



WALC 2024
Applied AI

Introduction to RNN - Time Series

Prof. Marcelo J. Rovai

rovai@unifei.edu.br

UNIFEI - Federal University of Itajuba, Brazil

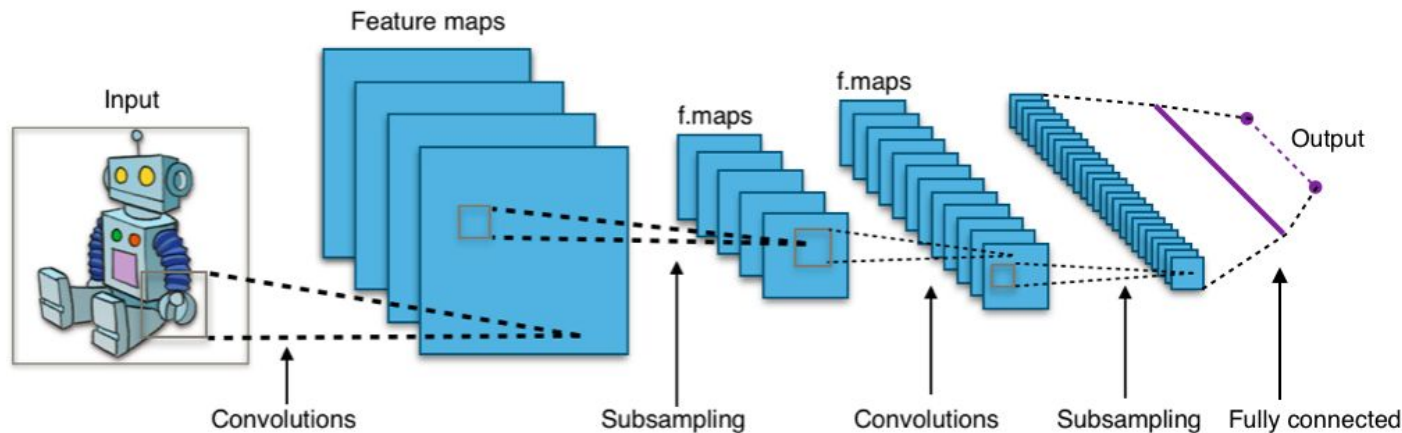
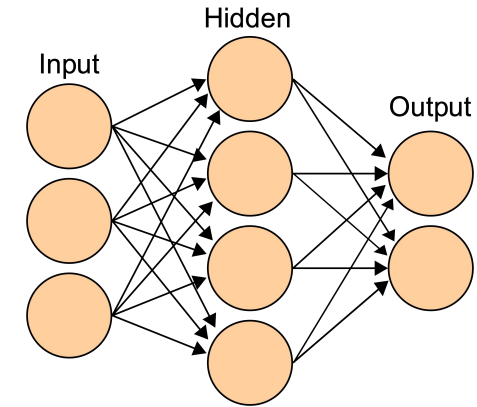
TinyML4D Academic Network Co-Chair



TINYML4D

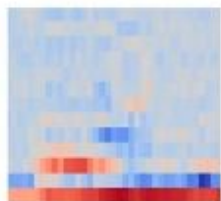
Deep Learning models (or artificial neural networks)

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.

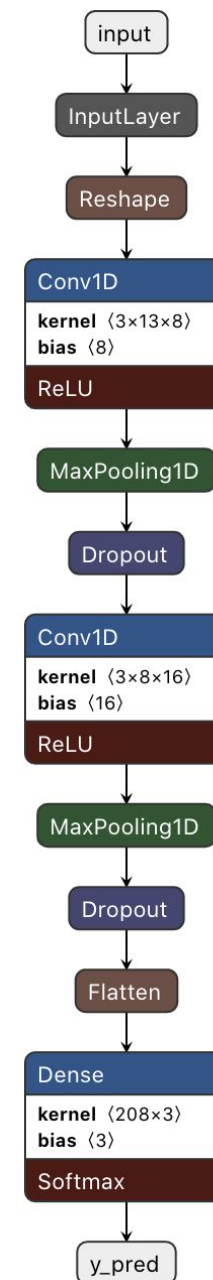
CNN for Audio:



Model: "sequential"

| Layer (type) | Output Shape | Param # |
|--------------------------------|----------------|---------|
| reshape (Reshape) | (None, 50, 13) | 0 |
| conv1d (Conv1D) | (None, 50, 8) | 320 |
| max_pooling1d (MaxPooling1D) | (None, 25, 8) | 0 |
| dropout (Dropout) | (None, 25, 8) | 0 |
| conv1d_1 (Conv1D) | (None, 25, 16) | 400 |
| max_pooling1d_1 (MaxPooling1D) | (None, 13, 16) | 0 |
| dropout_1 (Dropout) | (None, 13, 16) | 0 |
| flatten (Flatten) | (None, 208) | 0 |
| y_pred (Dense) | (None, 3) | 627 |

