

A 22-nm All-Digital Time-Domain Neural Network Accelerator for Precision In-Sensor Processing

Ahmed M. Mohey¹, Graduate Student Member, IEEE, Jelin Leslin², Gaurav Singh³, Member, IEEE, Marko Kosunen⁴, Member, IEEE, Jussi Ryyänen⁵, Senior Member, IEEE, and Martin Andraud, Member, IEEE

Abstract—Deep neural network (DNN) accelerators are increasingly integrated into sensing applications, such as wearables and sensor networks, to provide advanced in-sensor processing capabilities. Given wearables' strict size and power requirements, minimizing the area and energy consumption of DNN accelerators is a critical concern. In that regard, computing DNN models in the time domain is a promising architecture, taking advantage of both technology scaling friendliness and efficiency. Yet, time-domain accelerators are typically not fully digital, limiting the full benefits of time-domain computation. In this work, we propose an all-digital time-domain accelerator with a small size and low energy consumption to target precision in-sensor processing like human activity recognition (HAR). The proposed accelerator features a simple and efficient architecture without dependencies on analog nonidealities such as leakage and charge errors. An eight-neuron layer (core computation layer) is implemented in 22-nm FD-SOI technology. The layer occupies $70 \times 70 \mu\text{m}^2$ while supporting multibit inputs (8-bit) and weights (8-bit) with signed accumulation up to 18 bits. The power dissipation of the computation layer is $576 \mu\text{W}$ at 0.72-V supply and 500-MHz clock frequency achieving an average area efficiency of $24.74 \text{ GOPS}/\text{mm}^2$ (up to $544.22 \text{ GOPS}/\text{mm}^2$), an average energy efficiency of $0.21 \text{ TOPS}/\text{W}$ (up to $4.63 \text{ TOPS}/\text{W}$), and a normalized energy efficiency of $13.46 \text{ 1b-TOPS}/\text{W}$ (up to $296.30 \text{ 1b-TOPS}/\text{W}$).

Index Terms—Edge computing, human activity recognition (HAR), inertial measurement unit (IMU), in-sensor processing, multiply-and-accumulate multiply and accumulate (MAC), neural network accelerator, smart sensor interface, time-domain signal processing.

I. INTRODUCTION

THE development of embedded sensing applications is crucial to continue advancements in areas such as the Internet of Things (IoT), autonomous driving, and, more generally, new usage for connected objects around us. In particular, sensing and data processing based on inertial measurement units (IMUs) is a widely used principle in

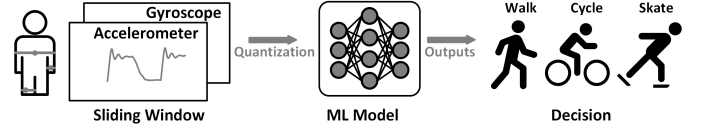


Fig. 1. ML applied to HAR [4].

various applications, for instance, in wearables for the general consumer market, in biomedical and health monitoring applications, or in automotive. Many of these applications are evolving toward sensory systems that combine accurate sensors with advanced signal processing or machine learning (ML) capabilities to develop more intelligent sensing systems, for example, motion or activity recognition. An example is human activity recognition (HAR), based on smart wearable sensors [1], [2], which is getting popular for the daily monitoring of health data (sleep and recuperation) and physical activities [3]. Fig. 1 shows an example of ML applied to HAR, to classify a range of activities performed by a user equipped with a wearable [4]. Here, the system starts by extracting features from the IMU sensors (a gyroscope and an accelerometer) and dividing them into sliding windows. These features can consist of, for example, statistical characteristics of the signal, such as mean, variance, or peak detection. Then, the obtained features are quantized and fed to a classifier that detects the current activity.

Optimizing the energy efficiency and the form factor of embedded sensing systems can be realized with in-sensor processing, where the sensor device is directly intercoupled with advanced data-processing capabilities offered by ML, particularly deep learning and deep neural network (DNN) models. However, the hard constraints required for the design of typical low-power tiny wearables necessitate utilizing highly efficient hardware to execute these DNNs. In this context, dedicated hardware accelerators are necessary to execute computation-hungry DNN models. These accelerators typically focus on executing efficient multiply and accumulate (MAC) operations, which are at the core of every DNN computation (i.e., every neuron performs a weighted sum of its inputs). The MAC operation realized by one neuron is expressed as

$$\text{MAC} = \sum_{i=1}^N X_i W_i \quad (1)$$

Received 8 April 2024; revised 15 August 2024 and 11 October 2024; accepted 6 November 2024. Date of publication 19 November 2024; date of current version 6 December 2024. This work was supported in part by the Academy of Finland Project WHISTLE under Grant 332218. (Corresponding author: Ahmed M. Mohey.)

Ahmed M. Mohey, Jelin Leslin, Gaurav Singh, Marko Kosunen, and Jussi Ryyänen are with the Department of Electronics and Nanoengineering, Aalto University, 02150 Espoo, Finland (e-mail: ahmed.mohey@aalto.fi).

Martin Andraud is with the Department of Electronics and Nanoengineering, Aalto University, 02150 Espoo, Finland, and also with UCLouvain, ICTEAM, 1348 Leuven-la-Neuve, Belgium.

Digital Object Identifier 10.1109/TVLSI.2024.3496090

where N is the number of connected neurons, X_i refers to each neuron's inputs, and W_i represents the synaptic weights. In classical processors, the MAC computation is heavily dominated by memory transfers, as these processors have separate calculation and memory units. Dedicated accelerators solve this bottleneck by bringing the MAC operation closer to memory, using Near-Memory Computing approaches, or directly inside the memory, using Compute-In Memory approaches [5]. As most DNN inference routines can be accommodated with low-resolution integer computation of 5–8 integer bits [6], [7], MAC operations can be realized either in the digital or analog domain, which implies design tradeoffs. Analog-domain MAC has shown better efficiency than digital implementations for low bit-width (up to 6–8 bits integer) [8] but falls steeply for higher bit-widths, due to noise constraints requiring four times more power per extra bit of precision [6]. Hence, digital-domain MAC is typically preferred for resolutions higher than 6–8 integer bits using their lower susceptibility to noise and better scaling properties. In this context, similar to other circuit blocks such as ADCs or sensor interfaces, time-domain implementations of MAC operations could offer a tradeoff between analog and digital computing methods. Indeed, time-domain MACs can be implemented using mostly digital circuits, benefiting from scaling; they consume less power than digital MACs due to lower switching activities (thanks to the reduction of digital buses) [9], [10]. Hence, multiple time-based computing methods have been proposed [11], [12], [13], [14], [15], [16]. Time-based MAC accelerators to attain energy efficiencies of 10–100 TOPS/W [11], [12], [16], which are among the best reported for AI accelerators [17]. However, the time-based nature limits the throughput of the system to 0.1–5 GOPS, hence time-based MACs are generally suited for applications requiring efficient computing but lower throughput requirements, for instance, HAR. Despite these good results, challenges remain in the development of time-based MAC accelerators. For instance, most existing implementations are not fully digital and rely on analog circuits for signal generation. Hence, these blocks suffer from analog nonidealities, in particular, leakage and charge errors, which can limit the accuracy of the system. In addition, they tend toward increasing complexity when multibit inputs/weights are enabled requiring a relatively complex design of the key components.

To solve these issues, this article presents an all-digital neural network accelerator that utilizes time-domain approaches to realize highly efficient multibit MAC operations with no dependency on analog nonidealities. It consists of an array of combined eight MAC units 8/8/18 (neurons) to act as the core processing layer in a multilayer perceptron (MLP) NN for HAR. The UCI HAR dataset [18] is used to verify the capability of the proposed architecture to process precision sensing data with a signed accumulation of up to 18 bits while achieving more than 90% classification accuracy. The main characteristics of this design are as follows.

- 1) A native precision-scalable and multibit support (up to 8 bits), allowing tailor design as per the application's needs, thereby saving area and power.

- 2) A sequential processing of input features enables the use of an arbitrarily large number of features for the system without redesigning the architecture. For this reason, the output can accumulate up to 18 bits.
- 3) A simple design and compact footprint to achieve a state-of-the-art area per MAC for time-based accelerators of $612.5 \mu\text{m}^2$. It relies on a simple design based on a single clock source in an event-triggered (on-need) activation. In addition, dynamic frequency scaling can be deployed for extra power saving at low computation demands. These features ease the integration of many cores for parallel in-sensor processing applications.

The proposed circuit is simulated using a 22-nm FD-SOI technology. The layer only occupies $70 \times 70 \mu\text{m}$, and it dissipates $576 \mu\text{W}$ at 0.72 V supply and 500-MHz clock frequency. It achieves an average operation rate of 0.12 GOPS (up to 2.67 GOPS), measuring an average energy efficiency of 0.21 TOPS/W (up to 4.63 TOPS/W), and an area efficiency of 24.74 GOPS/mm^2 (up to 544.22 GOPS/mm^2). It achieves also an average normalized energy efficiency [19] of 13.46 1b-TOPS/W (up to 296.30 1b-TOPS/W).

