# A Wearable Multi-Modal Edge-Computing System for Real-Time Kitchen Activity Recognition

Mengxi Liu[1], Sungho Suh[1,2], Juan Felipe Vargas[1], Bo Zhou[1,2], Agnes Grünerbl[1], and Paul Lukowicz[1,2]

[1] German Research Center for Artificial Intelligence (DFKI), Kaiserslautern, Germany
[2] Department of Computer Science, RPTU Kaiserslautern-Landau, Kaiserslautern, Germany
`firstname.lastname@dfki.de`

**Abstract.** In the human activity recognition research area, prior studies predominantly concentrate on leveraging advanced algorithms on public datasets to enhance recognition performance, little attention has been paid to executing real-time kitchen activity recognition on energy-efficient, cost-effective edge devices. Besides, the prevalent approach of segregating data collection and context extraction across different devices escalates power usage, latency, and user privacy risks, impeding widespread adoption. This work presents a multi-modal wearable edge computing system for human activity recognition in real-time. Integrating six different sensors, ranging from inertial measurement units (IMUs) to thermal cameras, and two different microcontrollers, this system achieves end-to-end activity recognition, from data capture to context extraction, locally. Evaluation in an unmodified realistic kitchen validates its efficacy in recognizing fifteen activities, including a null class. Employing a compact machine learning model (184.5 kbytes) yields an average accuracy of 87.83 %, with model inference completed in 25.26 ms on the microcontroller. Comparative analysis with alternative microcontrollers showcases power consumption and inference speed performance, demonstrating the proposed system's viability.

**Keywords:** First keyword · Second keyword · Another keyword.

## 1 Introduction

Human activity recognition (HAR) provides a promising method for personal health monitoring and is exploited in many daily scenarios like healthcare centers, sports monitoring, and smart home, and it has emerged as a paramount study direction for the pervasive community [16]. Although researchers have explored many application scenarios of HAR, kitchen activity recognition has not yet been widely studied in comparison with others. A kitchen is an important place that people visit almost every day. Human activities in the kitchen are more closely related to their eating habits than in other environments. For example, opening microwave ovens and refrigerators may reflect a person's long-term food consumption frequency. Unhealthy dietary habits such as overeating, eating too frequently, or at inappropriate times can lead to many diseases. Therefore, continuous monitoring of kitchen activity can help people build their model of dietary habits and encourage them to develop healthier dietary habits. At the same time, the physician can also benefit from the dietary data to diagnose possible diseases.

Many existing works have focused on extracting activity-related context using complex algorithms on public datasets [31,12,32] or developed new sensor modalities [6] to improve recognition accuracy, while few of them studies implemented real-time kitchen activity recognition on a low-power, low-cost edge computing device. Deep learning (DL) is the primary algorithm used in state-of-the-art research for HAR. DL models with millions of parameters bring challenges to the deployment in resource-constraint microcontrollers, resulting in data processing on the cloud, which requires additional data transmission causing drawbacks during the application in real life. For example, the risk of user privacy exposure, unexpected power consumption, and unnecessary latency. The edge-devices-centric architectures can be more power-efficient, provide better privacy, and reduce latency for inference [33]. With the help of the tiny machine learning framework, integrated sensors in the smartphone (inertial measurement unit (IMU), camera, barometer, and microphone) have turned the smartphone into a very competitive multi-sensor edge-computing platform [24]. However, limited sensor modalities make smartphones unsuitable for complex activity recognition in many applications, for example, food preparation activity in kitchens, which can be addressed by the smartphone camera though, privacy issues prevent it from wide utilization.

To address the problems mentioned above, in this work, we present a wearable multi-sensor edge-computing platform for real-time human activity recognition in the kitchen. The hardware contains six sensors, one low-power consumption microcontroller, and one high-performance microcontroller. Using two microcontroller units (MCUs) improves the system's data processing performance and extends multiple sensors' peripheral interface for data acquisition. Users can deploy the model to different microcontrollers to obtain high computational performance or energy efficiency. Besides, we deployed different machine learning models based on multi-channel time-series convolutional neural networks (MC-CNN) [30] and DeepConvLSTM [26] for kitchen activity recognition on these MCUs to study the hardware performance. The model inference process is implemented locally, by which the raw data from users can be well protected locally, and the power consumption of data transmitted to the cloud can be avoided.

Overall, in this paper, we present the following contributions:

– We developed a wearable multi-sensor edge-computing system to recognize fourteen human activities in a kitchen scenario, which consists of six sensors and two MCUs.
– We deployed two kinds of popular neural network models for HAR on our proposed hardware and two other microcontrollers and compared their performance.
– We demonstrated the performance in detecting fifteen activities, including a null class, by an 184.5-KByte neural network model with an average accuracy of 87.83% and an average inference time of 25.26 ms.

