
WaveSense: A Spiking Neural Network Inspired by the WaveNet Architecture for Keyword Spotting

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Ultra-low power local signal processing is a crucial aspect for edge applications on
2 always-on devices. Neuromorphic processors emulating spiking neural networks
3 show great computational power while fulfilling the limited power budget as
4 needed in this domain. In this work we develop a spiking neural network inspired
5 by the WaveNet architecture which is based on simple neural dynamics, fixed
6 time-constants and is strictly feed-forward and hence is particularly well suited for
7 a neuromorphic implementation. We test the capabilities of this model on several
8 datasets for keyword-spotting. The results show that the proposed network beats the
9 state-of-the-art of other spiking neural networks and reaches near state-of-the-art
10 performance of artificial neural networks such as CNNs and LSTMs.

11 **Keywords**— spiking, keyword spotting, temporal processing, streaming, audio, neuromorphic, wavenet,
12 wavesense, always-on, low-power

13 1 Introduction

14 Local signal processing is an important component of the computational pipeline for Internet-of-Things (IoT)
15 devices equipped with a range of sensors like audio, video, and motion sensing. A significant range of these
16 sensors capture signals comprising temporal features. Ideally these features need to be extracted by an on-board
17 processor before decision making or relaying the pre-processed information for further processing. Processing
18 temporal signals is often computationally challenging and requires large amounts of memory and power,
19 especially in always-on scenarios. Neuromorphic [13] processors with spiking-neural networks have shown
20 promise in this domain as ultra-low power compact solutions [10, 2, 14, 8, 6, 12].

21 In this work we propose WaveSense a Spiking Neural Network (SNN) model suitable for efficient neuromorphic
22 implementations while retaining high accuracy on temporal data streams. This work bridges the performance
23 gap between Artificial Neural Networks (ANNs) and SNNs for temporal tasks. Crucially, the proposed model

- 24 • accepts spike streams and not ‘buffered frames’ as input,
- 25 • requires no delays in its connectivity,
- 26 • utilizes a very simple spiking neuron model - Leaky Integrate and Fire (LIF) neuron - without the need
27 of any additional adaptive mechanisms,
- 28 • does not require recurrent connectivity (recurrent connectivity can often be difficult to tune or train).
- 29 • produces a high classification performance.

30 Recently several works have shown how to build efficient SNNs with accuracies equivalent to ANNs [7, 18]. In
31 these studies, spiking neurons are used in rate mode with equivalent response curves to ReLU activations, to
32 transfer weights from pre-trained ANNs to SNNs. This approach therefore completely neglects the temporal
33 capabilities of spiking neurons. On the other hand, surrogate gradient methods enable directly training SNN using
34 Back Propagation Through Time (BPTT) [15]. Shrestha and Orchard [21] for instance show temporal processing

35 on temporal gesture recognition task and show a good performance on a visual task. Similar approaches have
36 already been investigated on audio tasks [1, 26, 5]. Bellec et al. [1] demonstrate classification results on TIMIT
37 dataset using long time constants, a complex learning strategy (Deep-R) and require significant computational
38 resources to train. Wu et al. [26] train SNNs for automatic speech recognition tasks in a tandem approach with
39 an ANNs. This training pipeline integrated a language model and pronunciation model which goes beyond the
40 capabilities of the neuromorphic system. The same authors showed in an earlier study the capabilities of a SNN
41 in combination with a self-organized map to learn to recognize digits using the TIDIGIT dataset [25].

42 Blouw et al. [3] demonstrate high accuracy in a audio classification task using dense networks with spectrograms
43 as inputs. This approach requires passing the frequency data of previous time steps (defined by the spectrogram
44 time window) in every sample presentation to the network. Kugele et al. [11] show that by matching ANN
45 roll-out delays to the propagation delays in SNN, the resultant networks can demonstrate a high accuracy on
46 vision based spatio-temporal classification tasks. Implementation of delays in neuromorphic hardware requires
47 additional memory resources to store and deliver spikes in a delayed fashion and could be potentially quite
48 expensive. Yin, Corradi, and Bohté [27] use a Spiking Recurrent Neural Network (SRNN) architecture and
49 demonstrate that by utilizing *adaptive* LIF neurons and learning the time constants, these networks can perform
50 temporal classification tasks in a sequential manner fairly well. The authors demonstrate the effective use of
51 spiking neural networks and show a significant increase in power efficiency. Unfortunately, having fine tuned
52 time constants in low-power neuromorphic hardware can often be challenging especially while using fixed
53 precision numerical representations and computations.

54 We propose a novel network architecture for SNN that does not require buffering or delays and can directly
55 process spiking data from event-based sensor using simple LIF neurons. Our architecture is derived from first
56 principles and inspired by the WaveNet [16] architecture that does not necessitate learning the time constants of
57 the system but could be defined as the task demands. In addition we also propose an efficient training strategy
58 and a corresponding loss-function that is suitable for streaming based models, in particular models that could be
59 run in real-time with neuromorphic hardware.

60 The first key aspect of the WaveNet architecture is the use of a multi-layer causal dilated convolutions. The
61 *causal* refers to the use of data from the past, *dilated* refers to a sparse kernel and the *convolution* is along the
62 time axis. Stacking such convolutional operations along multiple layers enables the network to have a long
63 temporal memory. A second aspect is that it eliminates the need for sliding window based inference/prediction
64 and minimizes the number of computations within the network when operating on a continuous stream of data.