This article is organized as follows. Section II discusses the target application (HAR), including model training, quantization, and optimization. Section III reviews the recent development of time-domain neural accelerators. Section IV presents the design objectives and the proposed architecture. The circuit implementation and simulation results are discussed in Section V. A conclusion is given in Section VI.

II. TARGET APPLICATION: ONLINE ACTIVITY RECOGNITION

The proposed accelerator will be evaluated on online activity recognition applications. These applications require moderate inference speed (in the order of hundreds of milliseconds) but high energy efficiency to be run in the background of portable devices such as wearables or smartphones. Two different activity recognition benchmarks will be compared, the UCI HAR [18] containing six activities, and PAMAP containing 18 output activities [20]. In the following, we highlight the analysis for the HAR dataset, yet a similar analysis has been performed for PAMAP and illustrated in Figs. 2 and 3.

A. Human Activity Recognition

The HAR dataset was collected from experiments carried out with a group of 30 volunteers within an age bracket of 19–48 years. Each person performed six activities (walking, walking upstairs, walking downstairs, sitting, standing, and laying) wearing a smartphone (Samsung Galaxy S II) on the waist. The smartphone's embedded accelerometer and a gyroscope were used to capture three-axial linear acceleration and three-axial angular velocity at a constant rate of 50 Hz. From these acquired signals, a large number of features can be extracted with preprocessing blocks, typically involving noise filtering, applying various time- and frequency-domain transformations, and extracting statistics and other information like mean, variance, max, and entropy from these signals. These preprocessing steps are designed to capture the characteristics

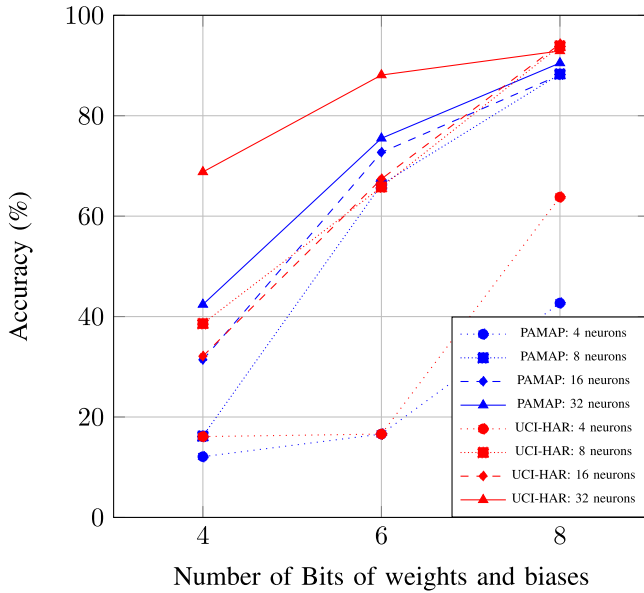


Fig. 2. Model accuracy versus number of bits in weights and bias for different neuron configurations.

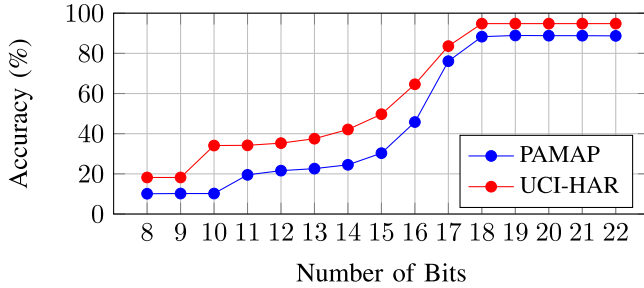


Fig. 3. Model accuracy versus number of bits in an HW accumulator.

of human activities in a form that is suitable for ML models. Although many ML models performed well on this dataset, enhancing the hardware efficiency typically requires heavily quantized data and computation with integer numbers, with the constraint of doing the inference in integer hardware. In that regard, ML models such as decision trees, SVMs, forests, and regression can have their performance drastically decreased under heavily quantized weights. On the other hand, neural networks have shown strong resilience and efficiency after quantization [21]. Given that the HAR dataset is composed of time-series data, which are essentially 1-D, classical MLPs can be aptly suited. MLPs can model relationships in sequential data effectively without necessitating the spatial context exploited by 2-D convolutional neural networks (Conv2D) or more advanced architectures like recurrent neural networks. In addition, the relative simplicity and lower computational demands of MLPs make them an efficient choice for HAR.

B. Model Training and Quantization

The HAR dataset is composed of 561 features extracted to describe the activity window. As detailed earlier, the chosen model is an MLP computing all 561 input features, where each neuron necessitates accumulating 561 multiplications. This accumulation of numerous products leads to substantial growth

in the dynamic range of intermediate values, up to 16 bits [7]. Various models address this accumulation problem within a limited-resources context (e.g., a wearable device) [4], [22]. These works highlight that quantizing weights and biases to lower bitwidths reduces the model's computational burden, yet the inference is performed on a 32-bit microcontroller, thereby setting the maximum dynamic range of intermediate values to 32 bits. This computation bottleneck is solved in the proposed architecture by performing a successive accumulation of all the features in a single high-resolution accumulation block, accommodating a large dynamic range. The accumulator is the only computation block with higher resolution, as detailed in Section IV.

The MLP model is composed of an input layer of 561 features, one hidden layer with a variable number of neurons (4, 8, 16, and 32), and one output layer of six neurons for HAR (one for each activity). The model is initially trained in Keras and compressed using the TFlite framework to perform quantization-aware training [23], which converts the pretrained 32-bit floating point (FP) weights to 8-bit integer (INT) values. In that regard, it has been proved that INT8 provides a superior efficiency than FP8 for DNN inference [7]. Quantization-aware training reduces the memory footprint of the model by up to 90%, depending on the final NN architecture. From there, a uniform quantization is done to further obtain the weights and biases fitting the proposed architecture (4–8 bits).

C. Model Selection

Fig. 2 shows the obtained classification accuracy values for both HAR and PAMAP benchmarks for various configurations of the hidden MLP layer (varying number of bits in the weights and biases for four different neuron configurations). First, we target to attain 90% accuracy on both benchmarks, which requires at least an 8-bit quantization of weights and biases. Then, we evaluate the minimum number of neurons in the hidden layer, common to both benchmarks, to find the best compromise between accuracy and computation efficiency around this baseline. For HAR, the best accuracy is for 16 hidden neurons with 94.3% (32 neurons have a 92.9% accuracy). Using eight neurons, the accuracy drop is only 0.4%, which we consider acceptable considering the substantial gain in terms of computations in the hidden layer (50%). Reducing to four neurons drops the accuracy to 63% which we considered too low for the application. For PAMAP, the maximum accuracy in our design space was 90.5% with 32 hidden neurons. This accuracy reduces to 88.3% with eight neurons, which we consider again acceptable considering 75% less computational load for this layer. Reducing to four neurons significantly degrades the accuracy to 42%. Fig. 3 shows how the two benchmarks perform with varying numbers of bits in the hidden layer's accumulator. The PAMAP dataset has 243 variables but it still requires a high accumulator precision, similar to the HAR which had 561 inputs, to prevent a significant accuracy loss. Hence, the chosen configuration for both datasets is to use 8-bit weights and biases, eight neurons in the hidden layer, and 18-bit resolution for the accumulator. It should be noted that the proposed hardware is fit for this

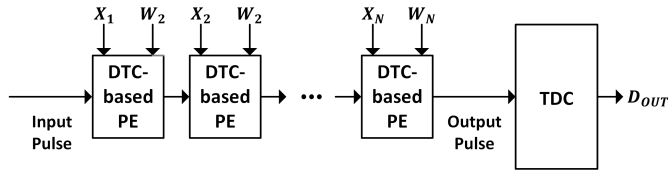


Fig. 4. DTC-based MAC.

application but could be made more generic by integrating spare hardware neurons.

III. BACKGROUND ON TIME-DOMAIN ACCELERATORS

Time-domain accelerators represent and compute data in the time domain. A time-based MAC operation can be encoded in the timing properties of a pulse, such as delay, phase, or frequency. As such, time-domain computing allows the representation of multibit data in a single wire, which can reduce the dynamic power dissipation. Fig. 4 depicts a baseline time-based architecture, based on a digital-to-time converter (DTC). The DTC converts multibit digital data (inputs X and weights W) into a single pulse, acting as an interface between the digital and time domain. To perform a MAC operation, the DTC is followed by an accumulation phase, governed by the output pulse. The signal is converted back to digital using a time-to-digital converter (TDC).