=====
Total params: 1,347
Trainable params: 1,347
Non-trainable params: 0
=====



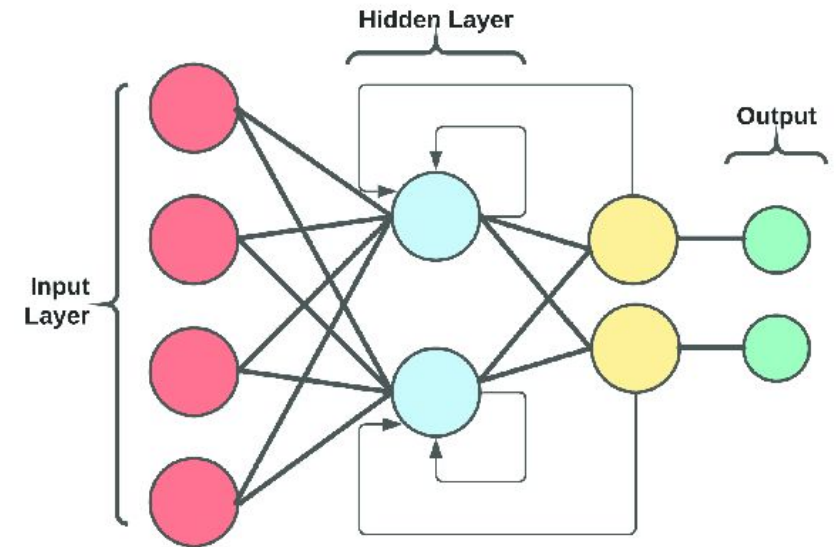
RNN (LSTM/GRU)

RNN (Recurrent Neural Network) architecture:

Think of RNN like a person reading a book - they remember what happened in previous pages to understand the current one. It's designed to work with sequences, like numbers and words in a sentence or notes in a melody.

Key components:

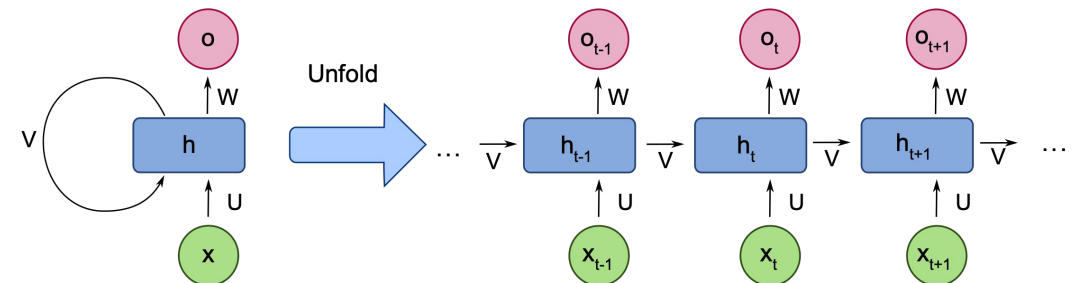
1. **Input Layer**: Receives data one element at a time (like words in a sentence)
2. **Hidden Layer**: Contains the "memory cells" that remember important past information
3. **Output Layer**: Produces predictions based on current input and stored memory



The unique feature is the "**feedback loop (v)**" - each step's output influences the next step's processing. This helps the network understand context and patterns over time.

Real-world applications include:

- Text prediction (like smartphone keyboards)
- Language translation
- Speech recognition
- Music generation
- Weather forecasting



LSTM (Long Short-Term Memory):

LSTM is an advanced version of RNN with three special "gates":

