

# Step Counting with Attention-based LSTM

Shehroz S. Khan and Ali Abedi

*KITE Research Institute, Toronto Rehabilitation Institute*

*University Health Network*

Toronto, Canada

{shehroz.khan, ali.abedi}@uhn.ca

**Abstract**—Physical activity is recognized as an essential component of overall health. One measure of physical activity, the step count, is well known as a predictor of long-term morbidity and mortality. Step Counting (SC) is the automated counting of the number of steps an individual takes over a specified period of time and space. Due to the ubiquity of smartphones and smartwatches, most current SC approaches rely on the built-in accelerometer sensors on these devices. The sensor signals are analyzed as multivariate time series, and the number of steps is calculated through a variety of approaches, such as time-domain, frequency-domain, machine-learning, and deep-learning approaches. Most of the existing approaches rely on dividing the input signal into windows, detecting steps in each window, and summing the detected steps. However, these approaches require the determination of multiple parameters, including the window size. Furthermore, most of the existing deep-learning SC approaches require ground-truth labels for every single step, which can be arduous and time-consuming to annotate. To circumvent these requirements, we present a novel SC approach utilizing many-to-one attention-based LSTM. With the proposed LSTM network, SC is solved as a regression problem, taking the entire sensor signal as input and the step count as the output. The analysis shows that the attention-based LSTM automatically learned the pattern of steps even in the absence of ground-truth labels. The experimental results on three publicly available SC datasets demonstrate that the proposed method successfully counts the number of steps with low values of mean absolute error and high values of SC accuracy.

**Index Terms**—step counting, attention mechanism, long short-term memory, variable-length sequences

## I. INTRODUCTION

As the fundamental unit of human locomotion, steps are the preferred method of quantifying ambulatory physical activity [1]. The association between steps per specific time periods and health variables has been reported in several cross-sectional studies [2]–[4]. To illustrate, a higher number of steps is inversely associated with the risk of cardiovascular events and premature death [5]. Step Counting (SC) is the automated counting of the number of steps an individual takes over a specified period of time and space. SC has applications for telemonitoring/telemedicine to measure the number of steps and monitor the daily physical activity of patients remotely [6]. There are many other applications for SC, including indoor navigation where global positioning systems are unreliable, and pedestrian dead reckoning [7].

The increasing ubiquity of smartphones and smartwatches equipped with a variety of built-in Inertial Measurement Unit (IMU) sensors, such as accelerometers, gyroscopes, and magnetometers, has led to the development of various SC methods

using the multivariate time series of sensor signals. The existing SC approaches are broadly categorized into non-machine-learning- and machine-learning-based approaches. The non-machine-learning-based approaches can be divided into time-domain and frequency-domain approaches, Fig. 1.

Time-domain approaches generally rely on thresholding or peak detection. In thresholding methods, a step is detected when sensor data satisfy predefined criteria. These methods are particularly effective when detecting movements at the foot, where heel strikes can cause large and short-lived accelerations [9]. Peak detection or zero-crossing methods usually work on low-pass filtered signals and detect the occurrence of steps according to the presence of peaks in the signal. Based on the peaks and the distance between the peaks, some methods try to find the inherent periods in the signal. Auto-correlation is another method to detect the period of signal and step accordingly [10]. Some methods work based on stride in which the stride template is formed offline and cross-correlated with the signal [8]. Frequency-domain approaches, on the other hand, generally use the Fourier transforms of the signals, such as short-term Fourier transform and wavelet transform, and utilize the features in the frequency domain to detect steps [8].

A major limitation of non-machine-learning-based approaches for SC is that they require careful tuning of several parameters [7], [8]. For instance, in thresholding, peak detection, and period detection approaches in the time domain, the main difficulty is finding an optimal threshold/criteria to detect a specific timestep of the signal as a peak or consider segments of the signal as a period. Optimizing the window length for short-term Fourier transform and the parameters of continuous or discrete wavelet transforms is the issue with frequency-domain approaches [8].

The machine-learning-based SC approaches are categorized into feature-based and deep-learning approaches. In the former, features such as mean, variance, standard deviation, energy, and entropy are first extracted from the windows of signal in the time/frequency domain and then classified into step versus non-step using traditional machine-learning techniques such as support vector machines and Hidden Markov Model (HMM) trained on sequences of features [11]. The traditional machine-learning-based approaches also suffer from the need for parameter tuning, e.g., the selection of most effective features and the length of windows in which the features are extracted [7], [8]. In deep-learning methods, neural networks

can learn features from raw sensor signals. The existing deep-learning-based SC methods are mostly based on Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) [12]–[17].

In most of the above-discussed approaches, the steps are first detected by thresholds, features, CNNs, or RNNs, and then the total number of steps is output [12]–[17]. One practical problem with such a setup is that to evaluate the algorithms, ground truth is needed for every step taken. Manually labelling such fine-grained labelled data for each step is arduous or time-consuming. Alternatively, additional hardware may need to be added (e.g., pressure sensors) on the heels of the shoes, which is infeasible in a real-world setting. Keeping these issues in mind, in this paper, we propose a novel SC method that overcomes the limitations mentioned thus far, namely the necessity to determine the window size and the need for ground-truth data for single steps. Our main contributions are as follows:

- For the first time, SC is formulated as a regression problem that is solved using an attention mechanism for many-to-one LSTMs capable of analyzing variable-length sequences.
- The SC problem is solved at the signal level that is capable of analyzing the entire input time-series signal as a whole.
- Extensive experiments conducted on three publicly available SC datasets [7], [8], [19] show the superiority of our approach to a variety of machine-learning and non-machine-learning SC methods.

