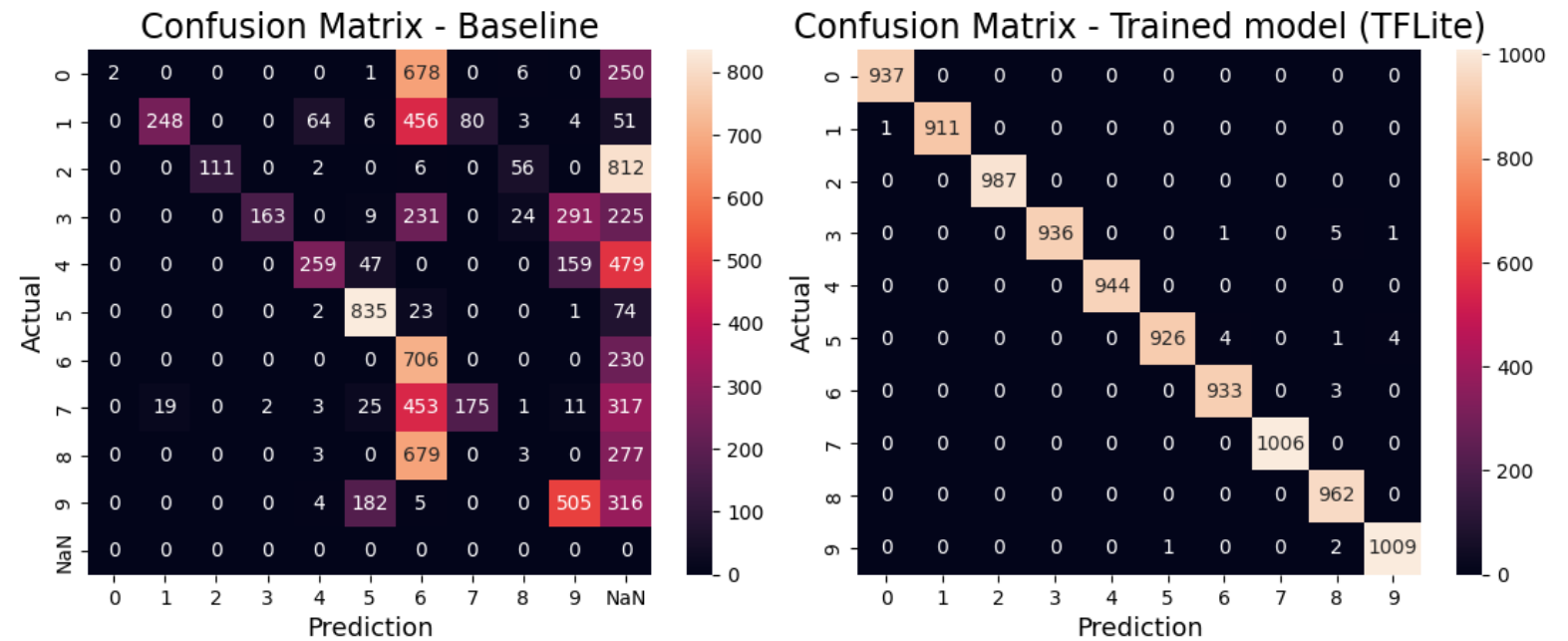
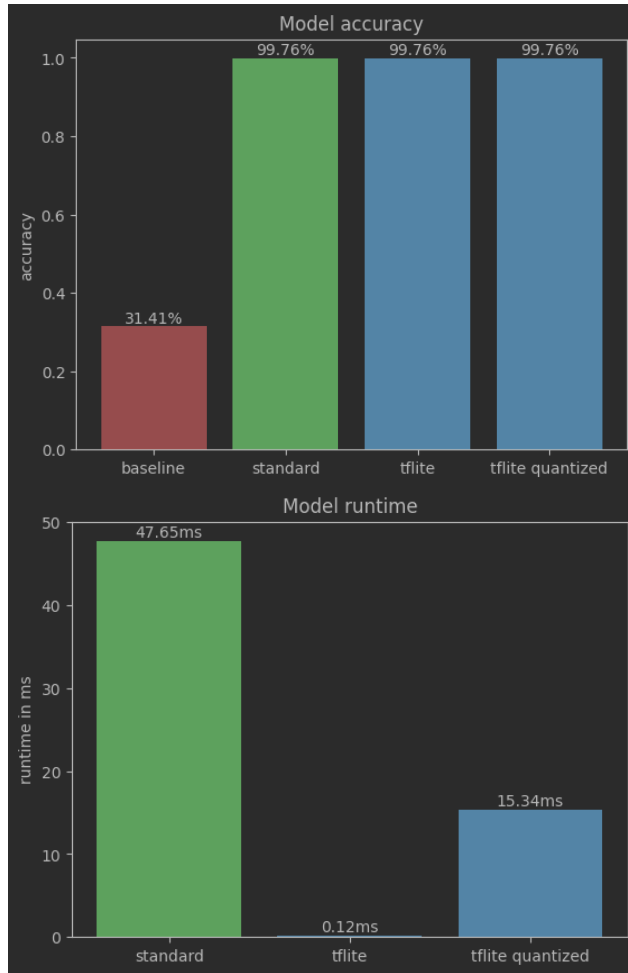
- Training auf dem Laptop
- Vorverarbeitung:
 - ROIs zu 32x20 Pixel ([Jupyter Notebook](#))
 - Augmentation
- Training
 - *Tensorflow* für den Aufbau + Training des Netzes
 - Train-Val-Split (70%/30%)
 - [Jupyter Notebook](#)
 - Umwandlung in TFLite
 - Post-training quantization