65 The WaveNet architecture is very amenable to general purpose micro-processors and micro-controllers but it still
66 requires a reasonable amount of memory and non-linear computations such as *tanh* and *sigmoid* which are
67 fairly complex (within the context of ultra-low power devices) and this inturn requires higher energy and power
68 requirements. Neuromorphic technology promises to bring the energy required for these tasks down by utilizing
69 SNNs to perform ultra-low power computation. So far, while neuromorphic devices have been demonstrated to
70 operate at an extremely low power, they have fallen short in demonstrating computational performance that is
71 on-par or comparable to state of the art ANNs for temporal tasks, most prominently in the audio domain.

72 WaveSense model proposed here aims to bridge this gap. For the results in this work, we focus on audio tasks as
73 spatio-temporal tasks without loss of generality. Sec. 2 details all the methods used for audio data pre-processing,
74 conversion to spikes and the details of the network architecture. In Sec. 3 we demonstrate the computational
75 capability of this network over several audio datasets for key-word spotting task. We compare our results to the
76 state of the art SNNs and ANNs. We conclude in Sec. 4 where we discuss the implications and impact of this
77 work and potential areas where this work could be utilized.

78 Most importantly all the code used to generate the reported results have been made open-source and can be found
79 at <http://XXXXXXXXXXXXXXXX>. We believe this will enable the research community to explore other avenues
80 that could take advantage of our work.

81 2 Materials and Methods

82 2.1 Neuron Model

83 In this work we used Leaky Integrate and Fire (LIF) neuron model with synaptic time constant τ_s and membrane
84 time constant τ_v . The sub-threshold dynamics of this neuron are described below.

$$\dot{v}(t) = -v(t)/\tau_v + i_s \quad (1)$$

$$\dot{i}_s(t) = -i_s(t)/\tau_s + \sum w_j s_j(t) \quad (2)$$

85 where v is the membrane potential, i_s is the synaptic current, w the synaptic weight and s is the input spike train.

86 In-order to optimize simulation time and computational efficiency, we make some alterations. Unlike a traditional
 87 LIF which resets to a resting potential upon reaching threshold θ , we *subtract* a fixed value θ from the membrane
 88 potential. Further more, if the membrane potential increases beyond $N\theta$, where N is an arbitrary positive integer,
 89 then this neuron produces N spikes, and proportionally $N\theta$ is *subtracted* from the membrane potential.

$$s(t) = \begin{cases} [v(t)/\theta], & \text{if } v(t) \geq \theta \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

90 The use of a mechanism for generating multiple spikes enables the computation to be more robust to the choice
 91 of simulation time steps. This enables us to choose a relatively large time step for our simulations. In practice
 92 we observe that some neurons occasionally *do* produce multiple spikes. (Producing multiple spikes in a single
 93 simulation time-step is not necessary if one chooses a small enough time step.)

94 For further simulation efficiency the LIF neurons were simulated using the Spike Response model (SRM) [9].

$$v(t) = \sum_j w_j (\epsilon * s_j)(t) + (\nu * s)(t) \quad (4)$$

95 where

$$\epsilon_s(t) = e^{-t/\tau_s} \quad (5)$$

$$\epsilon_v(t) = e^{-t/\tau_v} \quad (6)$$

$$\nu(t) = -\theta e^{-t/\tau_v} \quad (7)$$

$$\epsilon(t) = (\epsilon_s * \epsilon_v) \quad (8)$$

96 We use exponential kernels for synaptic $\epsilon_s(t)$ and membrane dynamics $\epsilon_v(t)$ and derive the Post Synaptic
 97 Potential (PSP) kernel $\epsilon(t)$. The refractory kernel $\nu(t)$ is also a negative exponential kernel with the same time
 98 constant τ_v as the membrane potential. The symbol $*$ denotes a convolution operation.

99 As spike generation is non-differentiable, we use a surrogate gradient [15]. Several profiles for the surrogate
 100 gradients have been proposed in literature. We use a modified exponential kernel for the surrogate gradient
 101 function. In order to accommodate the multi-spike behavior of the neuron, we choose a periodic exponential
 102 function (Figure 1) as the surrogate gradient. This function peaks as the membrane potential approaches multiples
 103 of neuron spiking threshold θ . Intuitively, this gradient function maximizes the impact of a parameter when the
 104 neuron is close to spiking or has just spiked and is a variant of the exponential gradient function. An extreme
 105 simplification of the periodic exponential would be a Heaviside function¹.

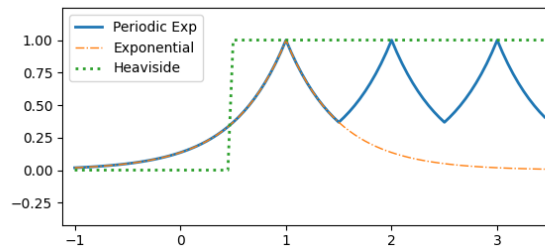


Figure 1: The surrogate gradient of the spiking neuron as a function of the (pre-spike) membrane potential.

106 The simulations were run with the library ‘sinabs’ [19]² using an adaptation of SLAYER [21].

¹The heaviside function like a Rectified Linear Unit (ReLU) has a range of membrane potentials where the gradient is 0 and could potentially prevent the network from learning at low activity levels.

