

Training neuromorphic neural networks for speech recognition and Bayesian optimisation of network hyperparameters

Rihards Klotiņš

Supervisor: Dorian Florescu

May 4, 2025

Abstract

Filler abstract - Spiking Neural Networks (SNNs), the third generation of neural networks, offer a promising alternative to traditional Artificial Neural Networks (ANNs) due to their energy efficiency, temporal processing capabilities, and potential for neuromorphic hardware implementation. Unlike ANNs, which rely on continuous-valued activations, SNNs process information through discrete spike events, closely mimicking biological neural systems. This characteristic enables them to efficiently handle sequential data, making them suitable for applications such as speech recognition, event-based vision, and edge computing. Training SNNs, however, remains a significant challenge due to the non-differentiability of spike events. While ANN-to-SNN conversion provides a workaround, it imposes computational costs and limits architectural flexibility. Direct training methods, such as Backpropagation-Through-Time (BPTT) with surrogate gradients, local learning rules, and biologically inspired approaches like e-prop and EventProp, have emerged as viable alternatives. Each method presents trade-offs in terms of computational complexity, biological plausibility, and performance. Notably, EventProp has demonstrated state-of-the-art results while reducing memory and computational overhead compared to BPTT. This work explores the theoretical advantages of SNNs, their energy-efficient processing, and recent advancements in training methodologies. It also highlights their potential impact on low-power computing applications, particularly in audio processing, where temporal encoding is crucial. By addressing current limitations and leveraging novel training strategies, SNNs could play a pivotal role in next-generation AI systems.

1 Introduction

Artificial Neural Networks (ANN) are computational models inspired by the brain; made up of layers of interconnected “neurons” which can learn patterns when given large amounts of data. After learning patterns, ANNs can be used to make inferences, predicting an output based on some given input. ANNs have long been used for application ranging from computer vision to financial forecasting. However the use of ANNs has surged recently, increasing fourfold since 2017, driven by the development of generative AI models like ChatGPT [1]. To compound this, models are growing in size and becoming more complex, consequently the global AI power consumption is increasing at an exponential rate [2]. Not only does this make models expensive to operate, but it also raises sustainability concerns.

A graphics processing unit (GPU) running a large language model (LLM) consumes hundreds of watts of power. On the other hand, the human brain processes sensory information, keeps involuntary biological systems functioning, runs its own language model, and enables conscious thought — all while using only about 20 watts. Such efficient information processing has motivated research into neuromorphic computing - a field aiming to emulate the brain’s neural architecture. Spiking Neural Networks (SNNs) are considered the third generation of neural network models [3], where ANNs are considered the second generation. Compared to neurons in ANNs, SNN neurons more closely model how biological neurons behave in the brain. For a second generation neuron, the output is a weighted sum of the inputs passed through an activation function (see Figure 6 (a)). For a third generation neuron, the output is a

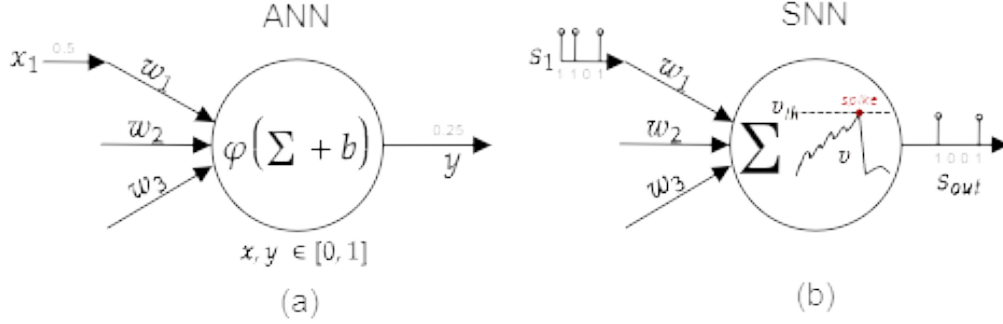


Figure 1: (a) 2nd generation neuron. (b) 3rd generation neuron.

spike, sometimes referred to as an event, which occurs when the voltage of a neuron surpasses a threshold (Figure 6 (b)), this results in neurons communicating with each other via series of spikes, which vary in timing and frequency. Since information is encoded in the timing of outputs, SNNs can inherently capture temporal patterns and dynamic behaviours in sequential data, making them inherently capable at processing tasks involving time-dependent signals, such as speech recognition, sensory processing, and event-based data streams. Moreover, it has been shown that SNNs theoretically possess higher computational power than the previous generation of ANNs [4]. Spiking neurons only activate when necessary, and don't require a clock signal, enabling massive power savings. For instance, the leading neuromorphic platform TrueNorth is capable of simulating a million spiking neurons in real-time while consuming 63 mW. The equivalent network executed on a high-performance computing platform was 100–200× slower than real-time and consumed 100,000 to 300,000× more energy per synaptic event [5]. The power savings of SNNs could be game-changing in edge computing devices that run on battery or have limited power budgets - like smartphones or remote environmental sensors. On top of this, their low-latency could benefit autonomous vehicles and robotics.

1.1 Problem to Solve

Spiking neural networks are inherently time dependent since information is encoded in the timing of spikes. This makes them naturally capable of being applied to temporal tasks - tasks where the data evolves with time. The problem with using ANNs for temporal data is that they require a fixed input size. Though they can be tweaked to deal with unspecified length temporal data using a paradigm called “recurrence”, this comes with its downsides, such as expensive training and large memory use [citation needed]. Temporal data can come in many forms, for instance video, audio, or even stock market information. The efficiency and potential effectiveness of SNNs to process temporal information could be game changing in edge devices such as mobile phones, wearable devices, and remote sensors which have small power budgets. Large language models are revolutionising how we interact with computers, they provide a way humans to interface with machines using natural language. As a result, companies are eagerly integrating LLMs into their products and services - e.g. Siri, Raybans [citation needed]. Accurate speech recognition is therefore posing to be a highly relevant application of SNNs, playing into their strengths of temporal processing and energy efficiency.

1.2 Spiking Heidelberg Digits Dataset

The Spiking Heidelberg Digits (SHD) dataset [citation needed] is a prominent spiking neural network benchmark. The wide use of SHD makes it good for fairly comparing different methods of training models. It is based on the Heidelberg Digits dataset, which is a collection of 10,000 high-quality recordings of spoken digits (0 to 9) in English and German. It is spoken by a relatively representative group of 6 males and 6 females of the ages 21 to 56 years old, with a mean age of 29. The HD dataset was converted into spike trains using a biologically realistic model of the human cochlea, outputting 700 channels of spikes representing different frequency bands picked up by the ear.

