

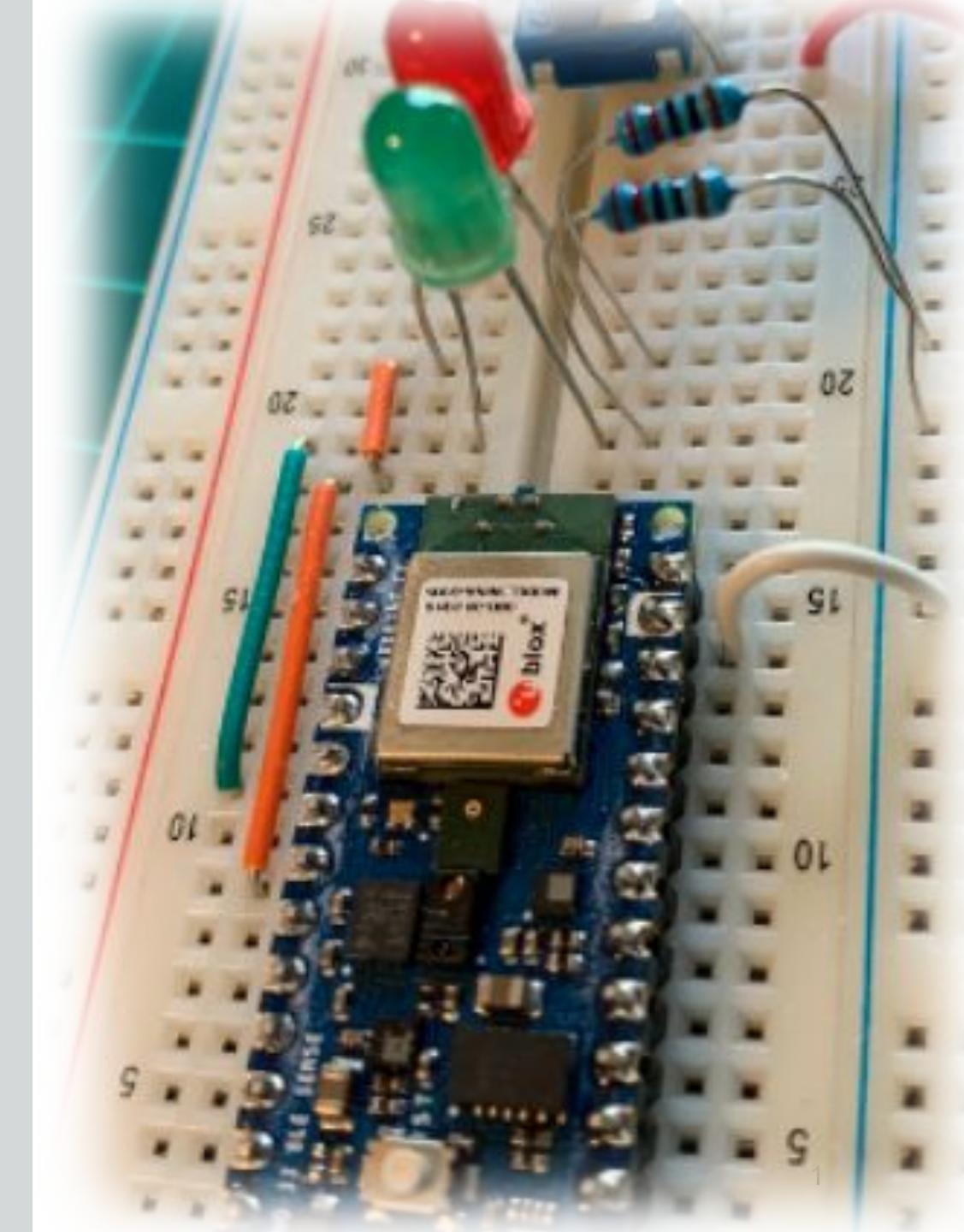
IESTI01 – TinyML

Embedded Machine Learning

29. EdgeAI Going Further



Prof. Marcelo Rovai
UNIFEI



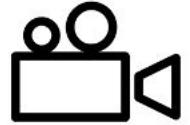
EdgeAI

Going Further

Definitions

- **Artificial Intelligence (AI)**: The broad discipline of creating **machines capable of performing tasks that typically require human intelligence**, including learning, reasoning, and problem-solving.
- **Machine Learning (ML)**: A subset of AI focused on building **systems that learn from data**, allowing computers to improve their performance on tasks without being explicitly programmed for each task.
- **Deep Learning (DL)**: An advanced **subset of machine learning involving neural networks** with multiple layers (deep networks), enabling powerful data modeling and decision-making based on vast amounts of data.
- **Edge AI**: The practice of **processing AI algorithms on local devices near the data source**, reducing latency and bandwidth use while enhancing privacy and speed in applications like real-time analytics and local data processing.
- **TinyML**: The field of **machine learning focused on developing low-power models** that can operate directly on small devices like microcontrollers, extending AI capabilities to edge devices.

Hardware



Anomaly Detection
Sensor Classification
20 KB



Rpi-Pico
(Cortex-M0+)



Arduino Nano
(Cortex-M4)



Arduino Pro
(Cortex-M7)



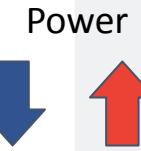
Wio

Image
Classification
250 KB+

KeyWord Spotting
Audio Classification
50 KB



TinyML



EdgeML (AI)

Video
Classification
2 MB+

Object Detection
Complex Voice
Processing
1 MB+



RaspberryPi
(Cortex-A)

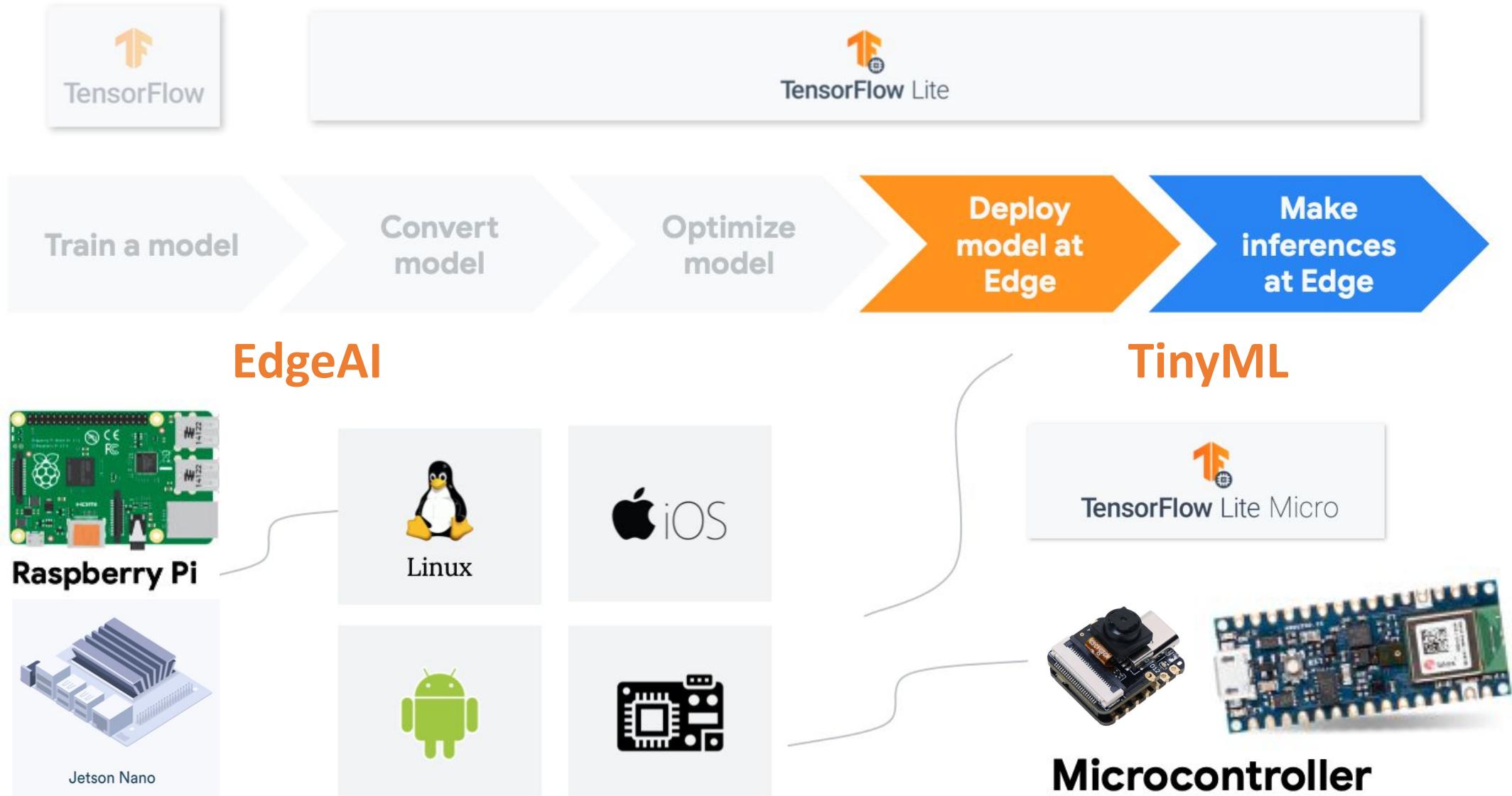


SmartPhone
(Cortex-A)



Jetson Nano
(Cortex-A + GPU)

Software

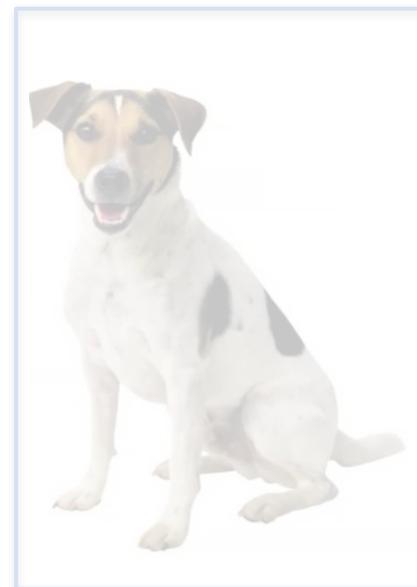


Computer Vision Main Types

Image Classification
(Multi-Class Classification)

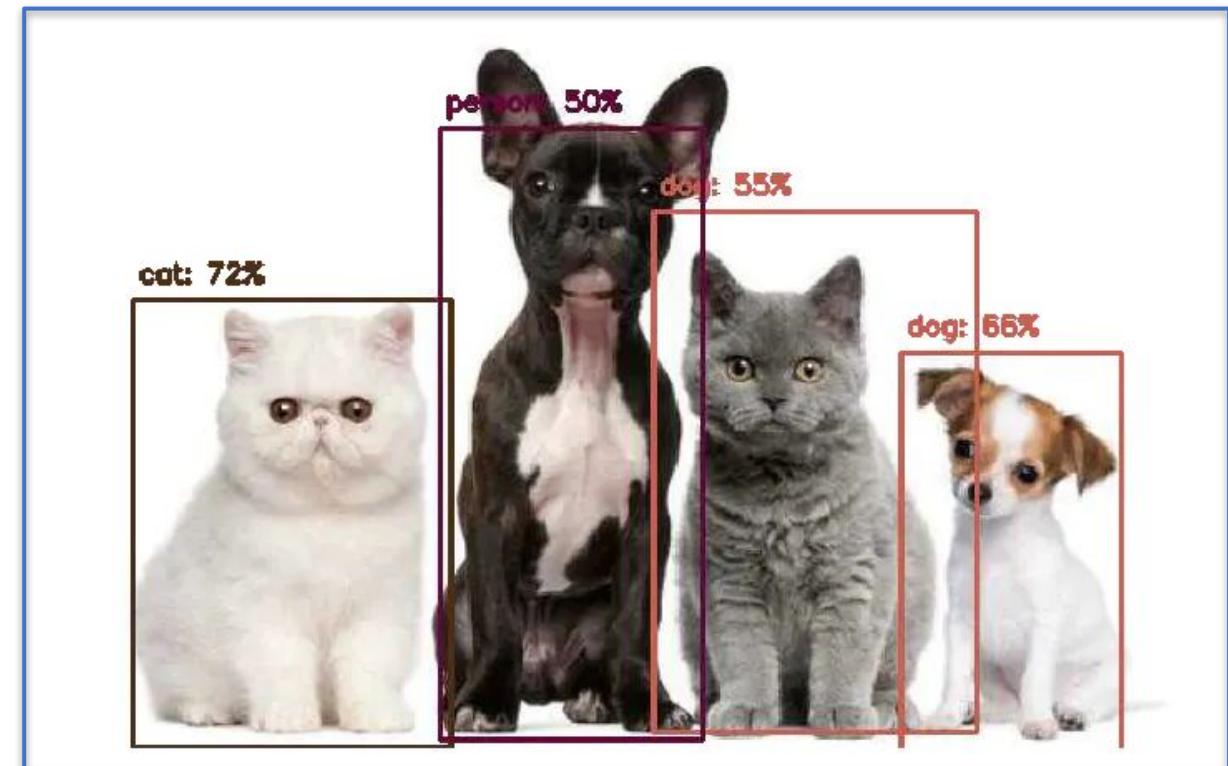


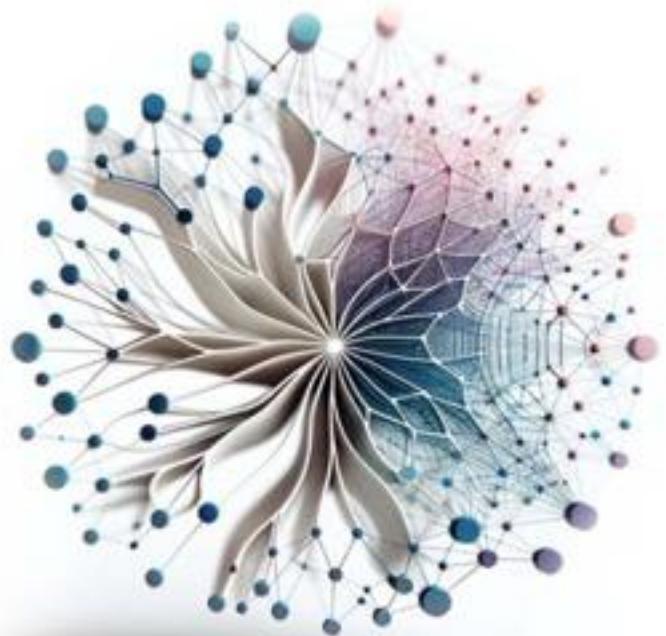
Cat: 70%



Dog: 80%

Object Detection
Multi-Label Classification + Object Localization





Machine Learning Systems

with TinyML

Written, edited and curated by
Prof. Vijay Janapa Reddi
Harvard University

With special thanks to the community for their contributions and support.