A. Baseline Time-Based Computing

In basic time-based computing, multibit data can be presented by pulsewidth modulation (PWM), such that the pulsewidth represents the magnitude of the data [15], [16], [24]. Hence, DTC-based processing elements (PEs) are used to output PWM pulses with a duration proportional to the inputs or the inputs-weights product (X, W). Thus, a single PE can perform time-domain multiplication, while cascading multiple PEs in a chain can carry out addition (accumulation) as the input pulse propagates through the chain. This process is expressed as

$$\begin{aligned} T_{out} &= T_{in} + (X_1 W_1 + X_2 W_2 + \dots + X_N W_N) \Delta\tau \\ &= T_{in} + \left(\sum_{i=1}^N X_i W_i \right) \Delta\tau \end{aligned} \quad (2)$$

where T_{in} is the duration of the input pulse and $\Delta\tau$ is the DTC resolution. The PE can be realized in different manners.

1) *Single-Bit PE Multiplication*: Fig. 5(a) illustrates a first DTC-based PE implementation [15]. It performs a single-bit multiplication using a voltage generator and a starved inverter with cascaded nMOS devices. The voltage generator produces four distinct voltage levels V_0, V_1, V_2 , and V_3 , each increasing the duration of the input pulse (T_{in}) by $1\Delta\tau, 2\Delta\tau, 3\Delta\tau$, or $4\Delta\tau$, respectively. One voltage level is selected based on the digital input (X_i) using a decoder. The delay path is only activated when the binary weight $W_i = '1'$; otherwise, the DTC is disabled by skipping the delay path. More delay levels can be generated by producing more distinct voltage levels or by cascading more PWM stages with binary weighted loads.

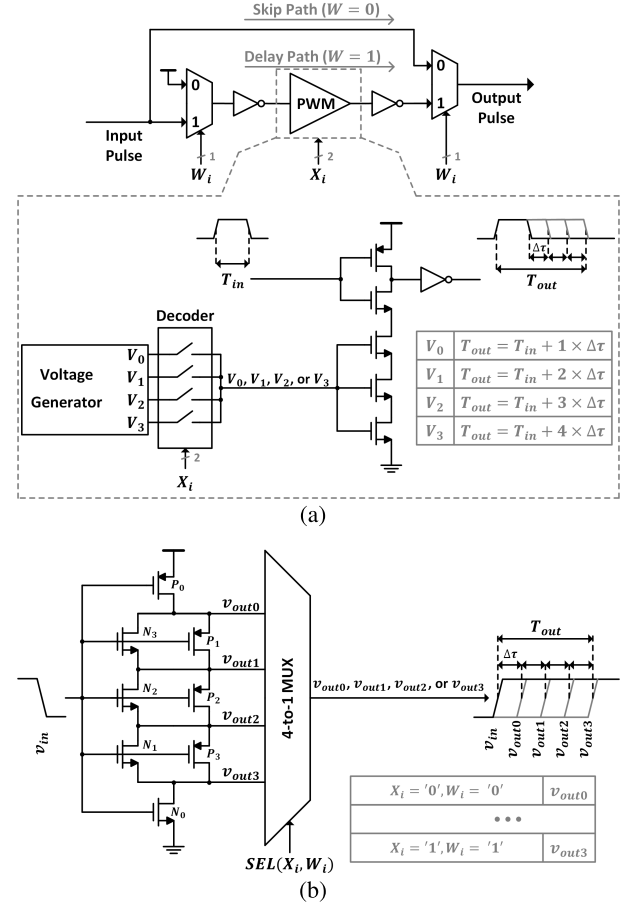


Fig. 5. DTC-based processing: (a) DTC-based PE implementation using PWM circuitry [15]. (b) DTC-based PE implementation using a ladder inverter [16].

2) *Multibit PE Multiplication*: To avoid using a voltage generator and allowing multibit multiplication, Fig. 5(b) shows a PE using ladder inverters [16], [25]. It relies on the sequential charging of v_{out0} to v_{out3} when v_{in} decreases. v_{out0} directly transits from low-to-high via the pMOS device P_0 . While v_{out1} transits via pMOS devices P_0 and P_1 , thus it can only transit after the transition of v_{out0} . Similarly, v_{out2} only transits after v_{out1} , and v_{out3} transits after v_{out2} . Following the ladder inverters, a multiplexer is used to map the input-weight product to one of the outputs and therefore set T_{out} with respect to the rising edge of v_{in} .

3) *Challenges of the basic architecture*: The architecture in Fig. 4 dedicates separate DTC-based PE for each input/weight pair (spatially unrolled). Hence, it relies on the ability of the DTC to output accurate absolute delays proportional to the input or the input-weight product. Utilizing this architecture for precision applications that require generating many delay levels (e.g., 255 levels for an 8-bit product) can be tedious because it necessitates using complex delay calibration and tuning and/or using larger devices to overcome mismatches between the PEs, reducing the overall area and power efficiencies. In a trial to overcome this limitation, more advanced architectures have been presented [11], [12].

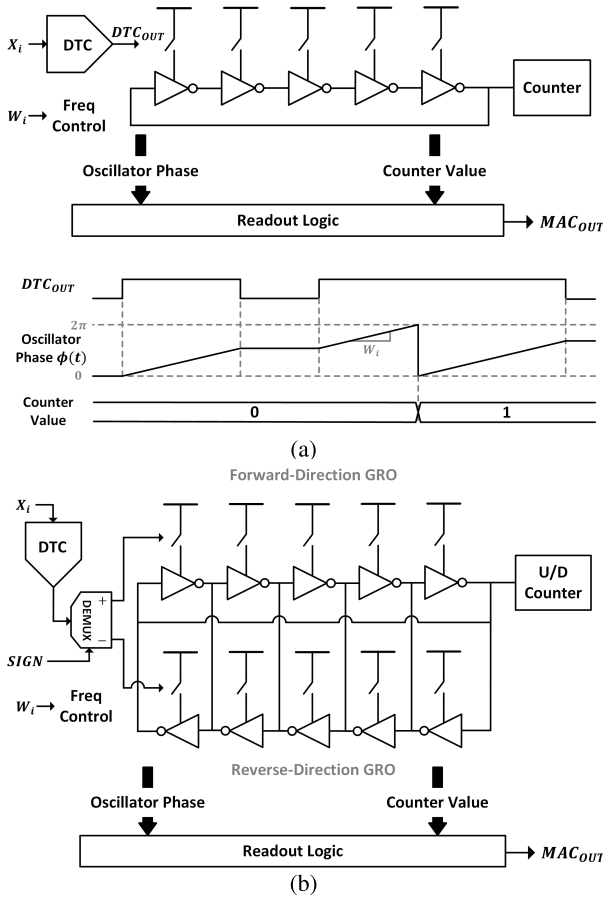


Fig. 6. GRO-based MAC [11]. (a) Phase accumulation using GRO. (b) Signed accumulation using bi-directional GRO.

B. Advanced Time-Based Computing

1) *Single DTC With GROs*: The area efficiency can be improved by employing a single DTC, with the accumulation performed in the phase domain using a gated ring oscillator (GRO) [11], as Fig. 6(a) shows, when the oscillator is enabled by the DTC output (DTC_{OUT}), its phase ϕ advances by $\phi[n] = \phi[n-1] + (2\pi/10)X_i W_i$ for a 5-stage oscillator, otherwise it holds its phase information. The pulsewidth of DTC_{OUT} is proportional to the digital input X_i , and the oscillation frequency is linearly controlled by the weight W_i . To realize continuous accumulation, a counter detects when the GRO phase returns to 0. The readout logic samples the GRO phase and the counter output to generate the MAC operation result. Fig. 6(b) shows how to realize a signed accumulation by utilizing bi-directional GRO. The output of the DTC is provided to the forward-direction GRO when the sign is positive and to the reverse-direction GRO when the sign is negative. This allows the phase to increment (resp. decrement) when the sign is positive (resp. negative). In this case, an up/down-counter detects when the oscillator returns to its initial state, such that it increments its value in the forward direction and decrements its value in the reverse direction. This architecture is very efficient (see Table II), but nonidealities must be carefully considered, specifically when scaling up. For instance, leakage and charge-injection errors [26] can degrade the oscillator phase information. Also, linearly controlling the oscillator frequency with high bit-width can be challenging.