[How are spikes encoded. Temporal or rate?] [Compare against other speech data sets.]

2 Spiking Neuron Model

In computational neuroscience, several neuron models have been developed to simulate how neurons behave, each offering a different trade-off between biological detail and computational efficiency. This section introduces three widely used models: the Hodgkin–Huxley model, the Leaky Integrate-and-Fire (LIF) model, and the Izhikevich model.

We will begin with the most influential model in neuroscience, the Hodgkin–Huxley model. Introduced in 1952 and awarded the Nobel Prize in Physiology or Medicine in 1963, this model provides a detailed mathematical description of how electrical signals—action potentials—are generated and propagated in neurons. While this model is highly accurate and biologically grounded, its complexity makes it computationally demanding, which limits its use in large-scale or real-time simulations.

Next we will discuss the Leaky Integrate-and-Fire model. To improve efficiency, the Leaky Integrate-and-Fire model simplifies neuron behaviour while retaining essential features. It treats the neuron like an electrical circuit that accumulates incoming signals over time. When the membrane voltage crosses a set threshold, the neuron emits a spike and resets. Though simplified, this model captures key spiking behavior and is highly efficient to compute, making it the most commonly used neuron model in neuromorphic engineering and machine learning applications.

Finally, we will discuss the Izhikevich model which offers a compromise between realism and efficiency. Using just two simple differential equations, it can reproduce a wide range of spiking and bursting patterns seen in real neurons.

2.1 Hodgkin–Huxley

The Hodgkin–Huxley is a detailed neuron model developed in 1952 by Alan Hodgkin and Andrew Huxley, was the first to quantitatively explain how neurons generate and propagate electrical impulses by describing the flow of sodium and potassium ions through voltage-gated channels. Its success in reproducing the full waveform of electrical impulses — including rapid depolarization, repolarization, and refractory behaviour — earned Hodgkin and Huxley the 1963 Nobel Prize in Physiology or Medicine. Due to the robustness and biophysical accuracy of this model it is a logical model to consider when looking to implement biologically inspired learning computationally.

A neuron is enclosed in a cell membrane, with gates for ion flow. We will consider positively charged potassium (K^+) and sodium (Na^+) ions as they are considered as the primary contributors to the electrical behaviour of neurons. The neuron membrane has channels which allow ions to flow between the inside and the outside of the neuron; these channels are selective, allowing only specific ions to pass through, e.g. only allowing Na^+ to flow. The amount of flow they allow depends on the membrane voltage V_m , which is the measure of voltage difference between the inside and outside of the neuron. Two factors govern the tendency of ions to flow from one area to the other - i.e. from inside to outside - diffusion and electric charge. Diffusion is the tendency of ions to flow from areas of high concentrations to low concentrations; in figure 2 Na^+ ions will tend to flow to the extracellular side of the membrane as the concentration of Na^+ ions is lower there. Particles which have the same electric charge repel, particles with opposite electrical charges attract. Therefore, Na^+ ions will be attracted towards the intracellular side of the membrane due to the inside of the neuron having a negative electric charge. The directions of these

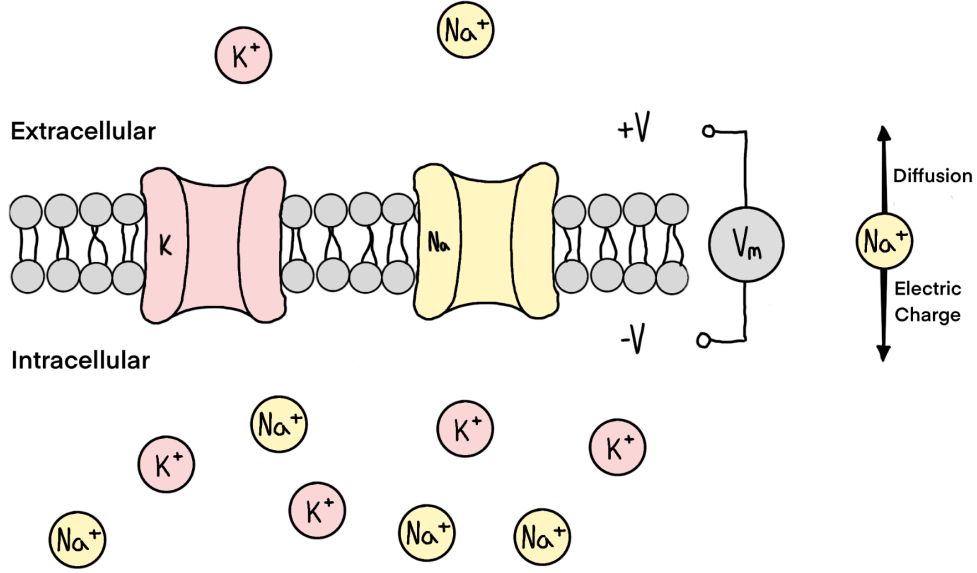


Figure 2: Diagram of neuron membrane.

opposing forces are displayed in figure 2. Naturally the cell reaches an equilibrium where the diffusion and electric charge forces are equal. When an external current is applied to the system, the equilibrium is lost and the system experiences transient behaviour defined by the mathematics of the Hodgkin-Huxley model.

The cell membrane acts as a dielectric separating two charged mediums, in other words it acts as a capacitor. According to the capacitor current-voltage equation, we know that $C \frac{dV}{dt} = I$. Applying this to the context of the neuron membrane, we get equation (8). Where V_m is the potential difference across the cell membrane, C_m is a constant representing the characteristic capacitance of the cell membrane, I_{Na} and I_K are the net charge flows of sodium and potassium ions respectively, I_L is the leakage current caused by the movement of other ions, and I_{ext} is the external current applied to the neuron. So we see that the sum of currents causes the voltage to change.

$$C_m \frac{dV_m}{dt} = -(I_{Na} + I_K + I_L - I_{ext}) \quad (1)$$

The currents - I_{Na} , I_K , and I_L - are determined by equations (2), (3) and (4). They can be examined like the classic $I = V/R = Vg$ equation (g being conductance). Where the left hand side is the current, the part in the brackets is the voltage, and the part outside the brackets is the conductance of the gates which the ions flow through.

$$I_{Na} = \bar{g}_{Na} m^3 h (V_m - E_{Na}) \quad (2)$$

$$I_K = \bar{g}_K n^4 (V_m - E_K) \quad (3)$$

$$I_L = \bar{g}_L (V_m - E_L) \quad (4)$$

As discussed before, when there is an imbalance in the force of *electric attraction* and *diffusion*, there is a *net flow of ions*, i.e. a current. This represented in equation (2) by the term $(V_m - E_{Na})$. Where E_{Na} is the equilibrium

