# Webinar Logistics

- Turn on your system sound to hear the streaming presentation.

- If you have technical problems, click "Help" or submit a question asking for assistance.

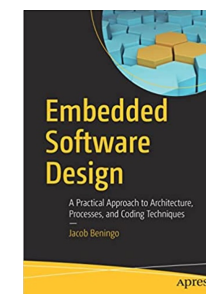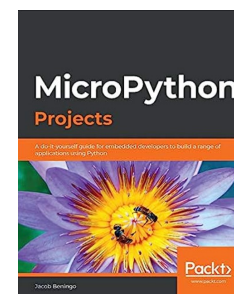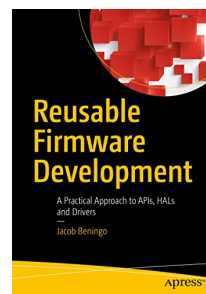- Participate in 'Group Chat' by maximizing the chat widget in your dock.

## THE SPEAKER

# Beningo Embedded Group - President

Focus: Embedded Software Consulting and Training

Specializes in <u>creating</u> and <u>promoting</u> embedded software **excellence** in businesses around the world.



### Jacob Beningo

Visit 'Lecturer Profile'

Blogs for:
- DesignNews.com
- Embedded.com
- EmbeddedRelated.com
- MLRelated.com

Visit www.beningo.com to learn more ...

# Course Sessions

- AI and ML for Microcontrollers
- Writing Embedded Software with ChatGPT and Open.AI
- Tools for Machine Learning in Microcontrollers
- Training a Model for the STM32
- **Deploying Machine Learning Models**

# 1 Preparing the Model for Export

# Exporting an STM32 Binary

# Exporting an STM32 AI pack

# Apply Optimization(s)



Available optimizations for NN Classifier

**Quantized (int8)** ⭐

Currently selected

This optimization is recommended for best performance.

| RAM USAGE | LATENCY | CONFUSION MATRIX | | | | |
|---|---|---|---|---|---|---|
| 1.5K | 1 ms | 87.8 | 7.0 | 0 | 0 | 5.2 |
| **ROM USAGE** | **ACCURACY** | 0 | 99.2 | 0 | 0 | 0.8 |
| 15.4K | 95.68% | 0 | 0 | 100 | 0 | 0 |
| | | - | - | - | - | - |

**Unoptimized (float32)**

Click to select

| RAM USAGE | LATENCY | CONFUSION MATRIX | | | | |
|---|---|---|---|---|---|---|
| 1.5K | 1 ms | 87.6 | 6.6 | 0 | 0 | 5.8 |
| **ROM USAGE** | **ACCURACY** | 0 | 100 | 0 | 0 | 0 |
| 17.7K | 95.88% | 0 | 0 | 100 | 0 | 0 |
| | | - | - | - | - | - |

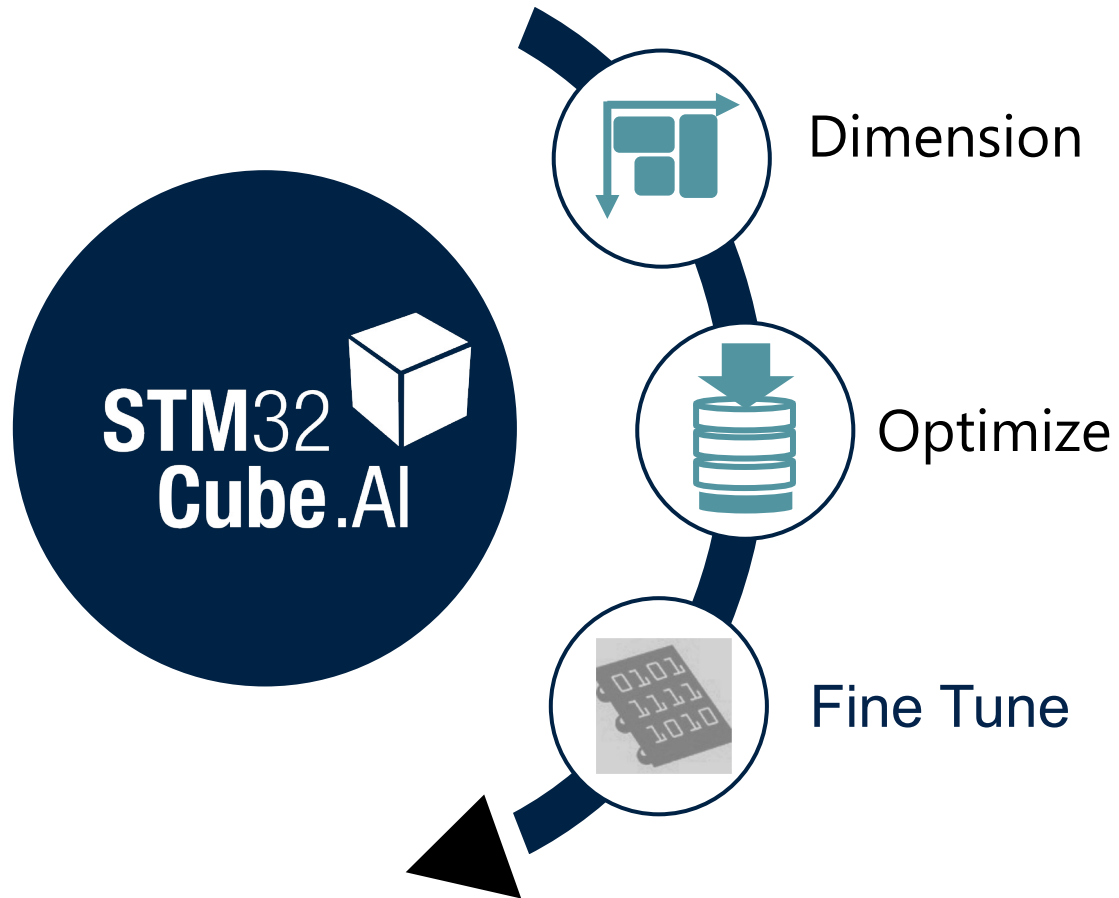Estimate for Cortex-M4F 80MHz (ST IoT Discovery Kit)