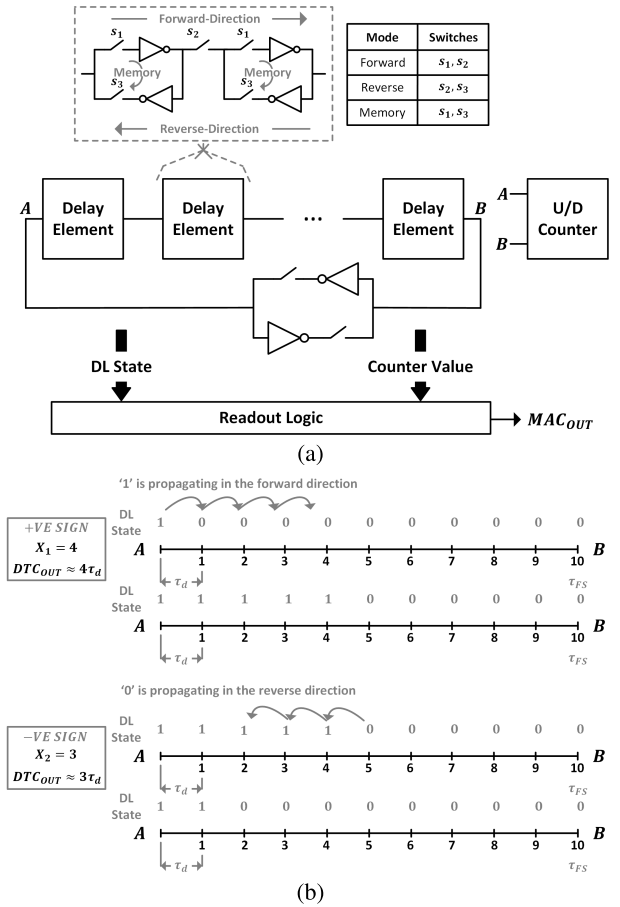


Fig. 7. GDL-based MAC [12]. (a) GDL block diagram and modes of operation. (b) GDL timing diagram.

2) *Single DTC With GDLs*: To provide a fully digital implementation, GROs can be replaced with a bi-directional gated delay line (GDL) [12]. Fig. 7(a) shows that the GDL is enabled by the DTC output, which allows a signal to propagate in the GDL (forward or reverse direction); otherwise, the DL state is preserved by the memory (latch) mode. The state of the delay line increases/decreases when the signal propagates in the forward/backward direction by the amount of the pulsewidth [27]. Fig. 7(b) shows a timing diagram example. When the sign is positive, the GDL state increases ("1" propagates in the forward direction). For $X_1 = 4$, the state (thermometer code) increases by the amount of the pulsewidth $DTC_{OUT} = 4\tau_d$: "0000000000" (initial state) to "1111000000." When the sign is negative, the GDL decreases ("0" propagates in the reverse direction). When the signal reaches the edge of the delay line, a full-scale signal is asserted to allow the signal's complement to propagate in a loop configuration through an inverter. An up/down-counter is used to detect the full-scale condition. This configuration enables long-duration time accumulation. In this architecture, no delay control is implemented. This means that multibit weights are realized either by utilizing a single delay line with configurable length sequentially or by utilizing multiple delay lines for each weight bit to allow parallel operation, which could impact the area efficiency.

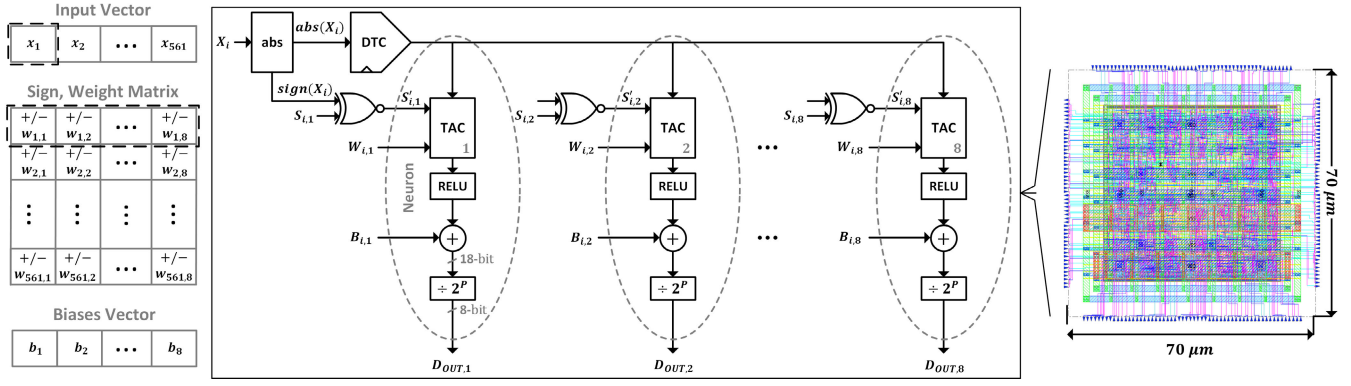


Fig. 8. Proposed eight-neuron computation layer.

C. Summary of the Current Challenges

Analog-based implementations are generally achieving competitive results for low precision of inputs, weights, and outputs (I/W/O), for instance, 1–8/1/8 precision with calibration in [15] or 1/3/4 with calibration in [16]. However, to the best of our understanding, these analog-based implementations follow a similar trend as described in [6]: when the precision increases, the efficiency falls steeply with an increasing implementation cost due to analog nonidealities (matching, leakage, noise, etc.). Thus, these implementations are typically targeting low-precision applications. Using analog blocks such as voltage generators or GROs increases the systems' susceptibility to leakage, noise, and variations, specifically considering modern CMOS technologies and high-precision applications. This may induce computation errors and/or highly increase the system complexity. Replacing GROs with delay lines allows for fully digital implementation, but maintains a relatively high complexity to build multiple proportional delay lines. Hence, this work targets a more elegant and efficient all-digital implementation for in-sensor processing applications.

IV. PROPOSED ARCHITECTURE

Fig. 8 depicts the proposed computation layer for the HAR application (which could be easily adapted to other applications integrating more neurons). In this section, we present the original time-based MAC array based on our previous work [28] and its extension toward a full system.

- 1) The precision of the TAC is upgraded to support up to 18-bit signed accumulation with 8-bit signed inputs and 8-bit weights (as decided in Section II).
- 2) Various postaccumulation functionalities integrated into the neurons including RELU activation, bias addition, and quantization.
- 3) The implementation of an SRAM memory, to store all necessary DNN weights and biases.

A. Time-Domain MAC Circuit

The proposed time-domain MAC architecture [28] is shown in Fig. 9. The topology has been chosen to enable a simple and fully digital implementation, which can easily be scaled according to the needs of in-sensor processing applications.

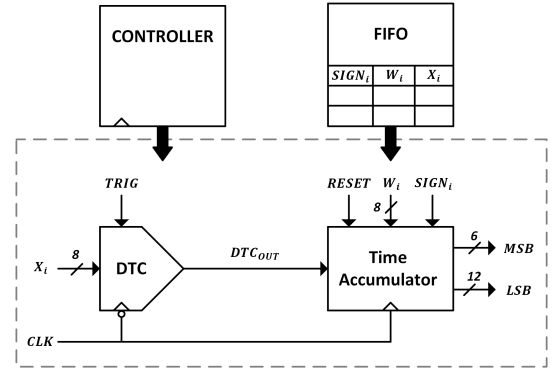


Fig. 9. Proposed time-based architecture.

It is composed of a DTC which converts 8-bit digital input X_i into a PWM signal DTC_{OUT} , and a time accumulator (TAC) which uses 8-bit weights W_i and a sign bit $SIGN_i = 1('1'), -1('0')$ to perform signed accumulation governed by the DTC. An 18-bit accumulation result D_{OUT} (six most-significant-bits MSB, and 12 least-significant-bits LSB) is produced by the TAC state machine. The data required by the DTC and the TAC (X_i , W_i , and $SIGN_i$) are provided by an FIFO (memory macro), while the required control signals (TRIG and RESET) are provided by a controller (sequencer).

1) *Digital to Time Converter (DTC)*: Fig. 10 shows the implementation of the proposed DTC and its timing diagram. As shown in Fig. 10(a), the DTC is composed of a falling-edge-triggered counter and a digital comparator. A DTC operation is triggered on need by the TRIG signal, which resets and enables a counter.

The digital comparator continuously compares the counter value CNT and the input value X_i . As long $CNT < X_i$, the comparator outputs "1" ($DTC_{OUT} = "1"$). Otherwise, it outputs "0" ($DTC_{OUT} = "0"$) and activates the HOLD signal to freeze the counter state and then wait for the trigger of a new operation.

The timing diagram of the DTC when a sequence of inputs equal 9, 5, and 7 is applied (these numbers are chosen for illustrative purposes) is shown in Fig. 10(b). The DTC uses the clock period T_{CLK} as its time reference, that is, the pulsedwidth (duration) of DTC_{OUT} equal to $X_i \times T_{CLK}$.