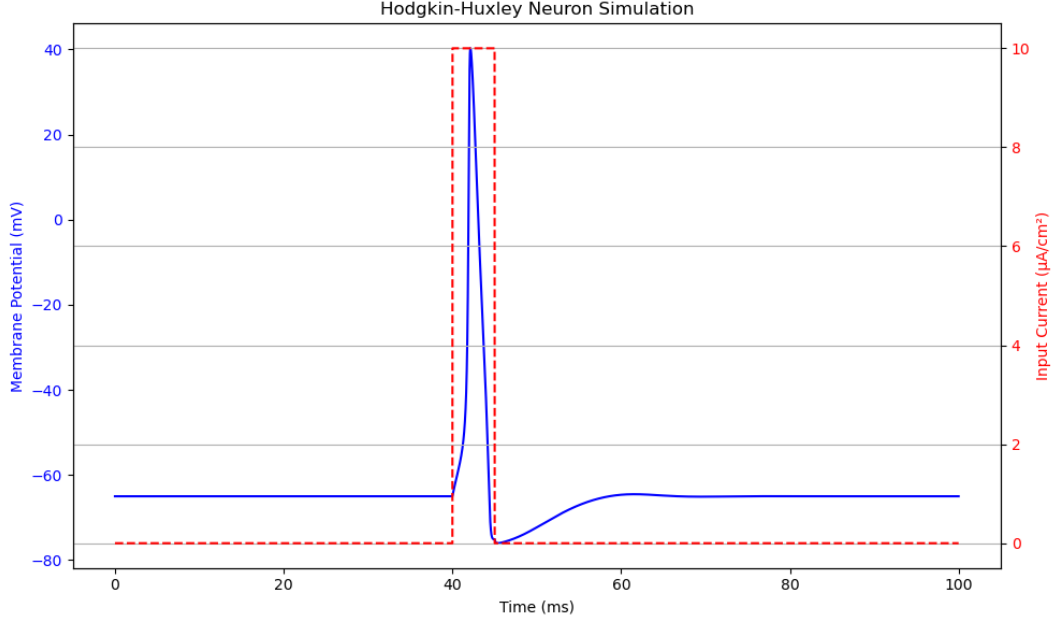


Figure 3: Response of the HH model (blue) to a 5ms step function (red).

potential due to diffusion forces. When the two variables are in balance their difference is 0, thus the current is also zero. The same also applies for equations (3) and (4).

As per equation (2) the current also depends on the *conductance* of the channels which the ions move through. There are specialised channels for Na⁺ ions and for K⁺ ions, which have different conductances. The conductance of sodium channels is given by $\bar{g}_{Na}m^3h$. The \bar{g}_{Na} term is the conductance of the channel when it is fully open, i.e. the channels maximum conductance. The m^3h term is the probability - 0 to 1 - that a sodium channel is open. This probability comes from the fact that each channel has 3 “m” gates and 1 “h” gate, ions can only flow through the channel when all gates are open. Due to the large number of channels in a neuron, the value of the probability will correspond to what proportion of the channels are open. The probabilities of m , h , n are given by equation (5), where α_x and β_x are rate constants at which speed the gates open and close respectively, they are voltage dependent.

$$\frac{dx}{dt} = \alpha_x(V_m)(1 - x) - \beta_x(V_m)(x), \text{ where } x \in \{m, h, n\} \quad (5)$$

Such simple first order ion flow equations characterise to high accuracy how the neurons in our brain function. Figures 3, 4, and 5 show how a neuron would react to a 5ms, 15ms, and 20ms pulse of current. While the HH neuron model has propelled our understanding of neurons due to it’s biophysical precision, it’s complexity - with its nonlinear differential equations and multiple gating variables - makes it computationally expensive to simulate. As a result researchers had to look elsewhere to find a model that could be used for training machines to think like humans.

2.2 Leaky Integrate-and-Fire Neuron

The leaky integrate-and-fire (LIF) neuron model is the most widely used neuronal model in SNNs due to its simplicity and efficiency which helps it scale for large networks. Simply put, the neuron receives spikes from “pre-synaptic” neurons (figure 6 (b)) - neurons which feed into the neuron in question; these spikes increase the “membrane potential” of the neuron - which is the voltage between the inside and outside of the neuron. When the membrane potential reaches a threshold voltage, the neuron in question fires a spike to the neurons that it is connected to, and it’s membrane voltage resets to a baseline value.

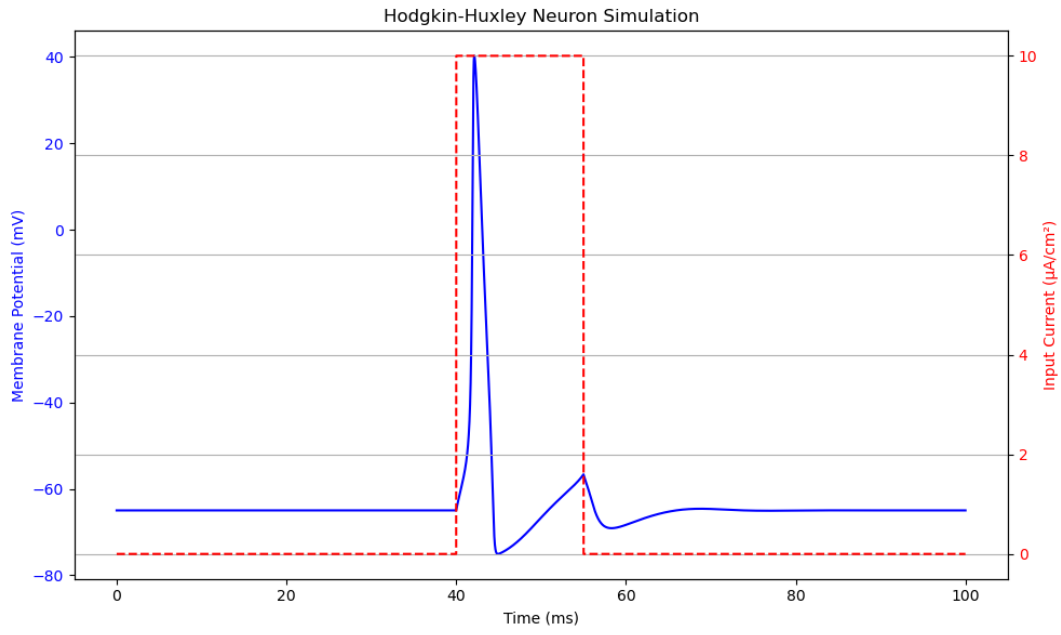


Figure 4: Response of the HH model (blue) to a 10ms step function (red).