## 2   Related Works

HAR benefits many areas, like healthcare monitoring, fitness tracking, and ambient-assisted living. Researchers have explored many methods to acquire data related to human activity. These methods can be divided into two groups: computer vision-based

methods and sensor-based methods. Computer vision-based methods have achieved high recognition accuracy in many application cases by using deep convolutional neural network (CNN) [31,12,32]. The sensor-based method is unsubstitutable in some special scenarios like physiological features detection [34,5,14,23]. Similar to our work, HAR in kitchen scenarios based on cameras or other sensors has been explored over the past years. For example, Bansal et al. [4] used a dynamic SVM-HMM hybrid model to predict nine cooking activities from video information. Lei et al. [19] proposed a study for fine-grained recognition of kitchen activities using RGB-D (Kinect-style) cameras. The proposed system can robustly track and accurately recognize detailed steps through cooking activities. Luo et al. [22] demonstrated a minimal and non-intrusive, low-power, low-cost radar-based sensing network system recognizing 15 kinds of activities. The related work contains solutions with remarkable recognition accuracy in recognizing kitchen activity and in complex kitchen environments. Another similar work for 15 kitchen activities recognition was presented in the work [21]. However, none of them have evaluated their model on edge devices in real-time, which prevents their widespread use in daily life.

Recently, researchers have studied the deployment of the machine learning model on embedded devices and the execution of real-time inference, which is crucial for a truly pervasive solution, thus bridging the gap between HAR research and commercial products. Edge devices are often limited by memory size and computational capacity, which implies CPU speed constraints, and some do not support floating-point operations. As a result, the machine learning model running on a workstation or cloud usually cannot be executed directly on the edge device. To address these problems, many studies have proposed software frameworks and tools for tiny machine learning, such as TensorFlow Lite Micro [11], MicroTVM [9], CMix-NN [8], CMSIS-NN [17] and STM X-Cube-AI [13]. In addition, many model size compression methods are also proposed, like pruning [3] and conversion/quantization [27]. With the help of tiny machine learning frameworks, machine learning models running on edge devices have become a reality. For example, Wan et al. [29] proposed DL models for real-time HAR with smartphones. Martin et al. [24] developed a lightweight algorithm for human activity detection based on Long short-term memory (LSTM) networks. The usability of the proposed algorithm is evaluated on real smartphone applications. Preetam et al. [1] proposed a machine learning pipeline to extract heart rate from pressure sensor data acquired on low-power edge devices, the size of their designed model was less than 40 kB, which was deployed on an ESP32 edge device. Bian et al. [7] demonstrated real-time hand gesture recognition based on capacitive sensing modality using tiny machine learning and deployed it on the Arduino Nano Sense platform. All these works have demonstrated the feasibility of HAR on edge-computing devices. However, most of the proposed tiny machine-learning models can only recognize a single category of human activities using multiple input channels from a single sensor modality.

Multi-modal wearable devices may offer advantages over uni-modal wearable devices for HAR, such as higher recognition accuracy, higher robustness (assuming uncorrelated error sources), and broader application scenarios. Thus, many works have explored multi-modal hardware. For example, Zhang et al. [34] designed a necklace with multiple embedded sensors with four types of sensors to detect eating-related ac-
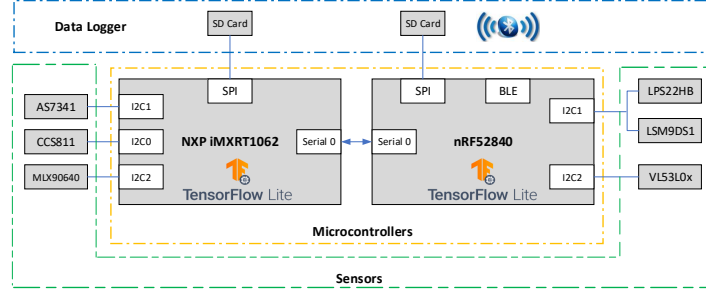
Fig. 1: Hardware Design of Wearable Multi-Modal Edge-Computing System for Real-Time Kitchen Activity Recognition

tivities. Bharti et al. [5] proposed a multi-modal and multi-positional system called "HuMan". The authors used five types of sensors (inertial and environmental) to recognize 21 complex activities in the home, with results up to 95%. Gravina et al. [14] presented a system based on body-worn inertial sensors combined with a pressure sensor to monitor in-seat activities, by which four ordinary basic emotion-relevant activities were recognized with high accuracy.

For all the above and to exploit the advantages of a multi-modal design, in this work, we propose a wearable multi-modal edge-computing system for real-time kitchen activity recognition locally.

## 3    Hardware Design and Implemented Neural Networks

As shown in Figure 1, the proposed wearable multi-modal edge-computing system consists of three main components; the MCU module, the sensor units, and the data logging module. All sensors are connected to two MCUs via the Inter-Integrated Circuit (I2C) interface. The storage capacity of each MCU is extended through the Serial Peripheral Interface (SPI) connected to an SD card. Data transmission between the two MCUs uses the universal asynchronous receiver/transmitter (UART) protocol. All components are integrated on a prototype PCB board with an overall dimension of 56×64 mm.

### 3.1    Sensors

Table 1 lists six sensors available on the proposed edge-computing system. As the sampling rate of each sensor varies, these six sensors are divided into two groups and connected to two MCUs, separately. A low sampling rate could lead to degradation of measurement accuracy, although it can also provide advantages such as low power consumption. The sampling rate range of all sensors is from 3 Hz to 12 Hz, and the final sampling rate of 6 Hz after synchronization is selected in this work.