²GNU AGPL v3 License

107 **2.2 Network description**

108 We take inspiration from the work of Coucke et al. [4] who uses the WaveNet architecture [16] for classification
 109 of continuous audio streams. The WaveNet architecture provides a prescription for distributing temporal memory
 110 and computation across layers without repeated presentation of previous input data.

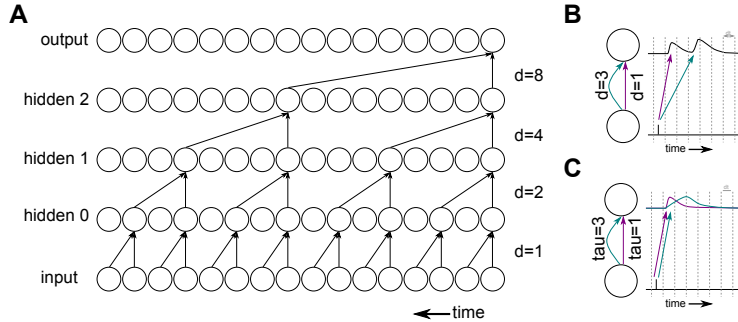


Figure 2: Relegating the job of delays in dilated convolutions to synaptic dynamics.

111 The original ANN WaveNet model comprises of a few different computational building blocks. We translate
 112 each of these building blocks to SNNs as follows.

113 **Dilated Temporal Convolutions** are the crux of what makes WaveNets efficient. The dilation enables convolu-
 114 tion over a large time window (or sample length) while still only using a small number of parameters. These
 115 convolutions perform a weighted-accumulation of information from different points in time, separated by the
 116 *dilation* parameter. This is done by storing previous activations in ANNs, and naively, this could be mimicked
 117 by replacing it by synaptic transmission delays. In this work, we observe that neuron and synaptic dynamics
 118 could be seen as proxies for temporal convolutional processing as shown in Figure 2. Figure 2B shows the
 119 contributions of each projection with a kernel size 2 would look like when realized with delays. Alternatively, if
 120 we choose an appropriate set of synaptic time constants τ_s , this is shown in Figure 2C. While quantitatively,
 121 these are different, qualitatively both these approaches provide the ability to transform and project information
 122 in the temporal domain [20]. This is the key insight we leverage to design our final model.

123 **Rectified Linear Unit (ReLU)** activations can be *approximated* by spiking neurons because a spiking neuron
 124 can only produce spikes if the membrane potential crosses a threshold and is silent otherwise [7].

125 **Non-linear activations** like *tanh* and *sigmoid* cannot be efficiently translated to SNNs. So we choose to
 126 replace these activations with SNNs activations (potentially ignoring the benefits of filtering and gating). In the
 127 original model, these two activations are preceded by two sets of weights, where as in our SNN we only use one
 128 weighted projection. (See Figure 3)

129 **Residual connections** and summation (+) are realized by a synaptic connection.

130 The WaveSense model is built upon these building blocks as shown in Figure 3. It comprises of several ‘blocks’,
 131 each of which comprises of three spiking neuron layers. The first spiking layer receives inputs filtered by two
 132 separate synapses with different time constants τ_s and weights. The time constants τ_s of the slow projections in
 133 each of these blocks are chosen such that they span a range of values relevant to the task. The number of blocks
 134 is chosen such that the sum of all these time constants is proportional to the temporal memory demanded by the
 135 task. This layer projects to the second spiking layer. Additionally a third spiking layer in each of these blocks
 136 projects to a ‘hidden’ layer followed by a non-spiking ‘low pass’ readout layer. The output of this block is the
 137 summation of its input (residual connection) and the output of the second layer. These ‘blocks’ are connected
 138 in a feed forward manner. The non-spiking ‘low-pass’(LP) layer simply acts as a weighted low pass filter on
 139 the spikes of the ‘hidden’ layer. This is equivalent to the synapse of a spiking neuron (without the neuron’s
 140 membrane potential or the spiking dynamics) and does not require any extra components unavailable to spiking
 141 neurons. The choice of leaving the output layer to be non-spiking is to enable a smooth, continuous valued
 142 readout useful for faster learning using BPTT.

143 **2.3 Datasets**

144 In order to evaluate the efficacy of the proposed model, we train and test it against several open-source publicly
 145 available audio datasets.