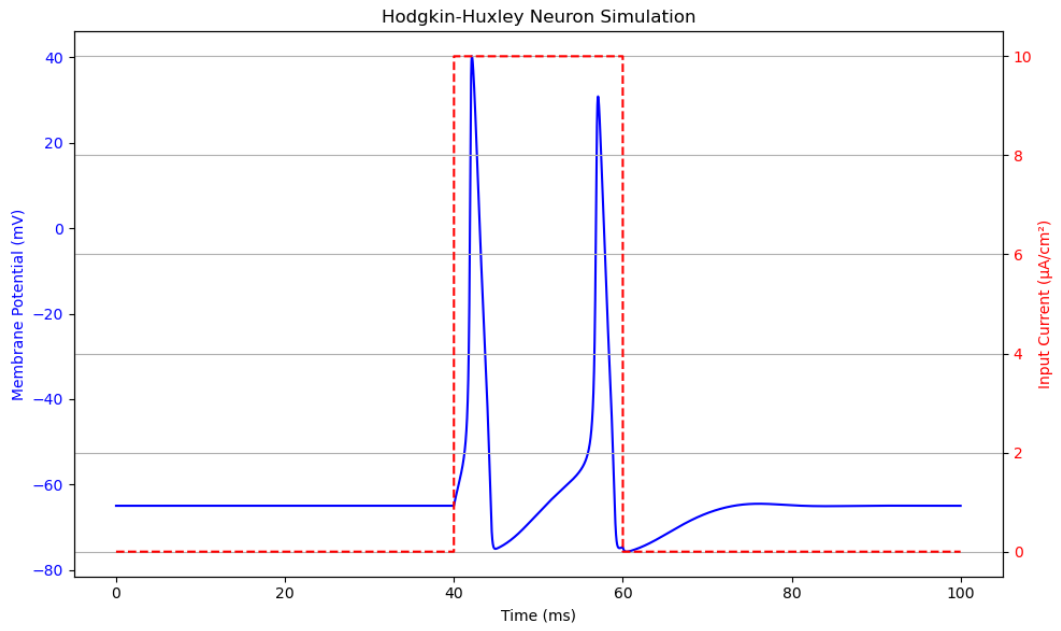


Figure 5: Response of the HH model (blue) to a 15ms step function (red).

LIF Neuron

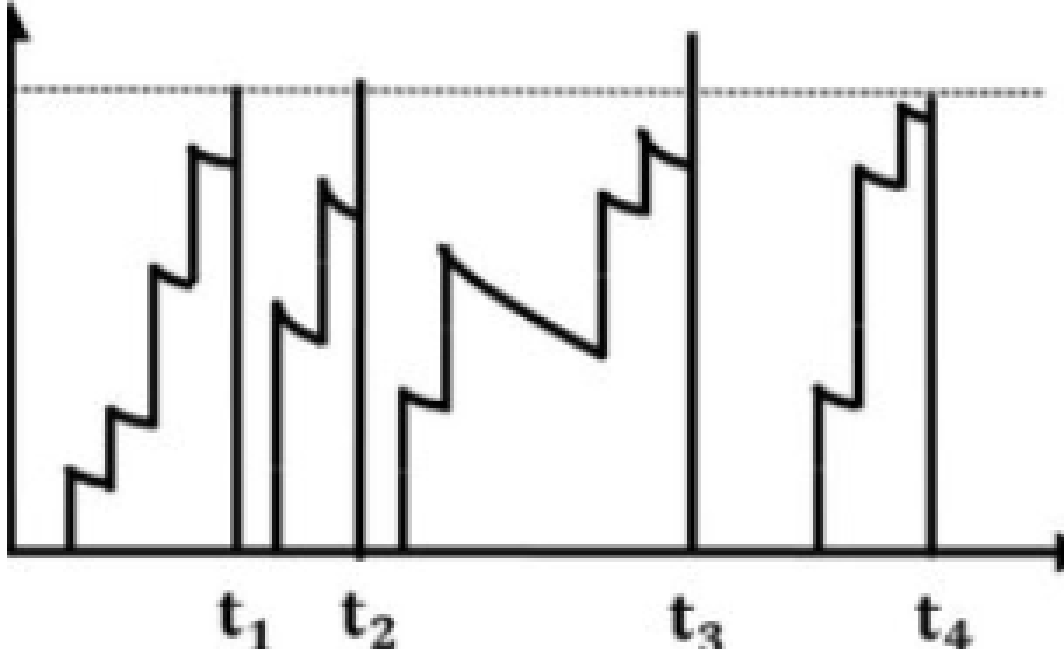


Figure 6: LIF neuron plot of V_{mem} .

In order to be modelled computationally, this behaviour is expressed in mathematical terms. Spikes received by a neuron at time t induce a current $X[t]$. This input current increases the membrane potential, $V[t]$, which is the voltage between the inside of the neuron and the outside. The voltage slowly decays over time, the speed of the decay depends on the membrane time constant, τ . This is the “Leaky Integrate” aspect of the neuron; $X[t]$ is integrated to give $V[t]$ at the same time as $V[t]$ leaks voltage. The neuron leaks voltage till it reaches its baseline voltage level, V_{reset} . This behaviour can be seen in equation (6), $H[t]$ is equal equation (6) to $V[t]$, unless a spike occurs during t . This spiking logic is defined in equation (7), where $\Theta(x)$ is a function that is 0 unless $x \geq 0$. In other words this function becomes 1 only when the threshold voltage V_{th} is reached. When the threshold voltage is reached, $V[t]$ is set to V_{reset} and a spike is released; if the threshold voltage is not reached then $V[t] = H[t]$. The membrane potential behaviour at spike time is defined in equation (8).

$$H[t] = V[t - 1] + \frac{1}{\tau}(X[t] - (V[t - 1] - V_{reset})), \quad (6)$$

$$S[t] = \Theta(H[t] - V_{th}), \quad (7)$$

$$V[t] = H[t](1 - S[t]) + V_{reset}S[t] \quad (8)$$

2.3 Izhikevich

Building on the simplicity of the leaky integrate-and-fire (LIF) neuron, which uses a single differential equation and a fixed threshold to generate spikes, the Izhikevich model strikes a balance between biological realism and computational efficiency [6]. Eugene Izhikevich observed that Hodgkin–Huxley-type models, while highly detailed, require integrating four stiff equations per neuron, making large-scale simulation infeasible, and that simple integrate-and-fire models cannot reproduce many cortical firing [6]. To address this, he derived a minimal two-variable model that can emulate diverse neuron behaviors—tonic spiking, bursting, chattering, and more—by tuning just four parameters [7].