Nicla Vision > Object Detection

Object Detection



Nicla Vision



XIAO ESP32S3 > Object Detection

Object Detection

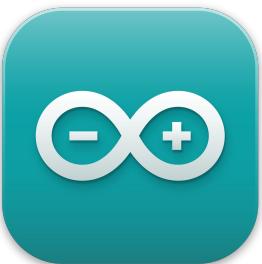


XIAO ESP32S3

FOMO

Object Detection model

Dataset



File > Examples > ESP32 >
Camera > CameraWebServer.ino



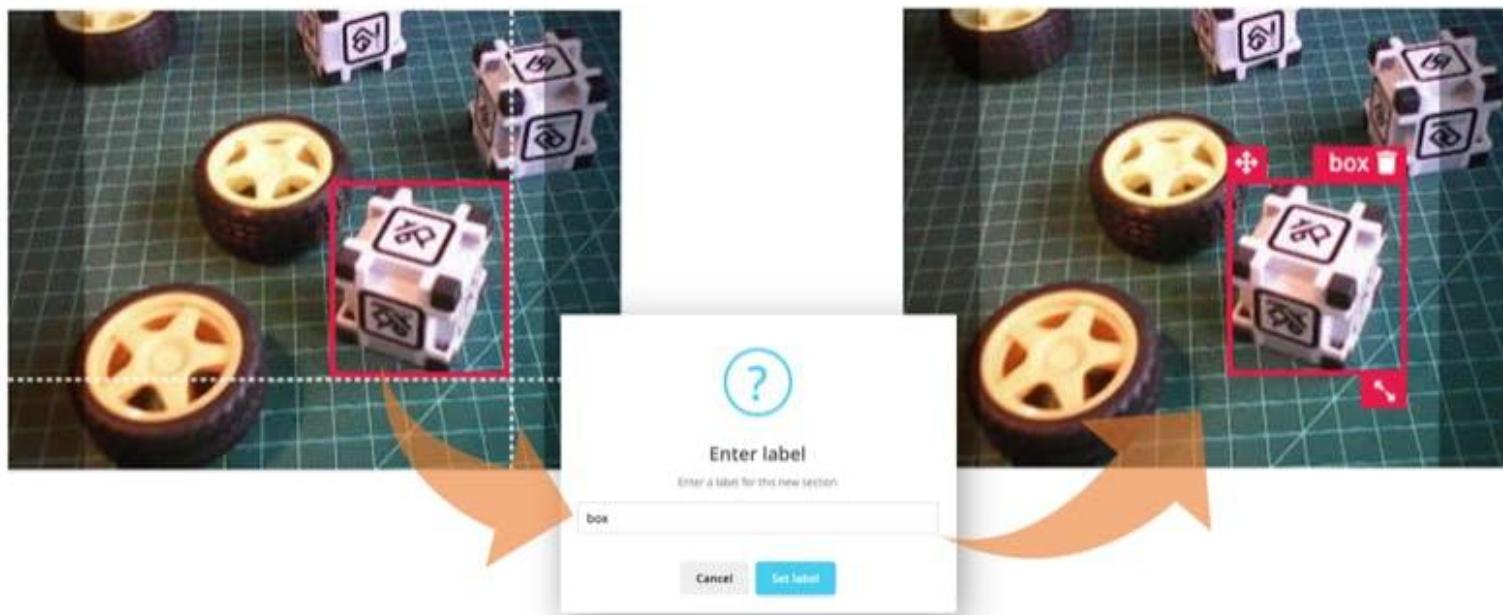
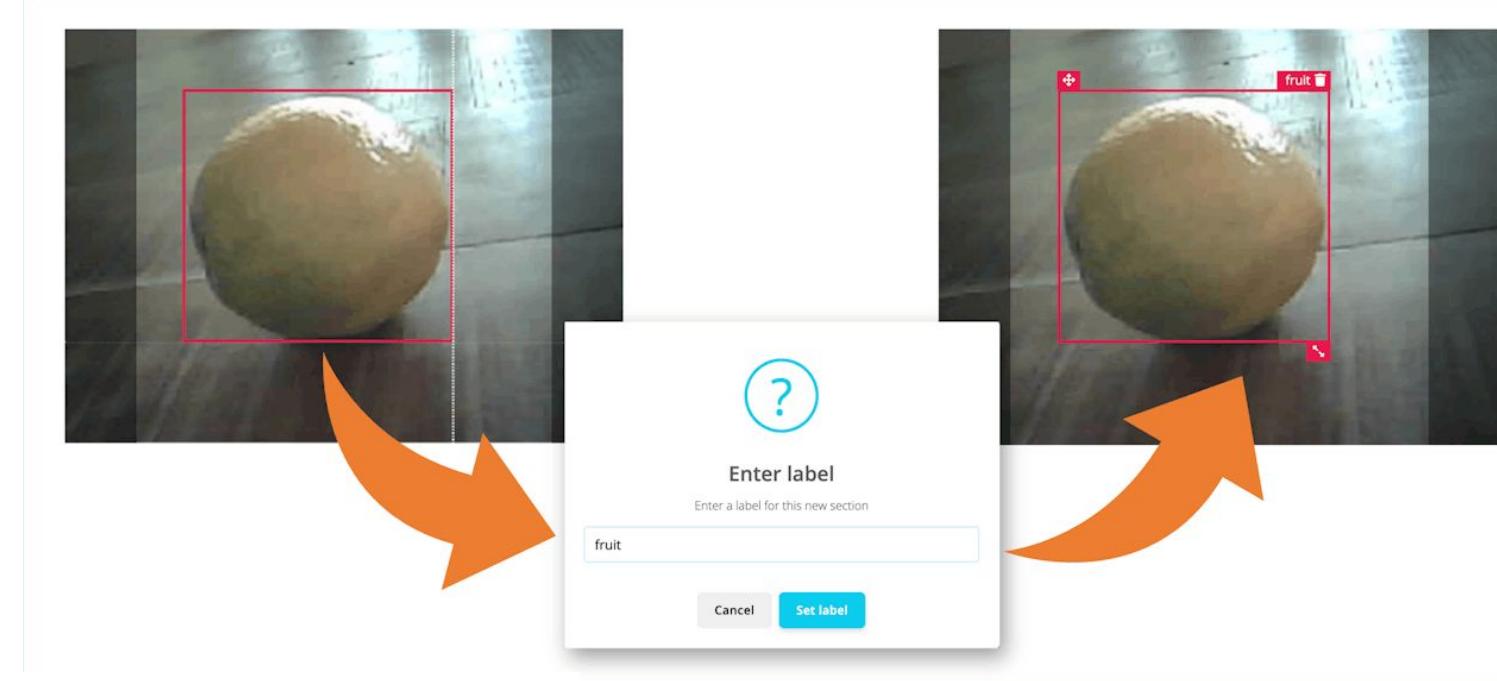
The screenshot shows the CameraWebServer interface running on an ESP32. The left side is a configuration menu with various camera parameters. The right side shows a live video feed of two green frogs. Two orange arrows highlight specific actions: one pointing to the 'Get Still' button at the bottom left and another pointing to the 'Save' button in the top right corner of the video frame.

Labeling

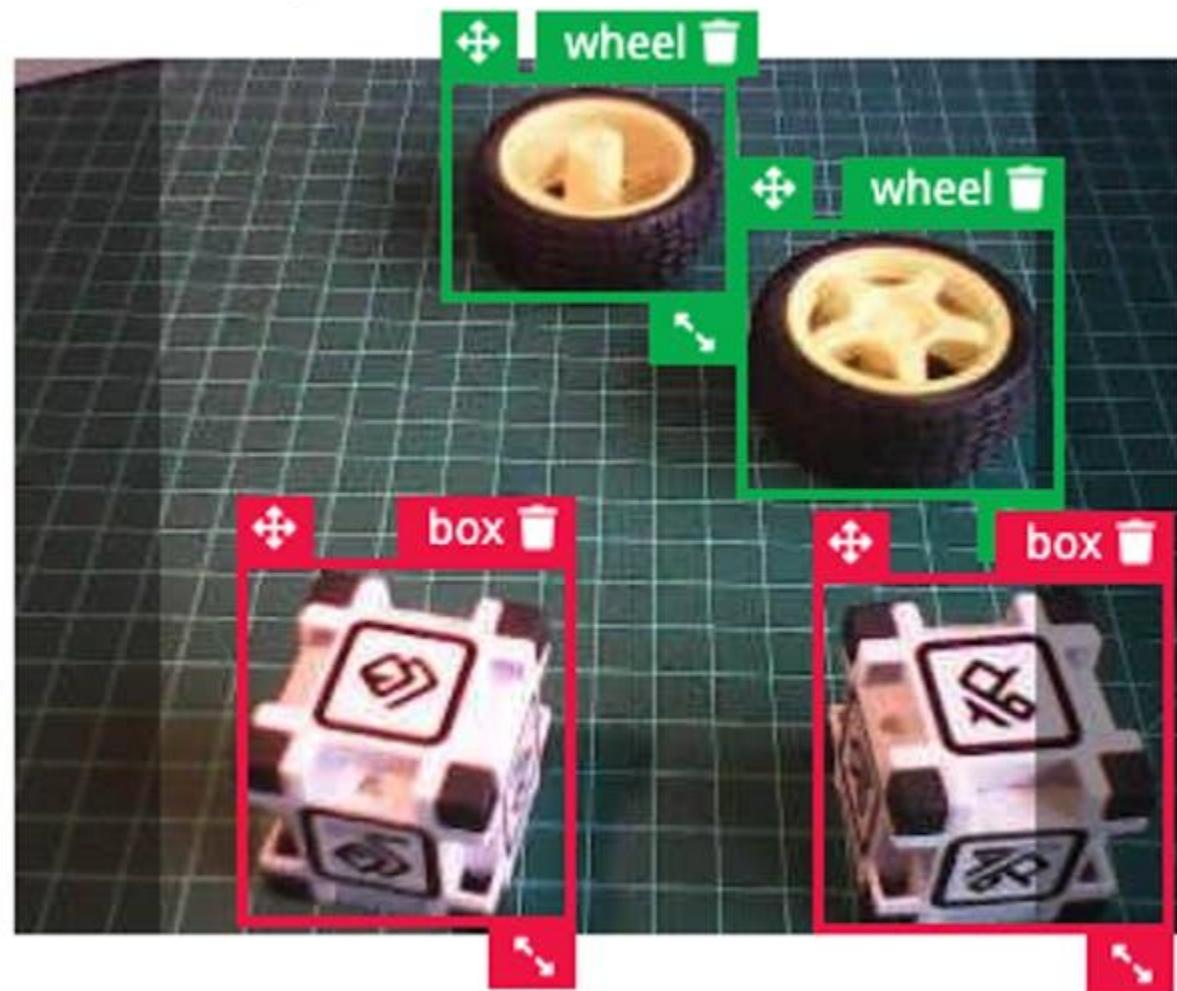
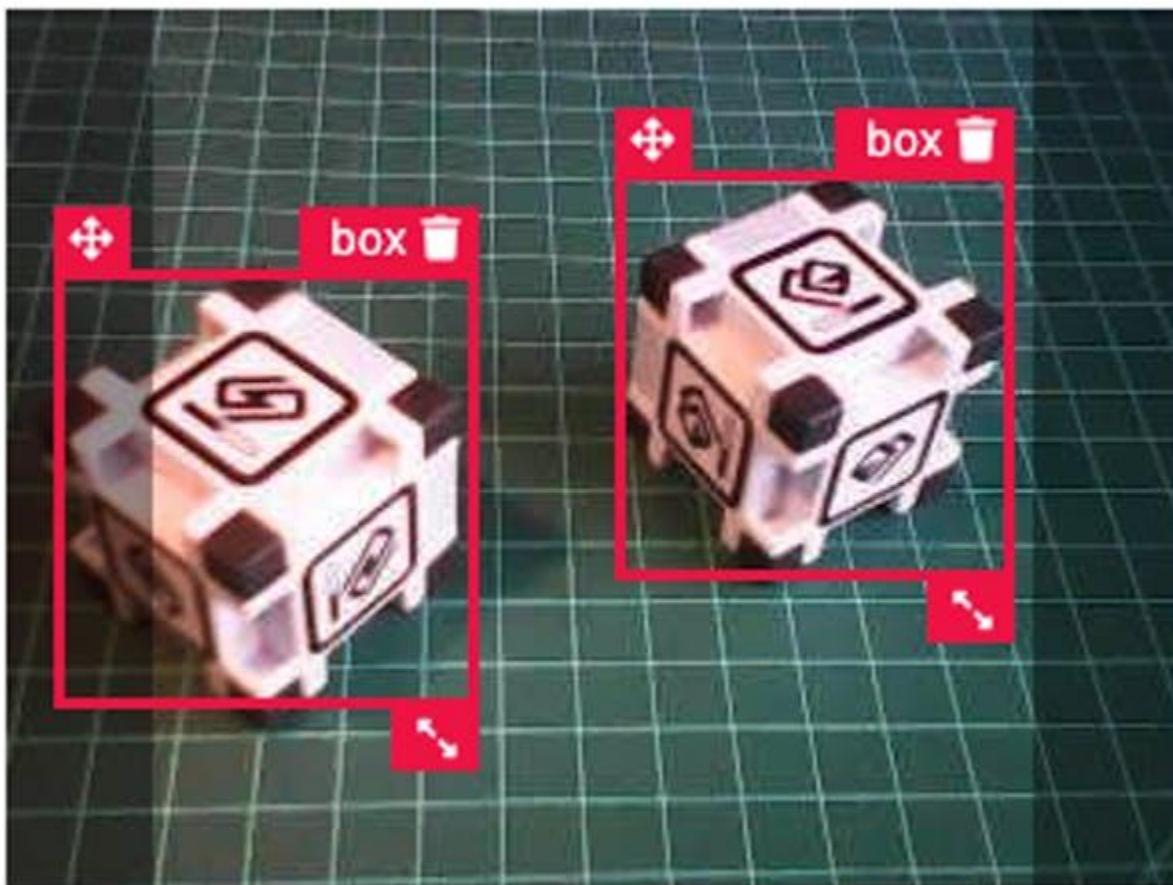
The screenshot shows the Edge Impulse web interface for dataset management. The main navigation bar includes 'Dataset', 'Data sources', and 'Labeling queue (47)', with the 'Labeling queue' tab highlighted by an orange arrow. The left sidebar lists various project management and documentation links. The central area displays a 'Dataset' table with columns for 'Training' (47) and 'Test' (0) samples, showing details like sample name, labels, add date, and length. To the right, a 'Collect data' section prompts to connect a device, and a 'RAW DATA' panel shows a timestamp '20231128151645' and a small image of two oranges and a green frog toy on a surface.

| SAMPLE NAME | LABELS | ADDED | LENGTH |
|----------------|--------|-----------------|--------|
| 20231128151645 | - | Today, 15:27:09 | - |
| 20231128150613 | - | Today, 15:27:09 | - |
| 20231128150604 | - | Today, 15:27:09 | - |
| 20231128150833 | - | Today, 15:27:09 | - |
| 20231128150600 | - | Today, 15:27:09 | - |
| 20231128150855 | - | Today, 15:27:09 | - |
| 20231128150458 | - | Today, 15:27:09 | - |
| 20231128150713 | - | Today, 15:27:09 | - |
| 20231128150908 | - | Today, 15:27:09 | - |

Labeling



Labeling



XIAO-ESP32S3-Sense-Objec X +

studio.edgeimpulse.com/studio/315759/create-impulse

EDGE IMPULSE

MJRoBot (Marcelo Rovai) / XIAO-ESP32S3-Sense-Object_Detection

An impulse takes raw data, uses signal processing to extract features, and then uses a learning block to classify new data.

Image data

Input axes
image

Image width 96 **Image height** 96

Resize mode
Squash

For object detection use a square image size, e.g. 96x96, 160x160 or 320x320.

Image

Name Image

Input axes (1)
 image

Object Detection (Images)

Name Object detection

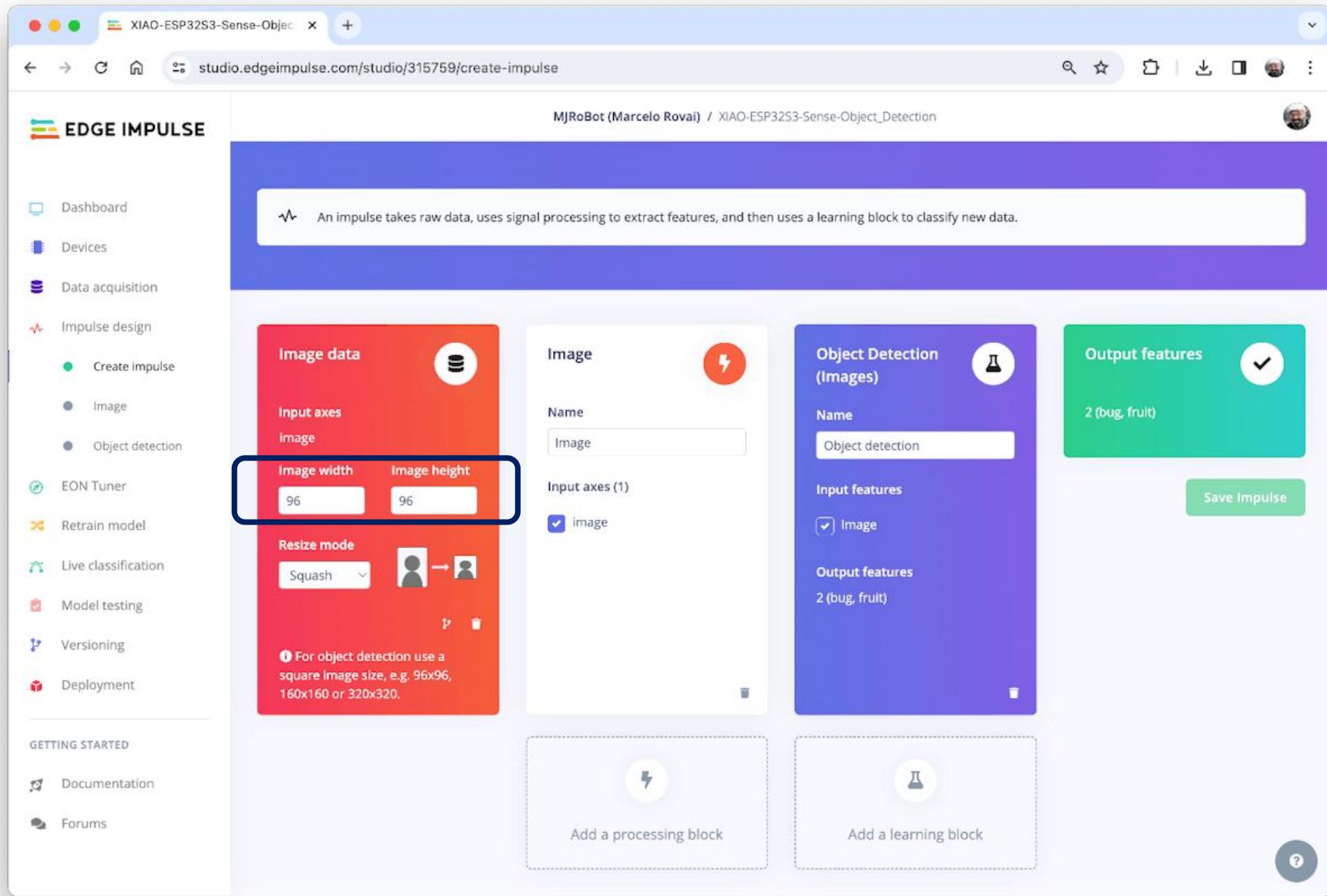
Input features
 Image

Output features
2 (bug, fruit)