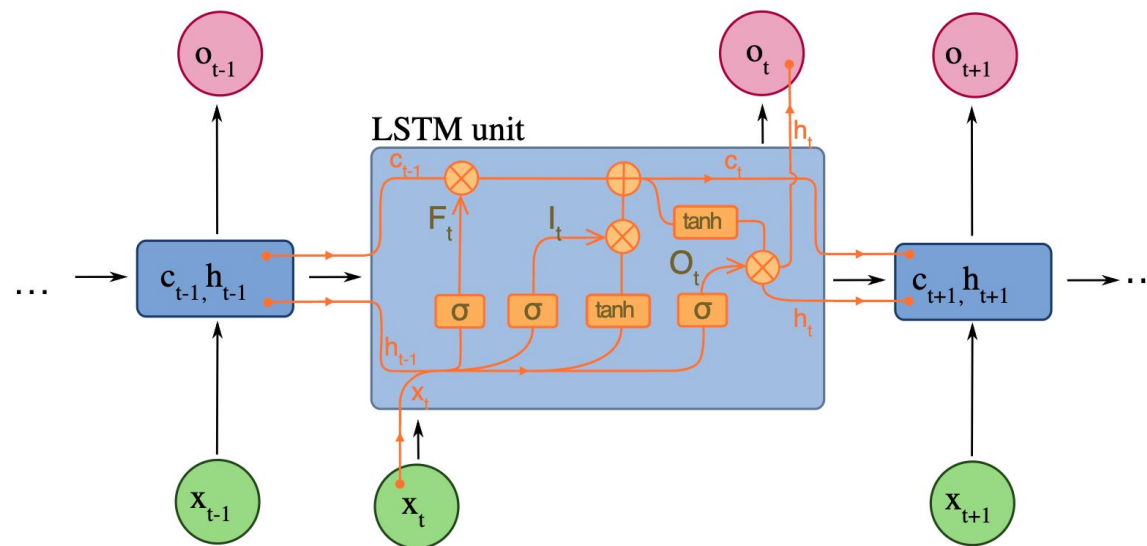
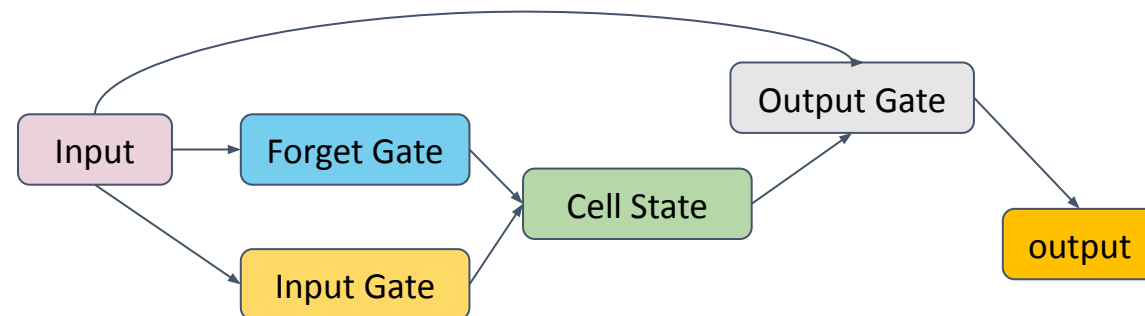
1. **Forget Gate**: Decides what information to discard
2. **Input Gate**: Decides what new information to store
3. **Output Gate**: Determines what parts of the cell state to output

Think of it like a smart note-taking system:

- Forget Gate: Erasing outdated notes
- Input Gate: Writing new important information