Mathematically, the model comprises these coupled ordinary differential equations:

$$\frac{dv}{dt} = 0.04v^2 + 5v + 140u + I(t), \quad (9)$$

$$\frac{du}{dt} = a(bvu), \quad (10)$$

where v is the membrane potential (in mV), u is a membrane recovery variable accounting for K^+ activation and Na^+ inactivation, and $I(t)$ is an external input current [6]. When v reaches the peak (typically 30 mV), it and u are reset:

$$\text{if } v \geq 30 \text{ mV, } \begin{cases} v \leftarrow c \\ u \leftarrow u + d \end{cases} \quad (11)$$

with (a,b,c,d) chosen to match specific neuron types. Despite its simplicity, this formulation reproduces over twenty cortical firing patterns by varying these four parameters, far surpassing the expressive capacity of the LIF model [7].

Critically, Izhikevich demonstrated that this model runs in real time for tens of thousands of neurons on a single CPU core - comparable to LIF performance -while maintaining Hodgkin–Huxley-level richness in spike dynamics. Thus, by augmenting the LIF framework with a single recovery variable and non-linear voltage dependence, the Izhikevich model provides an efficient yet versatile platform for large-scale SNN training and simulation.

2.4 Neural Networks

Alone, a neuron doesn’t do anything very impressive; in a network, intelligent behaviour can emerge. When many neurons are connected together, forming a neural network, they can collectively process information, detect patterns, and make decisions. Each neuron receives inputs from other neurons through connections called synapses (figure 7). These connections have strengths, or “weights”, which determine how much influence one neuron has on another; these strengths are visualised by the colour of the connections in figure 7. By carefully adjusting these weights, the network can learn to perform complex tasks such as image recognition, speech processing, or decision-making. This idea underpins both biological neural circuits and artificial neural networks used in machine learning.

3 Training Neural Networks

The objective of training a neural network for spoken word recognition is to construct a model that accurately maps audio input to the correct lexical output. Supervised learning remains the dominant training paradigm in this domain, consistently achieving state-of-the-art results in speech recognition tasks [8, 9]. Supervised learning

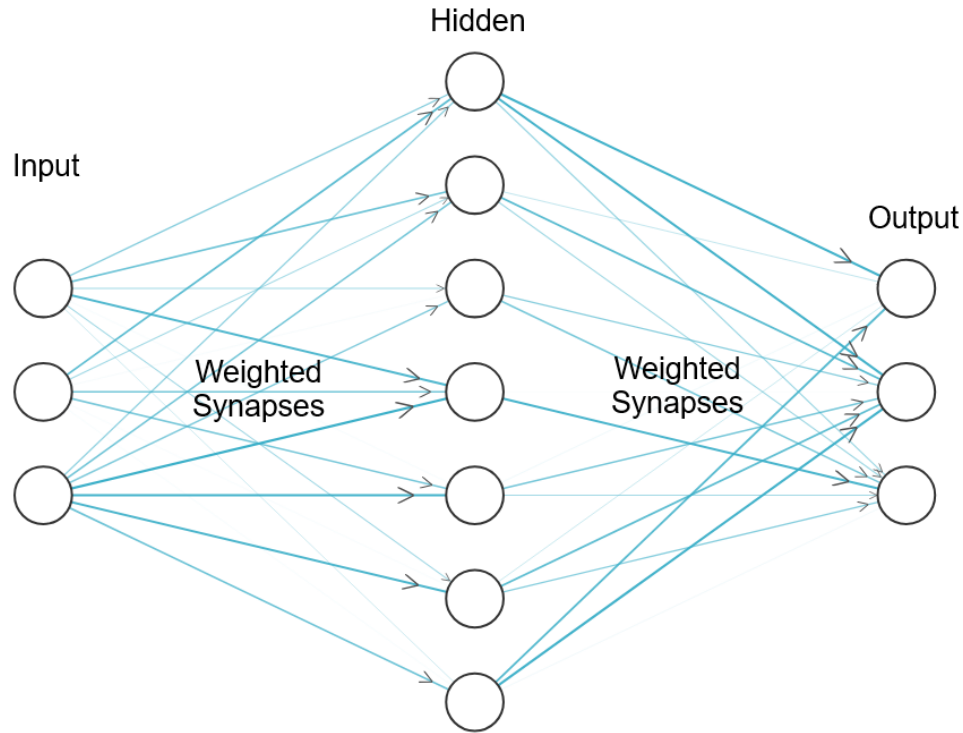


Figure 7: Diagram of a 3 layer neural network. The weights of the synapses are visualised by the colours of the connections.

relies on labelled datasets—typically composed of audio clips annotated with their corresponding transcriptions—providing explicit guidance for the model to learn the mapping between acoustic features and linguistic units. In contrast, unsupervised learning operates on unlabelled data, requiring the model to uncover inherent structure or patterns within the input without external annotation. Next I will discuss different ways that neural networks can be trained.

3.1 Backpropagation

Backpropagation - an efficient implementation of gradient descent - is the most effective and widely used method for training second generation artificial neural networks. It provides a systematic way to adjust the internal parameters - or weights - of a network to minimize the discrepancy between the network's predictions and the desired outputs. An overview of the backpropagation process is as follows:

1. Forward Pass: The model receives an input and processes it through its layers, generating an output.
2. Error Calculation: The output is compared to the desired target, and a loss function quantifies the error.
3. Gradient Computation: The derivative of the loss function is computed with respect to each weight in the network using the chain rule of calculus.
4. Weight Update: The weights are adjusted by subtracting a fraction of their corresponding gradient, typically scaled by a learning rate. This step moves the model toward a configuration that reduces error.

In supervised learning, the desired output is the label of the data. When a model predicts an output, the “loss” is calculated to quantify the error between the actual output and the desired output. For instance, a common way to calculate the loss is by calculating the mean squared error (MSE), where you find the sum of the squared differences between the actual output neuron values and the desired output neuron values (Equation (12)).

