

# Research Internship with RCS@TUM

Implementation of tiny machine learning models on microcontrollers

**Philipp van Kempen**

Department of Electrical and Computer Engineering, Technical University of Munich (TUM)

**24.08.2020 - 25.10.2020**



*TUM Uhrenturm*

# Motivation

- Machine learning in the past vs. today
- Compute-intensive algorithms
- State of the Art: ML/AI on smartphones
- Research topic: ML/AI on microcontrollers

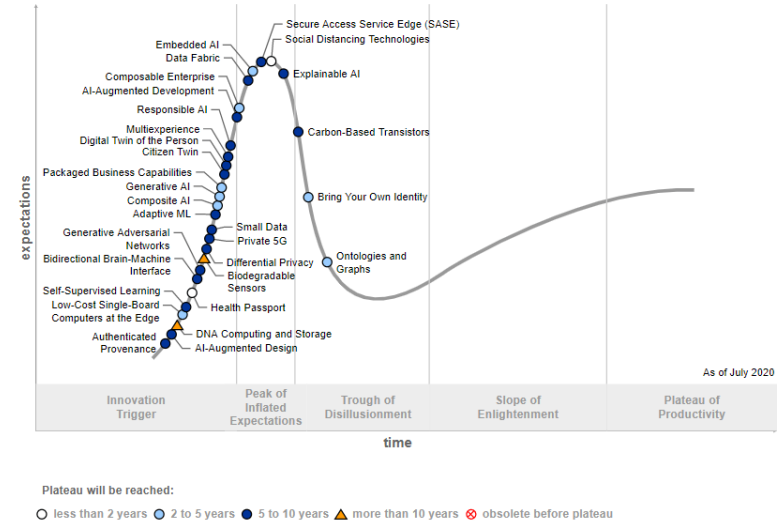
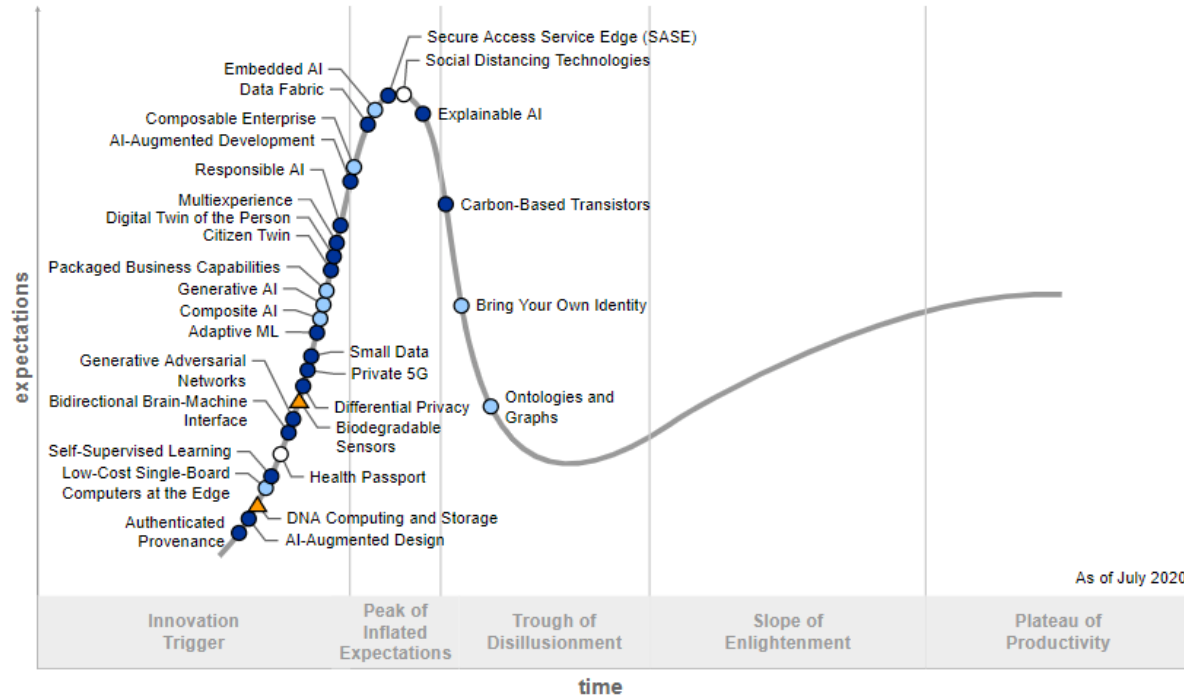