- Output Gate: Choosing which notes to share

This architecture helps solve the "vanishing gradient" problem of basic RNNs, allowing LSTMs to remember information for much longer periods. They're particularly good at:

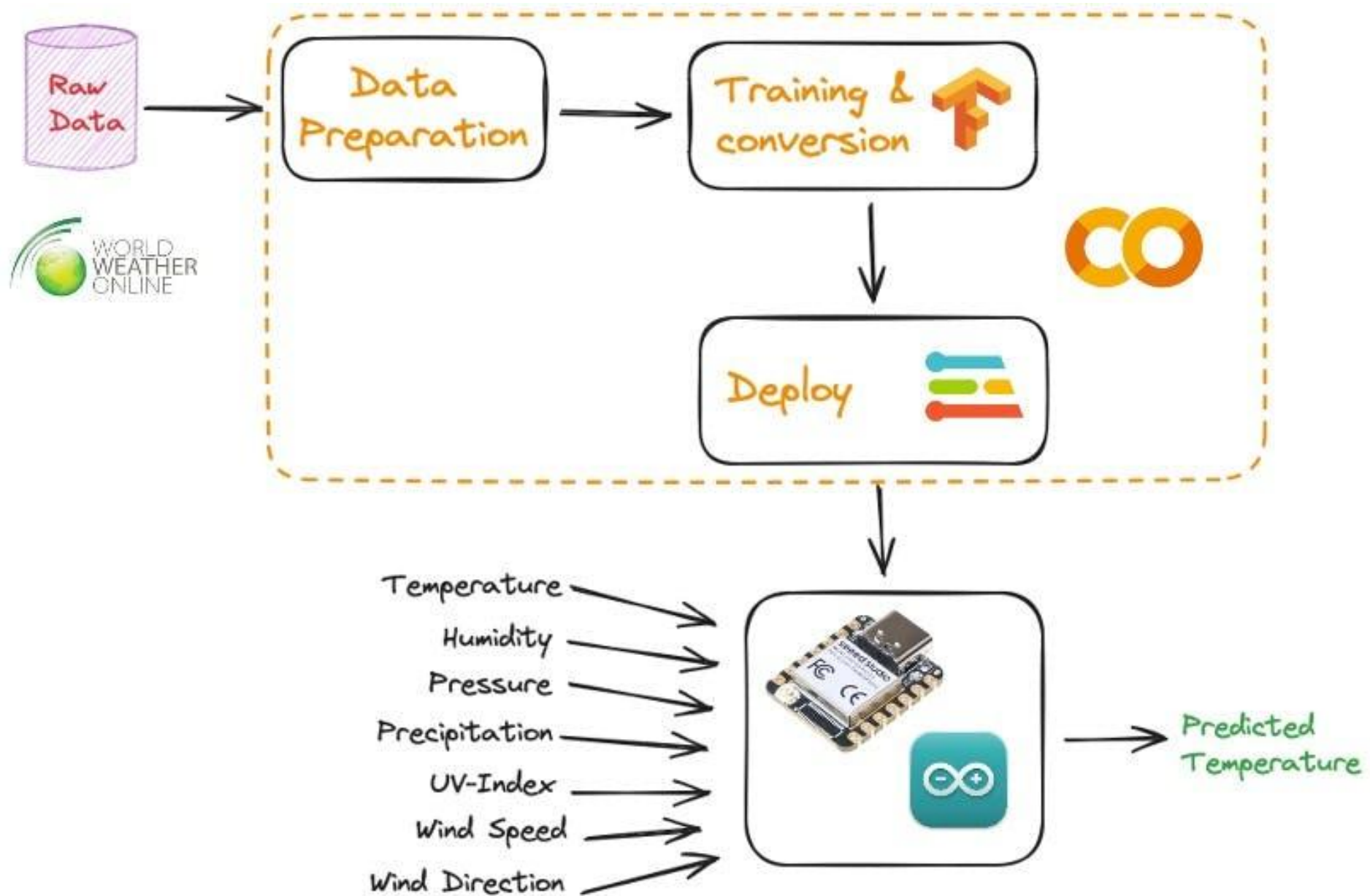
- Long text generation
- Time series prediction [Predicción Tráfico de Bicicletas usando DL/LSTM](#)
- Speech recognition
- Music composition