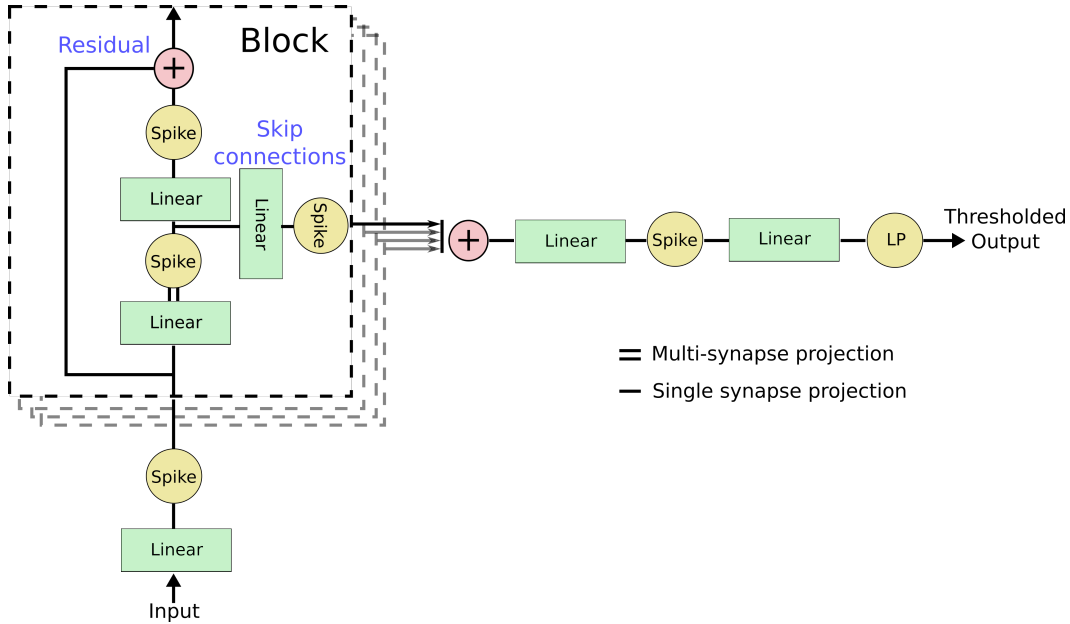


Figure 3: The WaveSense model prescribed in this work is a theoretical adaptation of the WaveNet architecture [16] based on first principles.

146 2.3.1 Aloha Dataset

147 The Aloha dataset [3] is a small collection of audio samples containing the keyword 'Aloha' and several
 148 distractors such as 'take a load off'. As the dataset is very small, only ~ 2000 samples, we augmented the
 149 samples using the MUSAN noise dataset [22]. For that end we standardized the sample length of each utterance
 150 in the training and validation set to five seconds and added randomly selected background noise data with a
 151 signal-to-noise ratio (SNR) of 5 dB to the training data.

152 2.3.2 Hey Snips Dataset

153 The 'Hey Snips' dataset [4] for wake phrase spotting distinguishes between two classes. The positive class
 154 contains 11'000 utterances from over 2'000 speakers of the wake phrase 'Hey Snips' while the negative (or
 155 distractor) class contains over 86'000 negative examples from more than 6'000 speakers. We split the data into
 156 a training-, validation- and test-set as provided by the authors of the dataset. We standardized the sample length
 157 of each utterance in the training and validation set to five seconds. As the dataset is already very large, we did
 158 not augment the data with noise.

159 2.3.3 Speech Commands Dataset

160 The Speech Commands [24] describe a dataset containing 35 keywords uttered in total 105'000 times from over
 161 2'600 speakers. The keywords contain the numbers 0 - 10, commands such as "stop", "go", "left" and "right" as
 162 well as other words like "Marvin", "Sheila", etc. This dataset was initially designed for keyword spotting in a
 163 limited vocabulary and the intended experiment is to detect 10 commands (plus silence) out of all 35 keywords (12 classes in total).
 164 Nevertheless, there are studies training models and showing results on all 35 classes [5].
 165 We augment the training set with noise data from the MUSAN dataset using a SNR of 5 dB just as we do for the
 166 Aloha dataset.

167 2.4 Pre-processing

168 All datasets we use in this work contain many wave audio files of a few seconds length. The data is pre-processed
 169 in several stages:

- 170 • **Noise augmentation** The training data is augmented with noise from the MUSAN noise dataset using
 171 a SNR of 5dB (except for the HeySnips data).
- 172 • **Length standardization and pre-amplification** The noise augmented waveform is cut into a standard
 173 length (dependent on the dataset) and the amplitude is normalized.

- 174 • **Band-pass filters** The audio is then passed through 64 Butterworth bandpass filters 2nd order. The
175 bandpass filters are distributed in Mel-scale between 100Hz and 8kHz.
- 176 • **Rectification** The response of the 64 bandpass filters is rectified using a full-wave rectifier.
- 177 • **Spike conversion and binning** The output of the rectifier is applied as direct input to the membrane
178 of 64 simplified LIF neurons resulting in a rate code. The spike trains are binned into 10ms timesteps
179 allowing multiple spikes per timestep.

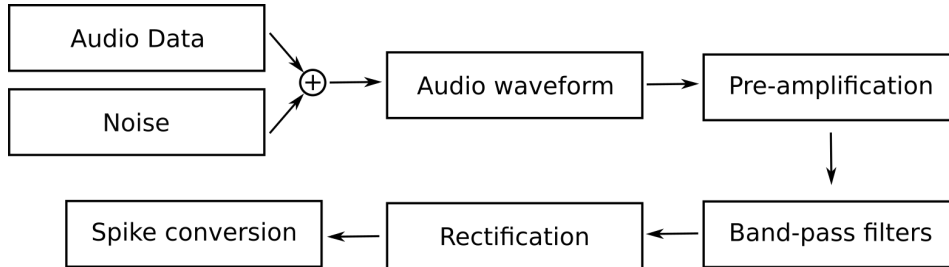


Figure 4: Data pre-processing pipeline figure.

180 2.5 Training Method

181 In order to train the parameters of SNN (See Sec. 2.2) we use BPTT. In particular we aim to be able to deploy
182 the network in streaming mode *ie.* the model receives the data stream directly generated from a sensor without
183 any frame based (sliding window) buffering. This requires us to employ an appropriate loss function.