Build

What method do you prefer for testing?

- Using the prebuilt binary

- Using the pack

- C++ library
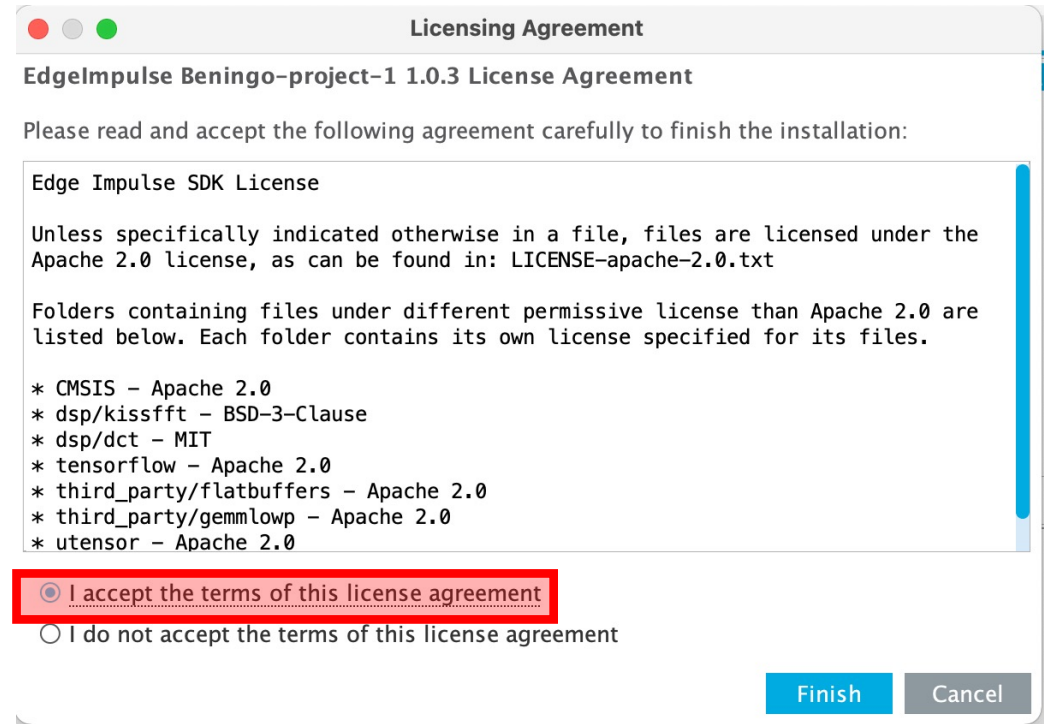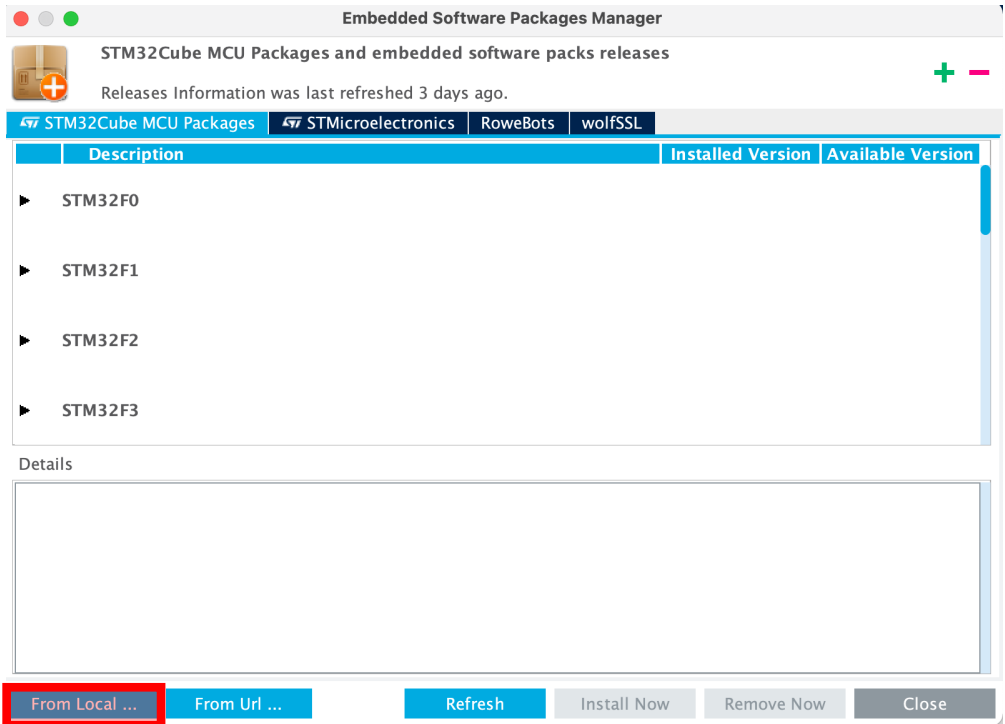
- other