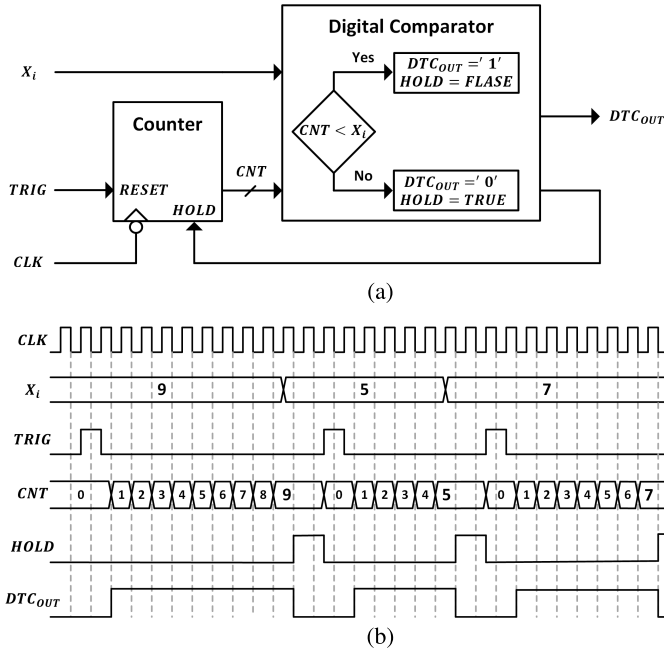


Fig. 10. Proposed DTC: (a) DTC block diagram. (b) DTC timing diagram.

2) *Time Accumulator*: The proposed TAC is implemented by the state machine shown in Fig. 11(a). The state machine is triggered by the clock's rising edge and it advances when the DTC_{OUT} is high. It can advance bi-directionally based on the sign of the accumulation $SIGN_i$, while the step size of each advancement is determined by the weight W_i .

If the accumulation contains more states than what is encoded in the state machine $LSB[11:0]$, the number of times the state machine returns to its initial state in both directions is tracked by $MSB[5:0]$ bits allowing long-time accumulation as needed (see Fig. 11(a)). A return while the state machine is advancing in the positive direction (+ve sign) leads to an increment ($MSB += 1$), and a return while the state machine is advancing in the negative direction (-ve sign) leads to a decrement ($MSB -= 1$). Both the MSB bits and the current state $LSB[11:0]$ represent the result of the MAC operation.

For example, the operations described in (3) are executed as shown in Fig. 11(b). First, the state machine advances from its initial state (0) in the positive direction (its value increments) by a step of 6 as defined by W_1 . The advancement continues for 9 clock cycles as defined by X_1 (DTC_{OUT} pulse duration). This operation results in an output equal to 54.

Following the first operation, the second operation ($X_2 = 5$ and $W_2 = 15$) is triggered with a negative sign. Therefore, the state machine advances in the negative direction (its value decrements) starting from the last state ($MSB = 0$ and $LSB = 54$). As shown in Fig. 11(b), the state machine returns to its initial state while it advances indicating that the output of this operation is negative. This return is tracked by decrementing the MSB ($MSB = -1$). In this case, the output of the operation is -21.

The last operation ($X_3 = 7$ and $W_3 = 12$) in Fig. 11(b) is triggered with a positive sign. The resultant advancement in the positive direction returns the state machine to its initial

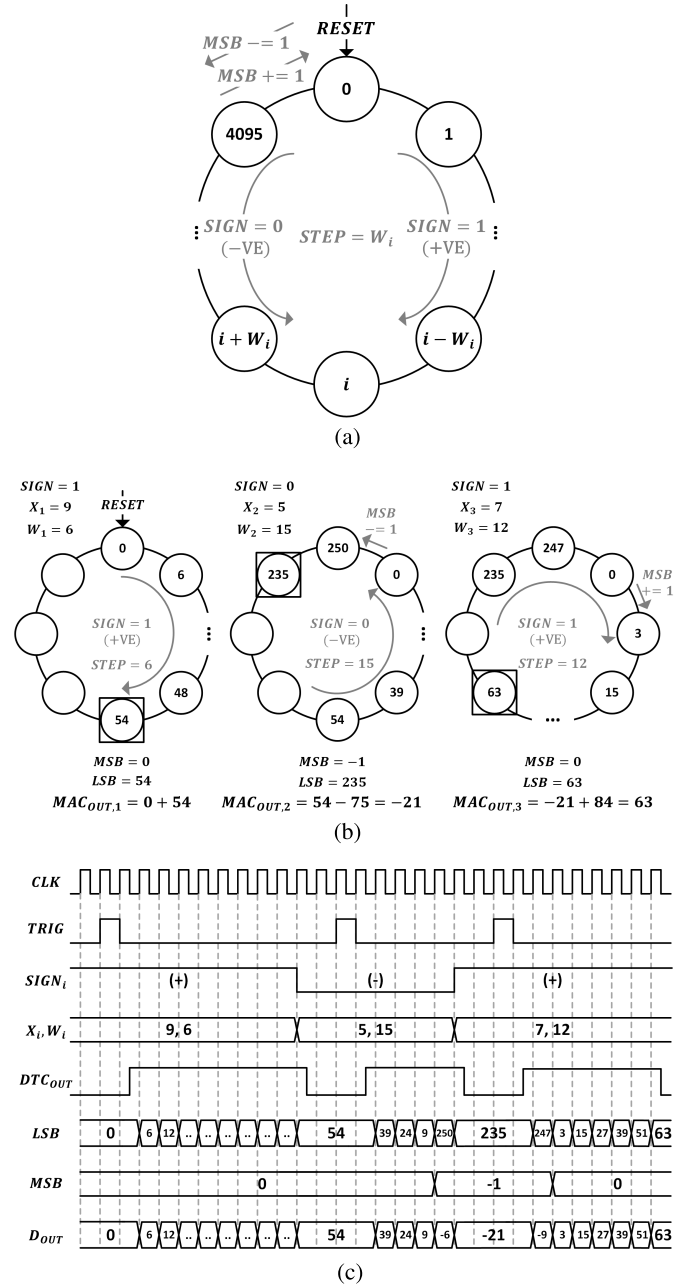


Fig. 11. Proposed time accumulator TAC. (a) TAC state machine. (b) Accumulation example. (c) MAC operation timing diagram.

state leading to a positive output equal to 63. This value is reserved until a new operation is triggered.

Fig. 11(c) shows the timing diagram of the MAC operation in (3) obtained by the simulation results utilizing both the proposed DTC and TAC. Note that having the DTC operating at the falling clock edge and the TAC operating at the rising clock edge, makes the architecture immune to time domain nonidealities like jitter and skew

$$\begin{aligned}
 X.(SIGN)W &= \begin{bmatrix} 9 \\ 5 \\ 7 \end{bmatrix} \begin{bmatrix} (+)06 \\ (-)15 \\ (+)12 \end{bmatrix} \\
 &= (54 - 75) + 84 = -21 + 84 = 63. \quad (3)
 \end{aligned}$$

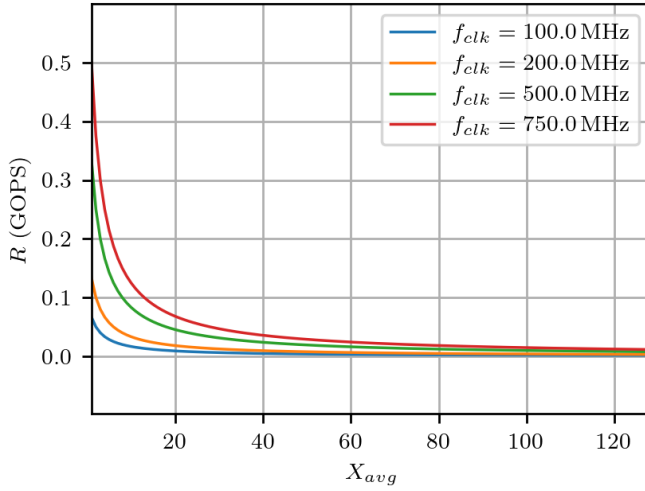


Fig. 12. Operation rate versus input and clock frequency.

B. Determination of the Operation Rate

In the proposed architecture, the necessary time T to perform $2N$ operations (MAC) depends on the clock frequency $1/T_{CLK}$ and the absolute values of the inputs $\text{abs}(X_i)$. T can be determined as

$$T = \sum_{i=1}^N \text{abs}(X_i) T_{CLK} + 2T_{CLK} \simeq (X_{avg} T_{CLK} + 2T_{CLK}) \times N. \quad (4)$$

Here, $X_{avg} \times T_{CLK}$ represents the average time required by the DTC (per MAC operation) while $+2T_{CLK}$ represents the controller overhead.

Similarly, the operation rate R is a function of the input X_i can be determined as follows:

$$R = \frac{2N}{T} \simeq \frac{2}{(X_{avg} T_{CLK} + 2T_{CLK})}. \quad (5)$$

Fig. 12 shows the operation rate R in GOPS as a function of the average input X_{avg} at different clock frequencies. A larger operation rate is obtained for lower input values and higher clock frequency. Also, employing more computation cores for the same number of operations increases the operation rate further.