Table 1: Sensor List

| Sensors | Description | Data Channels | Typical Application |
|---------|-------------|---------------|---------------------|
| AS7431 | Optical sensor | 10 | Food and beverages monitor [15] |
| CCS811 | Digital gas sensor | 2 | Electronic nose [2] |
| MLX90640 | Thermal IR array | 768 | Object detection [25] |
| LPS22HB | Air pressure sensor | 1 | Vertical movement [28] |
| LSM9DS1 | IMU (accelerometer, gyroscope, magnetic meter) | 9 | Fitness monitoring [10] |
| VL53L0X | Time-of-Flight ranging sensor | 1 | Environment recognition [18] |

Table 2: Important Parameters of Microcontrollers

| Microcontrollers | nRF52840 | MIMXRT1062 | STM32L4S5 | STM32F767 |
|------------------|----------|------------|-----------|-----------|
| Processor | ARM Cortex-M4 | ARM Cortex-M7 | ARM Cortex-M4 | ARM Cortex-M7 |
| Clock (Mhz) | 64 | 600 | 120 | 216 |
| Flash (Mbytes) | 1 | 8 | 2 | 2 |
| SRAM (Kbytes) | 256 | 1000 | 640 | 512 |
| Wireless connectivity | yes | no | yes (external) | no |

### 3.2   Microcontrollers (MCUs)

The microcontroller is an essential component of an edge computing system. The primary functions of the MCU in this work are the following three aspects: sensor data acquisition, data storage, and real-time context extraction. Table 2 presented some critical parameters of four MCUs with ARM Cortex-M series processor cores and with different memory (RAM) and storage (Flash) sizes as well as clock speeds. In DL-based models, SRAM size constrains the activation size (read and write) and Flash constrains the model size (read-only) [20]. Moreover, the inference time is determined by the neural network model and the MCU's clock speed. Therefore, the MCU nRF52840 with ARM Cortex-M4 core from Arduino Nano Sense board and the high-performance MCU MIMXRT1062 with ARM Cortex-M7 from the Teensy 4.1 board were selected for this work. The system clock of MIMXRT1062 is up to 600 MHz, by which application scenarios with high real-time requirements could be met at the expense of higher power consumption than nRF52840 MCU. The MCU nRF52840 has a slower system clock (64 MHz) and an on-chip Bluetooth module providing a wireless interface by which the extracted context and label can be transmitted. Besides, to investigate the performance of the tiny machine learning model for kitchen activity recognition, these models were also deployed on another two MCUs with similar processor cores (STM32L4S5 and STM32F767). All of them are supported by the TensorFlow Lite framework for tiny machine learning.

### 3.3   Neural Network Architecture

Seven hundred ninety-one channels data from six sensors were recorded, and two kinds of information fusion methods based on the neural network were utilized in our previous work [21]. The feature fusion method with 791 channels of input data has achieved the highest recognition accuracy and macro F1-Score, however, its size is the largest. The

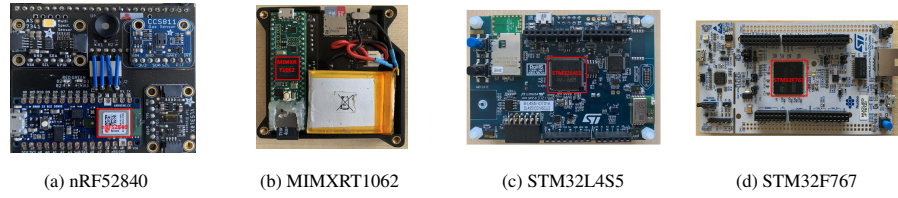(a) nRF52840      (b) MIMXRT1062      (c) STM32L4S5      (d) STM32F767

Fig. 2: Hardware Prototype and Commercial Board. (a) the front view and (b) the back view of the prototype of the proposed edge-computing device; (c) NUCLEO -767ZI board; (d) STM32L4S5 Discovery kit
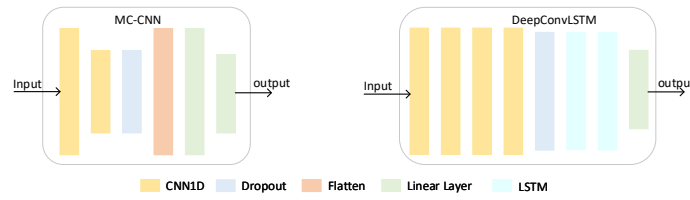


Fig. 3: Implemented Neural Network Architecture

data fusion method concatenating all the channel data before inputting to the neural network has the most straightforward architecture, and the best result of the data fusion method is very close to the feature fusion method. As the memory size is a limited resource of MCUs, the data fusion concatenating channels as input were selected to build a neural network running on the MCU in this work. Besides, one CNN layer was removed to reduce the model size compared to our previous work.