**2** Importing the Model

# STM32Cube.AI Overview

**Dimension**

- ✓ Quickly assess model footprint requirements
- ✓ Select and configure MCU in STM32CubeMX
- ✓ Review model layers in STM32Cube.AI

**Optimize**

- ✓ Generate C-code for pre-trained model
- ✓ Support quantized models to reduce RAM, flash and latency with minimal loss of accuracy
- ✓ Use light run-time libraries
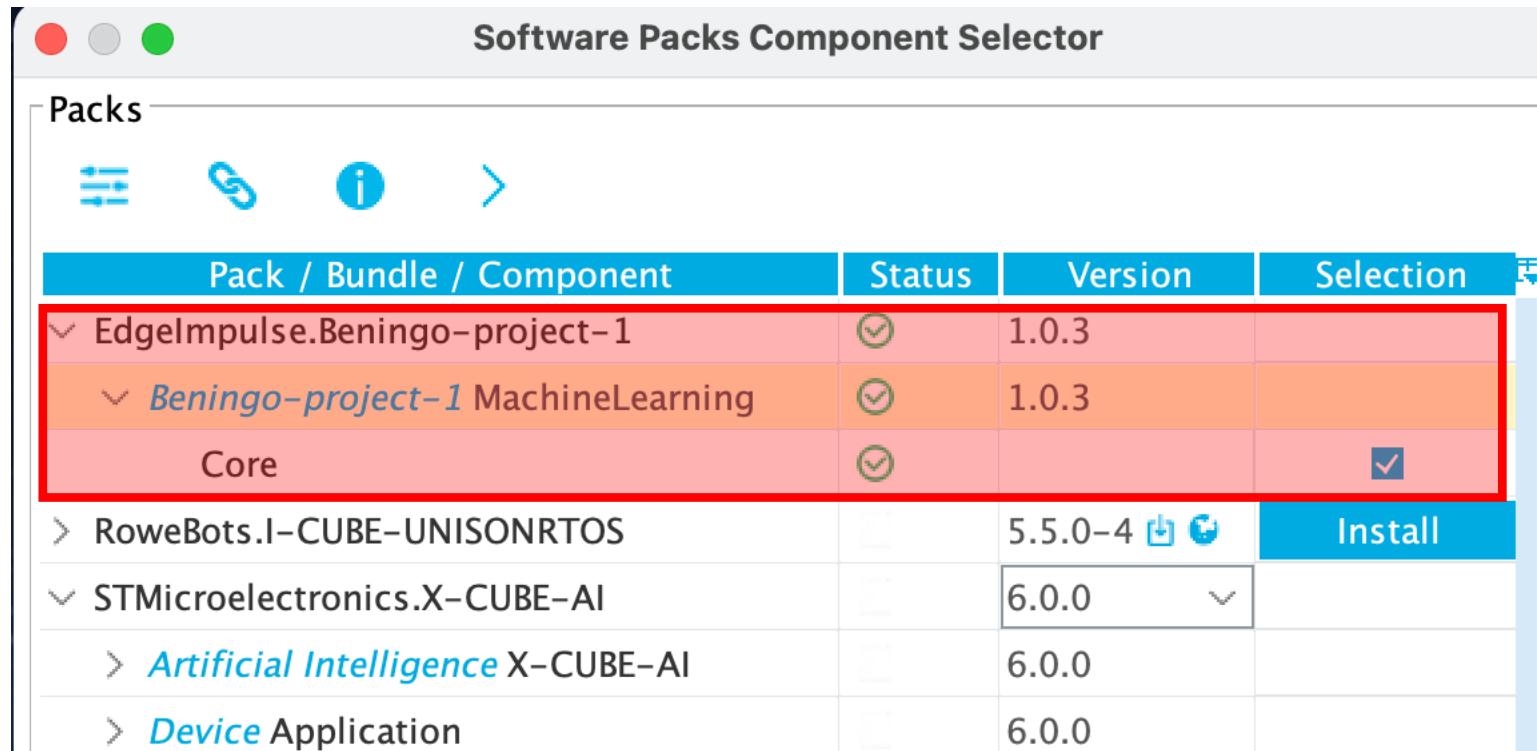- ✓ Optimize for performance