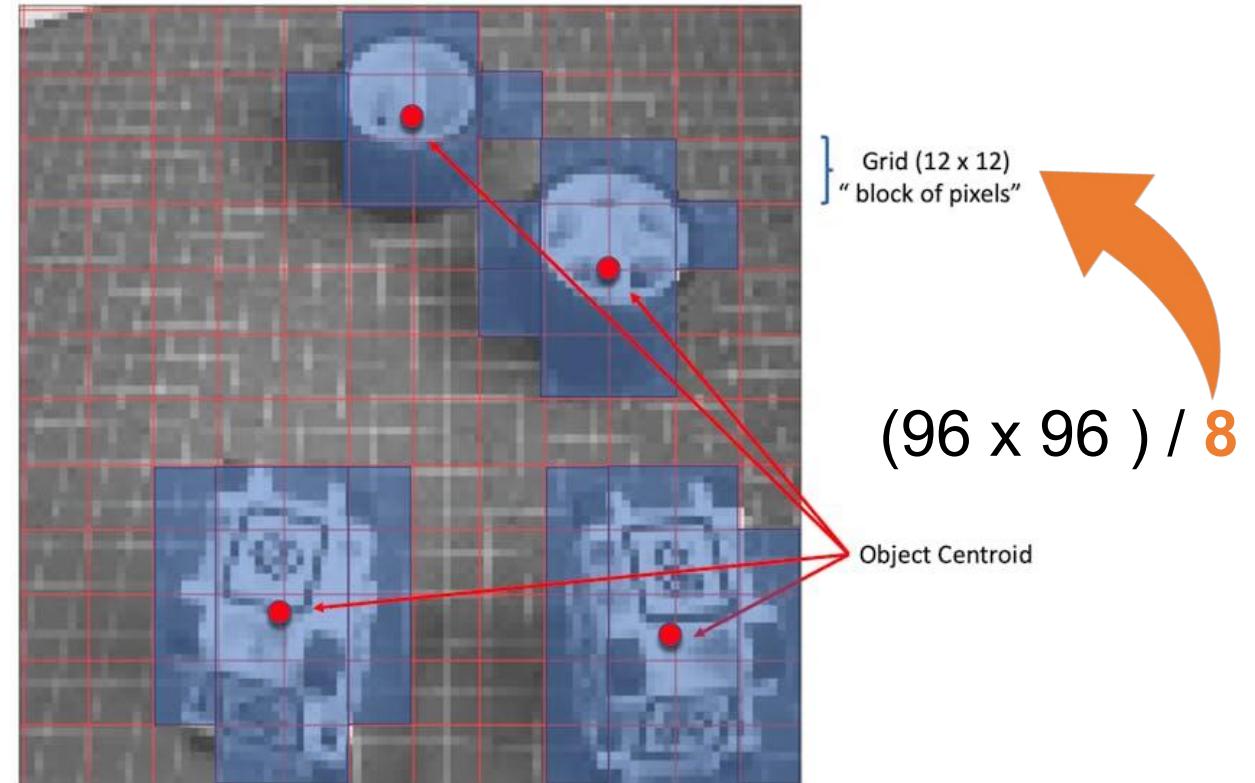
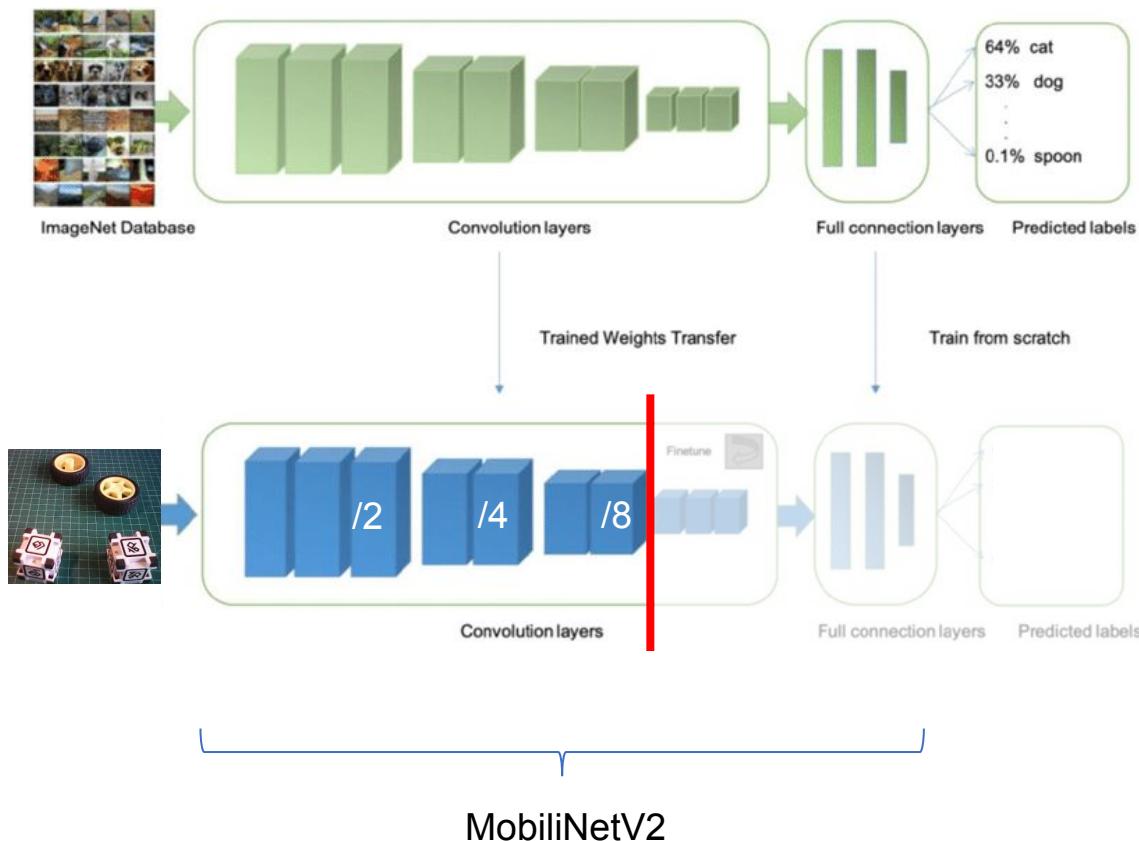
Save Impulse

Add a processing block

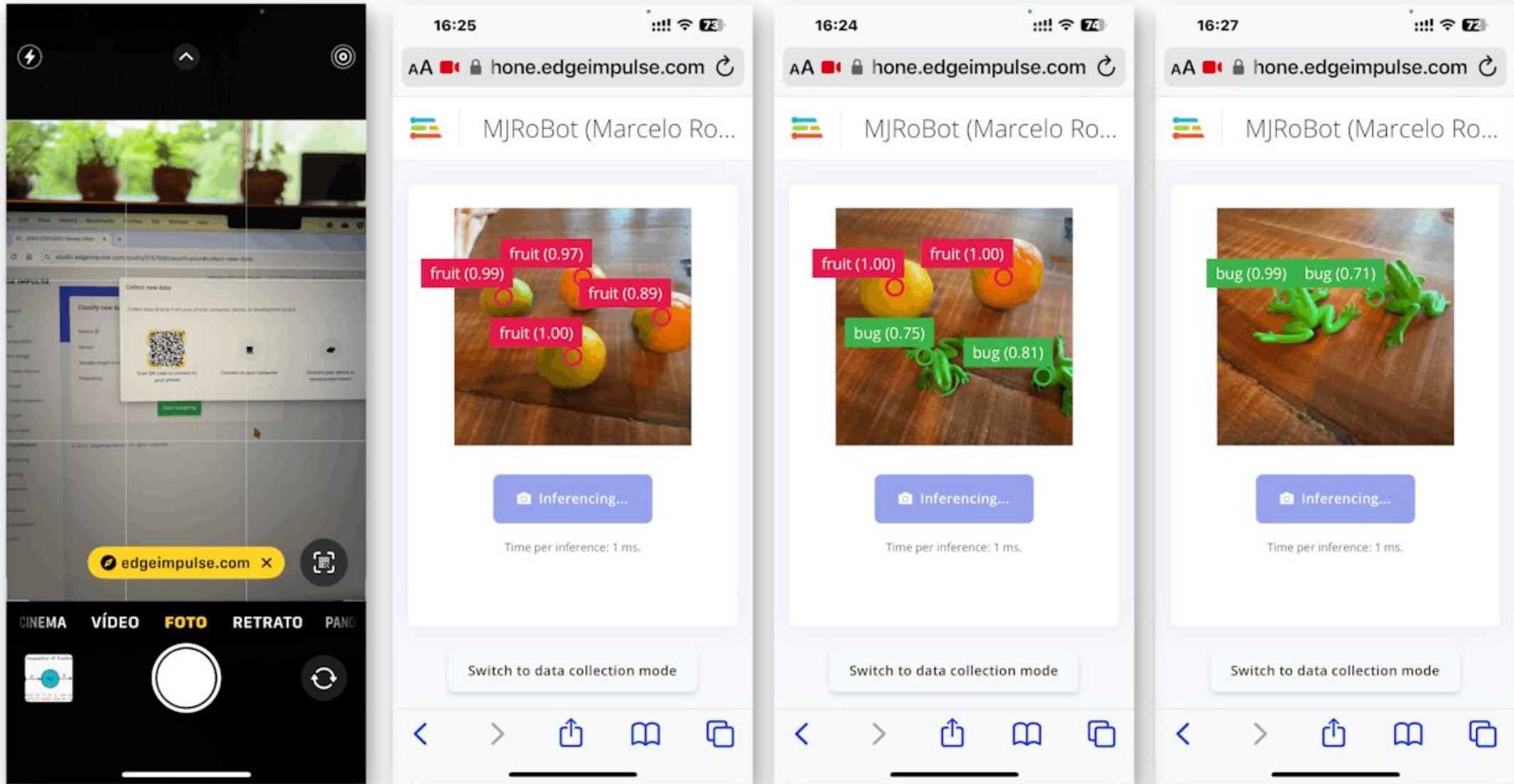
Add a learning block



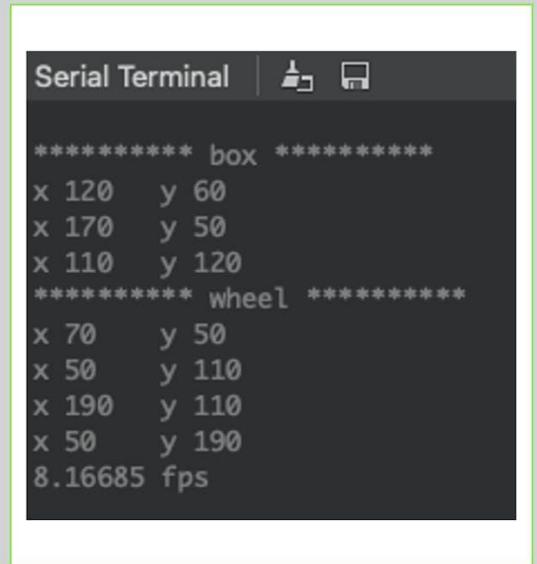
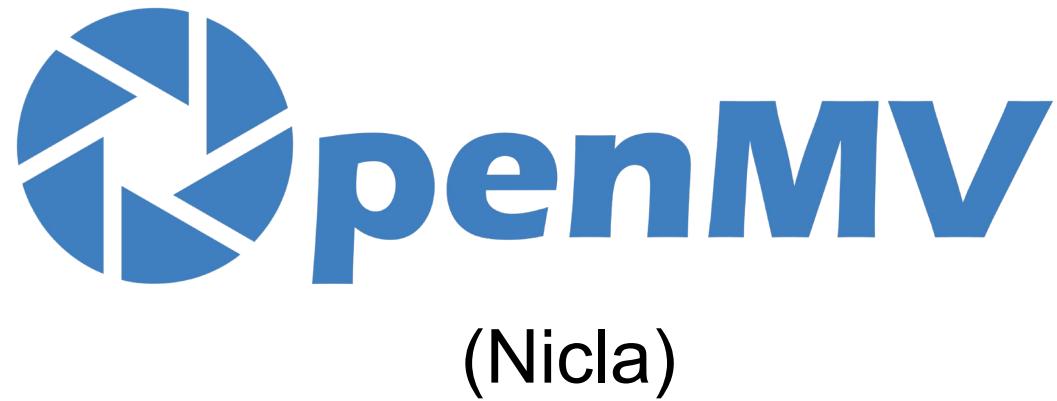
Model: FOMO



Inference Test



Deploy



Serial Terminal |  

```
***** box *****
x 120  y 60
x 170  y 50
x 110  y 120
***** wheel *****
x 70   y 50
x 50   y 110
x 190  y 110
x 50   y 190
8.16685 fps
```



[SenseCraft](#)
(XIAO)

Deploy

esp32_camera | Arduino IDE 2.2.1
XIAO_ESP32S3

```
18 * LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM,
19 * OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE
20 * SOFTWARE.
21 */
22
23 /* Includes -----
24 #include <XIAO-ESP32S3-Sense-Object_Detection_inferencing.h>
25 #include "edge-impulse-sdk/dsp/image/image.hpp"
26
27 #include "esp_camera.h"
28
29 // Select camera model - find more camera models in camera_pins.h file here
30 // https://github.com/espressif/arduino-esp32/blob/master/libraries/ESP32/examples/Camera/Camera
31
32 #define PWDN_GPIO_NUM      -1
33 #define RESET_GPIO_NUM     -1
34 #define XCLK_GPIO_NUM       10
35 #define SIOD_GPIO_NUM       40
36 #define STOC_GPIO_NUM       39
```

Serial Monitor Output

Message (Enter to send message to 'XIAO_ESP32S3' on '/dev/cu.usbmodem2101') Both NL & CR

```
fruit (0.566406) [ x: 56, y: 32, width: 8, height: 8 ]
Predictions (DSP: 4 ms., Classification: 143 ms., Anomaly: 0 ms.):
  No objects found
Predictions (DSP: 4 ms., Classification: 143 ms., Anomaly: 0 ms.):
  fruit (0.582031) [ x: 48, y: 32, width: 8, height: 8 ]
  fruit (0.773438) [ x: 80, y: 32, width: 8, height: 8 ]
Predictions (DSP: 4 ms., Classification: 143 ms., Anomaly: 0 ms.):
  fruit (0.550781) [ x: 64, y: 16, width: 8, height: 8 ]
Predictions (DSP: 4 ms., Classification: 143 ms., Anomaly: 0 ms.):
  fruit (0.722656) [ x: 64, y: 16, width: 8, height: 8 ]
```

Ln 48, Col 29 XIAO_ESP32S3 on /dev/cu.usbmodem2101



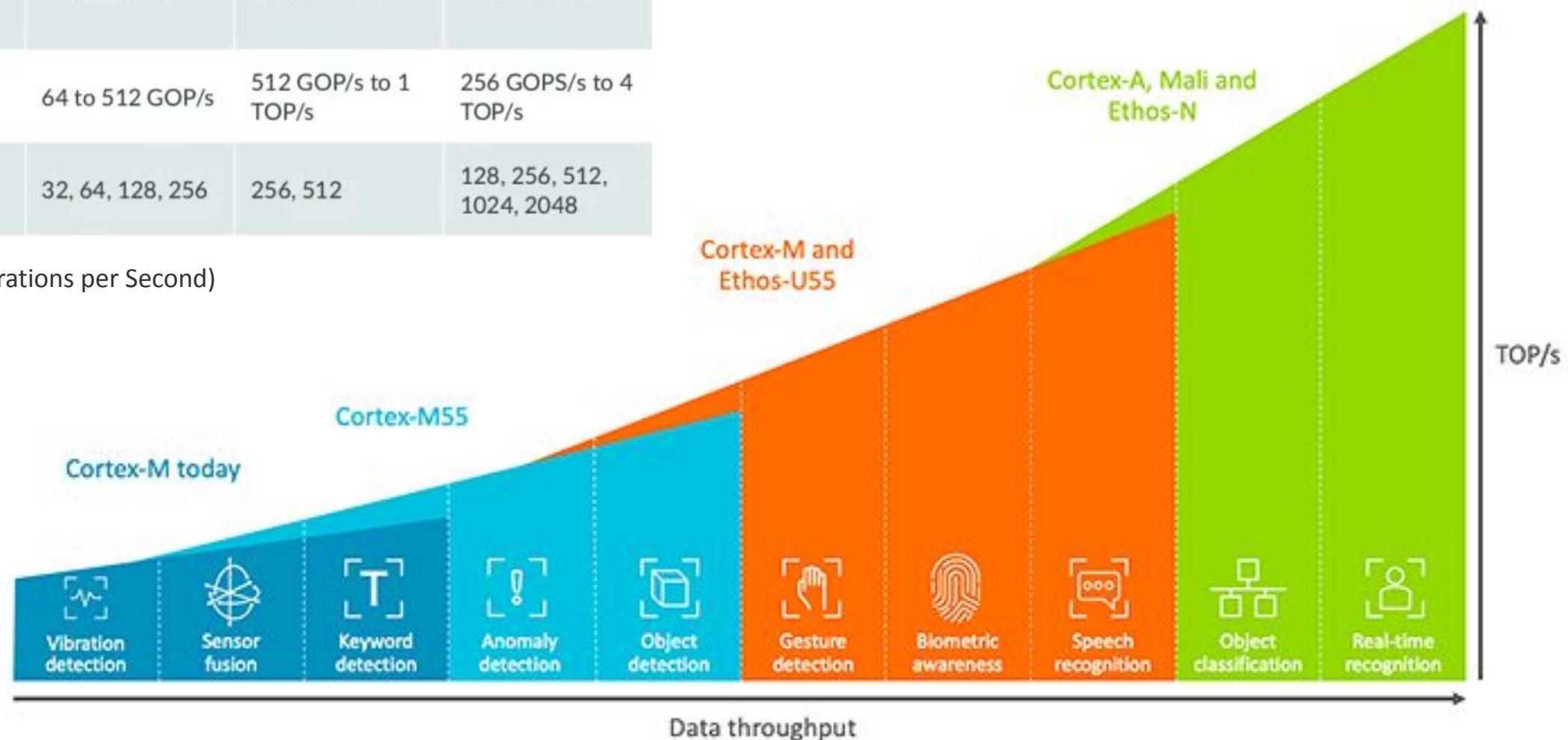
microNPU

a neural network unit for TinyML

ML- optimized Solutions

| | Ethos-U55 | Ethos-U65 | Ethos-U85 |
|---------------------------|------------------|----------------------|---------------------------|
| Performance (At 1 GHz) | 64 to 512 GOP/s | 512 GOP/s to 1 TOP/s | 256 GOPS/s to 4 TOP/s |
| MACs (8x8) | 32, 64, 128, 256 | 256, 512 | 128, 256, 512, 1024, 2048 |

TOPS (Tera Operations per Second)