Figure 3 shows the architecture of the neural network utilized in this work. At first, the selected channel data from sensors were concatenated directly. Then, the concatenated data were fed into two 1D-CNN layers with a ReLu activation function, a dropout, and 1D average pooling layers to extract features. At last, the extracted features were input into dense layers to get an activity recognition result. The first layer of the 1D CNN has four times the number of filters as the second 1D CNN layer. In order to achieve the best activity recognition results while keeping the model size small, the effect of hyperparameters, like filter number in CNN and input channel, on recognition results and model size was investigated. Theoretically, the number of filters of the 1D CNN has an essential effect on recognition results. Too many filters in CNN could result in over-fitting and a more extensive model size, however, fewer filters make the model size small and suitable for tiny machine learning, but it often causes under-fitting. Three filter size numbers of the first 1D CNN layer were selected: 128, 256, and 400 filters. The input data were also divided into four groups in terms of sensor features with different input channels, such as 791 channels (including all sensor data), 768 channels (only infrared array sensor data), 23 channels (including all sensor data except infrared array sensor data), 17 channels (including all sensor data except infrared array sensor data, accelerometer, and gyroscope data).

The neural network model utilized for kitchen activity recognition was built under **TensorFlow 2.10.0** framework and the model training process was performed on a laptop with the GeForce RTX 3080 Ti GPU. Then, the TensorFlow Lite model was generated by the TensorFlow Lite Converter. During conversion, optimization mechanisms like quantization, pruning, and clustering can be applied to reduce the model size and latency with minimal or no loss in accuracy. Here, the full integer quantization method was used, by which the latency and peak memory usage can be reduced. All model math is integer quantized so that this model can be executed on the hardware only supporting integer operation. It is worth noting that the range, i.e., (min, max) of all floating-point tensors in the model needs to be estimated. Therefore, a representative dataset is required by the converter. The training dataset was used as a representative dataset in this work.

To evaluate the performance of real-time recognition performed on MCU, two models with different precision representations were generated for this task as follows:

1. TFLite Model: converted TensorFlow Lite model without any optimization from TensorFlow Model by TensorFlow Lite converter, the model is represented with 32-bit precision float data
2. Full integer Model: converted TensorFlow Lite model with full integer quantization from TensorFlow Model by TensorFlow Lite converter, the model is represented with 8-bit precision integer data.

Since human activities are made of complex sequences of motor movements, and capturing these temporal dynamics is fundamental for successful HAR [26], a Deep-ConvLSTM model, which includes four convolutional layers, two LSTM layers, and one softmax layer, was also utilized to recognize these kitchen activities and deployed on the microcontrollers. A performance comparison for this task between DeepConvLSTM model and MC-CNN model was given.

## 4   Evaluation of Kitchen Activity Recognition in Real-Time

### 4.1   Dataset of Kitchen Activity

The dataset used to train this work's neural network is from previous work[21], 15 kinds of activity including the null class in the kitchen scenarios shown in Figure 4 were selected to recognize. Ten volunteers performed them in a realistic unmodified kitchen environment. The experiment was divided into five sessions, and the experiment lasted around one hour per volunteer. Most of the activities in this experiment strongly correlate with users' eating habits, e.g., recognizing beverage intake during consumption and in the case of food preparation, such as opening the microwave and boiling water.

### 4.2   Metrics of Comparison

To assess the performance of our proposed device in a kitchen environment, we chose five metrics: accuracy, macro F1 Score, model size, inference time, and power consumption. Given the varying durations of kitchen activities and a constant data slide
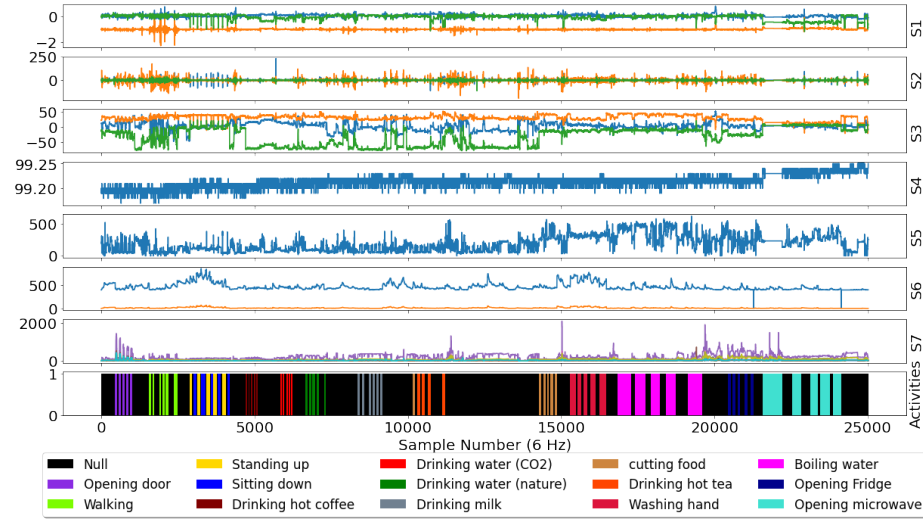
Fig. 4: Example of the measured multi-modal raw data from one volunteer and activity label. (**S1**: 3-channel accelerometer signal; **S2**: 3-channel gyroscope signal; **S3**: 3-channel magnetometer signal; **S4**: 1-channel Barometer signal; **S5**: 1-channel Distance signal; **S6**: 2-channel CO2 and TVOC gas signal; **S7**: 10-channel Optical Signal)

window size, the dataset suffers from imbalance. As a result, overall classification accuracy does not serve as a suitable performance indicator. Thus, the macro F1-score was also adopted as a benchmark. The MCU possesses significantly less memory and storage compared to computers and smartphones, making model size a vital metric for tiny machine learning applications. The main goal for deploying a machine model on an MCU is to facilitate local, real-time inferences. Nevertheless, due to constraints in computational resources such as the floating point unit (FPU) and clock speed, inference times can differ across MCUs. Additionally, power consumption is a also crucial metric for wearable devices as it affects battery life. Consequently, these five metrics were chosen as the evaluation criteria.