C. Postaccumulation Functionalities

As shown in Fig. 8, the input vector $[x_1, x_2, \dots, x_{561}]$ is applied sequentially (one-by-one) to the eight-neuron computation layer (the hidden layer in a two-layer neural network). For each signed input X_i , a vector of unsigned weights $[w_{i,1}, w_{i,2}, \dots, w_{i,8}]$ and a vector of sign bits $[s_{i,1}, s_{i,2}, \dots, s_{i,8}]$ is applied to neuron 1 to neuron 8, respectively. Here, the sign of the accumulation is determined by an XNOR gate (see Fig. 8). If the sign bit and the extracted sign of the input are the same (e.g., both are negative or both are positive), positive accumulation is carried out (the state machine advances in the positive direction). Otherwise, negative accumulation is carried out (the state machine advances in the negative direction).

Similarly, the biases $[b_1, b_2, \dots, b_8]$ are applied to the neurons in the last operation (zero biases are applied for other operations).

As shown in Fig. 8, the 18-bit accumulation result is quantized by level shifting ($\div 2^p = 10$). Then, the quantized 8-bit outputs of the hidden layer can be applied directly to the output layer which requires six neurons for the activities to be classified (walking, walking upstairs, walking downstairs, sitting, standing, and laying).

1) *FIFO Using SRAM Array*: As described in Fig. 9, the proposed architecture requires a first-in-first-out (FIFO) or a memory macro to store the application data as needed. An array of 6T SRAM is implemented for this purpose. Using SRAM leverages its unique architectural characteristics for high-speed data processing and low power consumption, for example, its quick access times, making it an ideal implementation choice for low-power applications that require fast data enqueueing and dequeuing. For our application, the SRAM occupies $180 \times 97 \mu\text{m}$, and it operates at frequencies up to 712MHz.

2) *SRAM and Input Control*: The memory and input sequence are controlled by an additional finite state machine (FSM). The FSM starts setting up the initial SRAM address in its idle state. Then, the controller generates the TRIG signal to trigger the DTC and perform a signed accumulation (governed by the DTC output pulse). The controller acknowledges the end of the ongoing accumulation by detecting the falling edge of the DTC pulse. A new operation is then initiated by generating the TRIG signal. This process is shown in Fig. 11(c), where the input values are applied sequentially, and the TRIG signal is generated following the falling edge of the DTC output (DTC_{OUT}) initiating a new operation.

For each operation, the FSM progresses through states to read the 8-bit input data, weights (64-bit, 8-bit for each neuron), and biases, incrementing the SRAM address for each step. The read data is then pushed to the accelerator for processing. Afterward, the FSM waits for the accelerator to complete the processing, to either loop back for more inputs or to proceed to write processed data back into the SRAM.

These cycles repeat until all data is processed, culminating in a cleanup phase that resets the system, making it ready for the next computation. In this way, the FSM ensures the sequential processing of data from reading, to processing, and writing back in a controlled manner.

V. SIMULATION RESULTS

The eight-neuron computation layer only occupies $70 \times 70 \mu\text{m}$, and it dissipates $576 \mu\text{W}$ at 0.72-V supply and 500-MHz clock frequency. In this case, the equivalent area of a single neuron is $612.5 \mu\text{m}^2$, which is the smallest compared to other time-based architectures. This allows the possibility of utilizing a large number of neurons even in area-constrained applications, to enable more parallelism and increase the system's throughput.

As the accumulation is time-based, the operation rate of the eight-neuron layer depends on the clock frequency and the value of the sum of the magnitude of the inputs. Higher clock

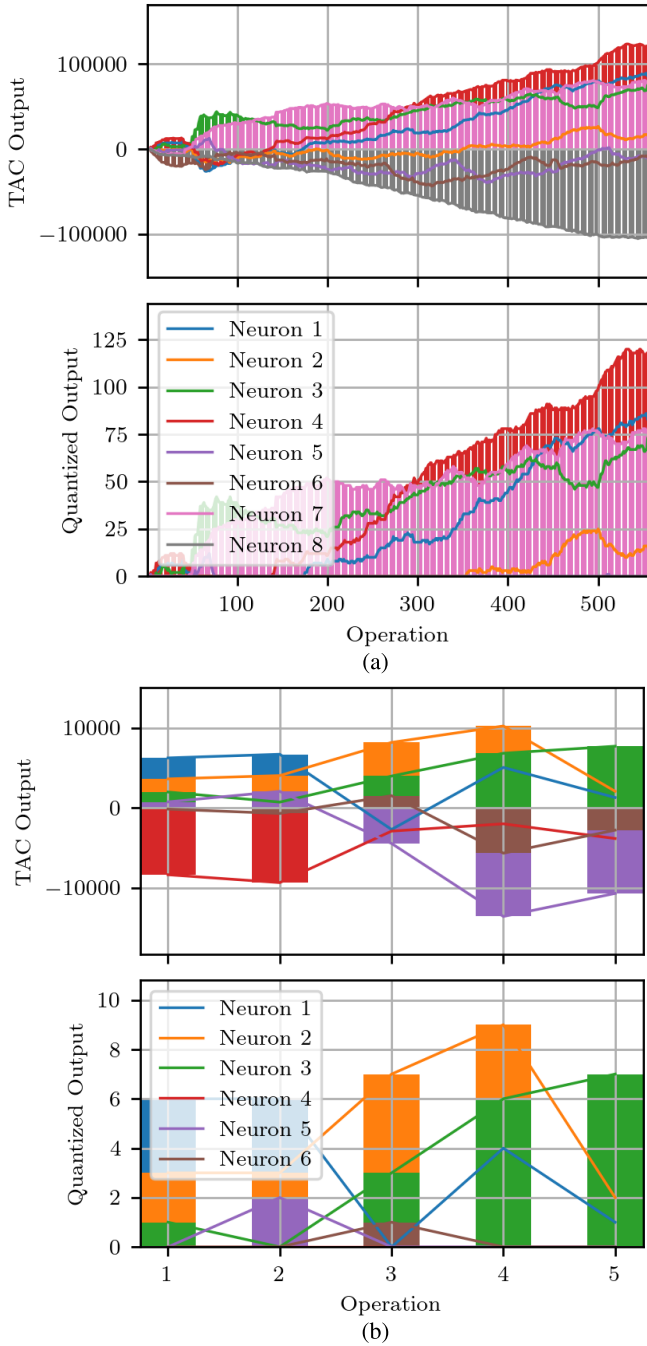


Fig. 13. Outputs of the computation layers. (a) Outputs of the hidden layer. (b) Outputs of the output layer.

frequency and low values can lead to higher operation rates (see Fig. 12).

A. Performances of the Proposed Architecture

The computation layer is characterized at 500 MHz and around the average value ~ 64 of the input (magnitude) range (0–128). Fig. 13 shows the RTL behavioral simulation of the hidden and output layers to verify the quantitative outputs of the TAC (signed 18-bit) and the final quantized outputs (signed 8-bit). The time taken for the hidden layer to output the quantized final results in this test scenario is $\sim 105 \mu s$ (52 813 clock cycles). While for the output layer, $\sim 746 ns$ (373 clock cycles)

TABLE I
POWER DISSIPATION VERSUS CLOCK FREQUENCY

Clock Frequency (MHz)	Power Dissipation (μW)
0.1	22.84
1	22.80
10	28.42
100	118.8
200	220.0
500	576.0
700	832.7

is required. The number of required clock cycles and operation rate given are extracted from the timing diagrams, and they can be determined from the formulas given in (4) and (5). The highest rate is recorded for binary inputs with a unity average. Note that the operation rate dependency on the clock frequency allows deploying dynamic frequency and voltage scaling at low computation demands for power savings. Table I shows the average power dissipation (tool extracted) at 0.72-V supply versus the clock frequency. We have included an extraction at 100 kHz operating frequency to match the settings used in [29] and provide a fair power consumption estimation. Note that below 10 MHz, the leakage dominates the power dissipation.

Since the performance is dominated by the hidden layer due to a large number of required operations, a single eight-neuron layer can be utilized for both the hidden and output layers in a temporally sequential manner. Only five operations are performed by the output layer with six neurons (two neurons will be idle).

Given this configuration, the computation layer achieves an average operation rate of 0.12 GOPS (up to 2.67 GOPS), average energy efficiency of 0.21 TOPS/W (up to 4.63 TOPS/W), and area efficiency of 24.74 GOPS/mm² (up to 544.22 GOPS/mm²). Its normalized energy efficiency [19] is 13.46 1b-TOPS/W (up to 296.30 1b-TOPS/W). While the normalized and scaled area efficiency is 1583.36 1b-GOPS/mm² (up to 18963.2 1b-GOPS/mm²).