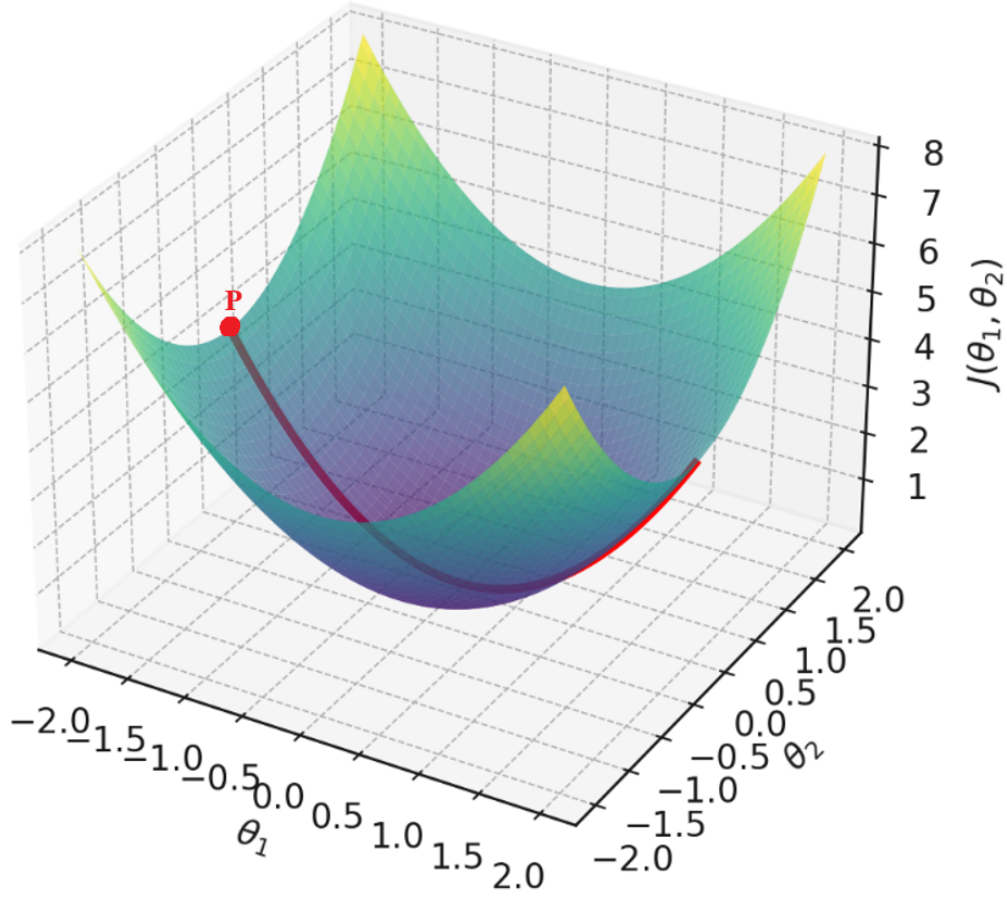


Figure 8: (a) ANN neuron (b) SNN neuron

$$L = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (12)$$

To visualise how the loss of a network is minimised, let us assume that a network has 2 weights, θ_1 and θ_2 . Figure 8 shows the loss plotted against the weights θ_1 and θ_2 . The optimal weight configuration occurs at the lowest point on the surface. If the weights of the neural network get initialised to the point P , we can find the gradient of the surface at that point which will indicate the direction of the steepest ascent. Going along the surface in the opposite direction of the gradient will follow the fastest path down the surface, minimising the loss function until the local minimum is reached where the gradient becomes 0. This process is known as *gradient descent*, backpropagation is an efficient implementation of gradient descent where the gradient of the network is calculated backwards layer by layer.

To demonstrate how backpropagation works, let's consider a simplified, fully-connected network consisting of a layer of input neurons, hidden neurons, and output neurons (Figure 9). Fully-connected means each neuron in a layer is connected to all neurons in the following layer.

The highlighted neuron, h_1 , is connected to all input neurons, x_1 , x_2 , and x_3 . Therefore h_1 is a weighted sum of all the input neurons (Equation 13).

$$h_1 = \sum_{i=1}^3 x_i \times W_i^1 \quad (13)$$

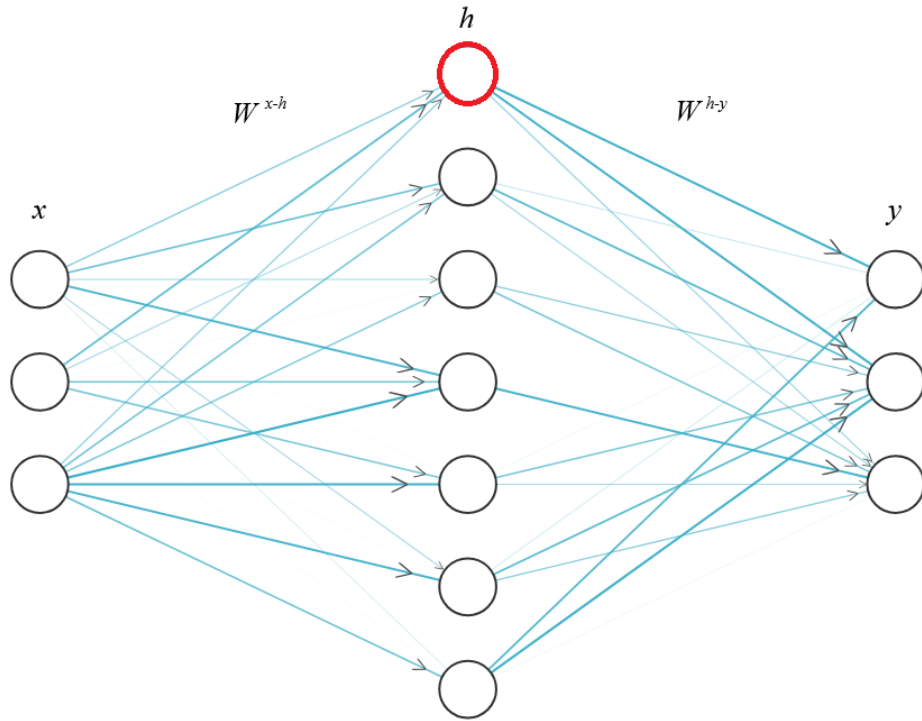


Figure 9: x is the input vector of length 3. h is the hidden layer of length 6. W^{x-h} is the matrix containing the values of each of the weights connecting x and h , thus it is of size 3 by 6. y is the output vector of length 3. W^{h-y} is the matrix containing the values of each of the weights connecting h and y . thus it is of size 6 by 3.

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \times \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} & w_{1,5} & w_{1,6} \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} & w_{2,5} & w_{2,6} \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & w_{3,5} & w_{3,6} \end{bmatrix} = \begin{bmatrix} h_1 & h_2 & h_3 & h_4 & h_5 & h_6 \end{bmatrix}$$

Figure 10: Matrix multiplication visualisation.

An efficient way to write these weighted summations for all of the neurons in a layer is by using matrix algebra.

$$\begin{aligned} h &= x \times W^{x-h} \\ y &= h \times W^{h-y} \end{aligned}$$

These two equations encompass all of the weighted summations that define the network. The relationship between the weighted summation equation (equation 13) and the matrix multiplication representation of the network is highlighted in figure 10. You can see that each input neuron x_i is multiplied by the corresponding weight $w_{i,1}$ and summed to give h_1 .