Our approach does not require windowing or annotation of individual steps - only the final count of steps is sufficient for training deep-learning models.

This paper is structured as follows. In Section II, we briefly study the existing deep-learning SC approaches. Section III, introduces the proposed method for SC. Section IV describes experimental settings and results on the proposed methodology. In the end, Section V presents our conclusions and directions for future works.

## II. RELATED WORK

In this section, we discuss some of the deep learning methods for SC from IMU sensors. Unfortunately, there are not many papers published on SC using deep learning. There are, however, some preprints available. Shao et al. [17] proposed a method for SC through step detection. A window is slid over the input signal, and a CNN classifies the signal in the window as a left step, right step, or no step. The CNN consists of a 1-dimensional convolutional layer followed by fully-connected layers. Chen [13] proposed a SC method using a many-to-many Long Short-Term Memory (LSTM), having the sequences of windows extracted from the signal as input and outputting step vs. non-step for each window. In the method proposed by Pillai et al. [15], the input signal is first segmented into windows, and the signal in each window is analyzed by an LSTM-based neural network. The last timestep of the LSTM is trailed by fully-connected layers to output left step or right step

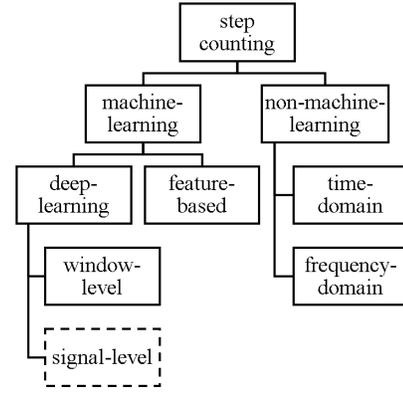


Fig. 1: Hierarchical classification of the existing SC approaches. The proposed method in this paper, dashed box, is the only deep-learning-based signal-level SC method that analyzes the entire signal as a whole and solves the SC as a regression problem (see Sections I, II, and III).

classes. Luu et al. [12] extended the work by Pillai et al. [15], where instead of the last timestep of the LSTM, the output of the LSTM over all timesteps followed by fully-connected layers are used to output step or non-step classes. In addition to LSTM, Luu et al. [12] have used WaveNet [18] and a 1D CNN for SC through step detection.

Even though the above deep-learning-based methods do not require tuning parameters such as a threshold on peak values to be considered steps as with the time-domain methods, the window size parameter still needs to be specified manually in these methods. The window size needs to be adjusted for different datasets using different sensors, with different sampling frequencies and with different sensor placements on the body of individuals, with different walking styles, walking speeds, ages, and disabilities. To illustrate, the window length used by Shao et al. [17] is “length of training step segments”, Pillai et al. [15] is 0.46 seconds “based on tuning experiments and literature”, and Luu et al. [12] is 2 seconds, and 4 seconds for LSTM, and CNN, respectively. None of the above works have examined the effect of different window sizes in different situations.

In the next section, we describe our approach for SC, which only needs a final count of the steps for training and evaluation of deep-learning models and does not need window size as a parameter to tune.

## III. STEP COUNTING THROUGH REGRESSION USING ATTENTION IN LSTM

Input data to the proposed attention-based many-to-one LSTM architecture for SC is time series of accelerometer sensor signals of variable length. The neural network solves a regression problem and outputs the number of specific learned patterns in the input time series corresponding to the number of steps in the accelerometer sensor signal. The raw  $x$ ,  $y$ , and  $z$  components of the accelerometer sensor signal or their  $l_2$  norm (as multivariate or univariate time series) can be input

to the neural network. The input time series is normalized into the range of 0 to 1 by subtracting the minimum value of the signal and dividing it by the range of the signal. We used LSTM as the main component of the proposed neural network architecture. It can either be the vanilla RNN, Identity RNN, or Gated Recurrent Unit [21]. These RNNs are capable of handling sequences of variable lengths and situations in which there are sequences of variable lengths in successive mini-batches of epochs.

The shape of the input tensor as a mini-batch is  $(N, L_i, input\_size), i = 0, 1, \dots, N$  in which  $N$  is the number of samples in the mini-batch,  $L_i$  is the length (number of timesteps) of the  $i$ -th signal sample, and  $input\_size$  is the dimensionality of the signal at each timestep. The LSTM has  $num\_layers$  layers where  $num\_layers = 2$  would mean stacking two LSTMs together to form a stacked LSTM, with the second LSTM taking in outputs of the first LSTM and computing the final results. The LSTM outputs a tensor of size  $(N, L_i, hidden\_size), i = 0, 1, \dots, N$ . To make the proposed method capable of analyzing variable-length sequences, the output tensor of LSTM for different samples in the current mini-batch are fed to the attention mechanism sample-by-sample. The  $(N, L_i, hidden\_size)$  output tensor of LSTM is divided into  $N$  tensors  $h_i$  of size  $(L_i, hidden\_size), i = 0, 1, \dots, N$ . A linear layer of size  $hidden\_size \times hidden\_size$  is used as an attention layer, takes the outputs of LSTM sample-by-sample  $h_i, i = 0, 1, \dots, N$ , and outputs energies  $e_i$  of the same size  $(L_i, hidden\_size)$ .  $e_i$  then will be multiplied by the summation of  $h_i$  over length dimension ( $s_i$ ) to output a vector whose softmax generates weights  $w_i$  of size  $L_i$ .  $w_i$ , as the result of training the attention layer, is multiplied by  $h_i$  to output context  $c_i$  of length  $hidden\_size$ .