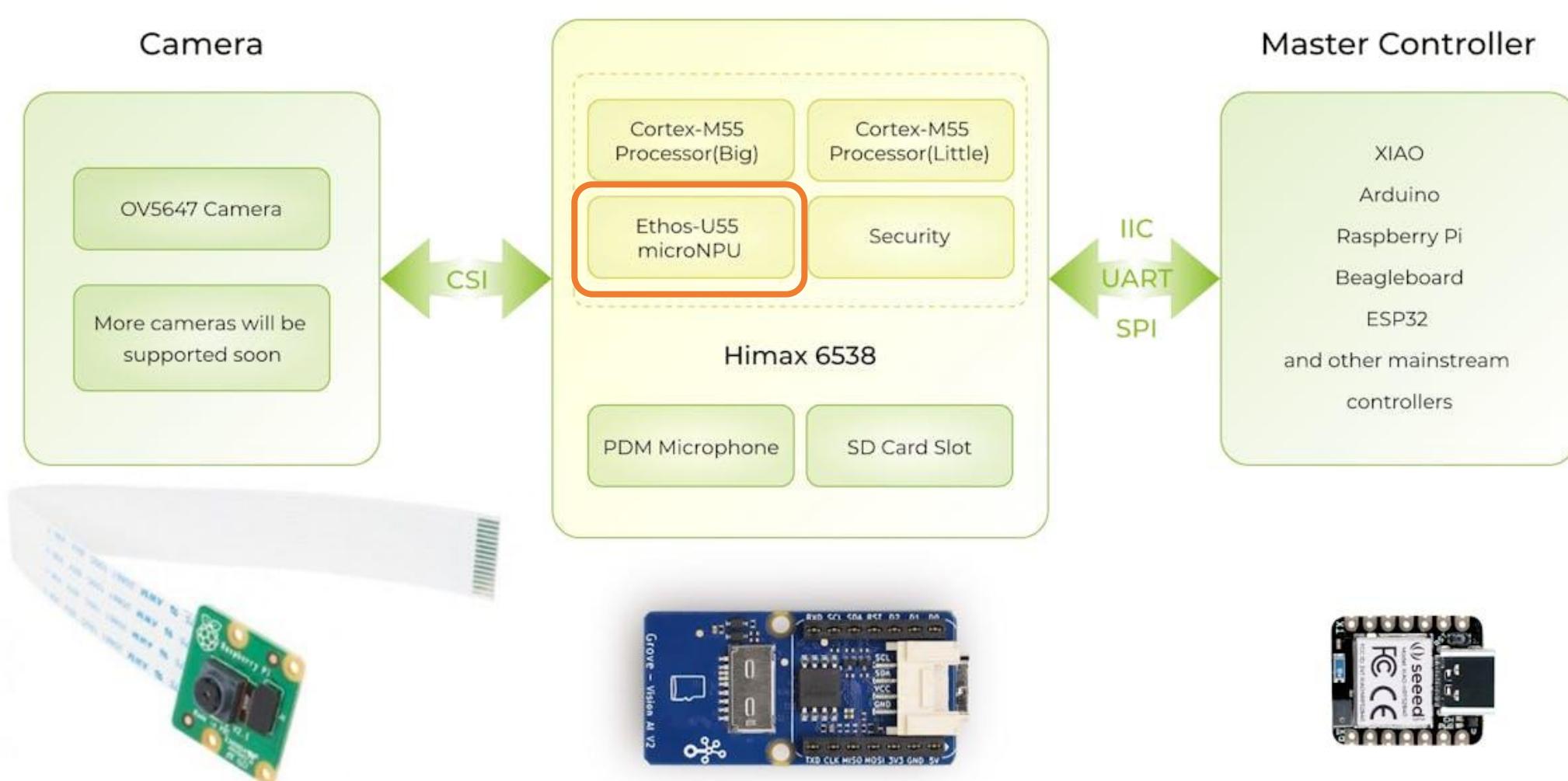


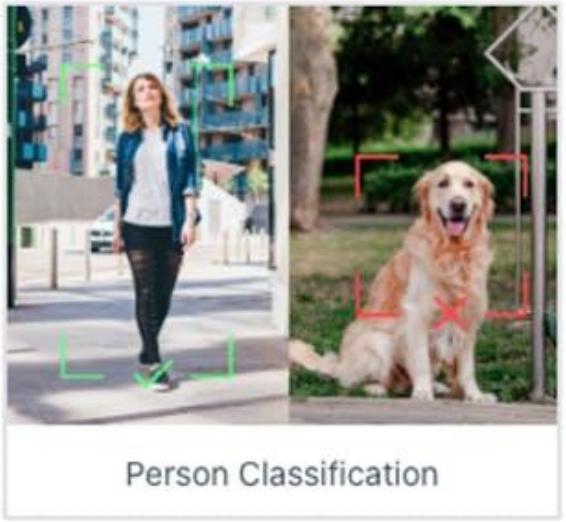


Computer Vision at the Edge with Grove Vision AI Module V2

Exploring Computer Vision applications such as Image Classification, Object Detection, and Pose estimation.

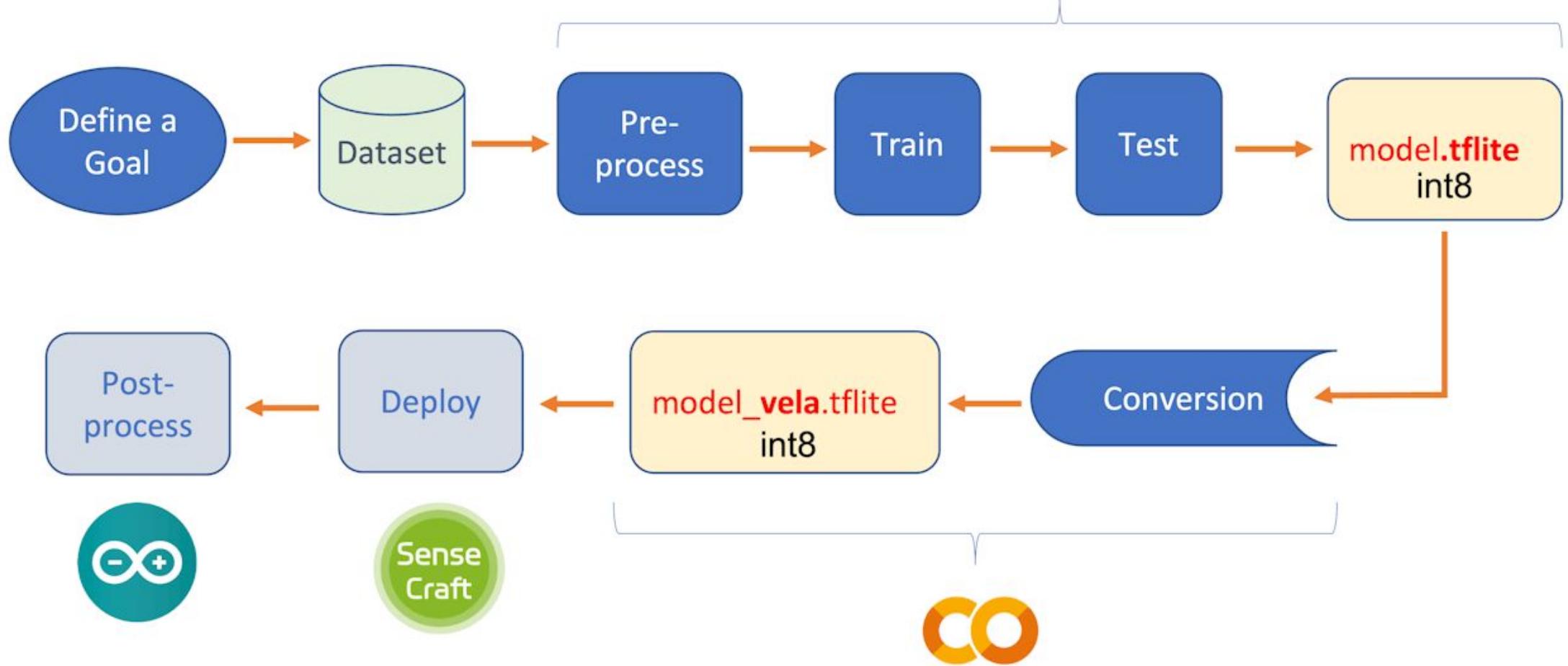
Grove Vision AI v2







EDGE IMPULSE





Classification: 687 ms

1.5 FPS



ESP - CAM
Xtensa LX6
240 MHz

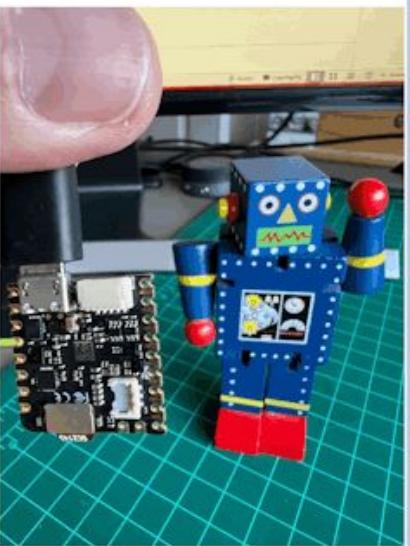


Classification: 142 ms

7.0 FPS



XIAO ESP32S3
Xtensa LX7
240 MHz

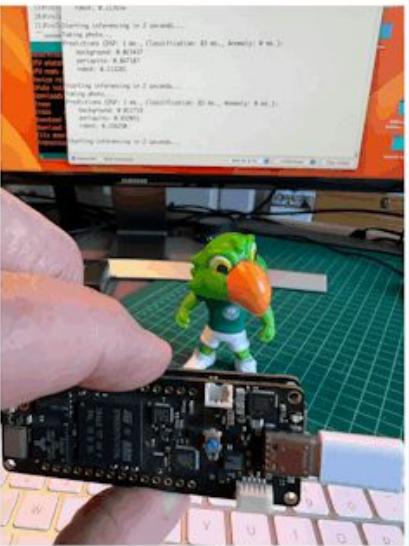


Classification: 86 ms

11.6 FPS



Nicla-Vision
ARM M7
480 MHz



Classification: 83 ms

12.0 FPS



Portenta H7
ARM M7
480 MHz



Classification: 6 ms

167 FPS

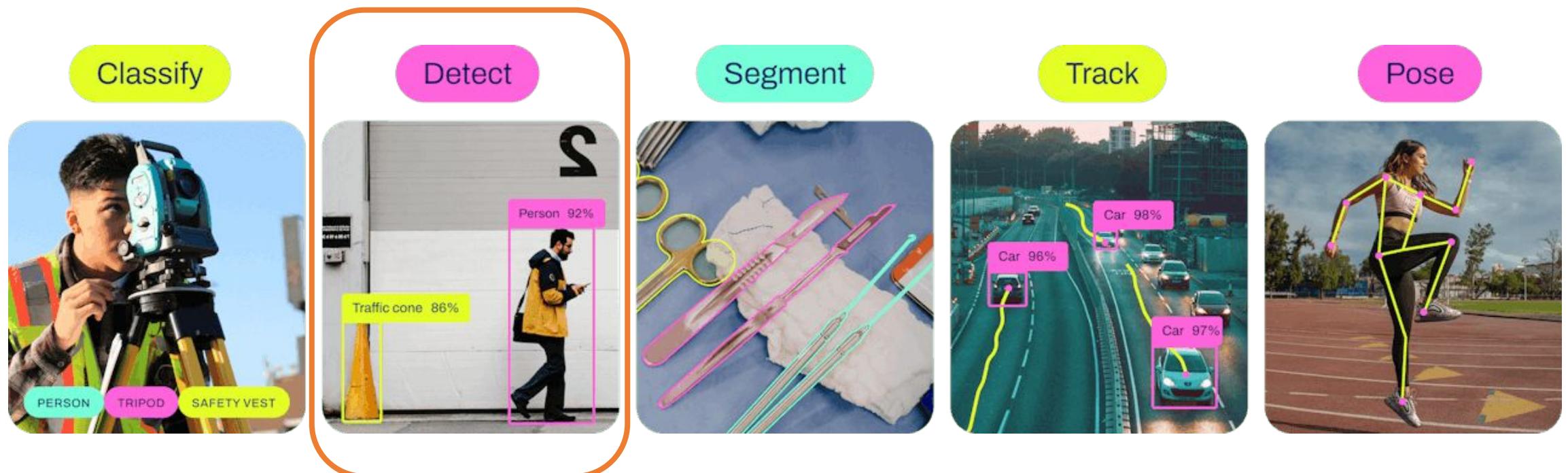


Grove Vision AI V2
ARM Ethus-U55
400 MHz

YOLO

Object Detection model

Ultralytics YOLO (You Only Look Once)



Real-time **object detection** systems that identify and classify many objects **very fast** in a single image pass.



BuzzTech: Machine Learning at the Edge

Deploying YOLOv8 on Raspberry Pi Zero 2W for Real-Time Bee Counting at the Hive Entrance.

Goal: Estimate the number of bees

Number of objects: 15 bees

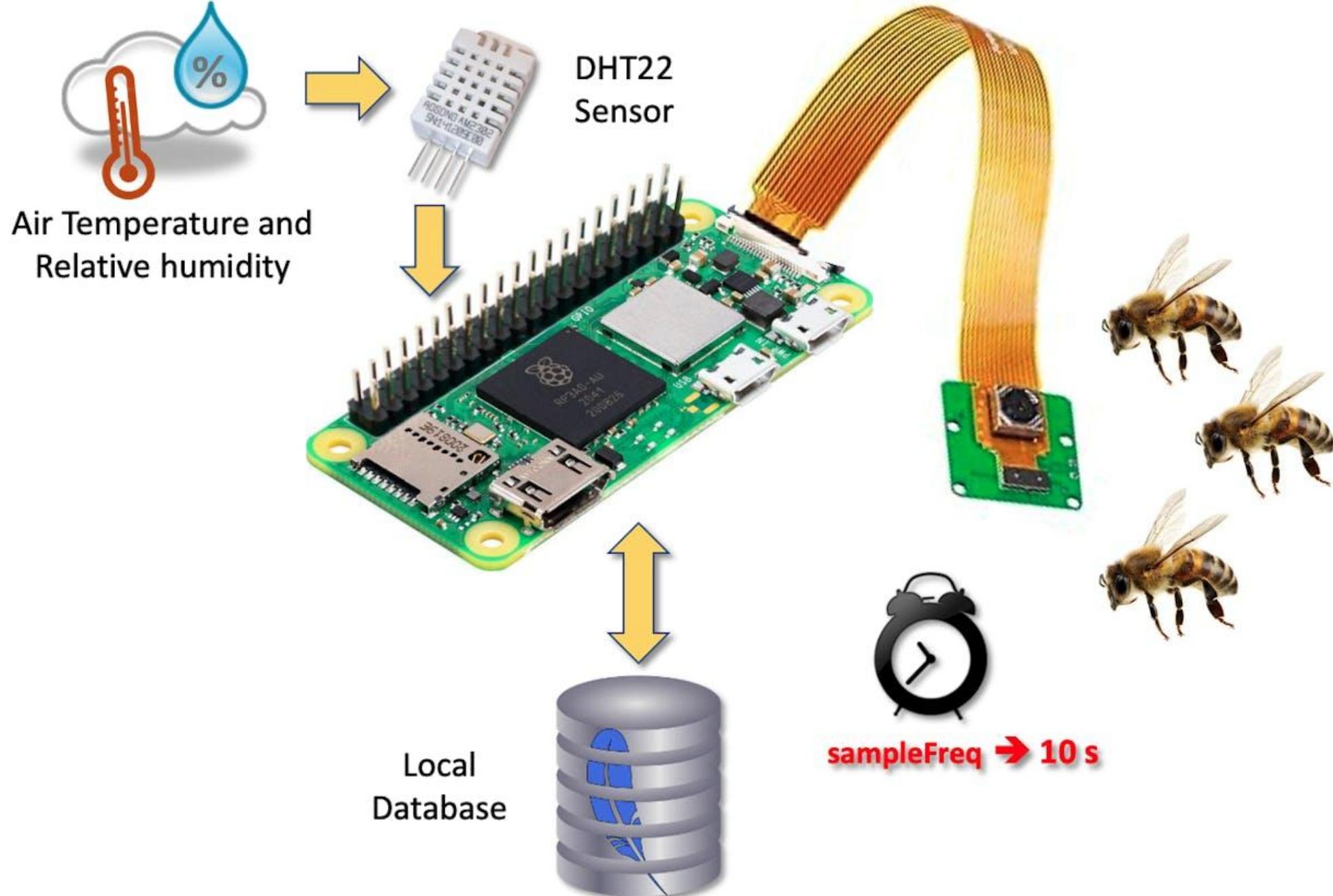


Number of objects: 36 bees



Number of objects: 28 bees





Create Project

https://app.roboflow.com/marcelo-rovali-riila/create

roboflow

Let's create your project.

Marcelo Rovai > New Public Project

Project Name: Bees_on_Hive_Landing_boards

License: CC BY 4.0

Annotation Group: bees

Project Type:

- Object Detection**: Identify objects and their positions with bounding boxes.
Best For: # Counting, % Tracking
- Classification**: Assign labels to the entire image.
Classification Type: Multi-Label (selected), Single-Label
Best For: Filtering, Content Moderation
- Instance Segmentation**: Detect multiple objects and their actual shape.
Best For: Measurements, Odd Shapes
- Keypoint Detection**: Identify keypoints ("skeletons") to subjects.
Best For: Pose Estimation

Show More ↓

Cancel Create Public Project

Bees_on_Hive_landing_board

https://app.roboflow.com/marcelo-rovali-riila/bees_on_hive_landing_boards/images/34IC4TuHjk5VtUNSKxC?queryText=&pageSize=50&startingIndex=0&browseQuer... ☆

BEES_ON_HIVE_LANDING_BOARDS > ANNOTATE
20230711b6510.jpg

Annotations
Group: bees-4uet
CLASSES LAYERS
bee 26

Annotation Editor
bee
Delete Save (Enter)

1 bee

Options ▾

Labels

Attributes

Comments

History

Raw Data

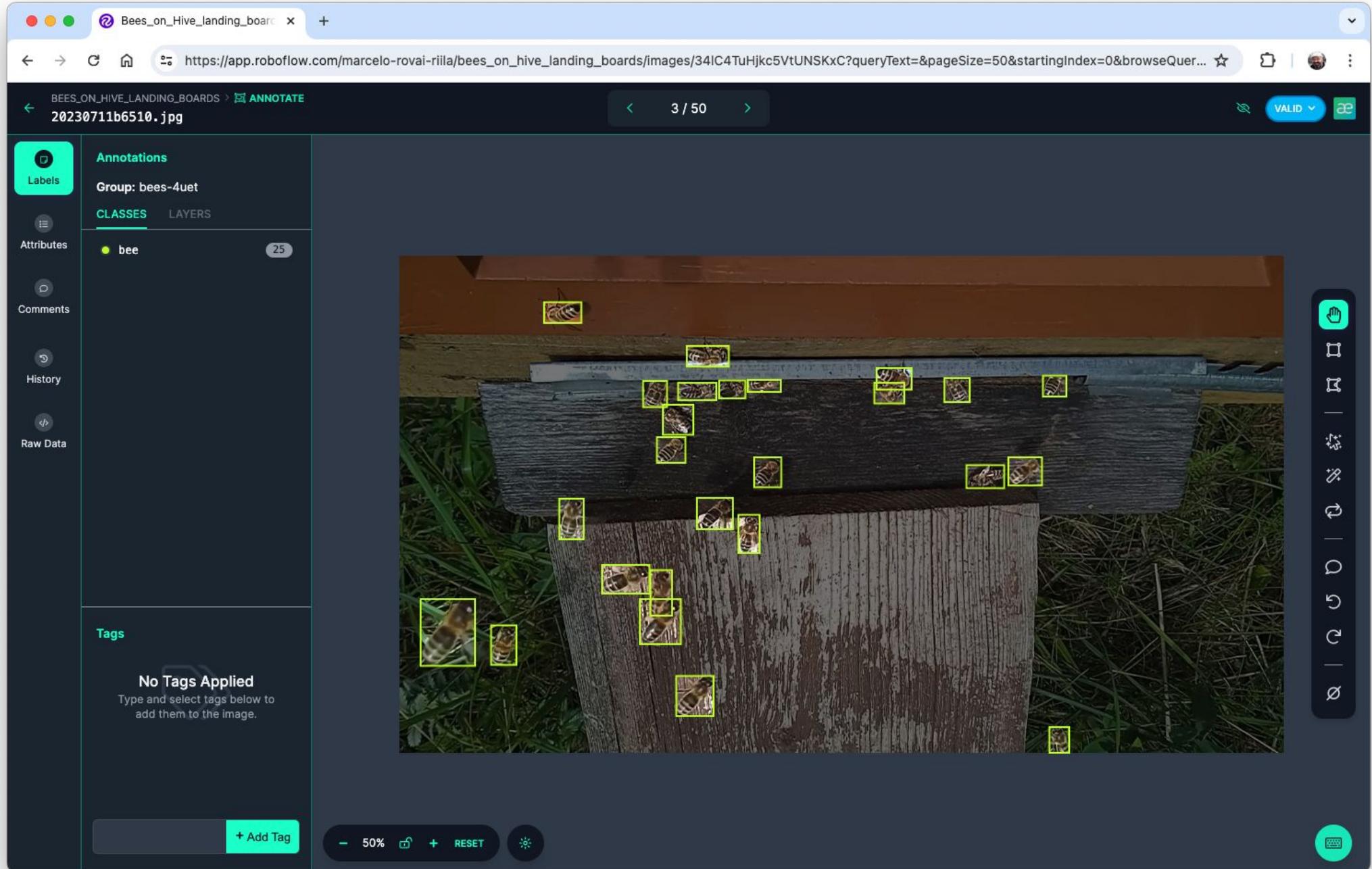
Tags

No Tags Applied
Type and select tags below to add them to the image.