Figure: 2020 Gartner Hype Cycle<sup>1</sup>

<sup>1</sup><https://www.gartner.com/smarterwithgartner/5-trends-drive-the-gartner-hype-cycle-for-emerging-technologies-2020/>



Plateau will be reached:

○ less than 2 years ● 2 to 5 years ● 5 to 10 years ▲ more than 10 years ✕ obsolete before plateau

# Goals

- 9 weeks
- Support Electronic System Level (ESL) research group (EDA+RCS)
- Implement reference implementations of TinyML models on STM32 Hardware
- Work based on previous attempts by Alex Hoffman
- Extend toolchain with new features and documentation
- Summarize results at the end of the internship

# Steps

## Major topics:

- Reading books and web pages
- Toolchain setup and extension
- Training of examples
- Implementation of examples
- Documentation and handover

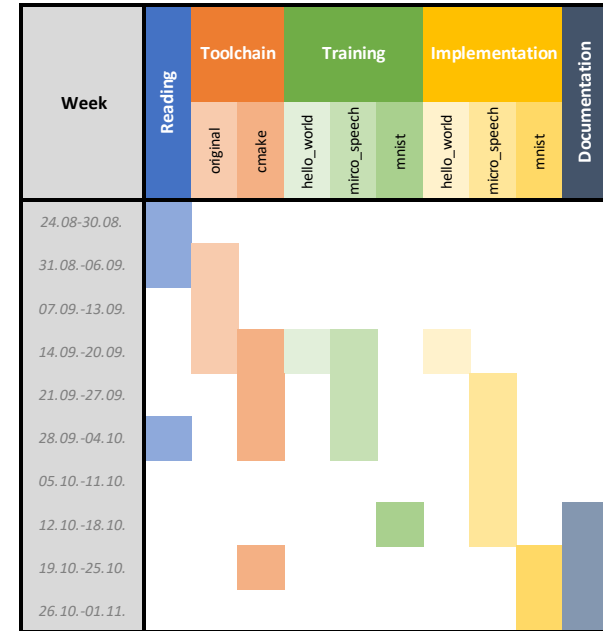


Figure: Schedule

# Reading

- Book by Pete Warden and Daniel Situnayake
- Introduction to Tensorflow Lite for Microcontrollers (TFLM) framework
- Referencing examples located in the Tensorflow source tree
  - **Hello World**
  - **Micro Speech**
- Official documentation of Tensorflow (v1.5 and latest)

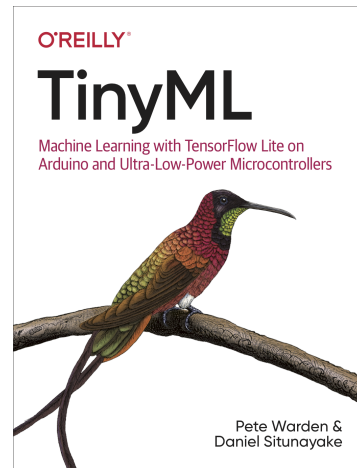


Figure: TinyML Book<sup>2</sup>

---

<sup>2</sup><https://tinymmlbook.com>

# Toolchain

- Original ARM mBed toolchain used as a reference
- CMake based - Initially developed by Konstantin Oblaukhov<sup>3</sup>, Improved and extended by Alex Hoffman
- Reference Application: STM3240G-EVAL-TensorFlow-MNIST
- Some workarounds required to fix upstream bugs