The concatenations of contexts  $c_i$  and summations  $s_i$  for all the samples form a tensor of size  $(N, hidden\_size \times 2)$ . The concatenation tensor is fed to a linear layer of size  $(hidden\_size \times 2) \times 1$  to generate the final output having a single real-valued number for each sample in the current mini-batch. The Mean Absolute Error (MAE) is used as the loss function of the neural network. During training, the attention layer and the corresponding weights (generated from the multiplication of the output of the attention layer and  $s_i$ ) learn to pay attention to the specific patterns (steps) and their summation in the input time series (accelerometer sensor signal), see Fig. 2 as an example.

The advantage of the proposed attention mechanism over the original versions of the attention mechanism for RNNs [22], [23] is that it can handle variable-length sequences. We provide the input to the attention layer, followed by multiplications and concatenations, in a sample-by-sample manner for individual training samples (with different lengths) successively.

#### IV. EXPERIMENTS

In this section, the performance of the proposed SC method is evaluated on three publicly available SC datasets using different evaluation metrics. We compare different settings of

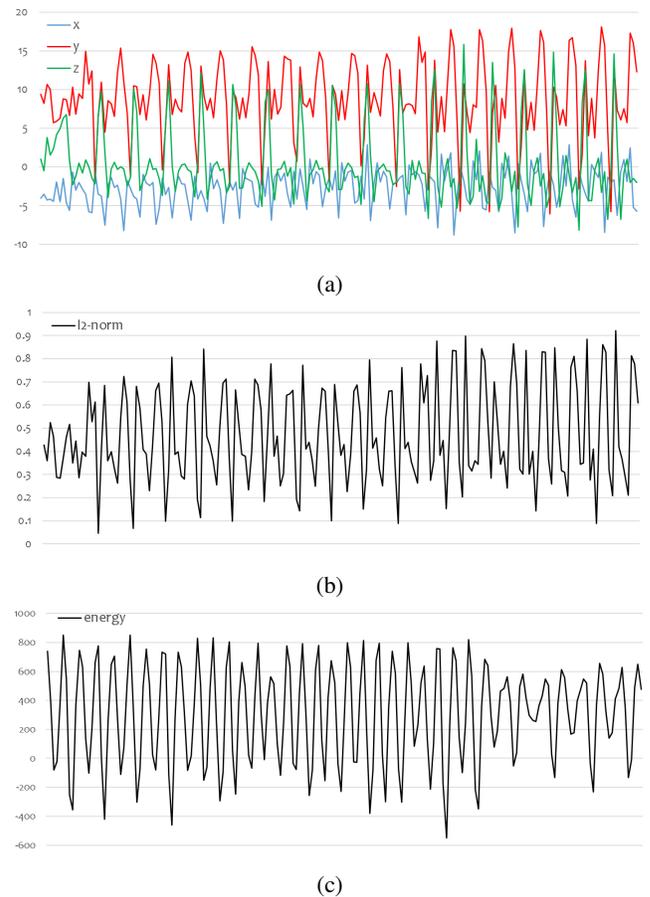


Fig. 2: (a) The raw  $x$ ,  $y$ , and  $z$  components, (b) the normalized  $l_2$ -norm of (a), and (c) the energy (corresponding to the input signal in (b)) of the attention mechanism in a trained attention-based LSTM, for the first half of an exemplary accelerometer signal in the WDSC dataset.

the proposed method with the previous machine-learning and non-machine-learning SC methods.

##### A. Datasets

**WDSC:** Brajdic and Harle [8] collected and annotated a dataset for walk detection and SC. The data was collected from 27 subjects of different ages and heights using the built-in accelerometer of an Android smartphone sampled at 100 Hz and under six different smartphone placements. There are 117 samples ( $x$ ,  $y$ , and  $z$  accelerometer signals of walks) in the dataset, each labelled based on the start and end time of the walk and the number of steps during the walk. The cropped signals from the start to end time of the walk are considered as the input, and the number of steps in the cropped signals is considered as the ground-truth labels. A major advantage of this dataset is the diversity of the location of the smartphone accelerometer sensors on the body. Table I presents the statistics of the samples in this dataset. Due to the lack of ground-truth labels for individual steps in the WDSC

TABLE I: The statistics, minimum, maximum, mean, standard deviation (STD), and skew of the ground-truth number of steps in the WDSC dataset [8], WeAllWalk dataset [7], and the regular and semi-regular parts of the Pedometer dataset [19].