+ Add Tag

- 50% + RESET

VALID ae



yolov8_beans_on_hive_landing_board.ipynb

All changes saved

Comment Share Gemini

RAM Disk Gemini

Files

content

datasets

Bees_on_Hive_landing_board

test

images

labels

train

images

labels

valid

images

labels

README.dataset.txt

README.roboflow.txt

data.yaml

databricks

dev

etc

Disk 46.44 GB available

+ Code + Text

```
[ ] 1 from ultralytics import YOLO  
2  
3 from IPython.display import display, Image
```

Dataset

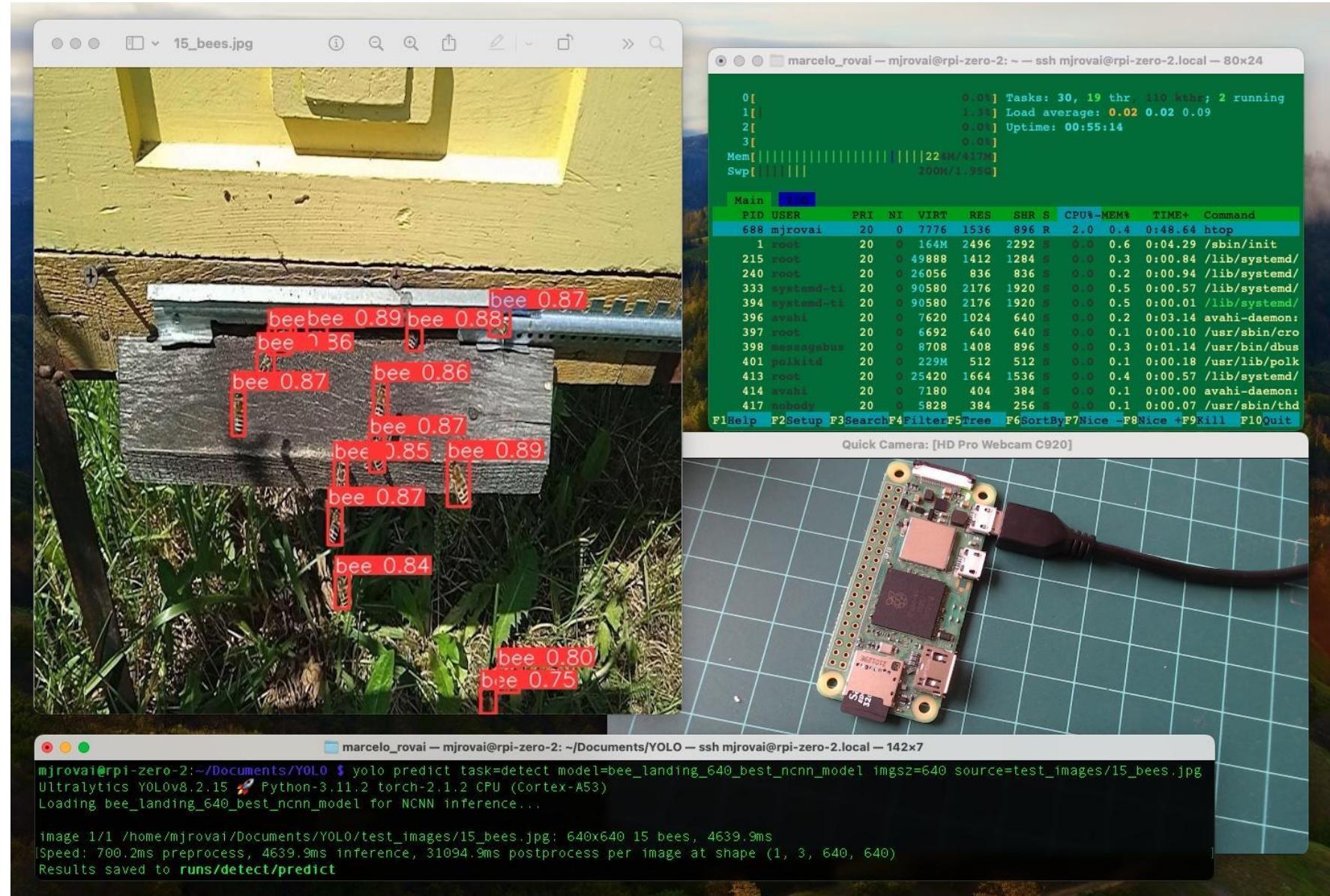
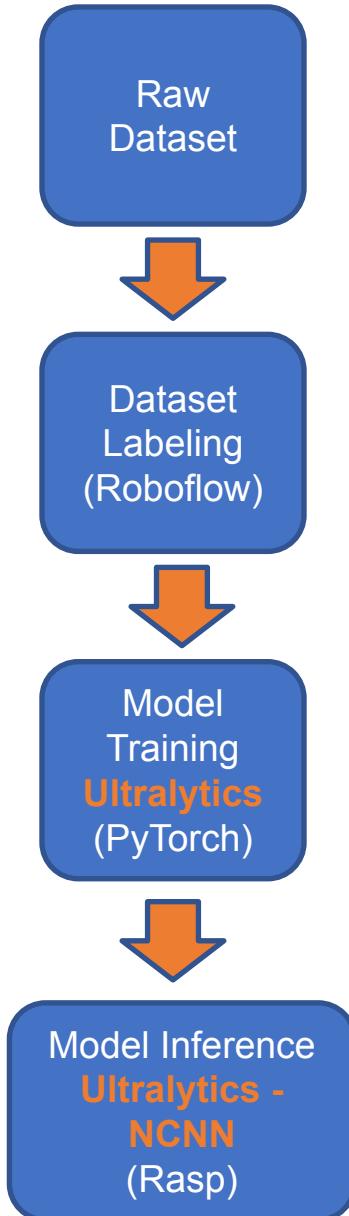
```
1 !mkdir {HOME}/datasets  
2 %cd {HOME}/datasets  
3  
4 !pip install roboflow --quiet  
5  
6 from roboflow import Roboflow  
7 rf = Roboflow(api_key="YOUR KEY HERE")  
8 project = rf.workspace("marcelo-rovai-riila").project("bees_on_hive_landing_boards")  
9 version = project.version(1)  
10 dataset = version.download("yolov8")  
11
```

/content/datasets

```
75.5/75.5 kB 3.6 MB/s eta 0:00:00  
158.3/158.3 kB 7.7 MB/s eta 0:00:00  
178.7/178.7 kB 8.3 MB/s eta 0:00:00  
58.8/58.8 kB 6.9 MB/s eta 0:00:00  
49.1/49.1 MB 16.7 MB/s eta 0:00:00  
54.5/54.5 kB 7.1 MB/s eta 0:00:00
```

loading Roboflow workspace...
loading Roboflow project...
Dependency ultralytics==8.0.196 is required but found version=8.2.23, to fix: 'pip install ultralytics==8.0.196'
Downloading Dataset Version Zip in Bees_on_Hive_landing_boards-1 to yolov8:: 100%|██████████| 1597328/1597328 [00:09<00:00]
Extracting Dataset Version Zip to Bees_on_Hive_landing_boards-1 in yolov8:: 100%|██████████| 32468/32468 [00:09<00:00]

50s completed at 8:46AM



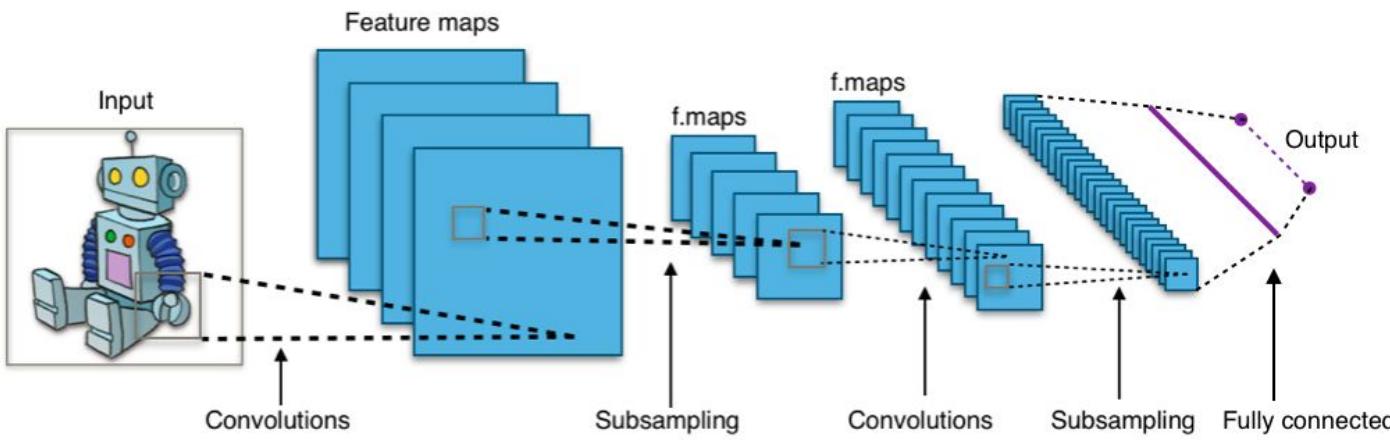
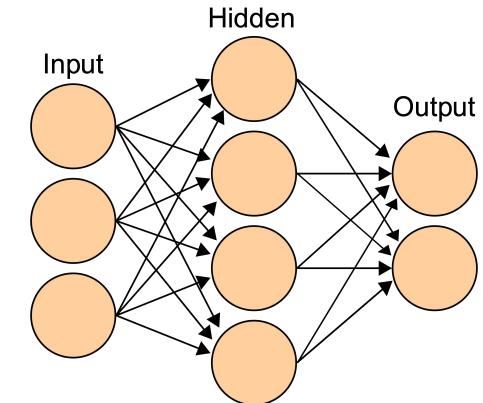
LSTM

Long Short-Term Model

*LSTM model is a type of recurrent neural network (RNN) that is well-suited for **sequence prediction** problems by effectively capturing long-term dependencies in data sequences.*

Deep Learning models (or artificial neural networks)

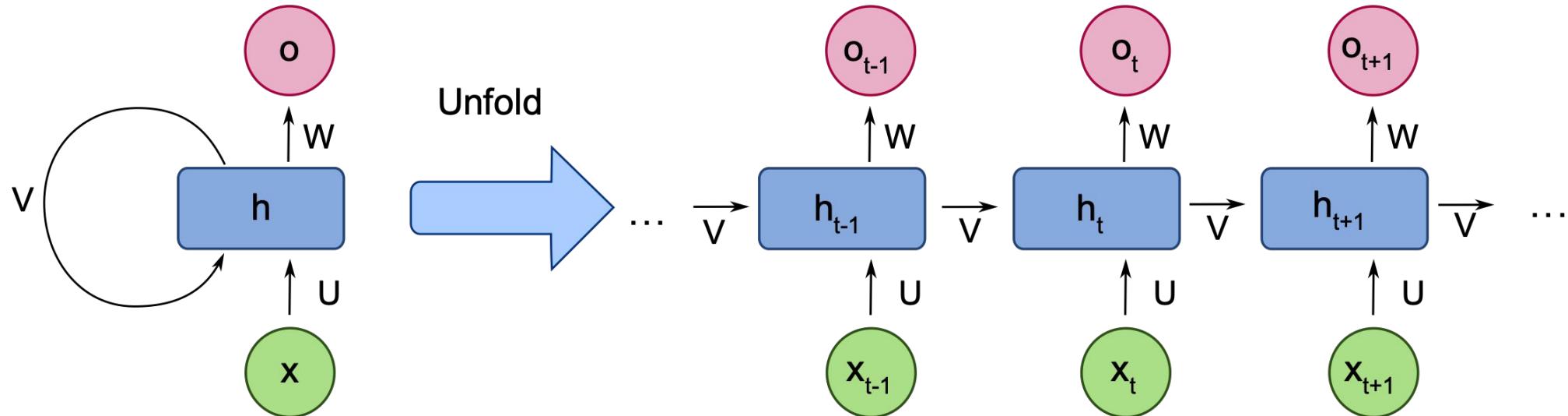
Fully Connected Neural Networks (FCNNs): Networks where **each neuron in one layer is connected to every neuron in the following layer**, useful for complex pattern recognition across diverse datasets.



Convolutional Neural Networks (CNNs): Specialized for **grid-like data such as images**, using convolutional layers to detect and learn spatial hierarchies of features.

Deep Learning models (or artificial neural networks)

Recurrent Neural Networks (RNNs): Designed for **sequential data like time series or text**, these networks use their internal state (memory) to process sequences of inputs.



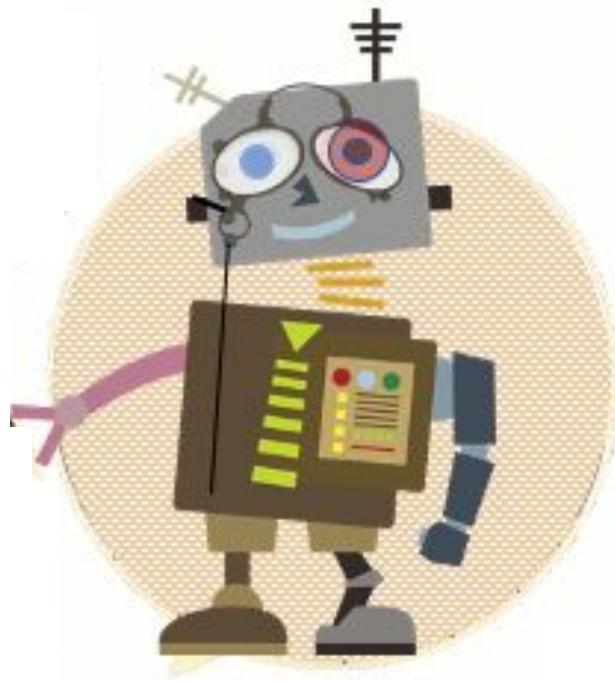
Machado de Assis Bot with RNN - GRU



The robot writer model is a Recurrent Neural network (RNN —GRU). To obtain the final AI model, 3.5 million parameters were trained with a **120-letter sequence** from seven of his books: *Memorias Posthumas de Braz Cubas*, *Dom Casmurro*, *Quincas Borba*, *Papeis Avulsos*, *A Mão e a Luva*, *Esaú e Jacob*, and *Memorial de Ayres*.