184 Often in a classification task, the output class can be determined by computing cross-entropy loss on the sum of
185 the outputs over the sample length for each output neuron. While this would yield a good classification accuracy,
186 the magnitude of the output trace at ‘a given point in time’ is not indicative of the network prediction. This is not
187 ideal for models being run in streaming mode.

188 2.5.1 Peak Loss

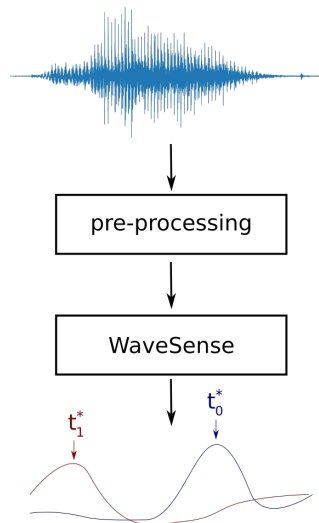


Figure 5: Visualization of the peak loss with peak times t_c^* .

189 Typically for streaming models, a signal is predicted as belonging to a certain class when the corresponding
190 output trace exceeds a ‘detection threshold’. This approach is also ideal for always-on neuromorphic systems.
191 We therefore design our loss function to reflect this detection mechanism and train our neural networks. We
192 determine the peaks of the output traces and use only the activation values at the peaks to compute the cross
193 entropy loss (see Figure 5) similar to *max-over-time loss* [5].

194 Consequently the loss is computed as follows:

$$L_{CE} = - \sum_c \lambda_c \log(p_c) \quad (9)$$

195 where λ_c yields 1 if class label c corresponds to the current input and 0 otherwise. p_c is the prediction probability
 196 by the neural network that the current input belongs to class c . It is calculated by a *softmax* operation as shown
 197 below.

$$p_c = \frac{e^{\hat{y}_c}}{\sum_i e^{\hat{y}_i}} \quad (10)$$

198 where \hat{y} are the 'logits' produced by the neural network.

199 For temporal tasks, the input $\mathbf{x} = x^T = x^1 \dots x^T$ and the output (logits) \hat{y} of the neural network are time-series
 200 over time T .

$$\hat{y}^t = f(x^t | \Theta, s^t) \quad (11)$$

201 where f is the transformation of the neural network, Θ are the network parameters and s^t is the internal state of
 202 the network at time t . In *peak-loss* we pass the peak of each output trace to the *softmax* function. The peaks are
 203 calculated as follows:

$$\hat{y}_c = \max(y_c^T) = y_c^{t_c^*} \quad (12)$$

204 where $t_c^* = \operatorname{argmax}(y_c^T)$ is the 'peak time', the time of maximal activation of output trace c (see Figure 5).

205 2.5.2 Spiking activity regularization

206 The activity of LIF neurons can change dramatically during the learning process. This is particularly true while
 207 using the multi-spike neuron model as we do in this WaveSense implementation. This could potentially eliminate
 208 the benefit of using spiking neurons by spiking at a high firing rate at every time step and consequently have no
 209 sparsity. This could also lead to high energy utilization of the network in a neuromorphic implementation.

210 In order to limit the activity of these neurons and maintain sparse activity, we include an activity regularizer term
 211 in our loss function [23].

$$L_{act} = (N_{spk}^\dagger / (T \cdot N_{neurons}))^2 \quad (13)$$

212 where the activation loss L_{act} is dependent on the total *excess* number of spikes N_{spk}^\dagger produced by the network
 213 with a population size $N_{neurons}$ in response to an input of length T time steps. N_{spk}^\dagger is given as:

$$N_{spk}^\dagger = \sum_i \sum_t N_i^t \Theta(N_i^t - 1) \quad (14)$$

214 is the sum of spikes from all neurons N_i exceeding 1 in each time bin t (Θ is a heaviside function).

215 Finally the loss function is given as:

$$L = L_{CE} + \alpha L_{act} \quad (15)$$

216 where α was chosen to be 0.01.

217 3 Results

218 In order to validate and verify that sufficient information from the input is retained after pre-processing and
 219 conversion to spikes, we train a state-of-the-art WaveNet classifier on the datasets considered in this work and
 220 check that we can obtain a high accuracy. We implement a non-spiking dilated Convolutional Neural Network
 221 (CNN) to replicate the WaveNet architecture very similar to that described in Coucke et al. [4] and Oord et al.
 222 [16] (see Section 2.2 for details).

223 We train this ANN on the HeySnips, Aloha and Speechcommands datasets and compare our results to those
 224 reported in literature [4, 3, 5]. The results obtained from this network are then used as baseline to evaluate the
 225 performance of the proposed SNN.

Table 1: Aloha result model size and resource comparison.

| Publication | #Neurons | #Parameters | Accuracy |
|------------------|----------|-------------|------------|
| Blouw et al. [3] | 541 | 172800 | 95.8 |
| This work | 864 | 18482 | 98.0 ± 1.1 |

226 **3.1 Aloha dataset**

227 In order to compare our model to other SNN implementations in the keyword spotting domain, we trained our
 228 WaveSense on the Aloha dataset [3]. Table 1 shows the network and memory resources the proposed model
 229 utilizes in comparison to the work demonstrated in Blouw et al. [3]. With an average accuracy of 98.0% with a
 230 standard deviation of 1.1% the model presented in this work performs significantly better while at the same time
 231 requiring a significantly smaller number of parameters. The best runs of the WaveSense model yielded 99.5%
 232 accuracy which is equal to the performance of the ANN model. It is important to note that the key focus of the
 233 work by Blouw et al. [3] is to benchmark energy and power consumption and not model performance.

