

**Implementation of a Quantized Convolutional
Neural Network for automatic detection of Atrial
Fibrillation in an 8-bit microcontroller**

**MAURICIO BAUTISTA PORRAS
LAURA CRISTINA MARTINEZ CRUZ**

**UNIVERSIDAD INDUSTRIAL DE SANTANDER
FACULTAD DE INGENIERÍAS FISICOMECÁNICAS
ESCUELA DE INGENIERÍA, ELÉCTRICA,
ELECTRÓNICA Y DE TELECOMUNICACIONES
BUCARAMANGA
2020**

**Implementation of a Quantized Convolutional
Neural Network for automatic detection of Atrial
Fibrillation in an 8-bit microcontroller**

**MAURICIO BAUTISTA PORRAS
LAURA CRISTINA MARTINEZ CRUZ**

Trabajo de grado presentado como requisito parcial para optar al título de
ingeniero electrónico

Director
Carlos Fajardo Ariza
Doctor en Ingeniería

**UNIVERSIDAD INDUSTRIAL DE SANTANDER
FACULTAD DE INGENIERÍAS FISICOMECAÑICAS
ESCUELA DE INGENIERÍA, ELÉCTRICA,
ELECTRÓNICA Y DE TELECOMUNICACIONES
BUCARAMANGA
2020**

AGRADECIMIENTOS

Agradezco a Dios por haber bendecido mi vida y guiado en este camino, a mis padres por su incondicional ayuda de todas las formas durante todo mi proceso universitario, igualmente a mi hermana por siempre estar ahí, a Duván por apoyar mis últimos días de estudiante universitaria y a las amistades que me deja la Universidad Industrial de Santander, porque el proceso fue mucho más provechoso con ellos, un café siempre pudo resolver cualquier situación. Gracias, a mi compañero de proyecto, al grupo de investigación CPS, especialmente a nuestro asesor privado por su gran aporte y de igual forma al director del proyecto, por su guía académica, profesional y anímica. Todo valió la pena, y hubiese sido más difícil sin ustedes. **Laura Cristina Martinez Cruz**

Doy gracias a Dios por guiar mi camino, por ser el faro en mi vida. A mis padres, porque gracias a ellos, soy quien soy el día de hoy y eternamente les estaré agradecido. A mis hermanos y familia, porque su granito de arena también cuenta. A mis amigos, la familia que uno escoge, les agradezco por estar a mi lado, por apoyarme en los momentos más tristes y difíciles, pero también, por ser partícipes de las mejores experiencias, momentos y vivencias que hoy puedo atesorar. Al equipo de la oficina, gracias, porque no hay problema que un café y una buena charla no puedan solucionar. A nuestro asesor privado de tesis (JDB) gracias. A mi querida colega, gracias porque logramos ser un excelente equipo. A la Universidad Industrial de Santander, a la rama IEEE UIS, al grupo de investigación CPS y en especial, a nuestro director, gracias por su consejo y apoyo. En fin, hay muchas más personas a las que quiero agradecer, y estas líneas no alcanzan para nombrarlas a todas, solo me queda decir gracias, gracias por apoyarme sin importar que. Y termino estos agradecimientos, con la mayor gratitud hacia la mujer que inició esto, que me llevó a mirar más lejos, a fijar una meta más grande y me apoyó en este camino, a la profe Cecilia ¡gracias!. Este es el inicio de un gran comienzo... **Mauricio Bautista Porras**

CONTENIDO

	pág.
Introduction	10
1 CONVOLUTIONAL NEURAL NETWORK	13
2 QUANTIZATION STRATEGY	15
2.1 TensorFlow Quantization	15
2.2 Heuristic Quantization	16
3 IMPLEMENTATION IN THE MICROCONTROLLER	21
4 RESULTS	23
5 CONCLUSIONS	24
BIBLIOGRAPHY	25

LISTA DE FIGURAS

	pág.
1 Histogram of the layer's output data Fully Connected 1	18
2 Convolutional Neural Network with fusion of some layers	21

LISTA DE TABLAS

	pág.
1 Convolutional Neural Network Architecture	14
2 Accuracy of Tensorflow methods	16
3 Layer output range on CNN	17
4 Neural Network accuracy for each value of Q	19
5 Outputs and parameters of each layer with the appropriate Q	20
6 Memory used on the MCU implementation	22
7 Accuracy, loss accuracy and supported for each quantization methods	23

RESUMEN

TÍTULO: : IMPLEMENTATION OF A QUANTIZED CONVOLUTIONAL NEURAL NETWORK FOR AUTOMATIC DETECTION OF ATRIAL FIBRILLATION IN AN 8-BIT MICROCONTROLLER.¹

AUTOR: MAURICIO BAUTISTA PORRAS Y LAURA CRISTINA MARTINEZ CRUZ ²

PALABRAS CLAVE: FIBRILACIÓN AURICULAR, MÉTODO HEURÍSTICO, MICROCONTROLADOR, RED NEURONAL CONVOLUCIONAL, CUANTIZACIÓN.