(a) STM32F413H-DISCOVERY



(b) STM32F769I-DISCOVERY

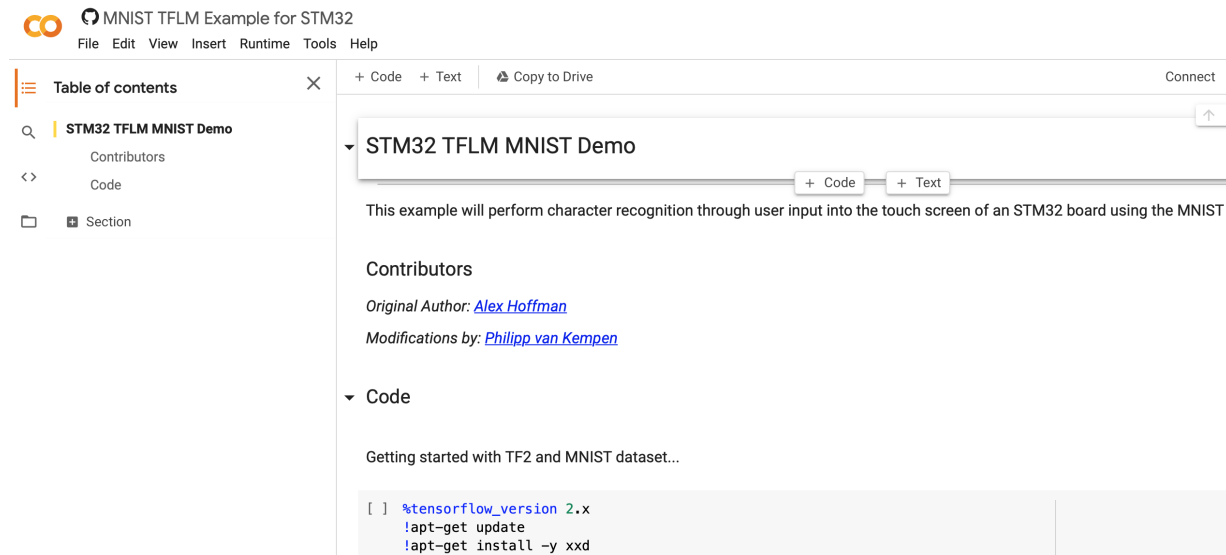
Figure: Target Boards<sup>4</sup>

<sup>3</sup><https://github.com/ObKo/stm32-cmake>

<sup>4</sup>Images taken from <https://www.st.com>

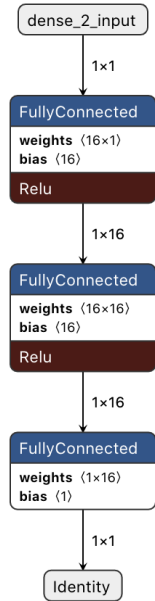
# Training

- Training scripts: Interactive notebooks hosted on Google Colaboratory<sup>5</sup>
- Deployment on Microcontrollers:
  - Quantization
  - Optimizations
  - Conversion to TFLite
  - Interpreter vs. Compiled offline model