234 **3.2 HeySnips dataset**

235 On the HeySnips dataset, our implementation of the WaveNet reaches an accuracy of 99.8% on the clean dataset.
 236 In Coucke et al. [4], the authors do not report any accuracy number but rather report the false rejection rate (FRR)
 237 of 0.12% for a fixed false alarm per hours (FAPH) of 0.5. In order to compare our results more accurately, we
 238 implement the same metrics; our WaveNet implementation reaches 0.95 FAPH and a 0.8% FRR on the test set.
 239 These results are slightly worse than the results reported by Coucke et al. [4] but that is expected as we do not
 240 apply the same specific methods to improve performance such as "End-Of-Keyword labeling" and "masking".
 241 Without those methods and without gating, the FRR reported by Coucke et al. [4] drops to 0.98%. On the other
 242 hand, our WaveNet implementation reaches similar or even better results than the CNN and LSTM reported by
 243 Coucke et al. [4]. This fact shows that our pre-processing method indeed extracts sufficient information from the
 244 input such that a neural network can reach very high accuracy. Hence, we train a spiking version of the WaveNet
 245 architecture (WaveSense), as described in 2.2, on the same data.

246 In the WaveSense model we do not use any gating mechanism, a kernel size of 2 and only 8 layers; much
 247 less compared to the 24 layers and kernel size of 3 as used in our WaveNet implementation and the WaveNet
 248 implementation by Coucke et al. [4]. The memory in our model is still long enough as the WaveSense implements
 249 the dilations using synaptic dynamics with long time constants but the number of parameters drops from 47090 to
 250 13042. Despite the low number of parameters and quantization from spiking activations, the WaveSense model
 251 achieves an average accuracy of 99.6% over 11 runs (only drops by 0.2%) . Our best run of the WaveSense
 252 model yielded the same accuracy (of 99.8%) as our WaveNet implementation. With an FRR = 1.0% and FAPH
 253 = 1.34 the performance is indeed lower than the WaveNet, but it is comparable to that of LSTM and CNN as
 254 reported by Coucke et al. [4].

255 **3.3 Speechcommands dataset**

256 We also trained WaveSense on the Speechcommands dataset. We evaluated our model by training it to classify
 257 all 35 classes in the dataset. In a study by Perez-Nieves et al. [17] in which the authors investigate the impact
 258 of heterogeneity of time constants on the performance, the best model reached ~ 57.3% accuracy on the same
 259 dataset. In Cramer et al. [5] the best performing SNN is a recurrent network which yields ~ 50.9% accuracy of
 260 all 35 classes. In the same study, also an LSTM and CNN are trained on the same data resulting in an accuracy
 261 of ~ 73% resp. ~ 77.7%. The WaveSense model reaches an average accuracy of 79.6% over 11 runs (best
 262 80.0%) which is significantly higher than the best SNN described in previous studies. Notably, WaveSense
 263 performs better than the reported LSTM and CNN [5].

264 **4 Discussion and Conclusion**

265 While the results demonstrated here are obtained using a fixed set of time constants, it is conceivable that
 266 according to the constraints of the neuromorphic hardware, an appropriate network could be trained to obtain
 267 qualitatively similar results. This holds true even for mixed-signal neuromorphic devices [10] with programmable
 268 weights and tune-able time constants. Because the algorithm provides a recipe for how to choose the time
 269 constants in the network, even if a neuromorphic substrate has a limited range of time constants, a number of
 270 layers with an appropriate combination (sum) of time constants can always be chosen to fit the temporal task.
 271 This is in stark contrast to recurrent neural networks that often require a tight balance between excitation and
 272 inhibition and long time constants [1, 27].

Table 2: A comparison of model performance for various datasets and network architectures.

| Publication | Dataset | Accuracy (%) | Architecture |
|--------------------------|--------------------|--------------------------------|--------------|
| Coucke et al. [4] | HeySnips | FRR 0.12 FAPH 0.5 | WaveNet |
| Coucke et al. [4] | HeySnips | FRR 2.09 FAPH 0.5 | LSTM |
| Coucke et al. [4] | HeySnips | FRR 2.51 FAPH 0.5 | CNN |
| This work | HeySnips | 99.8 (FRR 0.8 FAPH 0.95) | WaveNet |
| This work | HeySnips | 99.6 ± 0.1 (FRR 1.0 FAPH 1.34) | SNN |
| Cramer et al. [5] | SpeechCommands(35) | 50.9 ± 1.1 | SNN |
| Cramer et al. [5] | SpeechCommands(35) | 73 ± 0.1 | LSTM |
| Cramer et al. [5] | SpeechCommands(35) | 77.7 ± 0.2 | CNN |
| Perez-Nieves et al. [17] | SpeechCommands(35) | 57.3 ± 0.4 | SNN |
| This work | SpeechCommands(35) | 87.6 | WaveNet |
| This work | SpeechCommands(35) | 79.6 ± 0.1 | SNN |
| Blouw et al. [3] | Aloha | 93.8 | SNN |
| This work | Aloha | 99.5 | WaveNet |
| This work | Aloha | 98.0 ± 1.1 | SNN |