DESCRIPCIÓN:

La fibrilación auricular (FA) es una enfermedad silenciosa que es de difícil diagnóstico porque sus síntomas son esporádicos, tiene una alta tasa de mortalidad en el mundo cuando se diagnostica tarde. Actualmente, las redes neuronales convolucionales (CNN) son una herramienta importante utilizada para el diagnóstico de enfermedades como fibrilación auricular, cáncer de mama, entre otras. Sin embargo, las CNN tienen una alta demanda computacional y de memoria, lo que dificulta su implementación en dispositivos con bajos recursos computacionales. Un tema muy activo en la investigación son las redes neuronales cuantizadas ya que son una solución para reducir la cantidad de recursos informáticos y de memoria. Nuestro objetivo es implementar el proceso de inferencia de una CNN en un microcontrolador de 8 bits (ATMEGA2560) mediante el uso de estrategias de cuantización. Se probaron varias técnicas de cuantización de 8 bits antes de implementar la CNN en el microcontrolador. La implementación final se realizó mediante un método heurístico que llamamos cuantificación dinámica de capa. Este método nos permite lograr una forma efectiva de reducir la complejidad computacional de CNN y sus requerimientos de memoria. Nuestros resultados muestran una precisión del 89,48 %. Este trabajo es la primera etapa de un macroproyecto que tiene como objetivo construir un dispositivo portátil altamente confiable para la detección de fibrilación auricular.

¹Trabajo de Grado.

²Facultad: Ingenierías Fisicomecánicas. Director: Dr.Carlos Fajaro Ariza.

ABSTRACT

TITLE: Implementation of a Quantized Convolutional Neural Network for automatic detection of Atrial Fibrillation in an 8-bit microcontroller.³

AUTOR: MAURICIO BAUTISTA PORRAS Y LAURA CRISTINA MARTINEZ CRUZ.⁴

KEY WORDS: ATRIAL FIBRILLATION, HEURISTIC METHOD, MICROCONTROLLER, CONVOLUTIONAL NEURAL NETWORK, QUANTIZATION.

DESCRIPTION:

Atrial fibrillation (AF) is a silent disease that is difficult to diagnose because its symptoms are sporadic. This disease has a high mortality rate in the world when it is diagnosed late. Currently, convolutional neural networks (CNN) are an important tool used for the diagnosis of diseases as AF among others. However, CNN are both computationally and memory intensive, making them difficult to deploy in devices with low computational resources. Quantized networks are a solution to reduce the amount of computing and memory resources. We aim to implement the inference process of a CNN into an 8-bit microcontroller (ATMEGA2560) by using quantization strategies. Several 8-bits quantization techniques were tested before implement the CNN into the microcontroller. The final implementation was done by a heuristic method that we called Dynamic Layer Quantization. This method allows us to achieve an effective way to reduce the computational complexity of CNN and its memory requirements. Our results show an accuracy of 89.48%. This work is the first stage of a project that aims to build a highly reliable portable device for the detection of AF.

³Bachelor Thesis.

⁴Facultad: Ingenierías Fisicomecánicas. Director: Dr. Carlos Fajardo Ariza.

INTRODUCCIÓN

Atrial fibrillation (AF) is a type of cardiac arrhythmia, characterized by very rapid and uncoordinated atrial activity. This disorder of the electrical signals of the heart has important clinical implications. AF patients have a high risk of stroke and thromboembolism. Furthermore, this disease can be paroxysmal and asymptomatic, making its early diagnosis difficult⁵.

Several studies have proposed the convolutional neural networks (CNN) for the detection of atrial fibrillation achieving high levels of accuracy⁶⁷⁸. By achieving high performance, CNN-based methods demand high amount computation and memory resources. The demand of resources become challenging for the inference of CNN in integrated circuit applications as microcontrollers.

Quantization is an effective strategy to solve computing demand and memory resources. In the past, various CNN quantization methodologies have been studied to perform hardware implementation⁹. CNN quantization requires a set of methodologies to reduce the size of the architecture and the weights of the neural network¹⁰¹¹¹². Quantization seeks to maintain the original accuracy as much as possible¹³.

⁵Lip, Gregory Y H et al. “Atrial fibrillation”. In: *Nature Reviews Disease Primers* 2.1 (2016), p. 16016. ISSN: 2056-676X. DOI: 10.1038/nrdp.2016.16. URL: <https://doi.org/10.1038/nrdp.2016.16>.