Sichern des Netzes

- Standard Format (.keras)
 - Ermöglicht Training + Inferenz
- TFLite
 - Nur Training
 - Geringere Laufzeit
 - Geringere Größe
 - Weniger RAM-Verbrauch
- TFLite (Post-training quantization)
 - Wie TFLite
 - Alle Gewichte + Aktivierungsfunktionen werden in Integer umgewandelt (oder andere Variablentypen z.B. float64)
 - Geringere Genauigkeit

Laufzeit/Genauigkeit der Modelle



Ergebnisse

	Baseline (TFLite)	Standard model (keras)	TFLite	TFLite quantized
Size	349KB	692KB	215KB	62KB
Accuracy	31,43%	99,76%	99,76%	99,76%
Runtime	-	47,6 ± 7,9ms	0,12 ± 0,01ms	15,3 ± 3,3ms

Probleme bei der Klassifizierung

- Bei hoher Sonneneinstrahlung bilden sich Streifen auf den Bildern
-> Nicht mehr eindeutig



Inferenz auf anderen Geräten

z. B.: Raspberry Pi mit Python

- Modell im TFLite-Format
 - Installation von *tflite_runtime* (Python-Paket) mit *pip3 install tflite-runtime*
 - Siehe Code zum Ausführen der Klassifikation eines Bildes

Am PC kann auch direkt das *tensorflow.lite* Subpaket verwendet werden

```
import tflite_runtime.interpreter as tflite
import numpy as np
from PIL import Image

TFLITE_FILE_PATH = '../models/7seg2912.tflite'

# Load the TFLite model in TFLite Interpreter
interpreter = tflite.Interpreter(TFLITE_FILE_PATH)
interpreter.allocate_tensors()

# Get input and output tensors.
input_details = interpreter.get_input_details()
output_details = interpreter.get_output_details()

image_in = Image.open('../training/digits_resized/7_19_12_2023_18_41_26.jpg')
test_image = np.array(image_in, dtype="float32")
img = np.reshape(test_image, [1, 32, 20, 3])

# Test the model on random input data.
input_shape = input_details[0]['shape']
interpreter.set_tensor(input_details[0]['index'], img)

interpreter.invoke()

# The function `get_tensor()` returns a copy of the tensor data.
# Use `tensor()` in order to get a pointer to the tensor.
output_data = interpreter.get_tensor(output_details[0]['index'])
print(output_data)
print(output_data.argmax())
```

Quellen

- <https://www.computacenter.com/de-de/what-we-do/cloud-and-data-center/ki/ai-the-edge>
- <https://viso.ai/edge-ai/edge-ai-applications-and-trends/>
- <https://github.com/jomjol/AI-on-the-edge-device>
- <https://jomjol.github.io/AI-on-the-edge-device-docs/>
- <https://github.com/jomjol/neural-network-digital-counter-readout/>
- <https://www.tensorflow.org/lite/>
- https://www.tensorflow.org/lite/performance/post_training_quantization