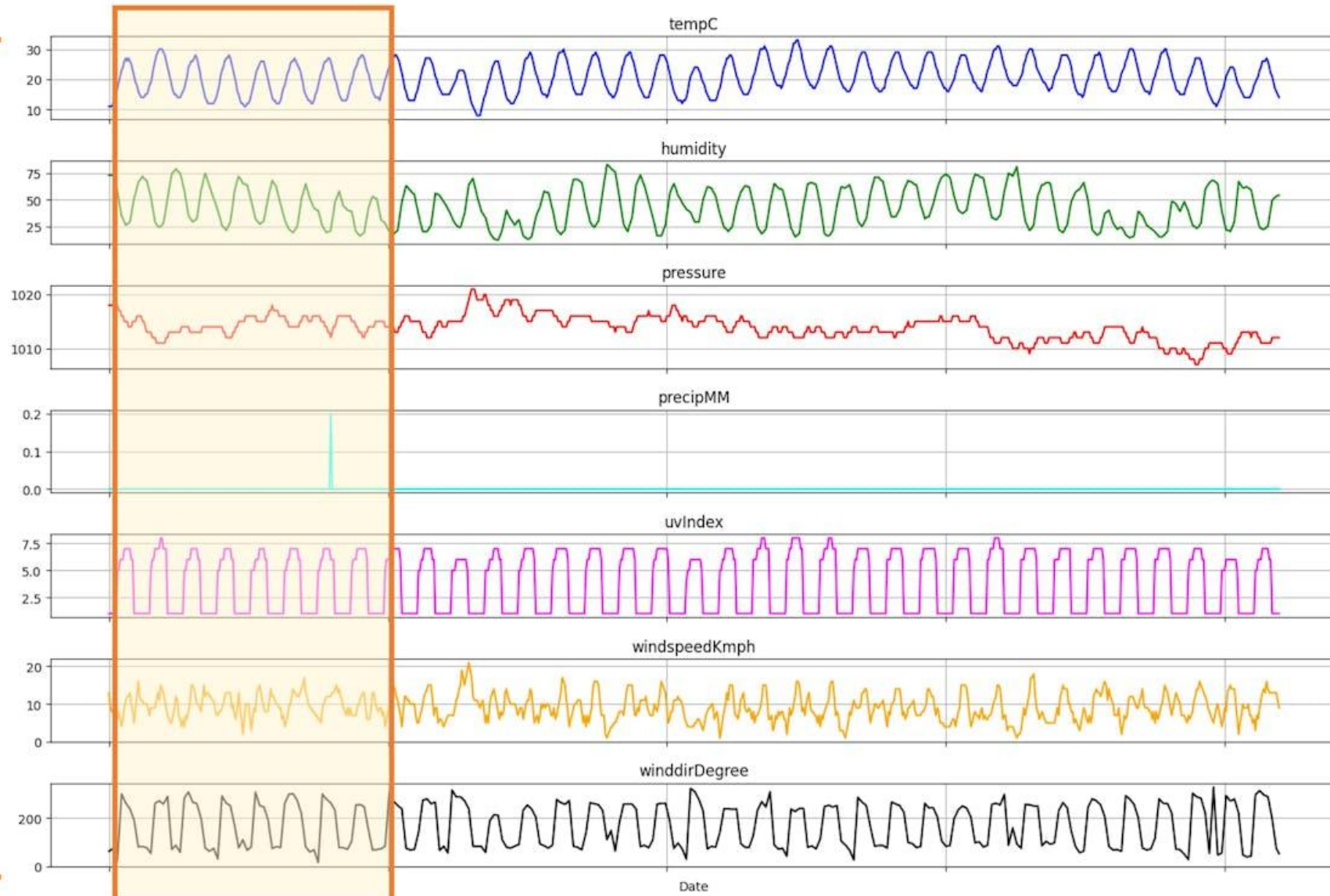


Temperature Prediction using an LSTM model

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

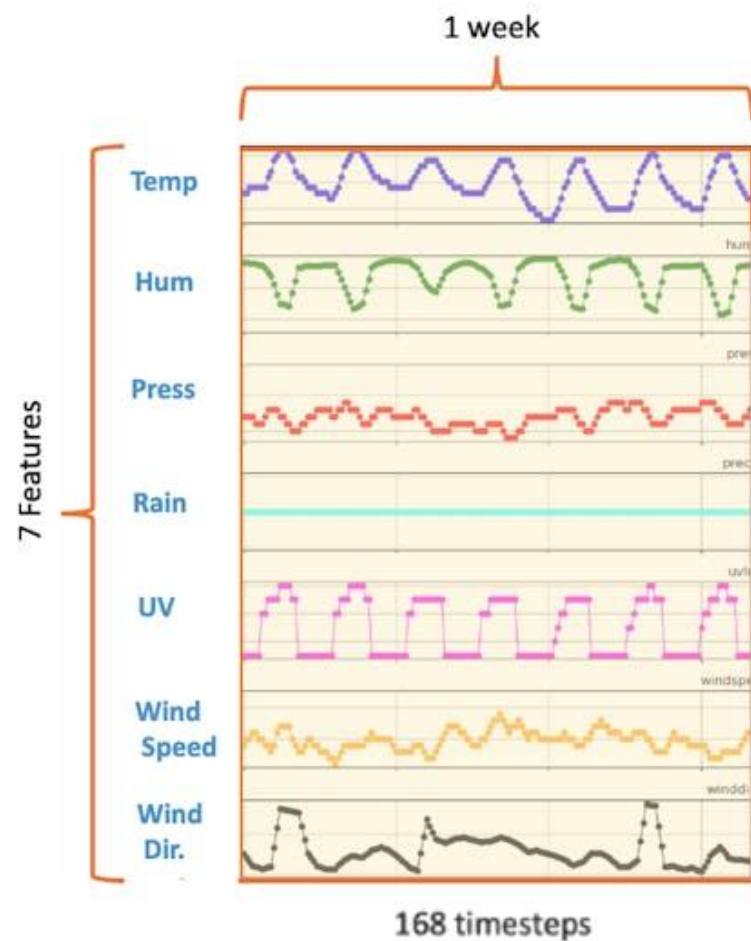


7 Features



168 timesteps

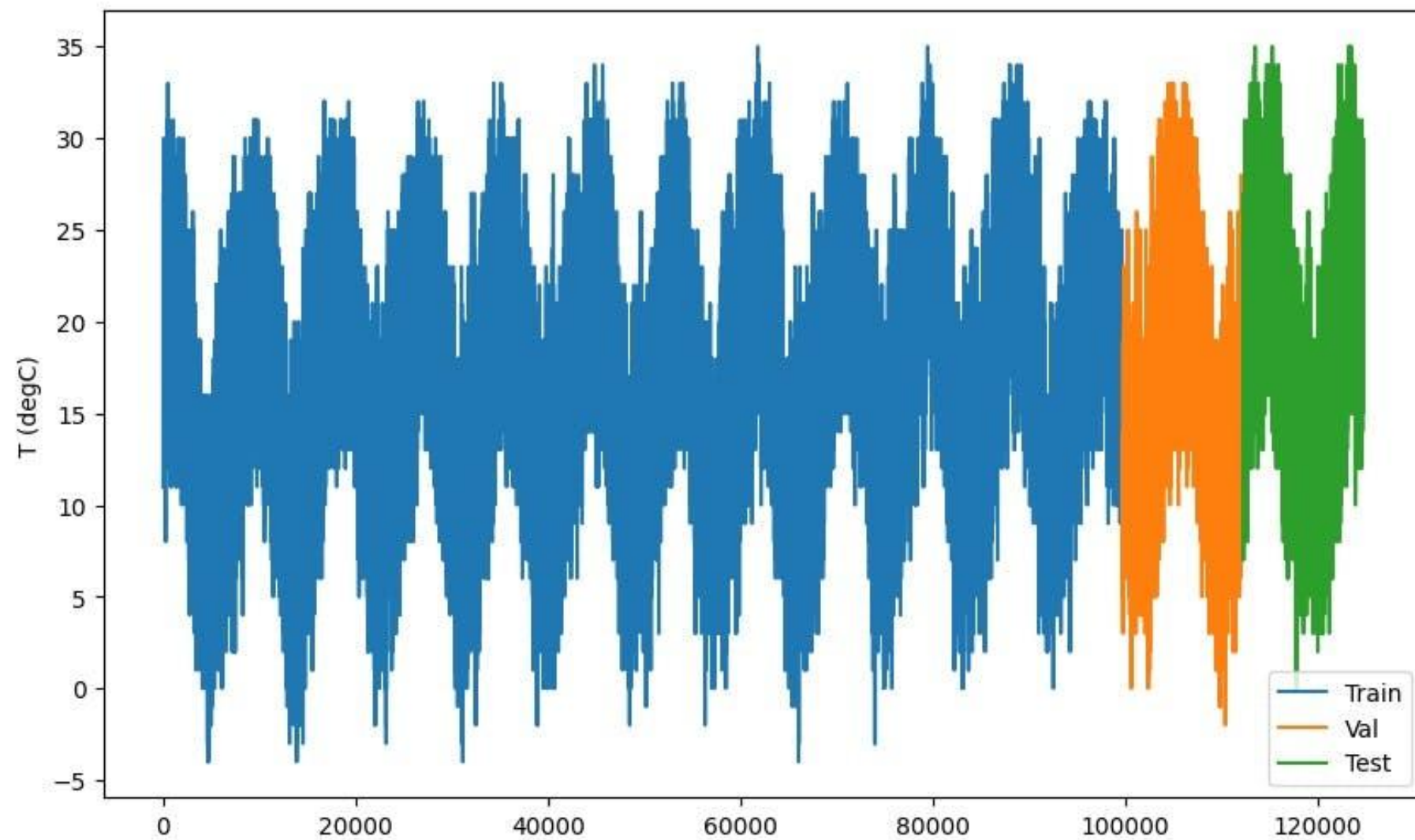




LSTM
Model



Temperature



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

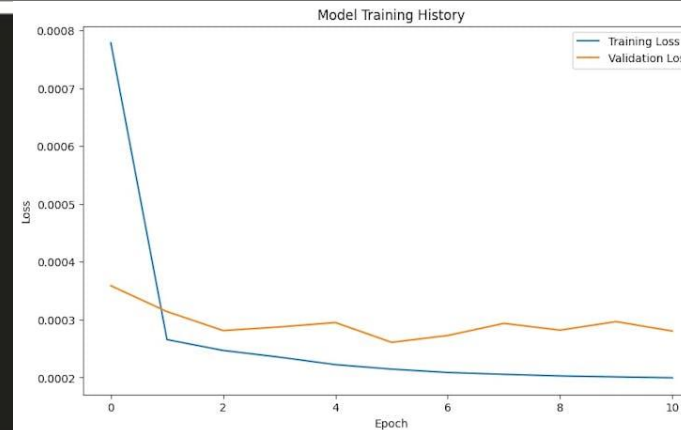
```
model.compile(optimizer='adam', loss='mse')