**Fine Tune**

- ✓ Optimize memory allocation
- ✓ Fine control of weight mapping
- ✓ Split between internal and external memory
- ✓ Update model without full FW update

# Install Pack into STM32Mx Project
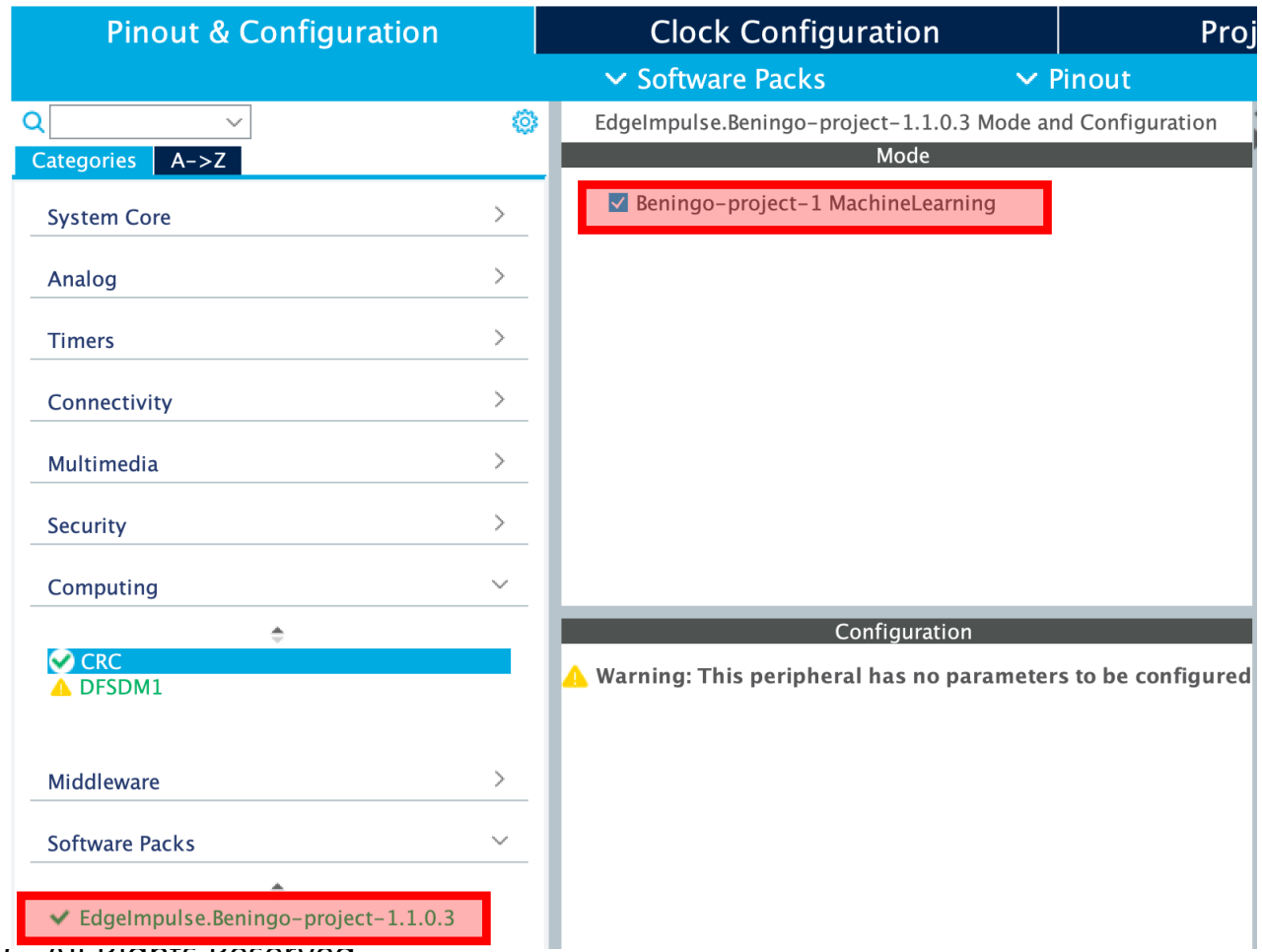
Help -> Manage Embedded Software Packages

# Install Pack into STM32Mx Project

# Install Pack into STM32Mx Project

# Install Pack into STM32Mx Project



Modify, Build, Deploy
https://bit.ly/32ESC3N

15

**3** Running the Model

# Running the Model

In a terminal, run the command: edge-impulse-run-impulse

```
Starting inferencing in 2 seconds...
Sampling... Storing in file name: /fs/device-classification.116
Predictions (DSP: 14 ms., Classification: 1 ms., Anomaly: 1 ms.):
    Circle: 0.99609
    Updown: 0.00000
    Wave: 0.00000
    anomaly score: -0.026
Finished inferencing, raw data is stored in '/fs/device-classification.116'. Use AT+UPLOADFILE to send back to Edge Impulse.


Starting inferencing in 2 seconds...
Sampling... Storing in file name: /fs/device-classification.121
Predictions (DSP: 15 ms., Classification: 0 ms., Anomaly: 2 ms.):
    Circle: 0.00000
    Updown: 0.00000
    Wave: 0.99609
    anomaly score: -0.132
Finished inferencing, raw data is stored in '/fs/device-classification.121'. Use AT+UPLOADFILE to send back to Edge Impulse.
```

# Running the Model

```
Starting inferencing in 2 seconds...
Sampling... Storing in file name: /fs/device-classification.118