Matrix multiplications perform the *forward pass* of the network - i.e. the inference. This inferred output value, y , is compared with the true “label” value, y_{label} . The comparison between the actual output and the desired output is done by a loss function, $l()$. There are many ways of implementing the loss function, one way by finding the mean squared error (MSE):

$$L = l(y, y_{label}) \quad (14)$$

The goal is to minimize this loss by updating the network’s weights in the direction of the negative gradient. This requires computing the gradient of the loss with respect to each weight w_i :

$$\frac{\partial L}{\partial w_i} \quad (15)$$

We first calculate the gradient of the layer closest to the output, W^{h-y} .

$$\frac{\partial L}{\partial W^{h-y}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial W^{h-y}} = \frac{\partial L}{\partial y} \cdot h \quad (16)$$

We then use the chain rule again to calculate the gradient of L with respect to W^{x-h} .

$$\frac{\partial L}{\partial W^{x-h}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h} \cdot \frac{\partial h}{\partial W^{x-h}} = \frac{\partial L}{\partial y} \cdot W^{h-y} \cdot x \quad (17)$$

This step-by-step application of the chain rule allows the error to be propagated backward through the network, enabling efficient computation of gradients at each layer - this is the essence of backpropagation.

Once the gradients are known, the weights are updated using the gradient descent rule:

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i} \quad (18)$$

where η is the learning rate. A higher η would mean faster learning but risks instability. A lower η would ensure stable convergence to the minima, but may make the training very slow.

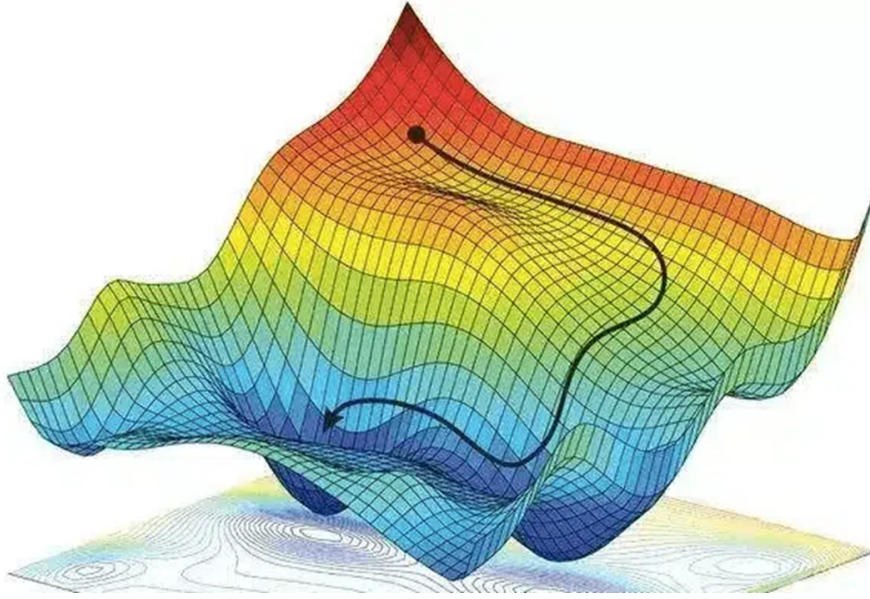


Figure 11: Loss function 3D surface being descended.

Repeating this backpropagation process across many training examples allows the network to gradually learn the desired input-output mapping by descending the loss function surface. A problem with this approach to be aware of is that the gradient descent algorithm may find a local minimum of the loss function instead of the global minimum. As a result the performance of the model may be significantly lower than it potentially could be. A popular way of dealing with this problem is by training the network several times with different random weight initialisations; if one network initialisation gets stuck in a local minima, others may not. If there are too many local minima, this approach may not be enough.

3.2 Spiking Neural Network Training

While backpropagation has proven to be a powerful algorithm for training artificial neural networks, it cannot be directly applied to spiking neural networks (SNNs). This is primarily due to the non-differentiable nature of spike events, which prevents the straightforward calculation of gradients required for weight updates. Additionally, SNNs operate in continuous time and rely on event-based communication, introducing complex temporal dynamics that further complicate gradient-based optimization. As a result, researchers have developed a range of alternative training algorithms tailored to the unique characteristics of SNNs. In this section, we introduce several of these methods, highlighting how they address the challenges of spike-based computation and enable effective learning in spiking systems.

3.2.1 ANN-to-SNN Conversion

Due to the popularity of ANNs, literature on training them is advanced, so a natural and popular choice for training SNNs has been by converting a trained ANN model into an SNN model. Usually a trained artificial neural network is transformed into a spiking neural network by substituting ReLU activation functions with spiking neuron models. This process often involves additional adjustments such as weight normalization and threshold balancing to maintain performance and stability. These ways of training have had good results for some tests [10, 11]. However, such a method incurs large computational costs during conversion and is limited by the architecture of ANNs which are less adaptable to dynamic data like audio [12]. Thus, to fully harness the benefits of SNNs — from energy efficiency to novel architectures — effective direct training methods are essential.

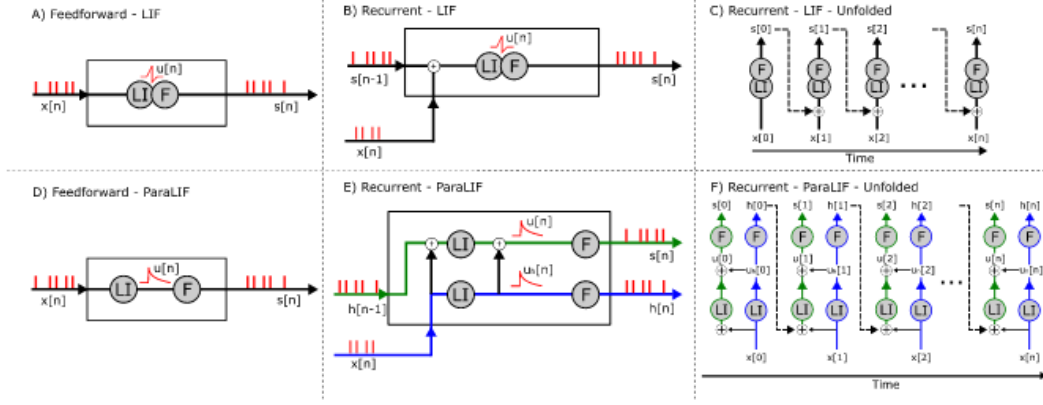


Figure 12: Matrix multiplication visualisation.