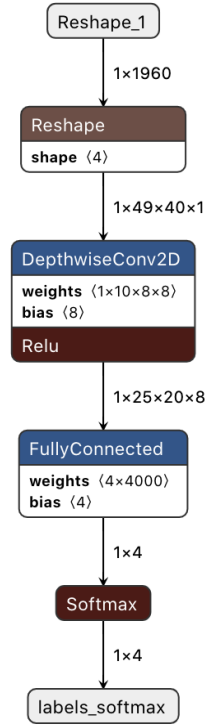


<sup>5</sup><https://colab.research.google.com>

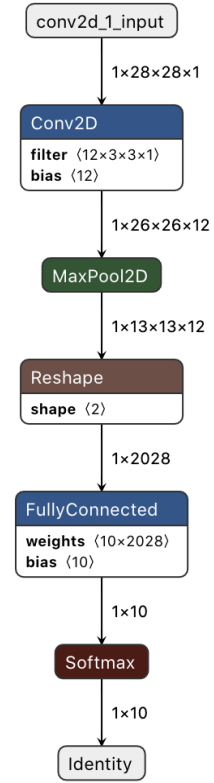




(a) Hello World



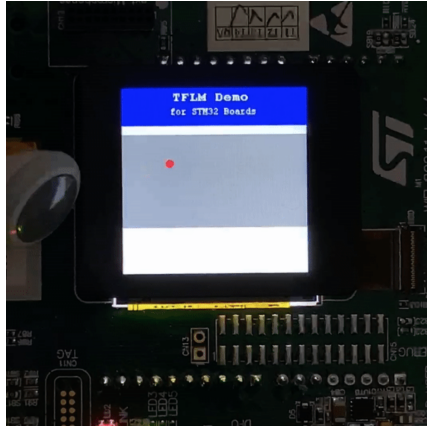
(b) Micro Speech



(c) MNIST

# Implementation

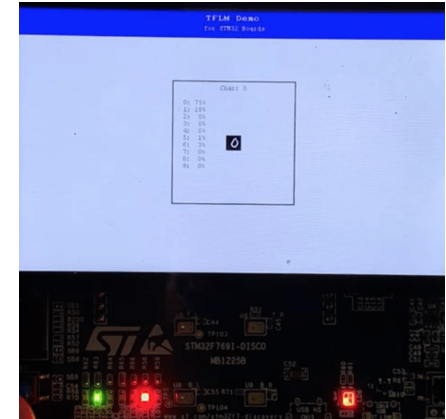
- **Hello World:** Ran almost out of the box
- **Micro Speech:**
  - Hardest challenge: Streaming real-time audio
  - Boards too slow?  $\Rightarrow$  Solution: Enable CMSIS-NN
  - Further Tuning of parameters was required to detect any voice commands  $\Rightarrow$  More false positives!
- **MNIST:**
  - Dataset: ten-thousands of handwritten digits
  - Quantization issue: Latest TFLM does not support unsigned uint8 inputs  $\Rightarrow$  transform to  $\{-128, \dots, 127\}$
  - Challenge: Transform touchscreen input into 28x28 pixel grayscale images
  - Bug in Alex's implementation: Inverted colors
  - Some digits can not be detected at all  $\Rightarrow$  Improve network architecture
- **Common:**
  - Added Benchmarking module
  - Added Memory Reporting support
  - Added possibility to choose TFLM interpreter or compiler
  - For easier Testing: SD-Card support to feed real samples to the network



(a) Hello World



(b) Micro Speech



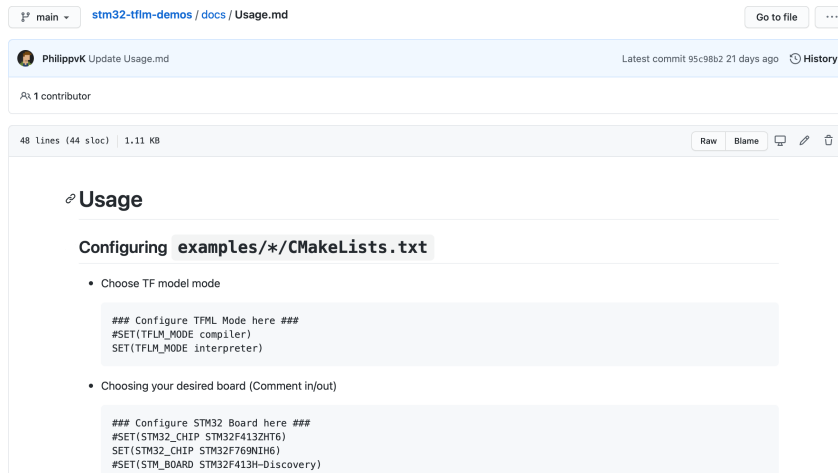
(c) MNIST

Figure: Examples running on the STM32 Boards<sup>6</sup>

<sup>6</sup>See <https://github.com/PhilippvK/stm32-tflm-demos> for GIFs!