⁶Yao, Zhenjie; Zhu, Zhiyong; Chen, Yixin. “Atrial fibrillation detection by multi-scale convolutional neural networks”. In: *2017 20th International Conference on Information Fusion (Fusion)*. IEEE, 2017, pp. 1–6.

⁷Xia, Yong et al. “Detecting atrial fibrillation by deep convolutional neural networks”. In: *Computers in biology and medicine* 93 (2018), pp. 84–92.

⁸Pourbabaee, Bahareh; Roshtkhari, Mehrosan Javan; Khorasani, Khashayar. “Deep convolutional neural networks and learning ECG features for screening paroxysmal atrial fibrillation patients”. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 48.12 (2018), pp. 2095–2104.

⁹Tang, Chuan Zhang; Kwan, Hon Keung. “Multilayer feedforward neural networks with single powers-of-two weights”. In: *IEEE Transactions on Signal Processing* 41.8 (1993), pp. 2724–2727.

¹⁰Kim, Doyun et al. “Convolutional Neural Network Quantization using Generalized Gamma Distribution”. In: *arXiv preprint arXiv:1810.13329* (2018).

¹¹Seo, Sanghyun; Kim, Juntae. “Efficient Weights Quantization of Convolutional Neural Networks Using Kernel Density Estimation based Non-uniform Quantizer”. In: *Applied Sciences* 9.12 (2019), p. 2559.

¹²Athar, Ali. “An Overview of Datatype Quantization Techniques for Convolutional Neural Networks”. In: *arXiv preprint arXiv:1808.07530* (2018).

¹³Kwasniewska, Alicja et al. “Deep Learning Optimization for Edge Devices: Analysis of Training Quantization Parameters”. In: *IECON 2019-45th Annual Conference of the IEEE Industrial Electronics Society*. Vol. 1. IEEE, 2019, pp. 96–101.

In this study, we aim to quantize the inference process of CNN Castillo-Granados¹⁴ and subsequently implement it in a microcontroller. We tested four different quantization methods. A heuristic method, that we called Dynamic Layer Quantization, was chosen because it provides the least loss of precision. With this method, we assign a different quantization factor to each layer. The Dynamic Layer Quantization method allow us to quantize CNN to 8-bit integers with an accuracy of of 89,48%. The implementation was carried out on the Atmega 2560 microcontroller using only 10% and 54% of the FLASH memory and SRAM respectively.

The next sections are structured as follows: Section 1 summarizes a brief description of the architecture of CNN Castillo-Granados¹⁵. Section 2 describes each of the quantization methods tested. Section 3 describes the implementation in the Atmega 2560 microcontroller. In Section 4 we present the most important results. Finally, this paper is closed with the conclusions in Section 5.

¹⁴Castillo, Jeyson A.; Granados, Yenny C.; Fajardo, Carlos A. “Patient-Specific Detection of Atrial Fibrillation in Segments of ECG Signals using Deep Neural Networks”. In: *Ciencia E Ingenieria Neogranadina* 30.1 (2020). DOI: <https://doi.org/10.18359/rcin.4156>.

¹⁵Ibid.

1. CONVOLUTIONAL NEURAL NETWORK

The artificial neural networks have an input layer, intermediate layers, and an output layer. The convolutional neural networks (CNN) have three types of layers between their intermediate layers. Those layers are: the convolutional layers that extract information by filters, the grouping layers that reduce the size of the input, and the fully connected layers that classify information extracted from previous layers.

In this work, we implemented the Castillo-Granados CNN¹ that detects atrial fibrillation (AF) with ECG signals. The network was trained using the MIT-BIH database². The ECG signals were stored in vectors of 500 samples at 250 [samples / s]. In³ the authors achieves an accuracy of 97,44 % by using the traditional 64-bit double-precision floating-point format.

The Castillo-Granados CNN is made up of 12 layers. An input layer, which is designed to process 500-data vectors. Four convolutional layers with max-pooling layers. Three fully connected layers with a single output data. The CNN has a total of 9.385 parameters. The Table 1 summarizes the characteristics of the layers.

¹Ibid.

²Goldberger, Ary L et al. "The MIT-BIH Atrial Fibrillation Database". In: (2000). URL: <http://physionet.incor.usp.br/physiobank/database/afdb/>.

³Castillo; Granados; Fajardo, "Patient-Specific Detection of Atrial Fibrillation in Segments of ECG Signals using Deep Neural Networks", op. cit.