273 The choice of time constants and number of layers is informed by the total temporal memory required by the task.
 274 We choose them in a similar fashion to that of WaveNet with time constants increasing with factors of 2 and
 275 such that the sum of all the time constants is proportional to τ_{task} . Typically we observe that a proportionality
 276 of 2.5 is suitable with a kernel size of 2. The proportionality factor is the length of time after which the effect of
 277 a PSP is negligible. This also translates to compact networks with fewer parameters for the same amount of
 278 temporal memory (at the same time resolution). In other words, given a network, the temporal memory of a
 279 given task can be computed as follows:

$$\tau_{task} \approx 2.5 \sum_i \tau_s^i \quad (16)$$

280 where i is the list of all the layers in the WaveSense network.

281 While the results reported here are significantly high, we believe this can be further improved by improving the
 282 loss function. For instance the *max-over-time* loss computed only during the presence of a keyword as opposed
 283 to the entire sample [4] has been shown to improve performance of such models. Furthermore, a thorough
 284 architecture search could potentially result in a better combination of time constants τ_m and τ_s , number of
 285 channels, kernel sizes etc.

286 A crucial factor in adopting a model is ease of training, deployment and power efficiency. By utilizing simple
 287 LIF neurons, we take full advantage of their computational efficiency [27] in addition to sparse computations
 288 afforded by SNNs. While training SNNs is relatively slow on CPUs and GPUs, utilizing the SRM in combination
 289 with the SLAYER algorithm [21], we are able to train at a relatively high speed. All experimental results reported
 290 in this manuscript were performed on a single NVIDIA 1080 Ti with a few hours of run time per experiment. We
 291 further improve upon this efficiency with a custom fork of the SLAYER implementation³. The resulting models
 292 while accurate within the SRM framework, are not *identical* to simulations based on LIF neurons, supported by
 293 most digital neuromorphic devices. But we find that they are a close approximation and a quick retraining can
 294 recover the model’s performance using LIF neurons.

295 The WaveNet architecture requires storing activations of each of its layers depending on their kernel size and
 296 dilation value: $N_{buf} \propto (k - 1) \cdot d + 1$. In contrast, WaveSense does not buffer any spikes(activations) from the
 297 past explicitly. Instead the information is retained in the neuron and synaptic states: $N_{buf} \propto k + 1$. This makes
 298 WaveSense extremely efficient in terms of memory utilization in contrast to WaveNet.

299 The results demonstrated here show that the WaveSense architecture is suitable for audio classification tasks and
 300 show a promising performance improvement in comparison to prior state-of-the-art. Audio signals, after they are
 301 pre-processed are equivalent to a population of neurons producing spike patterns with complex spatio-temporal
 302 correlations. We argue therefore, that the results presented here can be extended to other modalities of sensory
 303 data such as ECG, PPG, machine vibrations or DVS data.

304 This work we believe could contribute towards a future with a ubiquitous abundance of always-on audio and
 305 other sensory devices responding to user commands. This could lead to potential misuse of the technology
 306 for surveillance. Thankfully, neuromorphic algorithms such as the one proposed here require specialized
 307 neuromorphic hardware to take full advantage. If the availability of such hardware could be regulated, we hope
 308 that the society can benefit from this technology while protecting itself from misuse.