Model: "sequential"

| Layer (type) | Output Shape | Param # |
|-----------------------|-------------------|---------|
| embedding (Embedding) | (128, None, 64) | 7488 |
| gru (GRU) | (128, None, 1026) | 3361176 |
| dense (Dense) | (128, None, 117) | 120159 |
| Total params: | 3,488,823 | |
| Trainable params: | 3,488,823 | |
| Non-trainable params: | 0 | |



A LUVA DE CASMURRO II

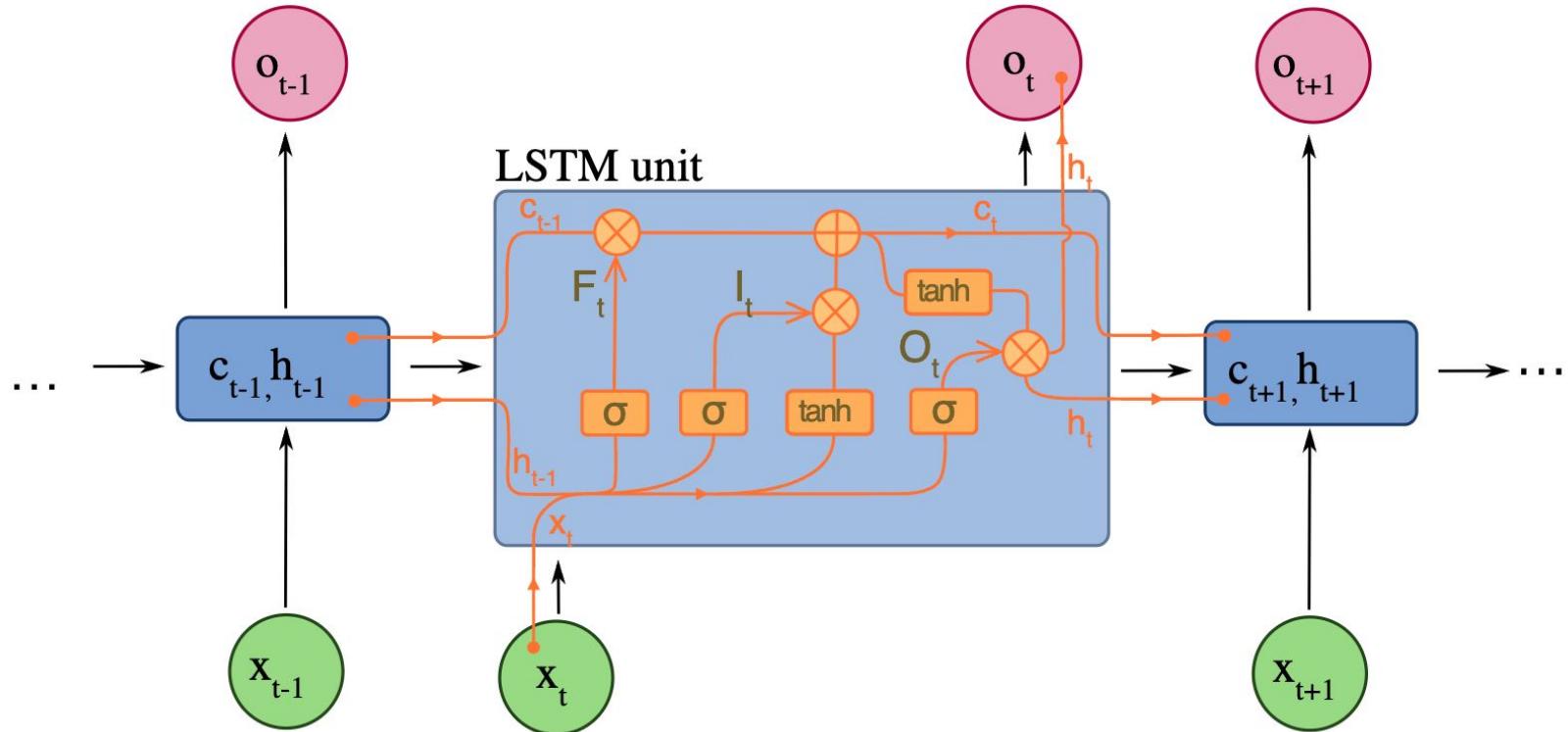
A missa do coupé e um presente e o governo devia cazar logo no papel, a morte do autor, e todos os seus considerados de alegria. Era um espirito de vinte e cinco annos, e eu não estou alguns passos no cerebro, como de outra cousa. Deus me disse:

--*Não digo que não. Se eu tivesse a intenção de um proboso. Palha acudiu a mulher, não havia nada. A noite vinha tambem para o seminario, tinha o aspecto do partido recto e de restaurar a minha mãe e do pae, pela primeira vez, a menor destinada a dispensar o chapéo, esperou que não vinhas com as suas mãos de creanças. A manhã della chegasse a baroneza e a maneira desta divida. Parece que é casada.*

--*Está bom, perdoa-lhe de todos os lados, a vida de que o comprar para o meu quarto de hora, e contavam com o fim de a anterior, e, a parede pouco tempo a alma de pessoas que definitivamente lhe interessam a menos para mim. De quando em quando, esses dous annos de conversação para o fim de deixar nenhuma pessoa que se dispersasse; mas não falo de uma cousa nem lhe pedia com a mão tremula, como se ella quizesse. Eu, apertando-lhe a mão, aliás o principio do governo, a proposito disso, com a desattenção de Estevão, e eu começo a aborrecel-o, e a solidão podia ser melhor, e a sympathia coloca da mãe, e não se sabe calar o enterro no meio do lagem, o que iam-se apanhados no chão, e para a mulher, não tendo visto, nem a mesma cousa.*

Deep Learning models (or artificial neural networks)

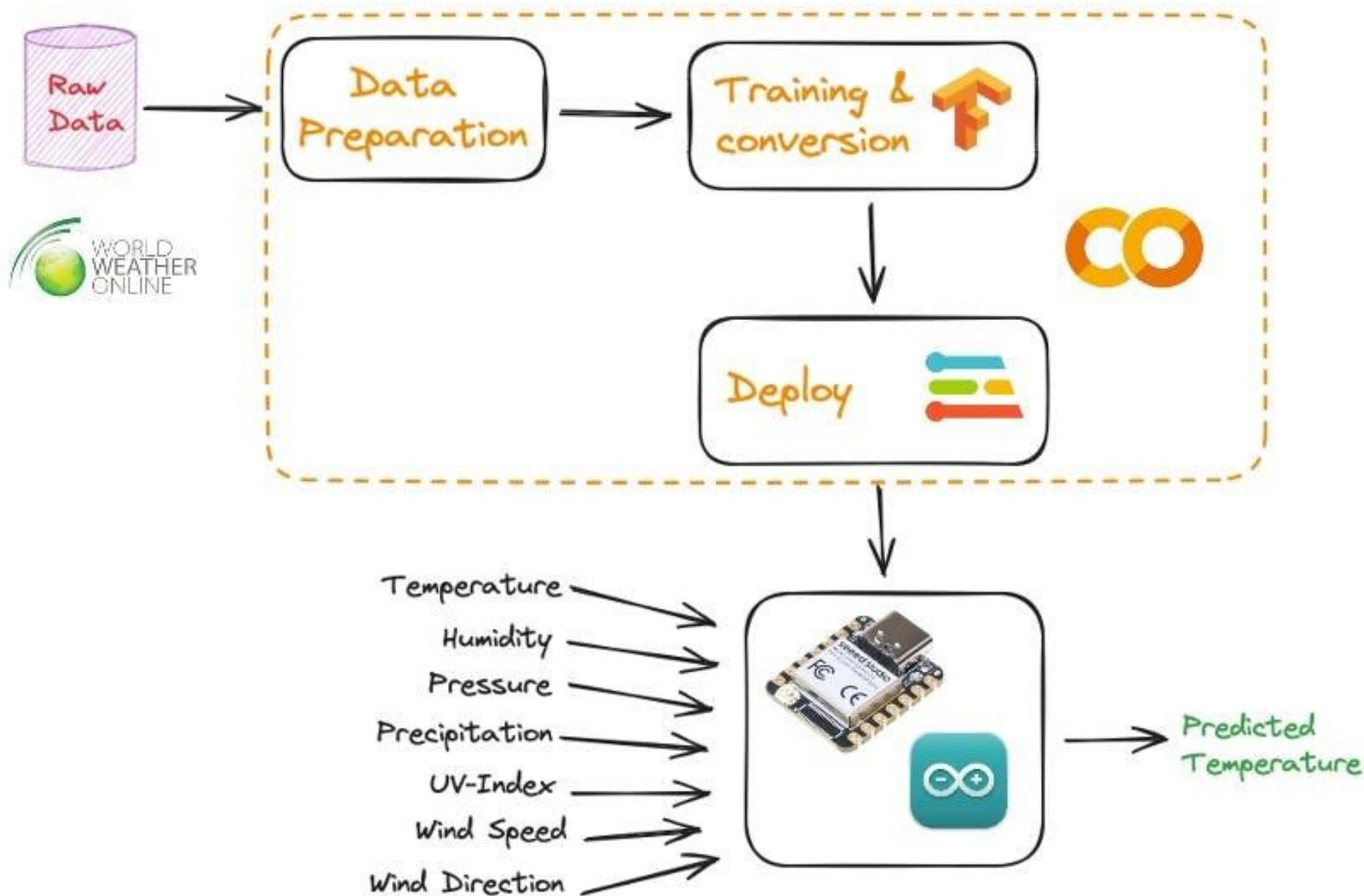
Long Short-Term Memory (LSTM): A type of RNN that can learn over long sequences without losing information, effectively managing long-term dependencies.

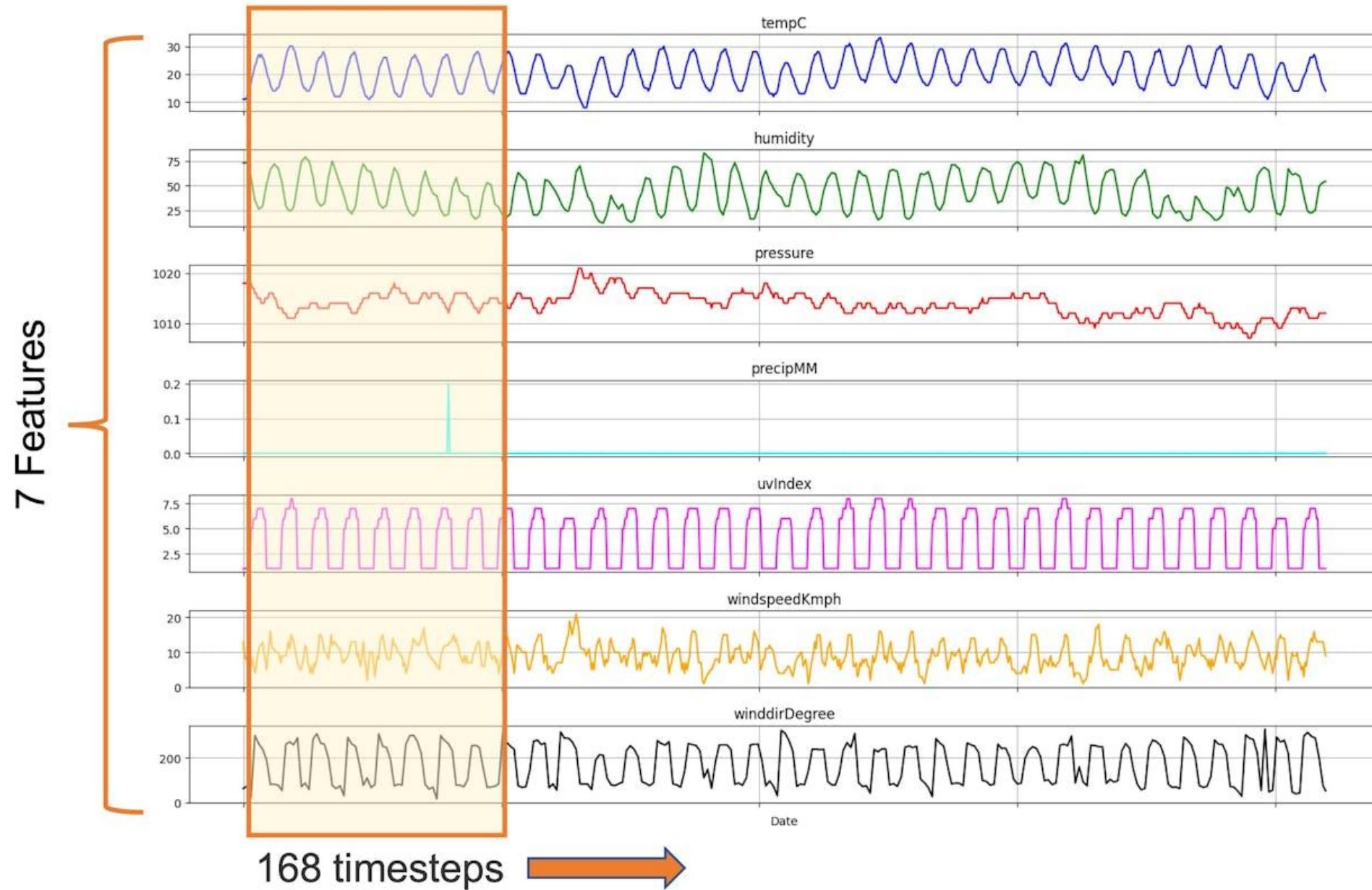


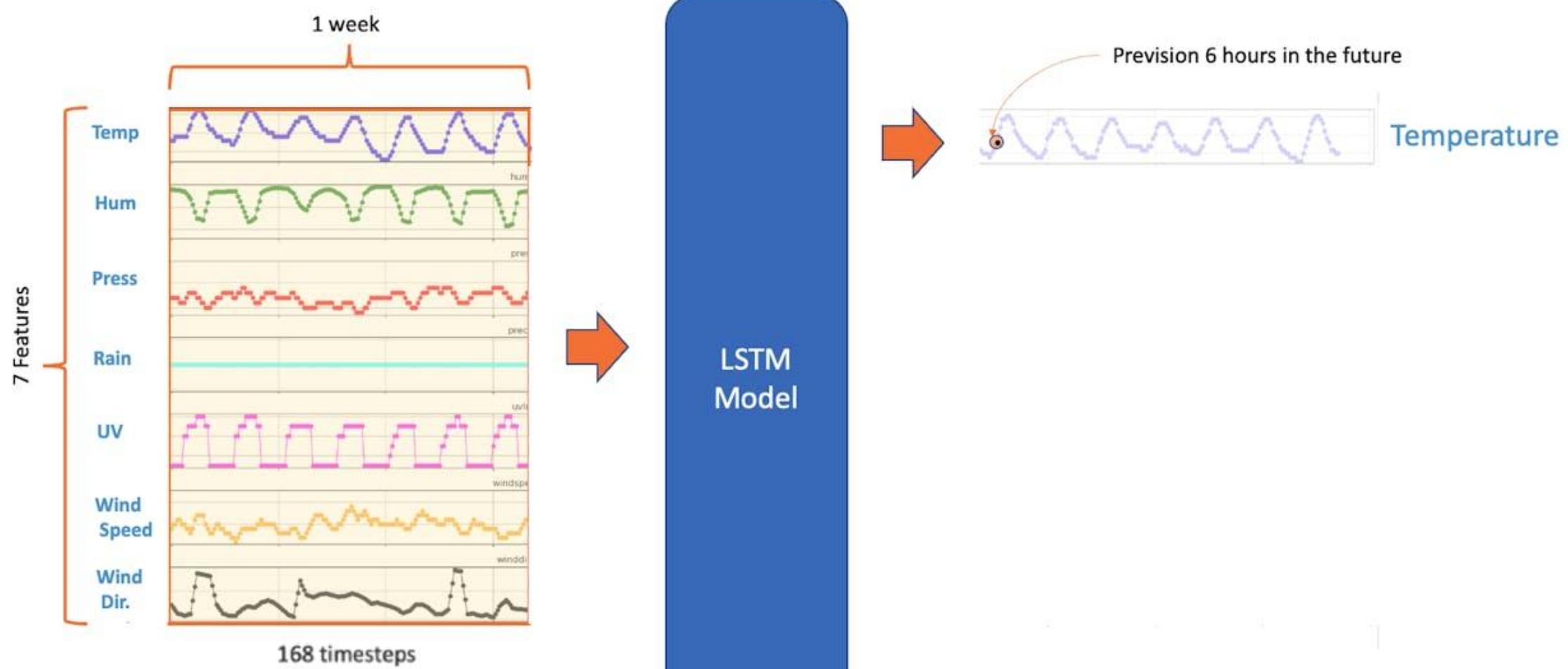


Temperature Prediction using an LSTM model

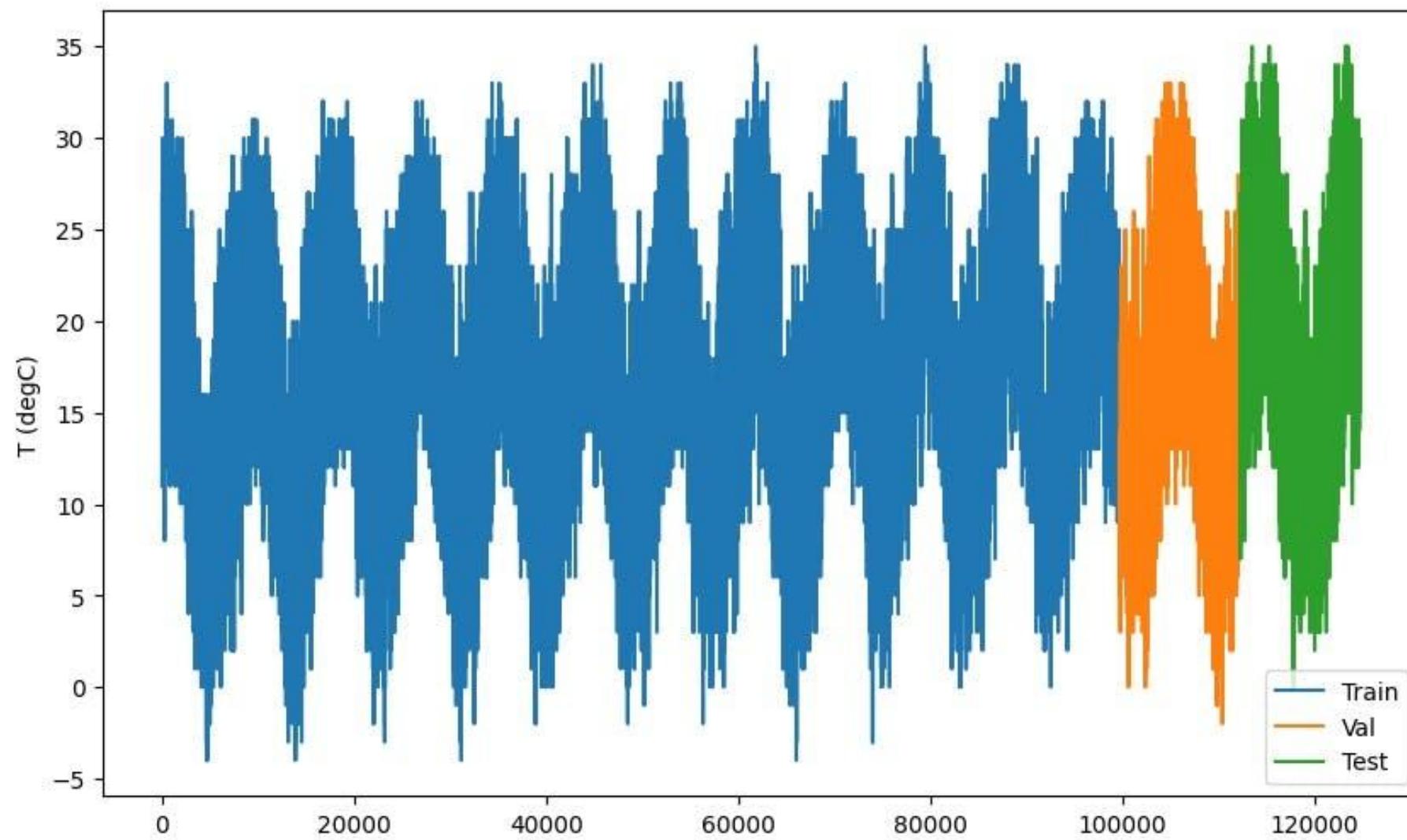
The trained LSTM model will be converted with TFLite-Micro and Edge Impulse Python SDK and deployed on an XIAO ESP32S3.