TABLE 1: Convolutional Neural Network Architecture

NEURAL NETWORK ARCHITECTURE		
LAYER TYPE	OUTPUT DIMENSION	PARAMETERS
INPUT	(500x1)	0
CONVOLUTION	(474x3)	84
MAX-POOLING	(237x3)	0
CONVOLUTION	(224x10)	430
MAX-POOLING	(112X10)	0
CONVOLUTION	(110x10)	310
MAX-POOLING	(55x10)	0
CONVOLUTION	(52x10)	410
MAX-POOLING	(26x10)	0
FLATTEN	260	0
FULLY-CONNECTED	30	7830
FULLY-CONNECTED	10	310
FULLY-CONNECTED	1	11

2. QUANTIZATION STRATEGY

Four different techniques were tested to quantize the CNN to 8 bits before implementing it in the microcontroller: First, we used two techniques contained in the TensorFlow backend¹, and then we used two heuristic methods based on Fake Quantization².

2.1. TENSORFLOW QUANTIZATION

The TensorFlow library offers a framework for deep learning called TensorFlow Lite³. This tool converts the model into a special storage format (FlatBuffer) with reduced computer resources and memory expenditure. TensorFlow Lite has different optimization methods available to quantize the parameters that are not quantized in the conversion. We apply Quantization for Size and Integer Quantization.

The Quantization for Size method optimizes the network depending on the weight of the operations of the model⁴. Once the Castillo-Granados CNN was quantized, we used the software Netron⁵ for the extraction and review of the quantized weights. This method converts all parameters to 32-bit floating-point, except for the first dense layer, that was converted to an 8-bit integer.

Integer quantization takes all weights and network activation functions to 8-bit integers⁶. This method applied in the Castillo-Granados CNN converted all parameters of all layers to

¹Google. *TensorFlow Lite guide*. Last accessed 2020-01-13. 2020. URL: <https://www.tensorflow.org/lite/guide>. (visited on 04/08/2020).

²Gupta, Suyog et al. "Deep Learning with Limited Numerical Precision". In: 37 (2015). ISSN: 19410093. DOI: 10.1109/72.80206. arXiv: 1502.02551. URL: <http://arxiv.org/abs/1502.02551>; Nagel, Markus et al. "Up or Down ? Adaptive Rounding for Post-Training Quantization". In: *arXiv preprint* (2020). arXiv: arXiv:2004.10568v1.

³Google, *TensorFlow Lite guide*, op. cit.

⁴Google. *Post-training Quantization for Size | TensorFlow Lite*. Last accessed 2020-01-13. 2020. URL: https://www.tensorflow.org/lite/performance/post%5C_training%5C_quant (visited on 04/08/2020).

⁵Roeder, Lutz. *Netron*. Last accessed 2020-01-21. 2020. URL: <https://pypi.org/project/netron/> (visited on 04/08/2020).

⁶Google. *Post-training integer quantization - TensorFlow Lite*. Last accessed 2020-01-13. 2020. URL: https://www.tensorflow.org/lite/performance/post%5C_training%5C_integer%5C_quant (visited on 04/08/2020).

8-bit integers.

Table 2 summarizes the results regarding the accuracy by using Quantization for Size and Integer Quantization.

TABLE 2: Accuracy of Tensorflow methods

TENSORFLOW	
METHOD	ACCURACY
QUANTIZATION FOR SIZE	92,77%
INTEGER QUANTIZATION	63,58%

2.2. HEURISTIC QUANTIZATION

We tested two heuristic methods, which are an adaption of the Stochastic Rounding Method proposed in⁷.

The first method is the static layer quantization, and the second method is the dynamic layer quantization.

The Static Layer Quantization (SLQ) applies the same quantization factor (2^Q) to all layers in the model, according the following equation,

$$B * 2^Q = A \tag{2.1}$$

Where B is the floating-point number to be quantized, Q is the number of bits after the point that will shift to the left. 2^Q is the quantization factor. Finally, A is a version of B shifted $Q - bits$ to the left.

In order to find the appropriate value of Q , we generate the histograms of the inputs and the activations of the network layers for the MIT signals (See section 1). Using the histogram, we analyze the dynamic range of the input data and intermediate values generated through the network. Figure 1 is an example of the histograms in this analysis, in this case, is showing

⁷Gupta et al., “Deep Learning with Limited Numerical Precision”, op. cit.

the distribution of the parameters in the FC1 layer. Note that the values are concentrated in the range of 0 to 6.

Table 3 shows the dynamic range for each layer in the network. In this table, $CONV_i$ and FC_i refer to convolutional and to full connected layers respectively. The data was analyzed after going through the ReLu activation function⁸ except for the output layer. The output layer uses the Sigmoid activation function⁹ to generate a probability distribution. For this reason, we analyze the data in the output layer before the Sigmoid function.