```

```
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

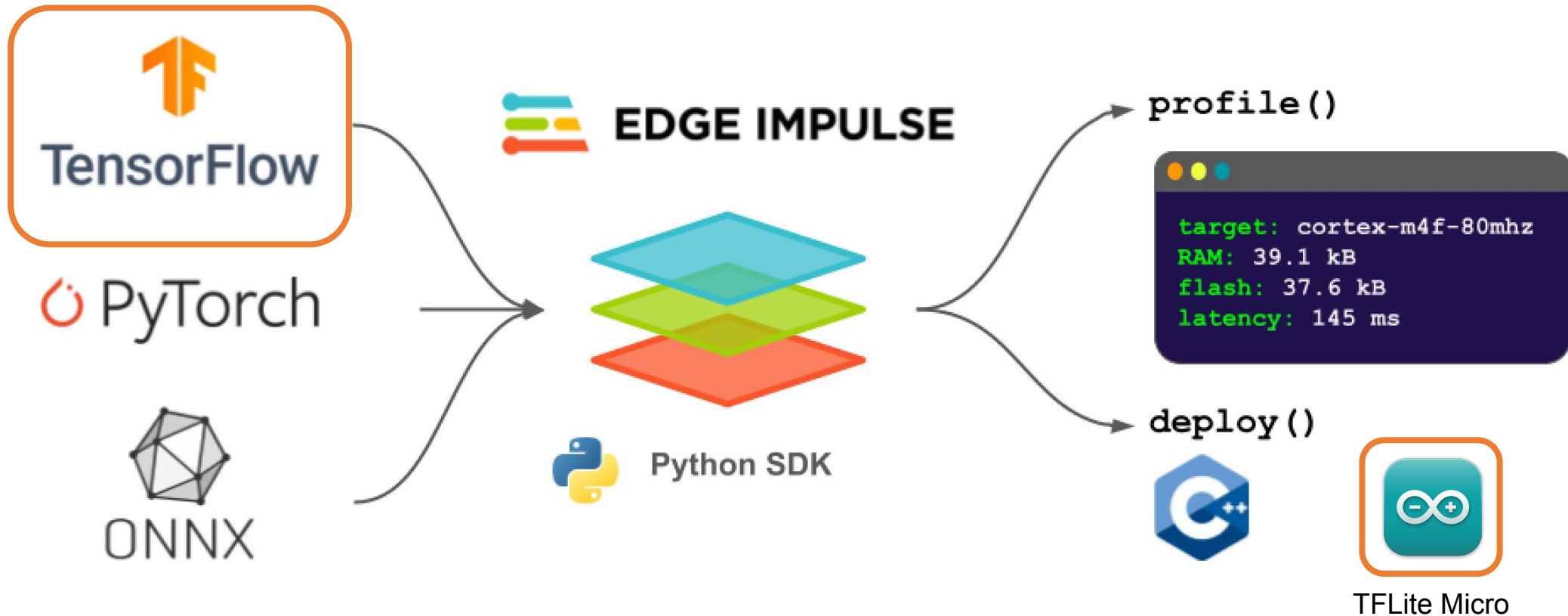


Model
Training
(TensorFlow)



Model
Conversion
(TFLite)

Edge Impulse Python SDK



converted_model.tflite

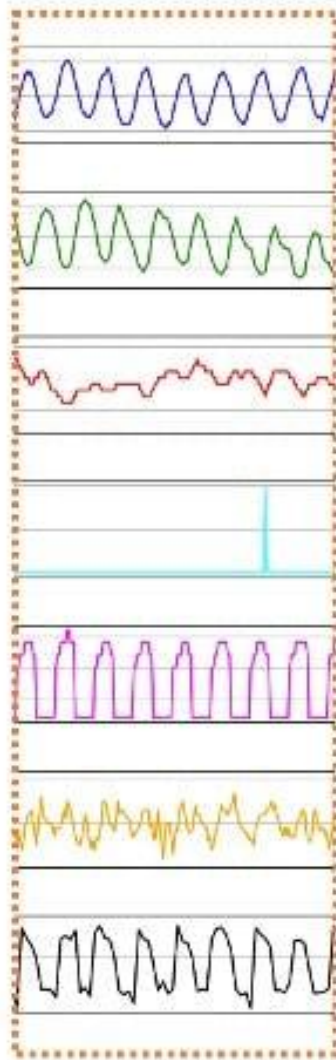


El Python SDK



lstm_float32_model.zip

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., 0.4516129, 0., 0.41666667, 0.41025641,  
0.5, 0.4516129, 0., 0.375, 0.,  
0.35555556, 0.46153846, 0.43877551, 0.4516129, 0.,  
0.375, 0.04347826, 0.46666667, 0.51282051, 0.37755102,  
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,  
0.60869565, 0.575, 0.74358974, 0.20408163, 0.38709677,  
0., 0.75, 0.73913043, 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.          , 0.575          , 0.2172012 , 0.65777778, 0.66666667
```