If the accumulation number of bits is considered in this metric, the normalized efficiency will increase 18 times. To the best of our knowledge, the proposed circuitry achieves one of the highest recorded accumulation ranges.

Moreover, thanks to the simplicity of the design, it can be easily tailored to offer precision scalability as per the application needs, thereby saving area and power. In addition, other configurations can be easily implemented to target various applications by integrating more neurons, or by utilizing more computation layers.

B. Comparison With Existing Time-Based Architectures

A comparison with the state-of-art is conducted in Table II. Based on the comparison, the key features of the proposed architecture can be summarized as follows.

- 1) The architecture has no dependency on analog and time-domain nonidealities, and it has an all-digital implementation allowing more reconfigurability and easier integration into other digital circuits.
- 2) The proposed architecture achieves good precision by supporting multibit inputs (8-bit), weights (8-bit), and

TABLE II
COMPARISON OF STATE-OF-THE-ART TIME-BASED MAC ARCHITECTURES

	JSSC'19 [11]	JSSC'20 [12]	ISSCC'19 [13]	JSSC'22 [14]	OJCAS'23 [16]	This work
Process	28 nm	40 nm	65 nm	65 nm	65 nm	22 nm
All-Digital Imp.	NO	YES	NO	NO	No	Yes
Domain	Phase	Time	Time	Hybrid	Time	Time
Clock frequency (MHz)	-	25	1×10^{-3} - 1.5	2.12-90	-	500
Power Dissipation (μ W)	60.73	30.17	0.3 - 3.4	126.72	697	576
Supply Voltage (V)	0.7	0.537	0.4 - 1	0.7 - 1.1	1.2	0.72
Precision (I/W/O)	8/8/8	4/1/8	3-8/3-8/-	4,7/4,7/16	1/3/4	8/8/18
Area per MAC (μm^2)	960	124×10^3	200×10^4	2000	~ 5025	612.5
Operation Rate (GOPS)	0.78	0.365	2.73×10^{-3}	5.98	4.03 ^a	0.12
Energy Effic. (TOPS/W)	12.4	12.08	9.1	47.19	116	0.21
Norm. Effic. (1b-TOPS/W)	793.6	48.32	81.9	755.04	208.8	13.46
Area Effic. (GOPS/ mm^2)	1300	2.94	1.36×10^{-3}	21.81	20	24.74
Norm. Area Effic. (1b-GOPS/ mm^2) ^b	125243.11	32.42	0.54	2191.42	376.79	1583.36

^aNormalized to 1 neuron throughput

^bNormalized 1-b Area Efficiency = Area Efficiency (TOPS/ mm^2) \times No. of input bits \times No. of weight bits $\times (1 + 80\%)^{\log_2 \frac{\text{node}^2}{(22\text{nm})^2}}$ assuming 80% area improvement per technology node

outputs (18-bit). Also, the precision can be tailored based on the application requirements.

- 3) The eight-neuron computation layer achieves an average operation rate of 0.12 GOPS and an average energy efficiency of 0.21 TOPS/W (13.46 1b-TOPS/W).
- 4) The equivalent area of a neuron is $612.5 \mu\text{m}^2$ which considered among the smallest implementations. The eight-neuron computation layer achieves an area efficiency of $24.74 \text{ GOPS}/\text{mm}^2$ ($1583.36 \text{ 1b-GOPS}/\text{mm}^2$). Note that, the normalized area efficiency in Table II is scaled with respect to the 22-nm node (see [30] for better insights).
- 5) The architecture utilizes a single-clock source, with an operation rate proportional to the clock frequency. Using high clock frequency leads to relatively high power dissipation. Therefore, the proposed architecture fits well with applications that do not require high operation rates like smart IMU sensors.

C. Comparison With Other Activity Recognition Implementations

To provide more comparison points with implementations using activity recognition, we evaluated our approach against state-of-the-art microcontroller units (MCUs) and custom ASIC implementations targeting HAR. Regarding accuracy, state-of-the-art MCU implementations reported for HAR are 92.3% for UCI HAR and 93.1% for PAMAP [4]. ASIC designs improved this accuracy using additional sensors and processing, achieving 95% on HAR [29]. Our implementation achieves 93.9% for UCI HAR and 88.3% for PAMAP, with standard data processing, which is competitive with state-of-the-art accuracies. Regarding power, the DNN classifier used in the digital ASIC implementation in [29] achieved $10 \mu\text{W}$ power consumption at a 100-kHz frequency and 1-V power supply. In comparison, Table I reports the estimated power dissipation for various clock frequencies, showing that under 1-MHz frequency, leakage dominates the overall power dissipation, at $22.8 \mu\text{W}$. This difference is due to the technology node (65 versus 22 nm). Yet, the proposed circuit can achieve

TABLE III
COMPARISON IN TERMS OF HAR ACCURACY

Ref.	[4]	[29]	This Work
Implementation	Microcontroller	ASIC	ASIC
PAMAP Accuracy	93.1%	-	88.3%
HAR Accuracy	92.3%	95%	93.9%

significantly higher frequency operation (700 MHz), opening the possibility of heavy duty-cycling to reduce the general power consumption.

Table III summarizes the achieved accuracies in HAR-related work.

VI. CONCLUSION

This article presents a time-domain neural network architecture that targets in-sensor processing applications. The proposed 8-neuron computation layer computes sequential inputs and supports a signed accumulation of up to 18 bits. A single clock is required allowing dynamic frequency and voltage scaling for configurable performance demands and power savings. Moreover, the architecture's small sizes and low complexity allow for utilizing a large number of neurons even in small chips to enable more parallelism and increase the throughput. The architecture has an all-digital implementation that benefits from technology scaling and it has no dependency on analog nonidealities allowing easy integration into other on-chip digital circuits. These features make the proposed accelerator suitable for precision sensing applications, adding advanced capabilities with low cost to the next-generation smart sensors.

REFERENCES

- [1] S. García-de-Villa, D. Casillas-Pérez, A. Jiménez-Martín, and J. J. García-Domínguez, "Inertial sensors for human motion analysis: A comprehensive review," *IEEE Trans. Instrum. Meas.*, vol. 72, pp. 1–39, 2023.
- [2] Z. Zhongkai, S. Kobayashi, K. Kondo, T. Hasegawa, and M. Koshino, "A comparative study: Toward an effective convolutional neural network architecture for sensor-based human activity recognition," *IEEE Access*, vol. 10, pp. 20547–20558, 2022.