	WDSC [8]	WeAll Walk [7]	Pedometer Regular [19]	Pedometer Semi-regular [19]
Minimum #steps	63	2	857	548
Maximum #steps	106	136	1100	814
Mean #steps	78	40.71	991.03	704.03
STD #steps	8.46	33.29	54.03	65.57
Skew #steps	0.56	0.81	-0.27	-0.21

dataset [8], the previous deep-learning SC approaches [12]–[17] cannot be trained and evaluated on this dataset. For the first time, we show the results of a deep-learning SC approach on this dataset.

**WeAllWalk:** Flores and Manduchi [7] collected and annotated a dataset from 15 subjects. The uniqueness of this dataset is due to the presence of sighted and blind subjects. Ten blind subjects (using a white cane or guide dog) and five sighted subjects contributed to the data collection. While walking ten different paths, the subjects carried two Apple iPhone 6S smartphones at two different body locations. Each smartphone recorded data from its accelerometer, gyroscope, and magnetometer at a rate of 25 Hz. The accelerometer signals are used in the experiments in this paper. There are 932 samples in the dataset categorized into three categories of 290 samples from sighted subjects, 468 samples from blind subjects with a white cane, and 174 samples from blind subjects with a guide dog. Table I presents the statistics of the samples in this dataset.

**Pedometer:** Mattfeld et al. [19] collected and annotated a dataset from 30 subjects for SC. The subjects wore three Shimmer3 sensors on their non-dominant wrist, hip, and non-dominant ankle. Each sensor recorded accelerometer and gyroscope data at 15 Hz. Unique to this dataset is the availability of different walking regularities, including regular (walking a path), semi-regular (conducting a within-building activity), and unstructured (conducting a within-room activity). There are a total of 270 samples in the dataset, nine ankle-, hip-, and wrist-placement signals collected during regular, semi-regular, and unstructured walking from each of the 30 subjects. According to the authors’ recommendations [20], 180 regular and semi-regular accelerometer signal samples of this dataset are used in our experiments, whose statistics are presented in Table I.

## B. Evaluation Metrics

The researchers who collected the above datasets also developed various evaluation metrics for SC. Considering  $y\_trues$ , and  $y\_preds$  as the ground-truth step count, and predicted step count, Brajdic and Harle [8] defined an Error Rate (ER) for each sample as follows. The mean and standard deviation of Equation 1 over all test samples are reported in this paper.

$$\frac{y\_preds - y\_trues}{y\_trues} \times 100 \quad (1)$$

Mattfeld et al. [20] defined a Running Count Accuracy (RCA) for each sample (defined below) and considered the prediction as undercount, or overcount if the accuracy is less than, or greater than 1. The mean and standard deviation of Equation 2 over all test samples are reported in this paper.

$$\frac{y\_preds}{y\_trues} \quad (2)$$

Flores and Manduchi [7] first calculated the number of samples with  $y\_preds - y\_trues < 0$  as undercount and the number of samples with  $y\_preds - y\_trues > 0$  as overcount, and defined their normalization as percentages of UnderCount (UC) and OverCount (OC) as follows.

$$\text{normalized undercount} = \frac{\text{undercount}}{\text{total number of samples}} \times 100 \quad (3)$$

$$\text{normalized overcount} = \frac{\text{overcount}}{\text{total number of samples}} \times 100 \quad (4)$$

Luu et al. [12] defined an ACCuracy (ACC) for each sample as follows. The mean and standard deviation of Equation 5 over all test samples are reported in this paper.

$$\left(1 - \frac{|y\_preds - y\_trues|}{y\_trues}\right) \times 100. \quad (5)$$

Since the proposed SC method solves a regression problem, the MAE between  $y\_trues$  and  $y\_preds$  is also reported.

## C. Experimental Settings

We implemented LSTM with and without attention mechanism to evaluate their performance. The proposed method is implemented with the 3-dimensional time series (x, y, and z components of the accelerometer signal)  $input\_size = 3$  or 1-dimensional time series (l2-norm of the x, y, and z components)  $input\_size = 1$  or their combination  $input\_size = 4$ . Compared to the previous methods, one major advantage of the proposed method is that it does not require many parameters that need to be changed depending on the dataset or signal. However, LSTM-specific parameters may still need to be set, which are as follows. The LSTM is unidirectional with  $num\_layers = 2$  and  $hidden\_size = 128$ . The attention layer is a linear layer of size  $hidden\_size \times hidden\_size$ . The final fully-connected layers, after the attention mechanism, contains two linear layers of size  $(hidden\_size \times 2) \times hidden\_size$  and  $hidden\_size \times 1$ . For the LSTM with no attention, the architecture of the LSTM and the final fully-connected layers is the same as the attention-based network. Due to not improving the results, no dropout is used in both the LSTMs with and without attention.

As discussed in Section IV-A, the frequency of WDSC, Pedometer, and WeAllWalk datasets are 100 Hz, 15 Hz, and 25 Hz. Before inputting the signals into the neural networks, the samples in the WDSC dataset are down-sampled by a factor of 4 to reduce the number of timestamps, and the samples in the

Pedometer and WeAllWalk datasets remain unchanged. Then, the signals are normalized as explained in section III.