Output Serial Monitor ✕

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

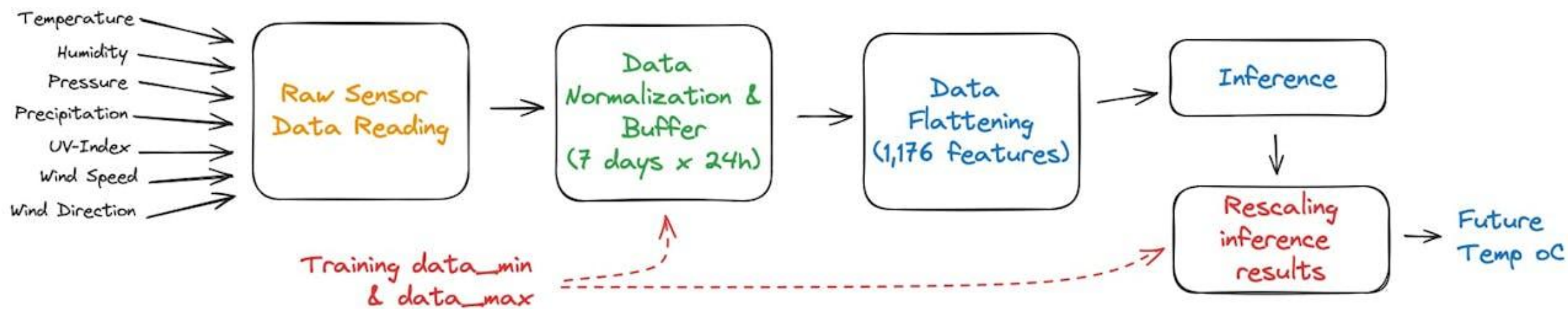
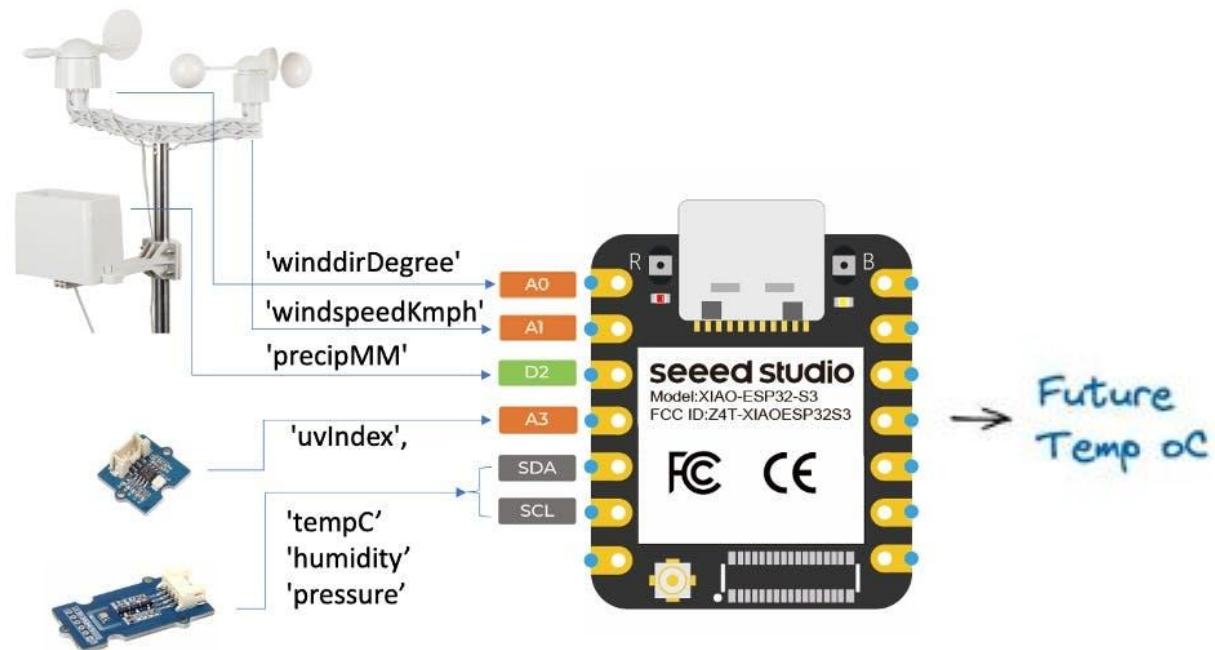
Both NL & CR ▾ 115200 baud ▾

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



GRU (Gated Recurrent Unit)

GRU (Gated Recurrent Unit) is a simplified version of LSTM with just two gates:

1. Reset Gate: Controls how much past information to forget
2. Update Gate: Decides how much past information to pass along

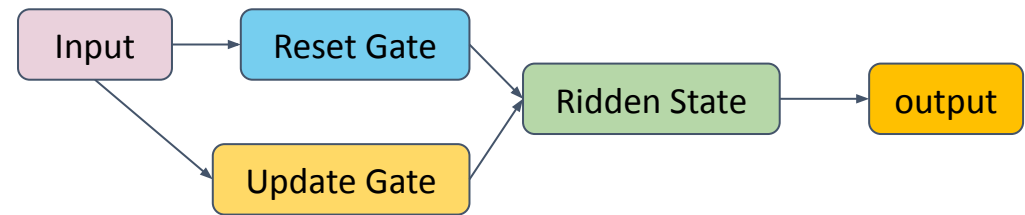
Key differences from LSTM:

- No separate cell state
- Fewer parameters, making it faster to train
- Often performs similarly to LSTM despite being simpler

Common applications:

- Machine translation
- Text generation
- Time series analysis
- Speech processing

The simpler architecture makes GRU a good choice when computational resources are limited or when working with smaller datasets.





VerneBOT: GENERATING TEXTS LIKE *JULES VERNE*

An Introduction to Language Models
Prof. Marcelo Rovai, UNIFEI

Questions?