### 4.3   Result of Recognition Accuracy, Macro F1 Score and Model Size

Table 3 shows the recognition results from two types of models with varying hyperparameters, such as the number of input features and the count of filters in the first layer of the 1D CNN. In the MC-CNN architecture, the initial CNN layer contains four times as many filters as its subsequent layer. In contrast, in the DeepConvLSTM architecture, the filter count remains consistent across all four layers. The highest accuracy achieved was 87.93% by the TFLite model with the MC-CNN architecture, configured with 23 input channels and 400 filters in its initial layer. The same configuration in the DeepConvLSTM architecture reached the best accuracy of 84.01%. Conversely, the lowest accuracy was observed in the Full Quantization Model using MC-CNN architecture with only 768 channel infrared sensor data, a similar result can be observed in

Table 3: Result Summery of Recognition Accuracy and Macro F1 Score (MC-CNN: N1=128, N2=256, N3=400. DeepConvLSTM: N1=32, N2=64, N3=100. The best accuracy results are denoted in red color, and the poorest results are denoted in bold font)

| Parameters | | | MC-CNN | | | DeepConvLSTM | | |
|---|---|---|---|---|---|---|---|---|
| Channels | Filters | Inference Model | Macro F1 Score | Accuracy | Model Size (KBytes) | Macro F1 Score | Accuracy | Model Size (KBytes) |
| 17 | N1 | TFLite Model | 67.94 | 84.31 | 107.51 | 62.64 | 77.84 | 935.38 |
| | | Full Quantization Model | **67.89** | 84.36 | **36.91** | **61.81** | 77.27 | **259.72** |
| | N2 | TFLite Model | 70.07 | 84.97 | 299.12 | 67.16 | 81.31 | 1118.55 |
| | | Full Quantization Model | 69.93 | 84.82 | 89.14 | 66.88 | 80.98 | 308.97 |
| | N3 | TFLite Model | 72.33 | 85.26 | 632.20 | 67.83 | 82.06 | 1412.74 |
| | | Full Quantization Model | 72.00 | 85.13 | 177.26 | 67.29 | 81.59 | 386.40 |
| 23 | N1 | TFLite Model | 72.53 | 86.56 | 116.72 | 65.65 | 79.96 | 937.69 |
| | | Full Quantization Model | 72.25 | 86.39 | 39.22 | 64.50 | 79.43 | 260.29 |
| | N2 | TFLite Model | 75.19 | 87.26 | 317.56 | 70.94 | 83.80 | 1123.16 |
| | | Full Quantization Model | 74.97 | 87.16 | 93.74 | 69.16 | 83.10 | 310.12 |
| | N3 | TFLite Model | 77.25 | <span style="color:red">87.93</span> | 661.00 | 72.10 | <span style="color:red">84.01</span> | 1419.94 |
| | | Full Quantization Model | 77.05 | 87.83 | 184.46 | 71.12 | 83.56 | 388.20 |
| 768 | N1 | TFLite Model | 71.95 | 82.24 | 1261.04 | 62.84 | 75.63 | 1223.77 |
| | | Full Quantization Model | 71.27 | 82.20 | 325.30 | 62.55 | **75.28** | 331.81 |
| | N2 | TFLite Model | 71.56 | **81.91** | 2606.20 | 66.37 | 77.26 | 1695.32 |
| | | Full Quantization Model | 71.74 | 82.30 | 665.90 | 64.49 | 76.23 | 453.16 |
| | N3 | TFLite Model | 71.60 | 82.03 | 4237.00 | 68.29 | 79.35 | 2313.94 |
| | | Full Quantization Model | 71.59 | 82.54 | 1078.46 | 66.92 | 78.53 | 611.70 |
| 791 | N1 | TFLite Model | 76.67 | 86.45 | 1296.37 | 67.51 | 78.88 | 1232.60 |
| | | Full Quantization Model | 76.53 | 86.35 | 334.13 | 66.46 | 78.23 | 334.02 |
| | N2 | TFLite Model | 77.68 | 86.75 | 2676.85 | 71.59 | 82.23 | 1712.98 |
| | | Full Quantization Model | 77.61 | 86.92 | 683.57 | 70.95 | 81.71 | 457.57 |
| | N3 | TFLite Model | <span style="color:red">77.96</span> | 87.05 | <span style="color:red">4347.40</span> | 72.48 | 82.81 | <span style="color:red">2341.54</span> |
| | | Full Quantization Model | 77.84 | 87.10 | 1106.06 | <span style="color:red">72.55</span> | 82.61 | 618.60 |

the model with the DeepConvLSTM architecture. Moreover, changing the number of filters in the MC-CNN models had minimal impact on accuracy, improving it by merely one percent when filters were increased from 128 to 400. However, in DeepConvLSTM models, increasing filters from 32 to 100, accuracy was improved by over three percent. Despite significant reductions in size, the Full Quantization Model maintained nearly the same accuracy and macro F1 score as the TFLite Model.