None of the datasets defined separate train and test sets. The non-machine-learning methods that used these datasets, e.g., [20], reported the results of their experiments on all the samples of the datasets. The machine-learning methods applied on these datasets, e.g., [12], used cross-validation and reported the results of their algorithms on all the samples in the dataset over all the folds. In this paper, cross-validation is implemented for evaluation, and the results are reported for all the samples in the datasets. The architecture and hyper-parameters of the models are the same in all the folds of cross-validation in all the datasets.

The MAE is the loss function, and Adam is the optimizer. The batch size and the number of epochs are 16, and 250, respectively. The learning rate starts from 0.001 and is scheduled to be reduced by a factor of 10 in each 75 epochs. The experiments were implemented in PyTorch [25] and scikit-learn [26] on a server with 64 GB of RAM and NVIDIA TeslaP100 PCIe 12 GB GPU. The code of our implementations is available at <https://github.com/abedICODES/stepcounting>.

#### D. Experimental Results

The proposed method is evaluated using various evaluation metrics described in Section IV-B. On the WDSC dataset [8], five-fold cross-validation is performed with the same network architecture across all folds. The average (and standard deviation) results of all samples across all folds using different evaluation metrics are presented in Table II. The results are shown for the LSTM without attention and the attention-based LSTM (described in Section IV-C) with different numbers of neurons in the hidden layers of the LSTM and correspondingly its following linear layers (including the attention layer), and different signal dimensionality, the raw  $x$ ,  $y$ , and  $z$  components of accelerometer signals (3-dimensional), the  $l_2$ -norm of the  $x$ ,  $y$ , and  $z$  components (1-dimensional), and both (4-dimensional). According to Table II, in all the configurations of LSTM and attention-based LSTM, a 2-layer LSTM with 128 neurons in the hidden layer has the best performance. The  $l_2$ -norm itself works better than the raw components and both the  $l_2$ -norm and raw components. In almost all the configurations of the vanilla LSTM, and attention-based LSTM, in the upper, and lower halves of Table II, respectively, adding attention mechanism to the vanilla LSTM results in significant improvements in different evaluation metrics. The best attention-based model achieves (in bold letters) very low values of MAE, UC and OC (Flores [7]), and ER (Brajdic [8]), and very close to one values of RCA (Mattfeld [20]) and very high values of ACC (Luu [12]).

Fig. 2 (a) illustrate the raw  $x$ ,  $y$ , and  $z$  components, 2 (b) the normalized  $l_2$ -norm of (a), and 2 (c) the energy (corresponding to the signal in 2 (b)) of the attention mechanism in a trained attention-based LSTM neural network (the best model in Table II), for the first half of an exemplary accelerometer signal in the WDSC dataset [8]. As shown in Fig. 2 (c), the energy, as the output of the attention layer, has a shape in which the

steps have been emphasized. As described in Section III, the energy and its corresponding weights will be multiplied by the output of the LSTM in the network. In this way, the attention mechanism learns and pays attention to the steps, modifies the output of the LSTM accordingly, and the neural network outputs the step count.

Table III presents the results of the proposed SC method using the LSTM without attention and the attention-based LSTM with 128 neurons in the two hidden layers of the LSTM and the  $l_2$ -norm of the accelerometer signal using different evaluation metrics for the WeAllWalk dataset [7]. Following the experimental settings in the original work introduced in the WeAllWalk dataset [7], leave-one-person-out cross-validation in different populations is implemented. The results are presented for sighted subjects, blind subjects with a white cane, and blind subjects with a guide dog to examine the robustness of the proposed method in different populations with different levels of walking regularity. As can be seen in Table III, in almost all the populations, in most of the evaluation metrics, adding attention significantly improves the SC performance. The proposed attention-based method is robust against irregular walking, i.e., in blind subjects with a white cane and blind subjects with a guide dog.

Table IV presents the results of the proposed SC method using the LSTM without attention and the attention-based LSTM with 128 neurons in the two hidden layers of LSTM and the  $l_2$ -norm of the accelerometer signal using different evaluation metrics (see Section IV-B) for the Pedometer dataset [19]. Following the experimental settings in the original work introduced the Pedometer dataset [20], leave-two-person-out cross-validation in two levels of walking regularity and three different sensor placements are implemented, and the average results are reported in six sections of Table IV. As can be seen in Table IV, in most of the sections using most of the evaluation metrics, adding attention improves the SC performance. However, compared to the LSTM, the performance deterioration (from regular to semi-regular walking) is more severe in the attention-based LSTM.

Table V shows the results of the proposed LSTM and attention-based LSTM SC methods with the hyper-parameters in the previous experiment on the Pedometer dataset [19], compared to the previous deep-learning-based methods on the whole regular walking samples in the Pedometer dataset [19] using leave-two-person-out cross-validation. Our proposed attention-based LSTM method competes with the CNN method [12]. As explained in Section II, the CNN method [12] is based on windowing and requires determining the window size. In addition, contrary to the proposed method, which only requires one single annotation data (the number of steps) for the entire signal, the CNN method [12] requires step annotation data at each timestep of the signal.