Temperature



```
model = Sequential([
    LSTM(128,
        input_shape=(n_steps, X_train.shape[2])),
    Dense(1)
])
```

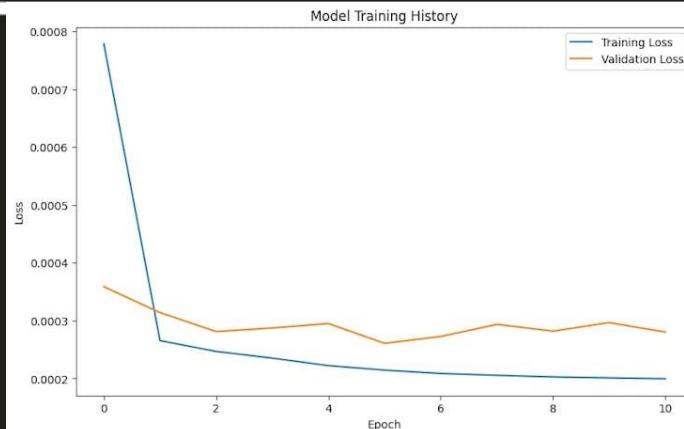
```
model.compile(optimizer='adam', loss='mse')
```

```
history = model.fit(
    X_train, y_train,
    validation_data=(X_val, y_val),
    epochs=20,
    batch_size=32,
    callbacks=[early_stopping]
)
```

```
prediccion = model.predict(X_test)
```

```
converter = tf.lite.TFLiteConverter.from_saved_model(MODEL_DIR)
tflite_model = converter.convert()
```

```
# Save the converted model to file
tflite_model_file = 'converted_model.tflite'
with open(tflite_model_file, 'wb') as f:
    f.write(tflite_model)
```



Time-Series
Dataset



Feature
Extraction

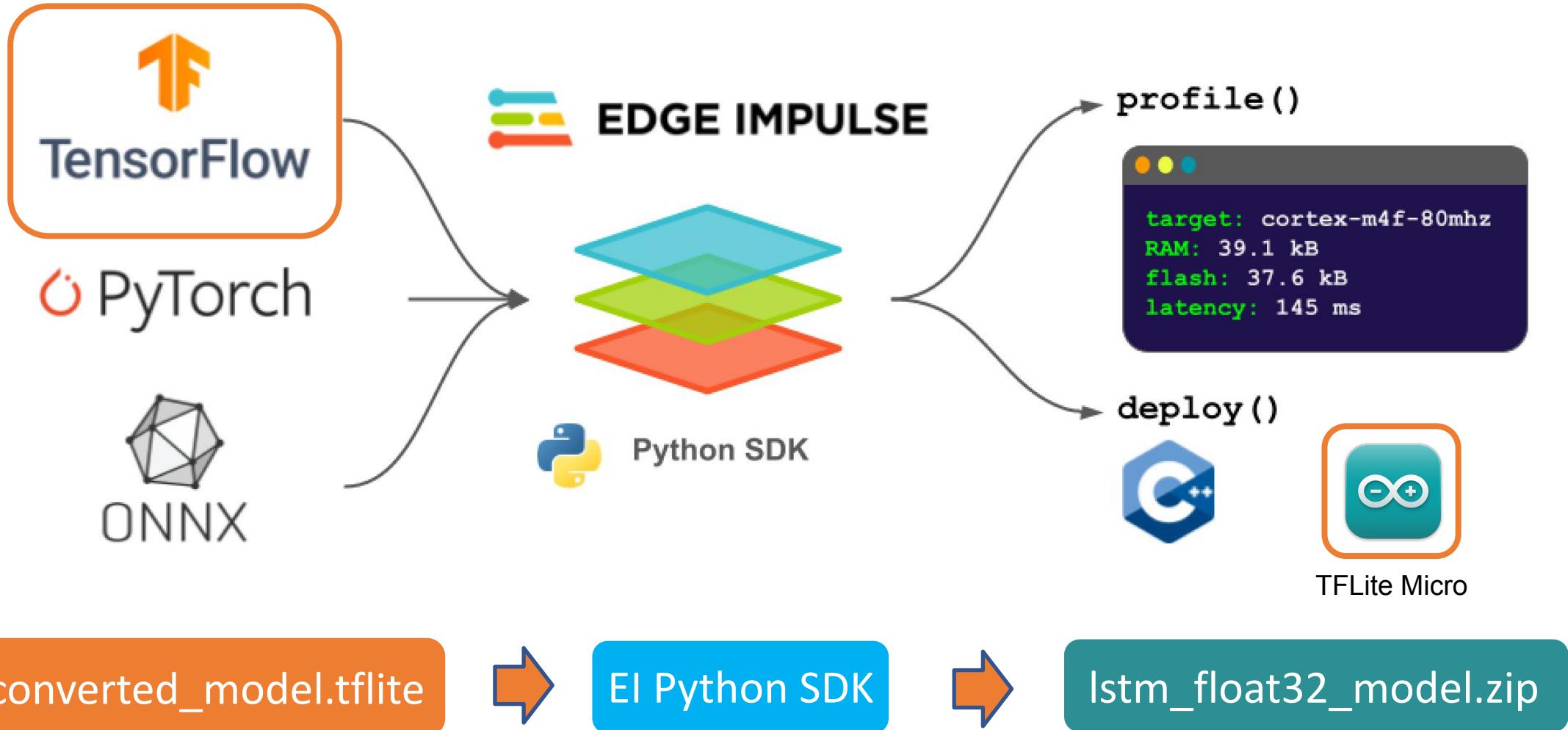


Model
Training
(TensorFlow)

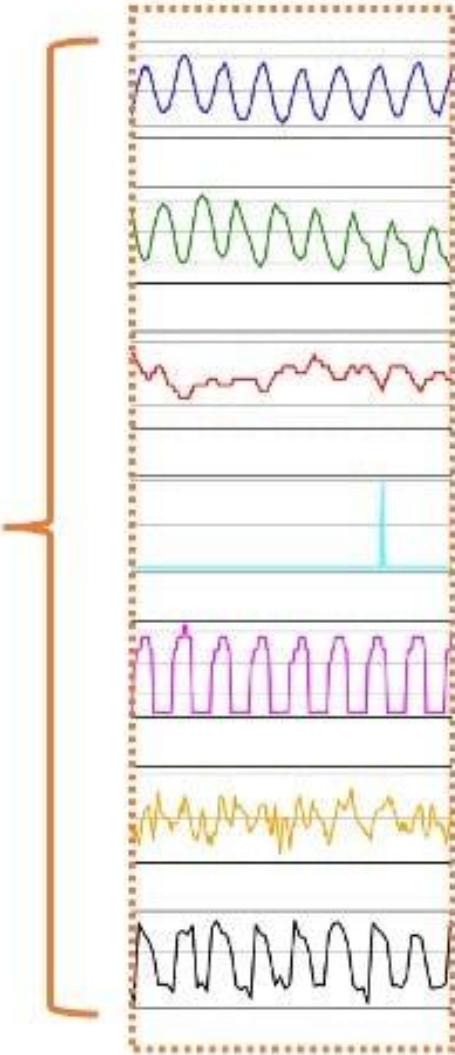


Model
Conversion
(TFLite)

Edge Impulse Python SDK



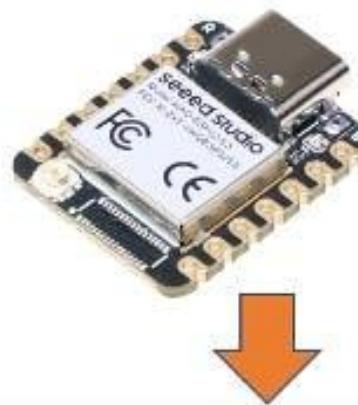
7 Features



168 timesteps

[0.38461538, 0.48979592, 0.4516129, 0., 0., 0.
0.04347826, 0.47777778, 0.38461538, 0.48979592, 0.4516129,
0., 0., 0., 0., 0.41666667, 0.41025641,
0.5, 0.4516129, 0., 0., 0.375, 0.,
0.35555556, 0.46153846, 0.43877551, 0.4516129, 0.,
0.375, 0.04347826, 0.46666667, 0.51282051, 0.37755182,
0.4516129, 0., 0.5, 0.13043478, 0.575,
0.58974359, 0.31632653, 0.4516129, 0., 0.5,
0.17391304, 0.68611111, 0.61538462, 0.28571429, 0.4516129,
0., 0.625, 0.2173913, 0.65277778, 0.66666667,
0.26530612, 0.4516129, 0., 0.625, 0.30434783,
0.61666667, 0.71794872, 0.24489796, 0.41935484, 0.,
0.625, 0.34782609, 0.58333333, 0.74358974, 0.2244898,
0.41935484, 0., 0.625, 0.47826087, 0.57777778,
0.76923077, 0.21428571, 0.41935484, 0., 0.75,
0.68869565, 0.575, 0.74358974, 0.28408163, 0.38789677,
0., 0.75, 0.73913843, 0.57222222, 0.69230769,
0.2244898, 0.41935484, 0., 0.625, 0.69565217,

1,176 Features



Edge Impulse Inferencing Demo:
Edge Impulse standalone inferencing (Arduino)
run_classifier returned: 0
Timing: DSP 0 ms, inference 2024 ms, anomaly 0 ms
Predictions:
value: 0.36065

static_buffer | Arduino IDE 2.3.2

XIAO_ESP32S3

static_buffer.ino

```
16
17  /* Includes */
18 #include <LoBa_Temp_Prediction_-_LSTM_inferencing.h>
19
20 static const float features[] = {
21     0.38461538, 0.48979592, 0.4516129 , 0.        , 0.        ,
22     0.04347826, 0.47777778, 0.38461538, 0.48979592, 0.4516129 ,
23     0.        , 0.        , 0.        , 0.41666667, 0.41025641,
24     0.5        , 0.4516129 , 0.        , 0.375        , 0.        ,
25     0.35555556, 0.46153846, 0.43877551, 0.4516129 , 0.        ,
26     0.375        , 0.04347826, 0.46666667, 0.51282051, 0.37755102,
27     0.4516129 , 0.        , 0.5        , 0.13043478, 0.575        ,
28     0.58974359, 0.31632653, 0.4516129 , 0.        , 0.5        ,
29     0.17391304, 0.68611111, 0.61538462, 0.28571429, 0.4516129 ,
30     0.        , 0.4725      , 0.7172012 , 0.65777779, 0.66666667
```

Output Serial Monitor X

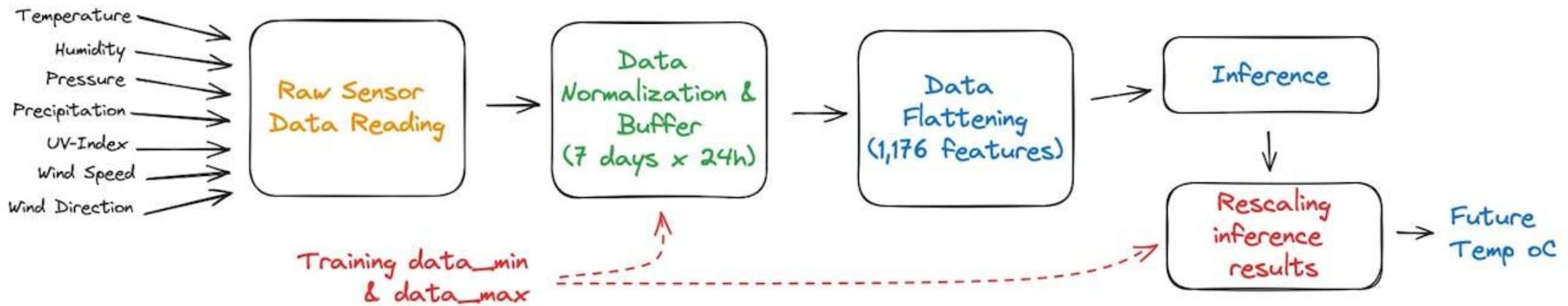
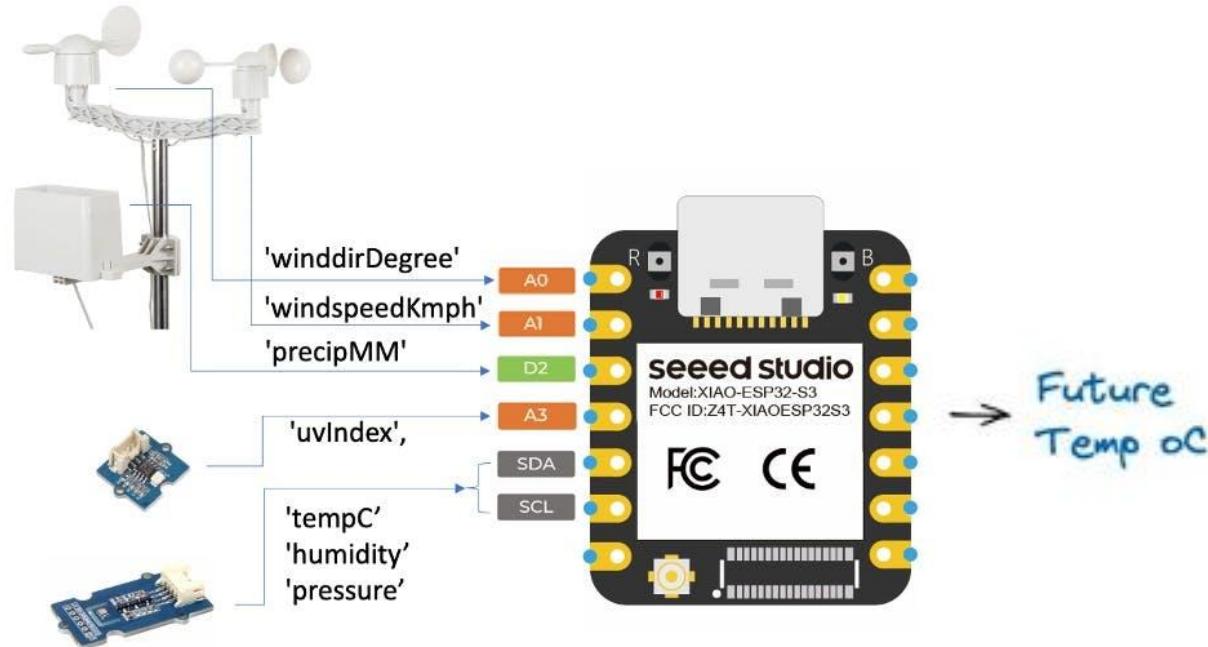
Message (Enter to send message to 'XIAO_ESP32S3' on '/dev/cu.usbmodem101')

Both NL & CR 115200 baud

```
Edge Impulse Inferencing Demo
Edge Impulse standalone inferencing (Arduino)
run_classifier returned: 0
Timing: DSP 0 ms, inference 2024 ms, anomaly 0 ms
Predictions:
    value: 0.36065
```

indexing: 14/49

Ln:256, Col:18 XIAO_ESP32S3 on /dev/cu.usbmodem101 2 49



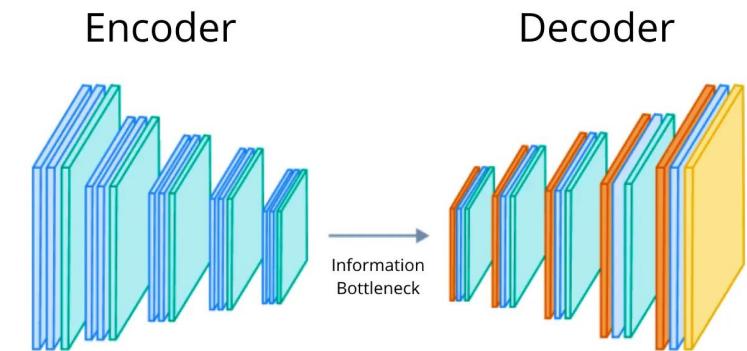
LLM / SLM

Large Language Model / Small Language Models

LLMs are **specialized deep learning models designed to understand and generate human language**, used for tasks like translation, summarization, and generating human-like text responses. SLMs are the same, but use a simpler, less resource-intensive approach (smaller in size).

Deep Learning models (or artificial neural networks)

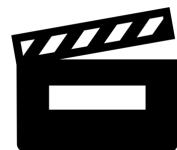
- **Autoencoders:** Used primarily for unsupervised learning tasks such as dimensionality reduction and feature extraction, autoencoders learn to compress data from the input layer into a shorter code and then reconstruct the output from this representation.
- **Transformer Models:** Highly effective in handling sequences, transformers use mechanisms like self-attention to weigh the importance of different words in a sentence, regardless of their position. The Transformer architecture, while innovative, can be seen as a derivative of earlier deep learning models, particularly those based on the concept of sequence modeling. However, the most direct lineage can be traced to the sequence-to-sequence (seq2seq) models that utilize **encoder-decoder** architectures. These earlier seq2seq models were often built using **recurrent neural networks (RNNs)** or their more advanced variants like **LSTMs (Long Short-Term Memory Networks)** or **GRUs (Gated Recurrent Units)**.



LLM/SLM – Large /Small Language Model

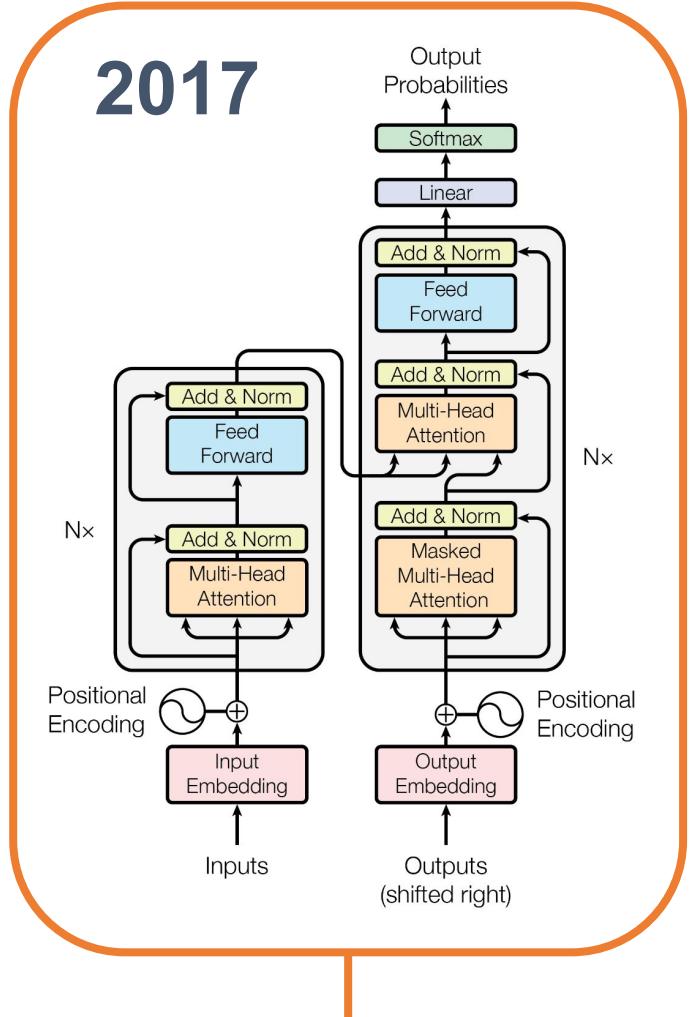
Large Language Models (LLMs) and SLMs are advanced neural networks based on the **Transformer architecture** that excel in understanding and generating human language. They represent a significant evolution from earlier sequence-based models like **LSTMs**, which surpass them in handling long-range dependencies and parallel processing efficiency.