The F1 score, which assesses model performance by balancing precision and recall, indicates that input sensor modalities greatly influence the macro F1 score; it rose by about ten percent with an increase in sensor types. The TFLite MC-CNN model achieved the highest macro F1 score of 77.96%. Given the variety of activities ranging from movement to drinking. While more sensor modalities and a greater number of convolutional filters can enhance feature extraction, improving the macro F1 score, it leads to the increase of the model size.

According to Table 3, the TFLite Model is approximately four times larger than the Full Quantization Model due to its use of 8-bit integers instead of 32-bit float data. Among all models, the largest is the TFLite MC-CNN architecture at 4346 KBytes

Table 4: Result of Inference Time and Power Consumption(the best accuracy results are denoted in red color, and the poorest results are denoted in bold font)

| Architecture | Metrics | Inference Model | nRF52840 | MIMXRT1062 | STM32L4S5 | STM32F767 |
|---|---|---|---|---|---|---|
| MC-CNN | Inference Time (ms) | TFLite Model | - | 169.63 | **4540.4** | 575.80 |
| | | Full Quantization Model | 394.49 | 25.26 | 195.21 | 31.57 |
| | Power Consumption (W) | TFLite Model | - | 0.78 | 0.67 | 1.13 |
| | | Full Quantization Model | 0.10 | 0.73 | 0.62 | 1.08 |
| DeepConvLSTM | Inference Time (ms) | TFLite Model | - | 253.84 | 2561.40 | 346.48 |
| | | Full Quantization Model | 685.95 | 56.04 | 516.07 | 70.15 |
| | Power Consumption (W) | TFLite Model | - | 0.72 | 0.68 | 1.13 |
| | | Full Quantization Model | 0.10 | 0.77 | 0.79 | **1.14** |

with 791 input channels, while the smallest is a Full Quantization Model at only 36.91 KBytes with 17 input channels. The presence of fully connected layers in the MC-CNN architecture causes a significant increase in model size with the number of input channels, though this effect is less pronounced in the DeepConvLSTM architecture. The smallest model size for DeepConvLSTM is about four times that of MC-CNN's smallest, but its largest model size is only half that of MC-CNN's largest. Enhancements in filters and input numbers improve recognition performance but also significantly increase the model size from kilobytes to megabytes. However, given that most MCUs, as shown in Table 2, have limited flash and SRAM sizes of up to 2 MBytes and 1 MByte respectively, many larger TFLite models may not be suitable. The largest Full Quantization Model size is around 1.1 MBytes, making full integer quantization a practical optimization strategy for these models, trading off minimal accuracy loss for substantial size reduction.

## 4.4   Result of Model Inference on Edge devices

To assess the real-time performance and power consumption of the model inference on the designated hardware, we opted for models with 23 input channels for implementation on microcontrollers due to their superior accuracy with a model size of approximately 1 Mbytes. Additionally, two other MCUs from development boards with comparable processors, as indicated in Table 2, were utilized to benchmark the real-time performance against our proposed hardware system. The primary testing sequence comprised reading sensor data, starting the timer counter, performing inference, predicting activities, measuring inference time, and handling outputs. Given that the six sensors for kitchen activity recognition were absent on the development boards, the sensor data, in 32-bit float format, was stored in these MCUs to evaluate model inference performance. Should the test model be a Full Quantization Model, the sensor data would be converted to an 8-bit integer prior to being fed into the model. A timer with microsecond accuracy was employed to record the inference time. The output is an array of 15 elements, each representing the probability of a specific activity. The activity predictor identifies and selects the element with the highest probability as the final output. Subsequently, the output handler activates the serial port to send the recognition result and inference time

to the laptop. To measure the power consumption during inference, an Oscilloscope was used along with a one Ohm Resistor placed in series in the 5V power supply line, with oscilloscope probes attached to both terminals of the resistor, allowing the inference current to be measured. The power consumption is determined by multiplying the measured current by the voltage.

Table 4 illustrates the hardware performance of selected models on four distinct microcontrollers. Given that nRF52840 possesses only 256 Kbytes of SRAM and 1 Mbyte of storage, it cannot support the TFLite Model due to its size. Thus, only the Full Quantization Model was tested on nRF52840. The fastest inference time was 25.26 milliseconds on the MIMXRT1062, a high-performance microcontroller with the fastest system clock. Conversely, the slowest inference time of 4540.4 milliseconds was realized on the STM32L4S5, which was running the TFLite Model with the MC-CNN architecture. Additionally, the inference time for the Full Quantization Model is significantly shorter than that of the TFLite Model. For instance, running the TFLite Model with MC-CNN architecture on STM32F767 takes about 18 times longer than the Full Quantization Model, primarily because floating-point operations consume much more time than 8-bit precision integer operations. The lowest power consumption during inference, at just 0.1 W, was achieved by nRF52840, whereas the highest was 1.14 W by STM32F767. It is also noted that model quantization has minimal impact on power consumption in this study.