Table VI shows the results of the proposed LSTM and attention-based LSTM SC methods with the hyper-parameters in the previous experiment on the Pedometer dataset [19] compared to the previous time-domain methods [20] on the whole regular walking samples and semi-regular walking sam-

TABLE II: The results of the proposed SC method on the WDSC dataset [8] using the LSTM with and without attention. In LSTM- $a \times b$ ,  $a$ , and  $b$  mean the number of layers in the LSTM, and the number of neurons in the layers, respectively. xyz means the raw  $x$ ,  $y$ , and  $z$  components of the accelerometer data and  $l2$  means the  $l2$ -norm of xyz as the inputs to the neural networks, using different evaluation metrics (see Section IV-B).

	MAE	UC, OC	ER	RCA	ACC
LSTM- $2 \times 64 - l2$	4.08	2.19, 3.03	0.78±6.32	<b>1.00±0.06</b>	0.95±0.04
LSTM- $2 \times 128 - l2$	2.83	1.56, 2.05	0.50±4.73	<b>1.00±0.05</b>	0.96±0.03
LSTM- $2 \times 256 - l2$	5.87	2.83, 4.70	2.00±8.97	0.98±0.09	0.93±0.06
LSTM- $2 \times 128$ -xyz	5.74	3.02, 4.32	1.56±8.69	0.98±0.09	0.93±0.05
LSTM- $2 \times 128 - l2$ and xyz	4.42	2.15, 3.52	1.45±7.13	0.98±0.07	0.94±0.05
LSTM- $2 \times 64 - l2$ (attention)	5.43	3.03, 3.94	1.12±8.34	0.99±0.08	0.93±0.05
LSTM- $2 \times 128 - l2$ (attention)	<b>2.33</b>	<b>1.39</b> , 1.60	<b>0.17±3.92</b>	<b>1.00±0.04</b>	<b>0.97±0.03</b>
LSTM- $2 \times 256 - l2$ (attention)	5.38	2.89, 4.01	1.29±8.1	0.99±0.08	0.93±0.05
LSTM- $2 \times 128$ -xyz (attention)	4.03	2.27, 2.90	0.61±6.71	0.99±0.07	0.95±0.04
LSTM- $2 \times 128 - l2$ and xyz (attention)	2.69	1.95, <b>1.50</b>	-0.55±4.80	<b>1.00±0.05</b>	0.96±0.03

TABLE III: The results of the proposed SC method on different populations of the WeAllWalk dataset [7] using the LSTM with and without attention with 2 layers of 128 neurons, and the  $l2$ -norm of the accelerometer signal (see Section IV-C) using different evaluation metrics (see Section IV-B).

	MAE	UC, OC	ER	RCA	ACC
Sighted Subjects					
LSTM- $2 \times 128$ - $l2$	2.51	2.80, 4.03	<b>3.67±10.76</b>	<b>0.97±0.1</b>	0.92±0.07
LSTM- $2 \times 128$ - $l2$ (attention)	<b>1.25</b>	<b>2.20</b> , <b>1.46</b>	-5.83±12.88	1.05±0.12	<b>0.93±0.11</b>
Blind Subject with a White Cane					
LSTM- $2 \times 128$ - $l2$	6.80	4.93, 9.49	7.6±23.32	0.92±0.24	0.83±0.18
LSTM- $2 \times 128$ - $l2$ (attention)	<b>4.09</b>	<b>3.50</b> , <b>5.10</b>	<b>-2.36±13.35</b>	<b>1.02±0.14</b>	<b>0.89±0.09</b>
Blind Subject with a Guide Dog					
LSTM- $2 \times 128$ - $l2$	4.64	5.53, 5.91	<b>4.16±26.25</b>	<b>0.96±0.27</b>	0.84±0.22
LSTM- $2 \times 128$ - $l2$ (attention)	<b>2.51</b>	<b>3.64</b> , <b>2.56</b>	-8.47±18.68	1.08±0.19	<b>0.88±0.17</b>

ples in the Pedometer dataset [19] using leave-two-person-out cross validation. Our proposed attention-based LSTM method outperforms the previous time-domain methods [20] with the advantage of not requiring many parameters and thresholds as in the time domain methods.

## V. CONCLUSIONS

This paper defined SC as a regression problem and used a many-to-one attention-based LSTM to solve it. Most of the previous methods work on windowed accelerometer signals in which SC results from step detection in individual signal windows. Our proposed method, working at the signal level, analyzes the entire accelerometer signal as a whole and outputs the number of steps. This signal level analysis eliminates the need for determining the window size and having ground-truth labels for every individual step. The proposed attention mechanism for RNNs that is capable of analyzing variable-length time series (signals) learns to pay attention to the steps and outputs their summation. The internal step identification is a consequence of applying the attention mechanism that is learned through the training of the neural network. The experimental results on three publicly available SC datasets demonstrated that the proposed method successfully counts the number of steps with low values of mean absolute error and high values of SC accuracy. Temporal Convolution Networks (TCNs) [27] are powerful neural networks for the analysis and modeling of sequences with extensive lengths. However, they are unable to handle variable-length sequences. Our future

work will investigate modifying TCNs [27] and attention-based TCNs [28] to handle variable-length sequences and applying them to signal-level SC. In addition, we plan to work on developing personalized SC models for specific populations and individuals.