TABLE 3: Layer output range on CNN

LAYER	LAYER OUTPUT RANGE	
CONV1	0	0,761
CONV2	0	1,103
CONV3	0	0,884
CONV4	0	2,196
FC1	0	7,691
FC2	0	11,357
FC3	-28,9	8,486

Then, we use the equation 2.2 to calculate the number of bits needed to represent the integer part.

$$C = \lceil \log_2(|OL_{max}|) \rceil \quad (2.2)$$

Where OL_{max} is the maximum value of the entire neural network, and C is the number of bits to represent the integer part of the number.

The equation 2.3 is used to calculate the number of bits needed to represent the decimal

⁸Agarap, Abien Fred. “Deep learning using rectified linear units (relu)”. In: *arXiv preprint arXiv:1803.08375* (2018).

⁹Zhang, Chao; Woodland, Philip C. “Parameterised sigmoid and ReLU hidden activation functions for DNN acoustic modelling”. In: *Sixteenth Annual Conference of the International Speech Communication Association*. 2015.

part.

$$Q = N_D - C \tag{2.3}$$

Where N_D is the number of bits to represent the data, in our case N_D is 8-bit. And Q is the number of bits to represent the decimal part of the number.

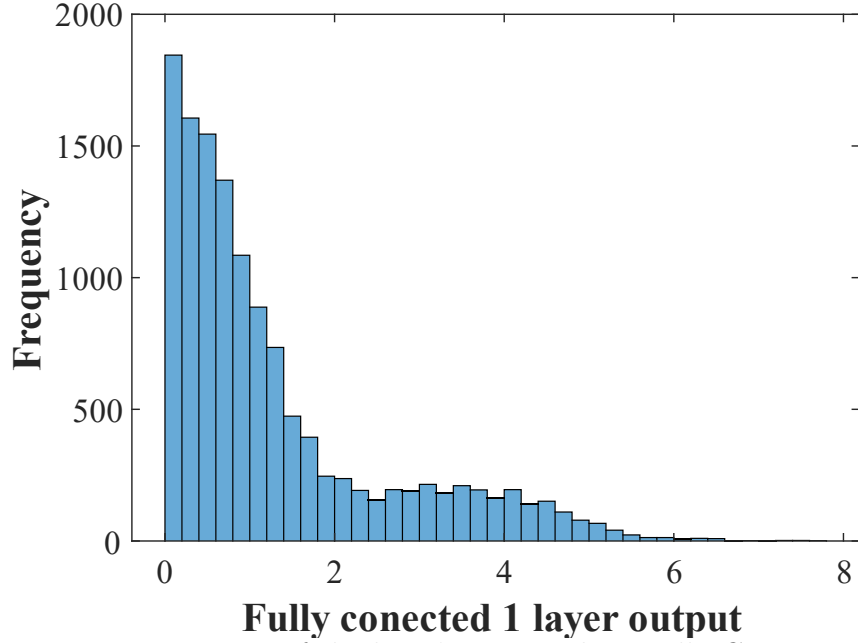


Figure 1: Histogram of the layer's output data Fully Connected 1

In Table 3, the absolute maximum value is found in the FC3 layer. Substituting this number in the equation 2.2 we obtain $C = 5$, then with the equation 2.3 we obtain $Q = 3$.

We tested this quantization factor on the network and the results showed a significantly reduction on the accuracy. Using $Q = 3$, we take the entire range of available values, including those that are far from the highest concentration of data (see Figure 1).

To improve the accuracy, we focus our analysis in the range to the highest data concentration intervals. Next, we apply equations 2.2 and 2.3 to calculate the value of the new Q . We repeat this process until we find the Q value for the maximum accuracy of the quantized CNN. Table 4 shows the results when we choose the highest data concentration to calculate quantization factor.

On the other hand, the quantization process generates an unquantized decimal part. We use truncation and rounding methods to remove the unquantized decimal part¹⁰¹¹. The rounding error is much less than the truncation error. In the Table 4, we show the results when applying the rounding method to eliminate the unquantized decimal part with all the Q found.

TABLE 4: Neural Network accuracy for each value of Q

QUANTIZATION FACTOR Q	ROUNDING METHOD ACCURACY
$Q = 3$	49,532 %
$Q = 4$	54,281 %
$Q = 5$	73,510 %
$Q = 6$	70,337 %
$Q = 7$	65,921 %

Dynamic Layer Quantization (DLQ) is an improvement to the Static Layer Quantization technique. The difference is that we include the dynamic range analysis of the values of the weights per layer as shown in Table 5. We analyze each layer individually and apply a different quantization factor to each one.