Predictions (DSP: 15 ms., Classification: 0 ms., Anomaly: 2 ms.):
    Circle: 0.01172
    Updown: 0.98828
    Wave: 0.00000
    anomaly score: -0.141
Finished inferencing, raw data is stored in '/fs/device-classification.118'. Use AT+UPLOADFILE to send back to Edge Impulse.
Starting inferencing in 2 seconds...
Sampling... Storing in file name: /fs/device-classification.119
Predictions (DSP: 14 ms., Classification: 1 ms., Anomaly: 1 ms.):
    Circle: 0.21094
    Updown: 0.78906
    Wave: 0.00000
    anomaly score: -0.164
Finished inferencing, raw data is stored in '/fs/device-classification.119'. Use AT+UPLOADFILE to send back to Edge Impulse.
```

What methods can be used to improve classifaction ?
- Running average on the output
- Monitor the anomaly value
- Set a minimum classification percentage
- All the above
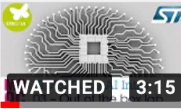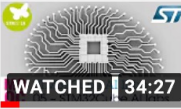- Other

**4** Going Further

# Next Steps

- Connect the output to a PWM LED channel
- Setup a DAC and drive an output voltage
- Configure the rate at which the inference runs (frequency control)
- Try and compare the Keras model behavior
- Improve the training model to provide a more accurate sine wave

# AI and ML Resources

- [Jacob's AI Blogs](#)
- [Jacob's CEC courses](#)
- [Jacob's ML Blogs](#)

- Embedded Bytes Newsletter
  - [http://bit.ly/1BAHYXm](http://bit.ly/1BAHYXm)

  [www.beningo.com](http://www.beningo.com)