## REFERENCES

- [1] Bassett, David R., Lindsay P. Toth, Samuel R. LaMunion, and Scott E. Crouter. "Step counting: a review of measurement considerations and health-related applications." *Sports Medicine* 47, no. 7 (2017): 1303-1315.
- [2] Pillay, Julian D., Hidde P. Van der Ploeg, Tracy L. Kolbe-Alexander, Karin I. Proper, Maartje Van Stralen, Simone A. Tomaz, Willem Van Mechelen, and Estelle V. Lambert. "The association between daily steps and health, and the mediating role of body composition: a pedometer-based, cross-sectional study in an employed South African population." *BMC public health* 15, no. 1 (2015): 1-12.
- [3] Paluch, Amanda E., Shivangi Bajpai, David R. Bassett, Mercedes R. Carnethon, Ulf Ekelund, Kelly R. Evenson, Deborah A. Galuska et al. "Daily steps and all-cause mortality: a meta-analysis of 15 international cohorts." *The Lancet Public Health* 7, no. 3 (2022): e219-e228.
- [4] Cuthbertson, Carmen C., Christopher C. Moore, Daniela Sotres-Alvarez, Gerardo Heiss, Carmen R. Isasi, Yasmin Mossavar-Rahmani, Jordan A. Carlson et al. "Associations of steps per day and step intensity with the risk of diabetes: the Hispanic Community Health Study/Study of Latinos (HCHS/SOL)." *International Journal of Behavioral Nutrition and Physical Activity* 19, no. 1 (2022): 1-14.
- [5] Sheng, Mingxin, Junyue Yang, Min Bao, Tianzhi Chen, Ruixue Cai, Na Zhang, Hongling Chen et al. "The relationships between step count and all-cause mortality and cardiovascular events: a dose-response meta-analysis." *Journal of sport and health science* (2021).
- [6] Thorup, Charlotte, John Hansen, Mette Grønkvær, Jan Jesper Andreassen, Gitte Nielsen, Erik Elgaard Sørensen, and Birthe Irene Dinesen. "Cardiac patients' walking activity determined by a step counter in cardiac telerehabilitation: Data from the intervention arm of a randomized

TABLE IV: The results of the proposed SC method on different walking regularities and sensor placements in the Pedometer dataset [19] using the LSTM with and without attention with 2 layers of 128 neurons, and the  $l_2$ -norm of the accelerometer signal (see Section IV-C) using different evaluation metrics (see Section IV-B).

	MAE	UC, OC	ER	RCA	ACC
Regular Walking—Sensor on Ankle					
LSTM-2 × 128-12	5.21	1.71, 1.87	-0.17±6.53	1.00±0.06	0.96±0.05
LSTM-2 × 128-12 (attention)	<b>2.00</b>	<b>1.04, 0.57</b>	<b>-0.68±3.77</b>	<b>1.00±0.037</b>	<b>0.98±0.03</b>
Semi-regular Walking—Sensor on Ankle					
LSTM-2 × 128-12	<b>13.87</b>	7.21, 8.39	-2.84±21.93	1.03±0.22	0.83±0.14
LSTM-2 × 128-12 (attention)	19.67	<b>5.44, 5.68</b>	<b>-1.75±14.69</b>	<b>1.02±0.15</b>	<b>0.88±0.09</b>
Regular Walking—Sensor on Wrist					
LSTM-2 × 128-12	5.65	2.20, 2.36	<b>-0.22±6.50</b>	<b>1.00±0.06</b>	0.95±0.05
LSTM-2 × 128-12 (attention)	<b>2.40</b>	<b>1.53, 0.40</b>	-1.54±7.70	1.02±0.080	<b>0.98±0.07</b>
Semi-regular Walking—Sensor on Wrist					
LSTM-2 × 128-12	<b>13.95</b>	9.26, 6.81	-7.46±25.43	1.07±0.25	0.81±0.19
LSTM-2 × 128-12 (attention)	14.59	<b>7.69, 8.82</b>	<b>-3.4±24.73</b>	<b>1.03±0.25</b>	<b>0.82±0.17</b>
Regular Walking—Sensor on Hip					
LSTM-2 × 128-12	5.37	2.11, 2.22	-0.26±6.40	1.00±0.07	0.96±0.05
LSTM-2 × 128-12 (attention)	<b>1.27</b>	<b>0.53, 0.50</b>	<b>-0.04±1.89</b>	<b>1.00±0.02</b>	<b>0.99±0.02</b>
Semi-regular Walking—Sensor on Hip					
LSTM-2 × 128-12	<b>9.29</b>	10.39, <b>10.68</b>	-7.72±35.81	1.07±0.35	0.75±0.27
LSTM-2 × 128-12 (attention)	9.38	<b>7.56, 13.35</b>	<b>-1.29±31.33</b>	<b>1.00±0.32</b>	<b>0.77±0.21</b>

TABLE V: The ACC [12] of the proposed LSTM and attention-based LSTM SC method compared to the previous deep-learning methods [12] on the whole (all the three sensor placements) regular walking samples in the Pedometer dataset [19] (see Sections IV-C and IV-B).

	ACC
LSTM [12]	0.9487
WaveNet [12]	0.9851
CNN [12]	<b>0.9872</b>
LSTM-2 × 128-12	0.9513
LSTM-2 × 128-12 (att)	<b>0.9868</b>

TABLE VI: The mean and standard deviation of the RCA [20] for the proposed LSTM with and without attention with 2 layers of 128 neurons compared to the previous time-domain methods [20] on the whole (all the three sensor placements) regular and semi-regular walking samples in the Pedometer dataset [19] (see Sections IV-C and IV-B).