To find the appropriate Q value in each of the layers, we analyzed the dynamic range of the parameters and the outputs for all layers, as shown in Table 5. We use equation 2.4 to calculate the number of bits needed to represent the integer part.

$$C_i = \lceil \log_2(\text{Max}\{|OL_i|, |P_i|\}) \rceil \tag{2.4}$$

Where C_i is the number of bits to represent the integer part of layer i_i , OL is the maximum value of the activations of layer i_i and P is the maximum value of the parameters of layer i_i .

¹⁰Sripad, Anekal; Snyder, Donald. “A necessary and sufficient condition for quantization errors to be uniform and white”. In: *IEEE Transactions on Acoustics, Speech, and Signal Processing* 25.5 (1977), pp. 442–448.

¹¹Schwarz, Eric M. “Rounding for quadratically converging algorithms for division and square root”. In: *Conference Record of The Twenty-Ninth Asilomar Conference on Signals, Systems and Computers*. Vol. 1. IEEE. 1996, pp. 600–603.

Then using equation 2.5, we calculate the number of remaining bits to represent the decimal part. The quantization factor Q represents the available decimal bits.

$$Q_i = 8 - C_i \quad (2.5)$$

Where Q_i is the number of bits to represent the decimal part of layer i .

Table 5 shows the dynamic range of the parameters and the outputs of the layers. The highlighted boxes represent the maximum values used to calculate the value of Q_i . The Q values shown in this table generated the best accuracy when quantizing the CNN when to using the DLQ method.

TABLE 5: Outputs and parameters of each layer with the appropriate Q

LAYER	LAYER OUTPUT RANGE		PARAMETERS RANGE		QUANTIZATION FACTOR Q
CONV1	0	0,761	-0,453	0,94	7
CONV2	0	1,103	-1,569	0,788	7
CONV3	0	0,884	-2,116	3	6
CONV4	0	2,196	-2,77	1,136	6
FC1	0	7,691	-2,072	0,99	5
FC2	0	11,357	-1,295	1,321	4
FC3	-28,9	8,486	-24	0,847	3

3. IMPLEMENTATION IN THE MICROCONTROLLER

The implementation in the microcontroller was done, describing the quantized CNN inference and verifying the accuracy of the network by sending the data from an SD module to the ATMEGA2560 microcontroller¹

CNN is made up of thirteen layers, including the input layer as shown in Table 1. The outputs of these thirteen layers must be stored in SRAM memory because their value is constantly changing. To reduce the amount of SRAM memory required, we merge the convolution and max-pooling layers. Figure 2 shows how was reduced the number of outputs from 13 to 9. We calculate all the output values of the convolutional layers, but we only store the outputs of the max-pooling layers.

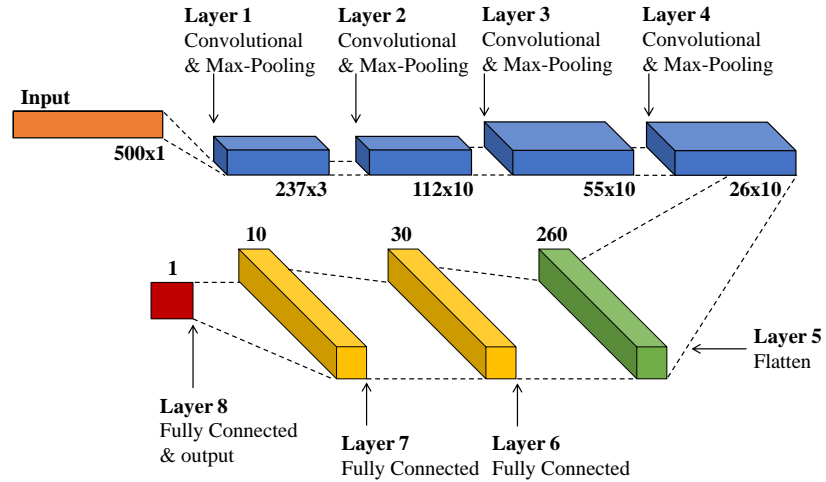


Figure 2: Convolutional Neural Network with fusion of some layers

At the output layer, negative values correspond to a non-fibrillated signal and positive values to fibrillated signals. We use this fact to change the activation function Sigmoid², in

¹Atmel. “ATmega 640/V-1280/V-1281/V-2560/V-2561/V - Datasheet”. In: (2014), p. 435. URL: https://ww1.microchip.com/downloads/en/devicedoc/atmel-2549-8-bit-avr-microcontroller-atmega640-1280-1281-2560-2561%5C_datasheet.pdf.

²Zhang; Woodland, “Parameterised sigmoid and ReLU hidden activation functions for DNN acoustic modelling”, op. cit.

the FC3 layer, by the function Hard-Limit³. The Hard-Limit function assigns 0 to negative values and 1 to positive values, so this function converts the CNN output into a binary output. This change reduces the number of bits needed to store CNN output data.