³<https://XXXXXXXXXXXXXXXXX>

309 References

- 310 [1] Guillaume Bellec et al. “Long short-term memory and learning-to-learn in networks of spiking
311 neurons”. In: *arXiv preprint arXiv:1803.09574* (2018).
- 312 [2] Ben Varkey Benjamin et al. “Neurogrid: A mixed-analog-digital multichip system for large-
313 scale neural simulations”. In: *Proceedings of the IEEE* 102.5 (2014), pp. 699–716.
- 314 [3] Peter Blouw et al. *Benchmarking Keyword Spotting Efficiency on Neuromorphic Hardware*.
315 2018. arXiv: 1812.01739 [cs.LG].
- 316 [4] Alice Coucke et al. *Efficient keyword spotting using dilated convolutions and gating*. 2018.
317 arXiv: 1811.07684 [cs.LG].
- 318 [5] Benjamin Cramer et al. “The Heidelberg Spiking Data Sets for the Systematic Evaluation of
319 Spiking Neural Networks”. In: *IEEE Transactions on Neural Networks and Learning Systems*
320 (2020).
- 321 [6] Mike Davies et al. “Loihi: A neuromorphic manycore processor with on-chip learning”. In:
322 *Ieee Micro* 38.1 (2018), pp. 82–99.
- 323 [7] Peter U Diehl et al. “Fast-classifying, high-accuracy spiking deep networks through weight
324 and threshold balancing”. In: *2015 International joint conference on neural networks (IJCNN)*.
325 ieee. 2015, pp. 1–8.
- 326 [8] Steve B Furber et al. “The spinnaker project”. In: *Proceedings of the IEEE* 102.5 (2014),
327 pp. 652–665.
- 328 [9] Wulfram Gerstner. “A framework for spiking neuron models: The spike response model”. In:
329 *Handbook of Biological Physics*. Vol. 4. Elsevier, 2001, pp. 469–516.
- 330 [10] Giacomo Indiveri et al. “Neuromorphic silicon neuron circuits”. In: *Frontiers in neuroscience*
331 5 (2011), p. 73.
- 332 [11] Alexander Kugele et al. “Efficient Processing of Spatio-Temporal Data Streams With Spiking
333 Neural Networks”. In: *Frontiers in Neuroscience* 14 (2020), p. 439. ISSN: 1662-453X. DOI:
334 10.3389/fnins.2020.00439. URL: <https://www.frontiersin.org/article/10.3389/fnins.2020.00439>.
- 335 [12] Qian Liu et al. “Live demonstration: face recognition on an ultra-low power event-driven con-
336 volutional neural network ASIC”. In: *Proceedings of the IEEE/CVF Conference on Computer*
337 *Vision and Pattern Recognition Workshops*. 2019.
- 338 [13] Carver Mead. “Neuromorphic electronic systems”. In: *Proceedings of the IEEE* 78.10 (1990),
339 pp. 1629–1636.
- 340 [14] Paul A Merolla et al. “A million spiking-neuron integrated circuit with a scalable communica-
341 tion network and interface”. In: *Science* 345.6197 (2014), pp. 668–673.
- 342 [15] Emre O. Neftci, Hesham Mostafa, and Friedemann Zenke. *Surrogate Gradient Learning in*
343 *Spiking Neural Networks*. 2019. arXiv: 1901.09948 [cs.NE].
- 344 [16] Aaron van den Oord et al. *WaveNet: A Generative Model for Raw Audio*. 2016. arXiv: 1609.
345 03499 [cs.SD].
- 346 [17] Nicolas Perez-Nieves et al. “Neural heterogeneity promotes robust learning”. In: *bioRxiv*
347 (2021), pp. 2020–12.
- 348 [18] Bodo Rueckauer et al. “Conversion of continuous-valued deep networks to efficient event-
349 driven networks for image classification”. In: *Frontiers in neuroscience* 11 (2017), p. 682.
- 350 [19] Sadique Sheik and Martino Sorbaro. *Project Title*. <https://sinabs.ai/>. 2013.
- 351 [20] Sadique Sheik et al. “Emergent Auditory Feature Tuning in a Real-Time Neuromorphic VLSI
352 System”. In: *Frontiers in Neuroscience* 6 (2012), p. 17. ISSN: 1662-453X. DOI: 10.3389/
353 fnins.2012.00017. URL: <https://www.frontiersin.org/article/10.3389/fnins.2012.00017>.
- 354 [21] Sumit Bam Shrestha and Garrick Orchard. “SLAYER: Spike Layer Error Reassignment in
355 Time”. In: *Advances in Neural Information Processing Systems* 31. Ed. by S. Bengio et al.
356 Curran Associates, Inc., 2018, pp. 1419–1428. URL: <http://papers.nips.cc/paper/7415-slayer-spike-layer-error-reassignment-in-time.pdf>.
- 357 [22] David Snyder, Guoguo Chen, and Daniel Povey. *MUSAN: A Music, Speech, and Noise Corpus*.
358 arXiv:1510.08484v1. 2015. eprint: 1510.08484.
- 359 [23] Martino Sorbaro et al. “Optimizing the energy consumption of spiking neural networks for
360 neuromorphic applications”. In: *Frontiers in neuroscience* 14 (2020), p. 662.
- 361
362
363

- 364 [24] Pete Warden. “Speech commands: A dataset for limited-vocabulary speech recognition”. In:
365 *arXiv preprint arXiv:1804.03209* (2018).
- 366 [25] J Wu et al. *A spiking neural network framework for robust sound classification*. *Front. Neurosci.*
367 *12* (2018). 2018.
- 368 [26] Jibin Wu et al. “A Tandem Learning Rule for Effective Training and Rapid Inference of Deep
369 Spiking Neural Networks”. In: *arXiv e-prints* (2019), arXiv–1907.
- 370 [27] Bojian Yin, Federico Corradi, and Sander M Bohté. “Effective and efficient computation
371 with multiple-timescale spiking recurrent neural networks”. In: *International Conference on*
372 *Neuromorphic Systems 2020*. 2020, pp. 1–8.

373 Checklist

- 374 1. For all authors...
- 375 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-
376 tions and scope? [Yes]
- 377 (b) Did you describe the limitations of your work? [Yes]
- 378 (c) Did you discuss any potential negative societal impacts of your work? [Yes] The potential
379 impact of this work extends from a broader discussion on the potential impact of low-power
380 keyword spotting technologies.
- 381 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 382 2. If you are including theoretical results...
- 383 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 384 (b) Did you include complete proofs of all theoretical results? [N/A]
- 385 3. If you ran experiments...
- 386 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
387 results (either in the supplemental material or as a URL)? [Yes]
- 388 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
389 [Yes]
- 390 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
391 multiple times)? [Yes]
- 392 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,
393 internal cluster, or cloud provider)? [Yes]
- 394 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 395 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 396 (b) Did you mention the license of the assets? [Yes]
- 397 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 398 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-
399 ing/curating? [N/A]
- 400 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-
401 tion or offensive content? [N/A]
- 402 5. If you used crowdsourcing or conducted research with human subjects...
- 403 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
404 [N/A]
- 405 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
406 approvals, if applicable? [N/A]
- 407 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
408 participant compensation? [N/A]