Prof. Jesus's Presentation about IA:



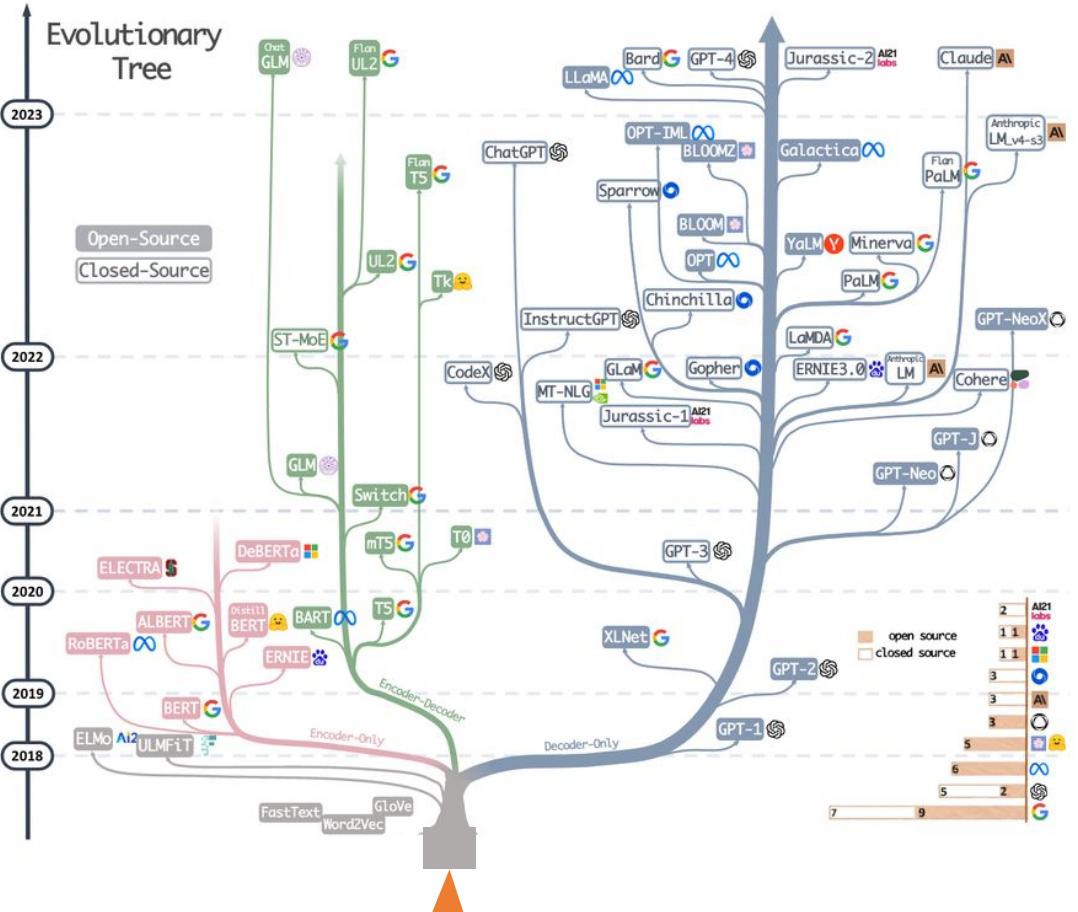
Transformers to LLMs and SLMs

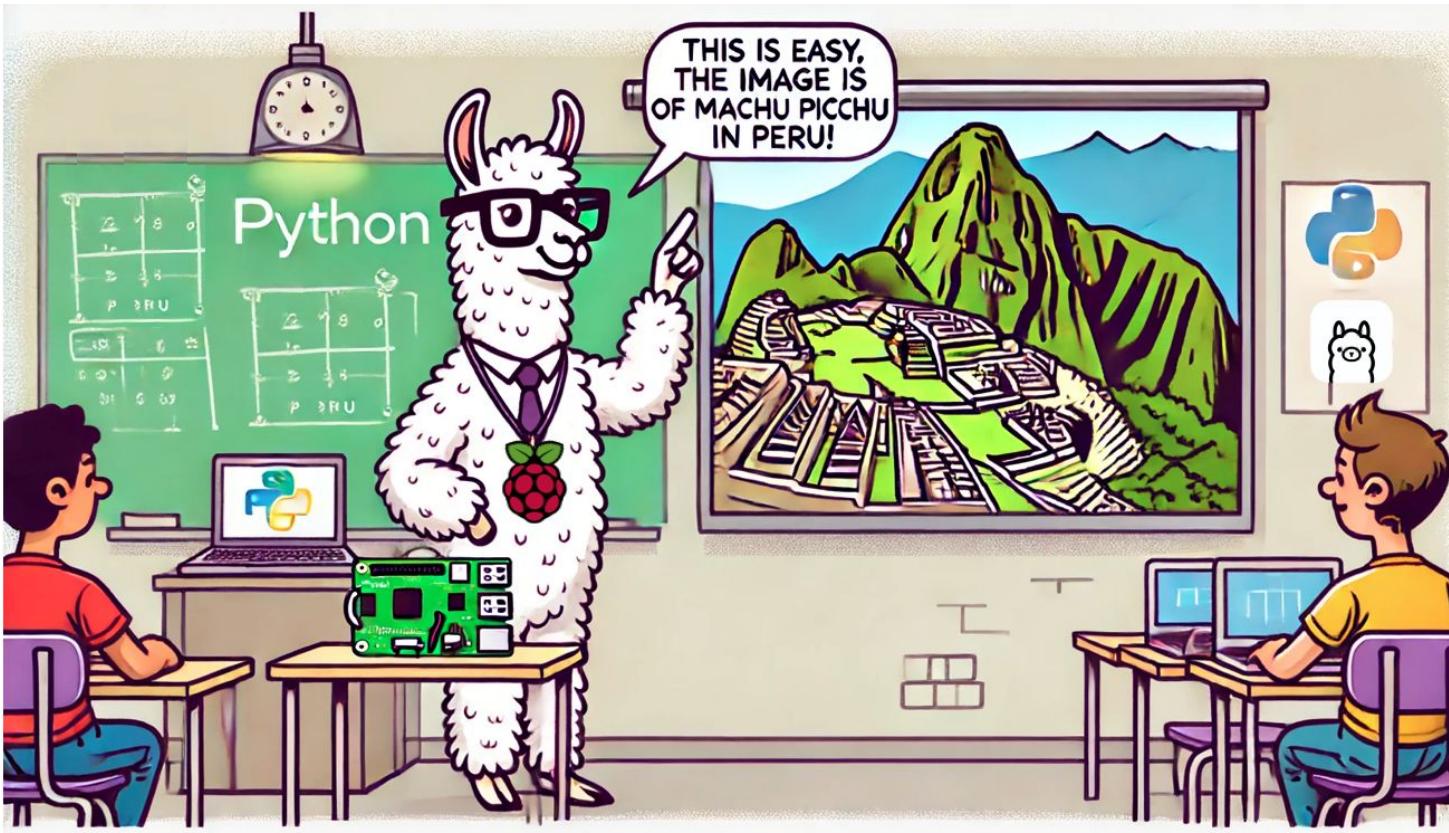
2017



Open

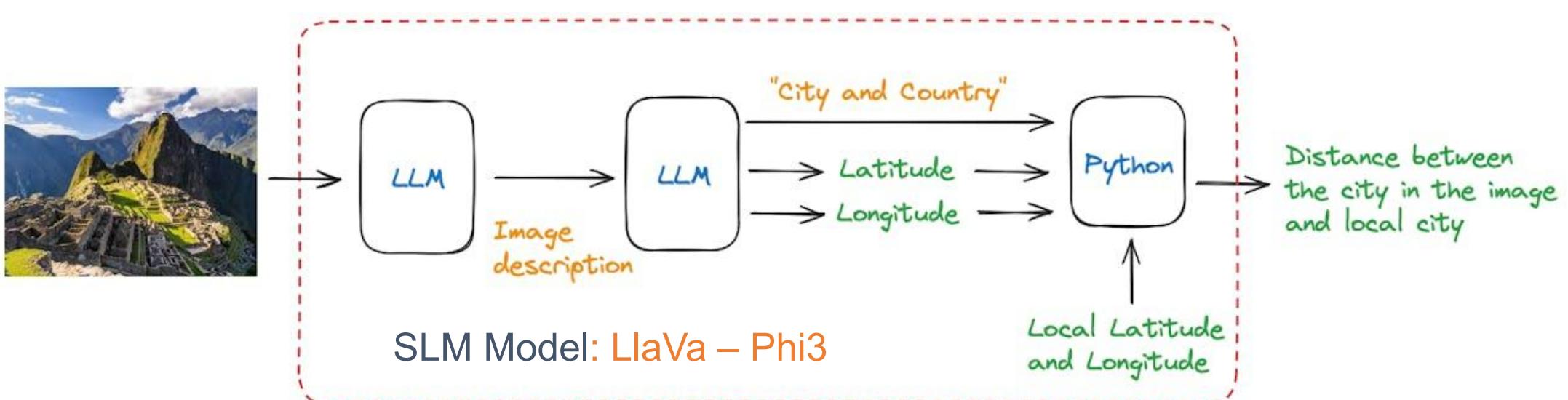
Closed



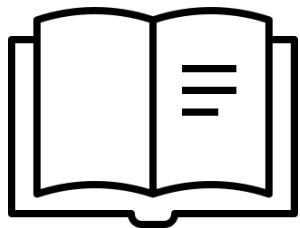


Running Large Language Models on Raspberry Pi at the Edge

Transform a Raspberry Pi into a powerful AI hub, running SLMs for real-time, on-site data analysis and insights using Ollama and Python.



llava-phi-3 is a LLaVA model (Large Language and Vision Assistant) fine-tuned from Microsoft Phi-3 mini



~ 350 pages



~ 300 words/page



1 word = ~ 1.4 token



A **4-bit** quantized **3.8 billion parameter *** language model trained on **3.3 trillion tokens****, whose overall performance, as measured by both academic benchmarks and internal testing, rivals that of models such as Mixtral 8x7B and GPT-3.5

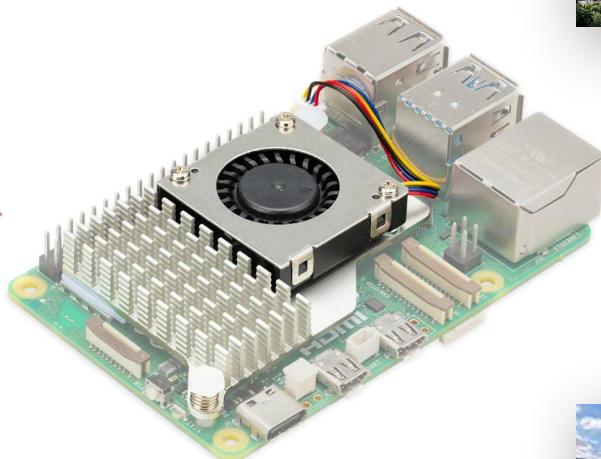
* 2.4 GB

** 22.5 Million books - 17% of all books written in the world

llava-phi-3 (2.9 GB)



Ollama



```
mjrovai@rpi-5:~\n\nFile Edit Tabs Help\n\n>>> Answer with one short sentence, what is the capital of France and its distance\n... in Km from Santiago, Chile\nThe capital of France is Paris and it is around 12,674 kilometers away\nfrom Santiago, Chile.\n\nTotal duration: 13.860074968s\nload duration: 1.537039ms\nprompt eval count: 27 token(s)\nprompt eval duration: 5.925386s\nprompt eval rate: 4.56 tokens/s\neval count: 26 token(s)\neval duration: 7.539223s\neval rate: 3.45 tokens/s\n>>> Send a message (/? for help)
```

(13 seconds)



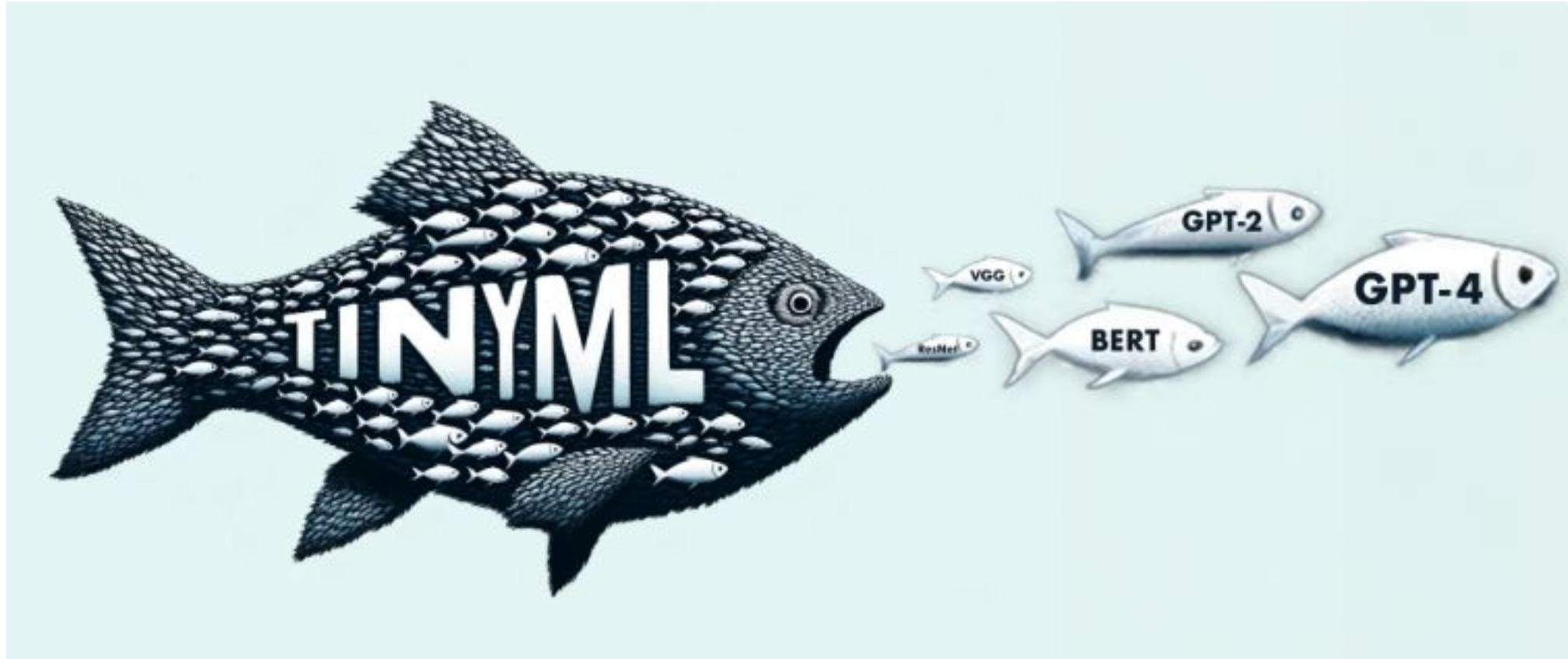
```
mjrovai@rpi-5:~/Documents/OLLAMA\n\nhelp\nroute.\n\n/Documents/OLLAMA $\n/Documents/OLLAMA $ python calc_distance_image.py /\n/home/mjrovai/Documents/OLLAMA/image_test_1.jpg\n\nThe image shows Paris, with lat:48.86 and long: 2.35, located in\nFrance and about 11,630 kilometers away from Santiago, Chile.\n\n[INFO] ==> The code (running llava-phi3), took 232.60845186299412\nseconds to execute.\n\nmjrovai@rpi-5:~/Documents/OLLAMA $
```



```
mjrovai@rpi-5:~/Documents/OLLAMA\n\nhelp\n\n/Documents/OLLAMA $\n/Documents/OLLAMA $ python calc_distance_image.py /\n/home/mjrovai/Documents/OLLAMA/image_test_3.jpg\n\nThe image shows Machu Picchu, with lat:-13.16 and long: -72.54,\nlocated in Peru and about 2,250 kilometers away from Santiago,\nChile.\n\n[INFO] ==> The code (running llava-phi3), took 267.579568572007\n7 seconds to execute.\n\nmjrovai@rpi-5:~/Documents/OLLAMA $
```

(4 minutes)

TinyML: Why the Future of Machine Learning is Tiny and Bright



Shvetank Prakash, Emil Njor, Colby Banbury, Matthew Stewart, Vijay Janapa Reddi

To learn more ...

Online Courses

[Harvard School of Engineering and Applied Sciences - CS249r: Tiny Machine Learning](#)

[Professional Certificate in Tiny Machine Learning \(TinyML\) – edX/Harvard](#)

[Introduction to Embedded Machine Learning - Coursera/Edge Impulse](#)

[Computer Vision with Embedded Machine Learning - Coursera/Edge Impulse](#)

[UNIFEI-HESTI01 TinyML: “Machine Learning for Embedding Devices”](#)

Books

[“Python for Data Analysis” by Wes McKinney](#)

[“Deep Learning with Python” by François Chollet - GitHub Notebooks](#)

[“TinyML” by Pete Warden and Daniel Situnayake](#)

[“TinyML Cookbook 2nd Edition” by Gian Marco Iodice](#)

[“Technical Strategy for AI Engineers, In the Era of Deep Learning” by Andrew Ng](#)

[“AI at the Edge” book by Daniel Situnayake and Jenny Plunkett](#)

[“XIAO: Big Power, Small Board” by Lei Feng and Marcelo Rovai](#)

[“MACHINE LEARNING SYSTEMS for TinyML” by a collaborative effort](#)

Projects Repository

[Edge Impulse Expert Network](#)

On the [TinyML4D website](#), You can find lots of educational materials on TinyML. They are all free and open-source for educational uses – we ask that if you use the material, please cite them! TinyML4D is an initiative to make TinyML education available to everyone globally.

Thanks



UNIFEI