Table 6 shows the memory usage in the ATMEGA2560 microcontroller. The network parameters were stored in FLASH memory because the numeric value does not change after stored. To store the network parameters we use the PROGMEM function⁴. The intermediate values are stored in SRAM memory because these values change with each inference. EEPROM memory was not used because it is small and has few write cycles. The implemented network has an inference time of 677[ms] using an oscillator of 12[MHz]. This inference time is enough for this application because the intervals of the ECG signal are collected each two seconds.

TABLE 6: Memory used on the MCU implementation

ATMEGA 2560			
MEMORY	AVAILABLE	USED	% USE
FLASH	256 KB	25,56 KB	10 %
SRAM	8 KB	4,44 KB	54 %

We developed a code in MATLAB to emulate the inference process carried out in the microcontroller. We use this code to validate the accuracy of the network for the quantized methods tested in this work.

³*Hard-limit transfer function - MATLAB hardlim.* Last accessed 2020-01-22. 2020. URL: <https://la.mathworks.com/help/deeplearning/ref/hardlim.html> (visited on 04/08/2020).

⁴*Arduino Reference - PROGMEM.* Last accessed 2020-02-11. 2020. URL: <https://www.arduino.cc/reference/tr/language/variables/utilities/proGMEM/> (visited on 04/08/2020).

4. RESULTS

Table 7 summarizes the results obtained in this work. In this table, we compare the accuracy of the quantization methods against the CNN at 64-bits double-precision floating-point format (Unquantized). We also show the percentage of accuracy loss regarding unquantized network. Finally, we indicate the possibility of implementing each quantized model in an 8-bit microcontroller.

TABLE 7: Accuracy, loss accuracy and supported for each quantization methods

METHOD	ACCURACY	% LOSS	SUPPORTED 8-BIT MCU
UNQUANTIZED	97,44 %	—	N
TF LITE QUANTIZATION FOR SIZE	92,77 %	4,67 %	N
TF LITE INTEGER QUANTIZATION	63,58 %	33,86 %	Y
SLQ	73,51 %	23,93 %	Y
DLQ	89,48 %	7,96 %	Y

The results show that the network quantized with the Dynamic Layer Quantization method has the best accuracy.

5. CONCLUSIONS

In this work, four quantization methods were tested to implement the Castillo-Granados CNN into an 8-bit microcontroller (ATMEGA2560). First, we used the TensorFlow backend to quantize the network. Our result showed that TensorFlow quantization reduces significantly the precision of the network. The best results regarding accuracy were obtained by a method that we called Dynamic Layer Quantization. Our results showed that the implementation achieves an accuracy of 89.48%, and only uses the 10% in the FLASH memory and 54% in the SRAM memory. A future job will focus on using post quantization strategies, which have proved to improve the accuracy of fake-quantized CNN¹. This work is part of a project that seeks the development of a low-cost portable device for the early diagnosis of AF.

¹Song Han. “EFFICIENT METHODS AND HARDWARE FOR DEEP LEARNING”. PhD thesis. STANFORD UNIVERSITY, 2017.

BIBLIOGRAPHY

AGARAP, Abien Fred. “Deep learning using rectified linear units (relu)”. In: *arXiv preprint arXiv:1803.08375* (2018).

Arduino Reference - PROGMEM. Last accessed 2020-02-11. 2020. URL: <https://www.arduino.cc/reference/tr/language/variables/utilities/progmem/> (visited on 04/08/2020).

ATHAR, Ali. “An Overview of Datatype Quantization Techniques for Convolutional Neural Networks”. In: *arXiv preprint arXiv:1808.07530* (2018).

ATMEL. “ATmega 640/V-1280/V-1281/V-2560/V-2561/V - Datasheet”. In: (2014), p. 435. URL: https://ww1.microchip.com/downloads/en/devicedoc/atmel-2549-8-bit-avr-microcontroller-atmega640-1280-1281-2560-2561%5C_datasheet.pdf.

CASTILLO, Jeyson A.; GRANADOS, Yenny C.; FAJARDO, Carlos A. “Patient-Specific Detection of Atrial Fibrillation in Segments of ECG Signals using Deep Neural Networks”. In: *Ciencia E Ingenieria Neogranadina* 30.1 (2020). DOI: <https://doi.org/10.18359/rcin.4156>.

GOLDBERGER, Ary L et al. “The MIT-BIH Atrial Fibrillation Database”. In: (2000). URL: <http://physionet.incor.usp.br/physiobank/database/afdb/>.

GOOGLE. *Post-training integer quantization - TensorFlow Lite*. Last accessed 2020-01-13. 2020. URL: https://www.tensorflow.org/lite/performance/post%5C_training%5C_integer%5C_quant (visited on 04/08/2020).