## 5    Discussion

As indicated by the results in Table 3, the presence of additional sensor modalities and convolutional layer filters enhances the accuracy of kitchen activity recognition, though their impact is minimal when compared to the increase in model size and the decrease in inference time. In practical settings, a neural network model operating on an edge device and processing data locally offers numerous benefits, such as enhanced privacy protection. Therefore, considerations of model size and inference time become essential for effective deployment on edge devices. In the context of our application, the MC-CNN architecture outperforms DeepConvLSTM in terms of recognition accuracy and inference time, the former also is approximately half the size and faster in inference than the latter. Consequently, the MC-CNN architecture is more suitable for our purposes. Theoretically, optimization methods for neural networks do influence power consumption, but this impact is negligible when comparing the power usage among four microcontrollers due to other dominant factors such as system clock speed and hardware architecture. In our experiments, the microcontroller MIMXRT1062, with the highest clock speed of 600 MHz, demonstrated superior real-time inference capabilities. In contrast, the microcontroller nRF52840 exhibited the lowest power consumption of the group. Thus, our proposed Wearable Multi-Modal Edge-Computing System is adaptable for various real-time HAR scenarios, allowing for the selection of a microcontroller based on desired performance or energy efficiency.

## 6   Conclusion

In this paper, we introduced a wearable multi-modal edge-computing system designed for real-time kitchen activity recognition. Initially, we explored how sensor channel inputs and convolutional layer filters impact both recognition accuracy and model size, examining the performance of 48 TensorFlow Lite models across various input channels, filter sizes, optimization techniques, and neural network architectures. The highest achieved recognition accuracy was 87.93% with a macro F1 score of 77.25%. Subsequently, the most accurate model was implemented on four microcontrollers to assess its real-time performance and energy efficiency. The MIMXRT1062 microcontroller demonstrated a fast inference time of 25.26 milliseconds for recognizing 15 types of activities, while the nRF52840 microcontroller exhibited the lowest energy consumption at 100 mW. These results demonstrate that the developed system holds significant promise for real-time HAR. Future work will focus on refining the tiny machine learning model and broadening its application contexts.

## References

1. Anbukarasu, P., Nanisetty, S., Tata, G., Ray, N.: Tiny-hr: Towards an interpretable machine learning pipeline for heart rate estimation on edge devices. arXiv preprint arXiv:2208.07981 (2022)
2. Arroyo, P., Meléndez, F., Suárez, J.I., Herrero, J.L., Rodríguez, S., Lozano, J.: Electronic nose with digital gas sensors connected via bluetooth to a smartphone for air quality measurements. Sensors **20**(3),  786 (2020)
3. Banik, D., Ekbal, A., Bhattacharyya, P.: Machine learning based optimized pruning approach for decoding in statistical machine translation. IEEE Access **7**, 1736–1751 (2018)
4. Bansal, S., Khandelwal, S., Gupta, S., Goyal, D.: Kitchen activity recognition based on scene context. In: 2013 IEEE International Conference on Image Processing. pp. 3461–3465. IEEE, Melbourne, VIC, Australia (2013)
5. Bharti, P., De, D., Chellappan, S., Das, S.K.: Human: Complex activity recognition with multi-modal multi-positional body sensing. IEEE Transactions on Mobile Computing **18**(4), 857–870 (2018)
6. Bian, S.: Human Activity Recognition with Field Sensing Technique. Ph.D. thesis, Technische Universität Kaiserslautern (2022)
7. Bian, S., Lukowicz, P.: Capacitive sensing based on-board hand gesture recognition with tinyml. In: Adjunct Proceedings of the 2021 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2021 ACM International Symposium on Wearable Computers. pp. 4–5 (2021)
8. Capotondi, A., Rusci, M., Fariselli, M., Benini, L.: Cmix-nn: Mixed low-precision cnn library for memory-constrained edge devices. IEEE Transactions on Circuits and Systems II: Express Briefs **67**(5), 871–875 (2020)
9. Chen, T., Moreau, T., Jiang, Z., Zheng, L., Yan, E., Shen, H., Cowan, M., Wang, L., Hu, Y., Ceze, L., et al.: {TVM}: An automated {End-to-End} optimizing compiler for deep learning. In: 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18). pp. 578–594 (2018)
10. Crema, C., Depari, A., Flammini, A., Sisinni, E., Haslwanter, T., Salzmann, S.: Imu-based solution for automatic detection and classification of exercises in the fitness scenario. In: 2017 IEEE Sensors Applications Symposium (SAS). pp. 1–6. IEEE (2017)