- [3] M. Rana and V. Mittal, "Wearable sensors for real-time kinematics analysis in sports: A review," *IEEE Sensors J.*, vol. 21, no. 2, pp. 1187–1207, Jan. 2021.
- [4] F. Daghero et al., "Human activity recognition on microcontrollers with quantized and adaptive deep neural networks," *ACM Trans. Embedded Comput. Syst.*, vol. 21, no. 4, pp. 1–28, Jul. 2022.
- [5] J.-S. Seo et al., "Digital versus analog artificial intelligence accelerators: Advances, trends, and emerging designs," *IEEE Solid State Circuits Mag.*, vol. 14, no. 3, pp. 65–79, Summer. 2022.
- [6] B. Murmann, "Mixed-signal computing for deep neural network inference," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 29, no. 1, pp. 3–13, Jan. 2021.
- [7] M. van Baalen et al., "FP8 versus INT8 for efficient deep learning inference," 2023, *arXiv:2303.17951*.
- [8] R. Sarpeshkar, "Analog versus digital: Extrapolating from electronics to neurobiology," *Neural Comput.*, vol. 10, no. 7, pp. 1601–1638, Oct. 1998.
- [9] H. A. Maharmeh, N. J. Sarhan, C.-C. Hung, M. Ismail, and M. Alhawari, "A comparative analysis of time-domain and digital-domain hardware accelerators for neural networks," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, May 2021, pp. 1–5.
- [10] P. S. Locatelli, D. M. Colombo, and K. El-Sankary, "Time-domain multiply-accumulate unit," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 31, no. 6, pp. 762–775, Jun. 2023.
- [11] Y. Toyama, K. Yoshioka, K. Ban, S. Maya, A. Sai, and K. Onizuka, "An 8 bit 12.4 TOPS/W phase-domain MAC circuit for energy-constrained deep learning accelerators," *IEEE J. Solid-State Circuits*, vol. 54, no. 10, pp. 2730–2742, Oct. 2019.
- [12] A. Sayal, S. S. T. Nibhanupudi, S. Fathima, and J. P. Kulkarni, "A 12.08-TOPS/W all-digital time-domain CNN engine using bi-directional memory delay lines for energy efficient edge computing," *IEEE J. Solid-State Circuits*, vol. 55, no. 1, pp. 60–75, Jan. 2020.
- [13] N. Cao, M. Chang, and A. Raychowdhury, "14.1 a 65nm 1.1-to-9.1TOPS/W hybrid-digital-mixed-signal computing platform for accelerating model-based and model-free swarm robotics," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2019, pp. 222–224.
- [14] S. Gweon, S. Kang, K. Kim, and H.-J. Yoo, "FlashMAC: A time-frequency hybrid MAC architecture with variable latency-aware scheduling for TinyML systems," *IEEE J. Solid-State Circuits*, vol. 57, no. 10, pp. 2944–2956, Oct. 2022.
- [15] J. Yang et al., "TIMAQ: A time-domain computing-in-memory-based processor using predictable decomposed convolution for arbitrary quantized DNNs," *IEEE J. Solid-State Circuits*, vol. 56, no. 10, pp. 3021–3038, Oct. 2021.
- [16] H. A. Maharmeh, N. J. Sarhan, M. Ismail, and M. Alhawari, "A 116 TOPS/W spatially unrolled time-domain accelerator utilizing ladder-inverter DTC for energy-efficient edge computing in 65 nm," *IEEE Open J. Circuits Syst.*, vol. 4, pp. 308–323, 2023.
- [17] P. Houshmand, J. Sun, and M. Verhelst, "Benchmarking and modeling of analog and digital SRAM in-memory computing architectures," 2023, *arXiv:2305.18335*.
- [18] D. Anguita et al., "A public domain dataset for human activity recognition using smartphones," in *Proc. ESANN*, vol. 3, 2013, p. 3.
- [19] N. R. Shanbhag and S. K. Roy, "Comprehending in-memory computing trends via proper benchmarking," in *Proc. IEEE Custom Integr. Circuits Conf. (CICC)*, Apr. 2022, pp. 1–7.
- [20] A. Reiss and D. Stricker, "Introducing a new benchmarked dataset for activity monitoring," in *Proc. Int. Symp. Wearable Comput. (ISWC)*, Newcastle, U.K., 2012, pp. 108–109. [Online]. Available: <https://api.semanticscholar.org/CorpusID:10337279>
- [21] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," 2015, *arXiv:1510.00149*.
- [22] K. Penner, F. Wittenfeld, B. Steinhagen, M. Hesse, and U. Rückert, "TinyML optimization for activity classification on the resource-constrained body sensor BI-vital," in *Proc. IEEE 19th Int. Conf. Body Sensor Netw. (BSN)*, Oct. 2023, pp. 1–4.
- [23] TensorFlow. (2024). *Quantization Aware Training*. Accessed: Mar. 10, 2024. [Online]. Available: https://www.tensorflow.org/model_optimization/guide/quantization/training
- [24] D. Miyashita, S. Kousai, T. Suzuki, and J. Deguchi, "A neuromorphic chip optimized for deep learning and CMOS technology with time-domain analog and digital mixed-signal processing," *IEEE J. Solid-State Circuits*, vol. 52, no. 10, pp. 2679–2689, Oct. 2017.
- [25] M. Alhawari and N. A. M. H. Perrott, "A 0.5 V< 4μW CMOS photoplethysmographic heart-rate sensor IC based on a non-uniform quantizer," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2013, pp. 1–3.
- [26] M. Z. Straayer and M. H. Perrott, "A multi-path gated ring oscillator TDC with first-order noise shaping," *IEEE J. Solid-State Circuits*, vol. 44, no. 4, pp. 1089–1098, Apr. 2009.
- [27] K. Kim, W. Yu, and S. Cho, "A 9 bit, 1.12 ps resolution 2.5 b/stage pipelined time-to-digital converter in 65 nm CMOS using time-register," *IEEE J. Solid-State Circuits*, vol. 49, no. 4, pp. 1007–1016, Apr. 2014.
- [28] A. M. Mohey, M. Kosunen, J. Ryyänen, and M. Andraud, "Toward all-digital time-domain neural network accelerators for in-sensor processing applications," in *Proc. IEEE Nordic Circuits Syst. Conf. (NorCAS)*, Oct. 2023, pp. 1–6.
- [29] G. Bhat, Y. Tuncel, S. An, H. G. Lee, and U. Y. Ogras, "An ultra-low energy human activity recognition accelerator for wearable health applications," *ACM Trans. Embedded Comput. Syst.*, vol. 18, no. 5s, pp. 1–22, Oct. 2019.
- [30] H. Wang, R. Liu, R. Dorrance, D. Dasalukunte, D. Lake, and B. Carlton, "A charge domain SRAM compute-in-memory macro with C-2C ladder-based 8-bit MAC unit in 22-nm FinFET process for edge inference," *IEEE J. Solid-State Circuits*, vol. 58, no. 4, pp. 1037–1050, Apr. 2023.



Ahmed M. Mohey (Graduate Student Member, IEEE) received the B.Sc. and M.Sc. degrees in electrical engineering from Ain Shams University, Cairo, Egypt, in 2012 and 2018, respectively. He is currently working toward the Ph.D. degree in electronics and nanoengineering from Aalto University, Espoo, Finland.

Since 2013, he held several industrial and research positions to develop software and electronic products. He is interested in analog and mixed-signal integrated circuit design with an emphasis on time-domain signal processing, sensing interfaces, power management, in-sensor/memory computing, and embedded systems.



Jelin Leslin received the bachelor's degree in electronics and communication engineering from Anna University, Chennai, India, in 2017, and the dual master's degree in electrical engineering from the Technical University of Eindhoven, Eindhoven, The Netherlands, and the KTH Royal Institute of Technology, Stockholm, Sweden, in 2020, through the EIT Digital dual university program. He is currently working toward the Ph.D. degree from Aalto University, Espoo, Finland, specializing in hardware-aware AI models. During his master's, he worked with Volvo Trucks on hardware acceleration for motion control in autonomous vehicles, which was also the focus of his thesis.

His research focuses on energy-efficient implementations of AI models through model compression, hardware-aware training, and custom hardware design.



Gaurav Singh (Member, IEEE) received the M.Sc. degree from VLSI from Amity University, Noida, India, in 2015. He is currently working toward the Ph.D. degree at the Department of Electronics and Nanoengineering, Aalto University, Espoo, Finland.

His research focuses on the development of low-power system-on-chip (SoC) solutions for sensor data processing. His work particularly explores microprocessors, low-power SRAM memories, and serial interfaces to enhance communication and debugging. His research work also includes integrating RISC-V processors with compute in-memory cores as hardware accelerators, aimed at improving power efficiency in AI applications.



Marko Kosunen (Member, IEEE) received the M.Sc., L.Sc., and D.Sc. (Hons.) degrees from Helsinki University of Technology, Espoo, Finland, in 1998, 2001 and 2006, respectively.

He is currently an Associate Professor at the Department of Electronics and Nanoengineering, Aalto University. From 2017 to 2019, he visited Berkeley Wireless Research Center, UC Berkeley, Berkeley, CA, USA, on a Marie Skłodowska-Curie grant from the European Union. In addition to his academic duties, currently, he is one of the three co-chairs of Microelectronics Finland, an academia–industry collaboration organization for microelectronics research and education. He has authored and co-authored more than a hundred journal and conference papers and holds several patents. His current research interests include programmatic circuit design methodologies, digital-intensive time-based data converters, transceiver circuits, and RISC-V microprocessor implementations with DSP accelerators.



Jussi Ryynänen (Senior Member, IEEE) was born in Ilmajoki, Finland, in 1973. He received the M.Sc. and D.Sc. degrees in electrical engineering from the Helsinki University of Technology, Espoo, Finland, in 1998 and 2004, respectively.

He is a Full Professor and the Dean of the School of Electrical Engineering, Aalto University, Espoo. He has authored or co-authored more than 200 refereed journal and conference papers in analog and RF circuit design. He holds seven patents on RF circuits. His research interests are integrated transceiver circuits for wireless applications.

Prof. Ryynänen is currently an SG Member for the European Solid-State Circuits Conference (ESSCIRC) and the IEEE Nordic Circuits and Systems Conference (NORCAS). He has served as a TPC Member of the IEEE International Solid-State Circuits Conference (ISSCC) and a Guest Editor for IEEE JOURNAL OF SOLID-STATE CIRCUITS.



Martin Andraud (Member, IEEE) received the Ph.D. degree from Grenoble University, Grenoble, France, in 2016.

He was a Post-Doctoral Researcher successively with TU Eindhoven, Eindhoven, The Netherlands, in 2016, and KU Leuven, Leuven, Belgium, from 2017 to 2019. From 2019 to 2024, he was an Assistant Professor at Aalto University, Espoo, Finland. He is currently an Assistant Professor in microelectronics at UCLouvain, Louvain-la-Neuve, Belgium, and a Visiting Professor at Aalto University. His research interests include the interface between edge AI, hardware/software co-design, testing, and reliability of custom ASIC for various AI accelerators, for instance, mixed-signal compute-in-memory architectures and alternatives to deep learning models (probabilistic reasoning or hybrid AI).