	RCA	
	Regular	Semi-regular
Peak [20]	0.92±0.11	1.30±0.21
Threshold [20]	1.03±0.17	1.34±0.17
Autocorrelation [20]	0.95±0.24	0.93±0.17
LSTM-2 × 128-12	<b>1.00±0.04</b>	1.07±0.29
LSTM-2 × 128-12 (att)	<b>1.00±0.06</b>	<b>1.03±0.21</b>

controlled trial." Journal of Medical Internet Research 18, no. 4 (2016): e5191.

- [7] Flores, German H., and Roberto Manduchi. "Weallwalk: An annotated dataset of inertial sensor time series from blind walkers." ACM Transactions on Accessible Computing (TACCESS) 11, no. 1 (2018): 1-28.
- [8] Brajdic, Agata, and Robert Harle. "Walk detection and step counting on unconstrained smartphones." In Proceedings of the 2013 ACM international joint conference on Pervasive and ubiquitous computing, pp. 225-234. 2013.
- [9] Scholkman, Felix, Jens Boss, and Martin Wolf. "An efficient algorithm for automatic peak detection in noisy periodic and quasi-periodic signals." Algorithms 5, no. 4 (2012): 588-603.
- [10] Jayalath, Sampath, Nimsiri Abhayasinghe, and Iain Murray. "A gyro-

scope based accurate pedometer algorithm." In International Conference on Indoor Positioning and Indoor Navigation, vol. 28, p. 31st. 2013.

- [11] Mannini, Andrea, and Angelo Maria Sabatini. "A hidden Markov model-based technique for gait segmentation using a foot-mounted gyroscope." In 2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society, pp. 4369-4373. IEEE, 2011.
- [12] Luu, Long, Arvind Pillai, Halsey Lea, Ruben Buendia, Faisal M. Khan, and Glynn Dennis. "Accurate Step Count with Generalized and Personalized Deep Learning on Accelerometer Data." Sensors 22, no. 11 (2022): 3989.
- [13] Chen, Ziyi. "An LSTM recurrent network for step counting." arXiv preprint arXiv:1802.03486 (2018).
- [14] Bagui, Sikha, Xingang Fang, Subhash Bagui, Jeremy Wyatt, Patrick Houghton, Joe Nguyen, John Schneider, and Tyler Guthrie. "An improved step counting algorithm using classification and double autocorrelation." International Journal of Computers and Applications 44, no. 3 (2022): 250-259.
- [15] Pillai, Arvind, Halsey Lea, Faisal Khan, and Glynn Dennis. "Personalized Step Counting Using Wearable Sensors: A Domain Adapted LSTM Network Approach." arXiv preprint arXiv:2012.08975 (2020).
- [16] Ren, Peng, Fatemeh Elyasi, and Roberto Manduchi. "Smartphone-based inertial odometry for blind walkers." Sensors 21, no. 12 (2021): 4033.
- [17] Shao, Wenhua, Haiyong Luo, Fang Zhao, Cong Wang, Antonino Crivello, and Muhammad Zahid Tunio. "DePedo: Anti periodic negative-step movement pedometer with deep convolutional neural networks." In 2018 IEEE international conference on communications (ICC), pp. 1-6. IEEE, 2018.
- [18] Oord, Aaron van den, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. "Wavenet: A generative model for raw audio." arXiv preprint arXiv:1609.03499 (2016).
- [19] Mattfeld, Ryan, Elliot Jesch, and Adam Hoover. "A new dataset for evaluating pedometer performance." In 2017 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 865-869. IEEE, 2017.
- [20] Mattfeld, Ryan, Elliot Jesch, and Adam Hoover. "Evaluating Pedometer Algorithms on Semi-Regular and Unstructured Gaits." Sensors 21, no. 13 (2021): 4260.
- [21] Sherstinsky, Alex. "Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network." Physica D: Nonlinear Phenomena 404 (2020): 132306.
- [22] Luong, Minh-Thang, Hieu Pham, and Christopher D. Manning. "Effective approaches to attention-based neural machine translation." arXiv preprint arXiv:1508.04025 (2015).
- [23] Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." arXiv preprint arXiv:1409.0473 (2014).

- [24] Coxe, Stefany, Stephen G. West, and Leona S. Aiken. "The analysis of count data: A gentle introduction to Poisson regression and its alternatives." *Journal of personality assessment* 91, no. 2 (2009): 121-136.
- [25] Paszke, Adam, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen et al. "Pytorch: An imperative style, high-performance deep learning library." *Advances in neural information processing systems* 32 (2019).
- [26] Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel et al. "Scikit-learn: Machine learning in Python." *the Journal of machine Learning research* 12 (2011): 2825-2830.
- [27] Bai, Shaojie, J. Zico Kolter, and Vladlen Koltun. "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling." *arXiv preprint arXiv:1803.01271* (2018).
- [28] Hao, Hongyan, Yan Wang, Yudi Xia, Jian Zhao, and Furoo Shen. "Temporal convolutional attention-based network for sequence modeling." *arXiv preprint arXiv:2002.12530* (2020).