GOOGLE. *Post-training Quantization for Size | TensorFlow Lite*. Last accessed 2020-01-13. 2020. URL: https://www.tensorflow.org/lite/performance/post%5C_training%5C_quant (visited on 04/08/2020).

— *TensorFlow Lite guide*. Last accessed 2020-01-13. 2020. URL: <https://www.tensorflow.org/lite/guide>. (visited on 04/08/2020).

GUPTA, Suyog; AGRAWAL, Ankur; GOPALAKRISHNAN, Kailash; NARAYANAN, Pritish. “Deep Learning with Limited Numerical Precision”. In: 37 (2015). ISSN: 19410093. DOI: 10.1109/72.80206. arXiv: 1502.02551. URL: <http://arxiv.org/abs/1502.02551>.

Hard-limit transfer function - MATLAB hardlim. Last accessed 2020-01-22. 2020. URL: <https://la.mathworks.com/help/deeplearning/ref/hardlim.html> (visited on 04/08/2020).

KIM, Doyun; YIM, Han Young; HA, Sanghyuck; LEE, Changgwun; KANG, Inyup. “Convolutional Neural Network Quantization using Generalized Gamma Distribution”. In: *arXiv preprint arXiv:1810.13329* (2018).

KWASNIEWSKA, Alicja; SZANKIN, Maciej; OZGA, Mateusz; WOLFE, Jason; DAS, Arun; ZAJAC, Adam; RUMINSKI, Jacek; RAD, Paul. “Deep Learning Optimization for Edge Devices: Analysis of Training Quantization Parameters”. In: *IECON 2019-45th Annual Conference of the IEEE Industrial Electronics Society*. Vol. 1. IEEE. 2019, pp. 96–101.

LIP, Gregory Y H et al. “Atrial fibrillation”. In: *Nature Reviews Disease Primers* 2.1 (2016), p. 16016. ISSN: 2056-676X. DOI: 10.1038/nrdp.2016.16. URL: <https://doi.org/10.1038/nrdp.2016.16>.

NAGEL, Markus; AMJAD, Rana; BAALEN, Mart; LOUIZOS, Christos; BLANKEVOORT, Tijmen. “Up or Down ? Adaptive Rounding for Post-Training Quantization”. In: *arXiv preprint* (2020). arXiv: arXiv:2004.10568v1.

POURBABAEE, Bahareh; ROSHTKHARI, Mehrsan Javan; KHORASANI, Khashayar. “Deep convolutional neural networks and learning ECG features for screening paroxysmal atrial fibrillation patients”. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 48.12 (2018), pp. 2095–2104.

ROEDER, Lutz. *Netron*. Last accessed 2020-01-21. 2020. URL: <https://pypi.org/project/netron/> (visited on 04/08/2020).

SCHWARZ, Eric M. “Rounding for quadratically converging algorithms for division and square root”. In: *Conference Record of The Twenty-Ninth Asilomar Conference on Signals, Systems and Computers*. Vol. 1. IEEE. 1996, pp. 600–603.

SEO, Sanghyun; KIM, Juntae. “Efficient Weights Quantization of Convolutional Neural Networks Using Kernel Density Estimation based Non-uniform Quantizer”. In: *Applied Sciences* 9.12 (2019), p. 2559.

SONG HAN. “EFFICIENT METHODS AND HARDWARE FOR DEEP LEARNING”. PhD thesis. STANFORD UNIVERSITY, 2017.

SRIPAD, Anekal; SNYDER, Donald. “A necessary and sufficient condition for quantization errors to be uniform and white”. In: *IEEE Transactions on Acoustics, Speech, and Signal Processing* 25.5 (1977), pp. 442–448.

TANG, Chuan Zhang; KWAN, Hon Keung. “Multilayer feedforward neural networks with single powers-of-two weights”. In: *IEEE Transactions on Signal Processing* 41.8 (1993), pp. 2724–2727.

XIA, Yong; WULAN, Naren; WANG, Kuanquan; ZHANG, Henggui. “Detecting atrial fibrillation by deep convolutional neural networks”. In: *Computers in biology and medicine* 93 (2018), pp. 84–92.

YAO, Zhenjie; ZHU, Zhiyong; CHEN, Yixin. “Atrial fibrillation detection by multi-scale convolutional neural networks”. In: *2017 20th International Conference on Information Fusion (Fusion)*. IEEE. 2017, pp. 1–6.

ZHANG, Chao; WOODLAND, Philip C. “Parameterised sigmoid and ReLU hidden activation functions for DNN acoustic modelling”. In: *Sixteenth Annual Conference of the International Speech Communication Association*. 2015.