# Documentation

- Common toolchain parts as submodules  $\Rightarrow$  less redundant code
- Wrapper repository<sup>7</sup>
- Common documentation at a single place
- Great Extend-ability



The screenshot shows a GitHub repository page for 'stm32-tflm-demos' with the file 'Usage.md' selected. The file is 48 lines long, 1.11 KB, and was last committed 21 days ago. The content of the file is as follows:

```
Usage

Configuring examples/*/CMakeLists.txt

• Choose TF model mode

### Configure TFML Mode here ###
#SET(TFLM_MODE compiler)
SET(TFLM_MODE interpreter)

• Choosing your desired board (Comment in/out)

### Configure STM32 Board here ###
#SET(STM32_CHIP STM32F413ZHT6)
SET(STM32_CHIP STM32F769NIH6)
#SET(STM_BOARD STM32F413H-Discovery)
SET(STM_BOARD STM32F769NIH6-Discovery)
```

<sup>7</sup><https://github.com/PhilippvK/stm32-tflm-demos>

# Evaluation

See report PDF or Documentation<sup>8</sup> for details! Only showing tables here...

	Boards		
Metrics	STM32F413HDISCOVERY	STM32F769IDISCOVERY	Units
Clock Frequency	100	216	<i>MHz</i>
Special Features	-	Double Issue, I/D-Cache	-
Flash Memory	1.5	2	<i>MB</i>
SRAM Memory	256	512	<i>kB</i>

Figure: Board Metrics

<sup>8</sup><https://github.com/PhilippvK/stm32-tflm-demos/blob/main/docs/Metrics.md>

# Evaluation - Memory Usage & Number of Ops

	Examples			
Type	hello_world	mirco_speech	mnist	Units
Model Size (FLASH)	2	18	23	<i>kB</i>
TensorArena Size (SRAM)	1	7	11	<i>kB</i>

Figure: Memory Usage (approx.)

Examples		
hello_world	mirco_speech	mnist
41 <i>FLOPS</i>	689980 <i>FLOPS</i>	202810 <i>FLOPS</i>

Figure: Number of Ops (after quant.)

# Evaluation - Runtime Measurements & CMSIS-NN

		<b>Examples</b>			
<b>Section</b>	<b>CPU</b>	hello_world	mirco_speech	mnist	<b>Units</b>
Populate	F4	~0	38	132	<i>ms</i>
	F7	~0	11	88	
Invoke	F4	~0	49	34	<i>ms</i>
	F7	~0	52	13	
Respond	F4	~0	~0	125	<i>ms</i>
	F7	~0	~0	93	

Figure: Runtime Measurements (approx.)

		<b>Examples</b>			
<b>Settings</b>		hello_world	mirco_speech	mnist	<b>Units</b>
CMSIS_NN	OFF	~1	413 (unusable!)	52	<i>ms</i>
	ON	~0	52	13	<i>ms</i>
Difference		(?)	(-87)	(-75)	%

Figure: CMSIS-NN Improvements (approx.)

# Conclusions and Outlook

- Main goal fulfilled
- MNIST was optional but also implemented successfully
- Models need more tuning
- Demonstrated capabilities of TinyML on microcontroller platforms

---

<sup>9</sup>[https://github.com/munober/thesis/blob/master/digital\\_edition.pdf](https://github.com/munober/thesis/blob/master/digital_edition.pdf)



# Conclusions and Outlook

- Main goal fulfilled
- MNIST was optional but also implemented successfully
- Models need more tuning
- Demonstrated capabilities of TinyML on microcontroller platforms

## Future Work:

- Support USB-Storage as alternative to SD-Card
- Add TinyFace Model<sup>9</sup>
- Merge Deployment and Evaluation Flow with EDA RISCv-toolchain
- Enable usage of FreeRTOS instead of baremetal
- Add possibility to parse results via UART
- Feed samples via UART for automated testing

⇒ Extract data on the accuracy of the models

---

<sup>9</sup>[https://github.com/munober/thesis/blob/master/digital\\_edition.pdf](https://github.com/munober/thesis/blob/master/digital_edition.pdf)