11. David, R., Duke, J., Jain, A., Janapa Reddi, V., Jeffries, N., Li, J., Kreeger, N., Nappier, I., Natraj, M., Wang, T., et al.: Tensorflow lite micro: Embedded machine learning for tinyml systems. Proceedings of Machine Learning and Systems **3**, 800–811 (2021)
12. Du, Y., Wang, W., Wang, L.: Hierarchical recurrent neural network for skeleton based action recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1110–1118. IEEE, Boston, MA, USA (2015)
13. Falbo, V., Apicella, T., Aurioso, D., Danese, L., Bellotti, F., Berta, R., Gloria, A.D.: Analyzing machine learning on mainstream microcontrollers. In: International Conference on Applications in Electronics Pervading Industry, Environment and Society. pp. 103–108. Springer (2019)
14. Gravina, R., Li, Q.: Emotion-relevant activity recognition based on smart cushion using multi-sensor fusion. Information Fusion **48**, 1–10 (2019)
15. Huang, H., Yu, H., Xu, H., Ying, Y.: Near infrared spectroscopy for on/in-line monitoring of quality in foods and beverages: A review. Journal of food engineering **87**(3), 303–313 (2008)
16. Kumrai, T., Korpela, J., Maekawa, T., Yu, Y., Kanai, R.: Human activity recognition with deep reinforcement learning using the camera of a mobile robot. In: 2020 IEEE International Conference on Pervasive Computing and Communications (PerCom). pp. 1–10. IEEE (2020)
17. Lai, L., Suda, N., Chandra, V.: Cmsis-nn: Efficient neural network kernels for arm cortex-m cpus. arXiv preprint arXiv:1801.06601 (2018)
18. Laković, N., Brkić, M., Batinić, B., Bajić, J., Rajs, V., Kulundžić, N.: Application of low-cost vl53l0x tof sensor for robot environment detection. In: 2019 18th International Symposium INFOTEH-JAHORINA (INFOTEH). pp. 1–4. IEEE (2019)
19. Lei, J., Ren, X., Fox, D.: Fine-grained kitchen activity recognition using rgb-d. In: Proceedings of the 2012 ACM Conference on Ubiquitous Computing. pp. 208–211. ACM, Pittsburgh, USA (2012)
20. Lin, J., Chen, W.M., Lin, Y., Gan, C., Han, S., et al.: Mcunet: Tiny deep learning on iot devices. Advances in Neural Information Processing Systems **33**, 11711–11722 (2020)
21. Liu, M., Suh, S., Zhou, B., Gruenerbl, A., Lukowicz, P.: Smart-badge: A wearable badge with multi-modal sensors for kitchen activity recognition. In: Adjunct Proceedings of the 2022 ACM International Joint Conference on Pervasive and Ubiquitous Computing and the 2022 ACM International Symposium on Wearable Computers. pp. 356–363 (2022)
22. Luo, F., Poslad, S., Bodanese, E.: Kitchen activity detection for healthcare using a low-power radar-enabled sensor network. In: ICC 2019-2019 IEEE International Conference on Communications (ICC). pp. 1–7. IEEE (2019)
23. Mehrang, S., Pietila, J., Tolonen, J., Helander, E., Jimison, H., Pavel, M., Korhonen, I.: Human activity recognition using a single optical heart rate monitoring wristband equipped with triaxial accelerometer. In: EMBEC & NBC 2017, pp. 587–590. Springer (2017)
24. Milenkoski, M., Trivodaliev, K., Kalajdziski, S., Jovanov, M., Stojkoska, B.R.: Real time human activity recognition on smartphones using lstm networks. In: 2018 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO). pp. 1126–1131. IEEE (2018)
25. Naser, A., Lotfi, A., Zhong, J., He, J.: Human activity of daily living recognition in presence of an animal pet using thermal sensor array. In: Proceedings of the 13th ACM International Conference on PErvasive Technologies Related to Assistive Environments. pp. 1–6 (2020)
26. Ordóñez, F.J., Roggen, D.: Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition. Sensors **16**(1), 115 (2016)
27. Pourghasemi, H.R., Gayen, A., Lasaponara, R., Tiefenbacher, J.P.: Application of learning vector quantization and different machine learning techniques to assessing forest fire influence factors and spatial modelling. Environmental research **184**, 109321 (2020)

28. Vanini, S., Faraci, F., Ferrari, A., Giordano, S.: Using barometric pressure data to recognize vertical displacement activities on smartphones. Computer Communications **87**, 37–48 (2016)
29. Wan, S., Qi, L., Xu, X., Tong, C., Gu, Z.: Deep learning models for real-time human activity recognition with smartphones. Mobile Networks and Applications **25**(2), 743–755 (2020)
30. Yang, J., Nguyen, M.N., San, P.P., Li, X.L., Krishnaswamy, S.: Deep convolutional neural networks on multichannel time series for human activity recognition. In: Twenty-fourth international joint conference on artificial intelligence (2015)
31. Yang, X., Tian, Y.: Super normal vector for human activity recognition with depth cameras. IEEE transactions on pattern analysis and machine intelligence **39**(5), 1028–1039 (2016)
32. Yang, X., Zhang, C., Tian, Y.: Recognizing actions using depth motion maps-based histograms of oriented gradients. In: Proceedings of the 20th ACM international conference on Multimedia. pp. 1057–1060. ACM, Nara, Japan (2012)
33. Zaidi, S.A., Hayajneh, A.M., Hafeez, M., Ahmed, Q.: Unlocking edge intelligence through tiny machine learning (tinyml). IEEE Access (2022)
34. Zhang, S., Zhao, Y., Nguyen, D.T., Xu, R., Sen, S., Hester, J., Alshurafa, N.: Necksense: A multi-sensor necklace for detecting eating activities in free-living conditions. Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies **4**(2), 1–26 (2020)