3.2.2 Backpropagation-Through-Time

Similar to the ANN-to-SNN method, Backpropagation-Through-Time (BPTT) attempts to carry over the effectiveness of the backpropagation algorithm to SNNs. However, this time the SNNs get trained directly which allows them to learn dynamic patterns from temporal data such as speech. Familiar gradient-based methods are applied in tandem with “surrogate gradients” (SG) which are spike gradient approximations used to overcome the fact that spikes are non-differentiable. You can think of the network’s activity over time as a very deep chain of simple processing steps; BPTT “unrolls” this chain so that the error at the end can be traced back step by step to adjust every connection [13]. Because a spike is a discontinuous event (it either happens or it doesn’t), we replace its true derivative—which is zero almost everywhere and jumps to infinity at spike times - with a smooth “surrogate” function during training. This surrogate lets us compute approximate gradients so that standard optimisers like gradient descent can still work [12].

When applied to speech-recognition benchmarks—such as the Spiking Speech Commands (SSC) and Spiking Heidelberg Digits (SHD) datasets—this method achieves accuracy on par with conventional neural networks while operating in a sparse, event-driven fashion that can be more energy-efficient at inference time [11, 14]. Researchers have even swapped out recurrent layers in end-to-end speech models for surrogate gradient trained spiking modules, showing only small drops in word-error rate and offering a path toward low-power, real-time processing [11].

However, BPTT comes with two major downsides. First, it requires storing every intermediate state over the entire duration of an input—meaning memory usage grows with the length of the audio clip, which can quickly exceed hardware limits for long recordings [14]. Second, because the learning rule relies on a global error signal propagated across many time steps and layers, it differs starkly from the local, synapse-by-synapse learning observed in biological brains—undermining some of the potential efficiency gains of neuromorphic hardware.

Parallelisable LIF The major bottleneck in training spiking neural networks (SNNs) using BPTT is the strictly sequential nature of classic Leaky Integrate-and-Fire (LIF) neurons, which update their membrane potential step by step in time. The Parallelizable LIF (ParaLIF) model overcomes this by decoupling the linear integration of inputs from the spiking (thresholding) operation and executing both across all time steps in parallel. This reorganization leverages highly optimized matrix operations on modern accelerators to deliver dramatic speed-ups in training, without altering the fundamental membrane-and-spike dynamics that give SNNs their event-driven efficiency [15, 16].

In a standard LIF neuron, the membrane potential $V(t)$ at time t depends on its previous value $V(t1)$ plus any new inputs, and a spike is emitted once V crosses a threshold. ParaLIF rewrites this process as two separate GPU kernels. The first kernel computes, for every neuron, the entire sequence of membrane-potential updates in one batched matrix multiplication; the second applies the threshold-and-reset rule simultaneously at all time points. By removing the need for “time-step loops,” ParaLIF converts a fundamentally serial simulation into a fully vectorized parallel computation [15].

When benchmarked on neuromorphic speech (Spiking Heidelberg Digits), image and gesture datasets, ParaLIF achieves up to $200\times$ faster training than conventional LIF models, while matching or exceeding their accuracy with comparable levels of sparsity [15, 16]. Compared to other parallel schemes—such as the Stochastic Parallelizable Spiking Neuron (SPSN) approach - ParaLIF maintains similar speed-ups on short sequences and far greater scalability on very long inputs [16].

While parallelising neurons has achieved impressive results on the SHD benchmark this method comes with notable disadvantages. By reorganizing the time dimension, ParaLIF departs from the continuous, step-by-step integration that real neurons exhibit, reducing its biological plausibility [17]. This parallel update can also undermine the network’s ability to capture fine temporal dependencies, since precise spike timing and sequential context are approximated rather than explicitly modelled [18, 19]. Moreover, the specialized GPU kernels and data-layout transformations needed for ParaLIF introduce implementation complexity and may not map efficiently to more constrained neuromorphic hardware, limiting its applicability in low-power edge scenarios, a key potential application of SNNs.

3.2.3 Eligibility Propagation

Eligibility propagation, or e-prop, is a method for training spiking neural networks (SNNs) that aligns more closely with how learning is believed to occur in the brain while still taking inspiration from the clearly effective backpropagation algorithm. Unlike traditional training methods like backpropagation-through-time (BPTT), which require storing the entire history of neuron activities and propagating errors backward through time, e-prop simplifies this process by using two key components: eligibility traces and a learning signal. Eligibility traces act like short-term memories at each synapse, recording recent activity patterns. They capture how the timing of spikes affects the potential for learning. The learning signal is a global factor that represents the overall error or feedback from the network’s output. Instead of sending detailed error information back through every layer and time step, as in BPTT, e-prop uses this single signal to modulate the eligibility traces. When the network makes a mistake, the learning signal adjusts the synapses with high eligibility traces, effectively correcting the connections that contributed most to the error. By updating synaptic weights immediately based on recent pre- and post-synaptic activity, e-prop reduces memory requirements compared to BPTT [20] and can dramatically lower energy consumption on event-driven hardware [21, 22]. However, because it uses approximate gradients, e-prop-trained models typically exhibit lower accuracy than fully BPTT-trained networks, reflecting a trade-off between biological plausibility and performance.

In its original demonstration, Bellec et al. applied e-prop to train spiking recurrent networks on the TIMIT speech corpus, showing that eligibility traces derived from slow neuronal dynamics could capture phonetic temporal dependencies without backward passes [20]. Subsequent work has enriched e-prop with spike-timing-dependent plasticity (STDP)-like eligibility decay and local random broadcast alignment to improve phoneme classification accuracy. Van der Veen demonstrated that modulating eligibility traces according to precise spike timing and using randomized local error broadcasts allowed spiking networks to approach conventional LSTM performance on phonetic labels, all while preserving the sparse activity characteristic of SNNs [23].

E-prop has been implemented on neuromorphic hardware for keyword spotting. On the SpiNNaker 2 system, Frenkel and Indiveri trained spiking recurrent networks on the Google Speech Commands dataset, achieving over 91% accuracy with only 680KB of training memory—over $12\times$ lower energy consumption than GPU-based BPTT solutions [24].

Despite its advantages, e-prop also has notable drawbacks. First, it requires maintaining multiple eligibility traces per synapse (e.g., for membrane potential and adaptive threshold), as well as optimizer state such as moment vectors, resulting in significant memory overhead for large networks [22?]. Second, because it employs approximate surrogate gradients rather than true backpropagation, e-prop-trained models typically achieve lower accuracy than their BPTT-trained counterparts [20]. Third, although e-prop avoids backward error propagation through time, it still depends on a global learning signal to modulate local eligibility traces, introducing communication overhead and deviating from strictly local synaptic updates—factors that can limit its energy efficiency on distributed neuromorphic hardware [21].

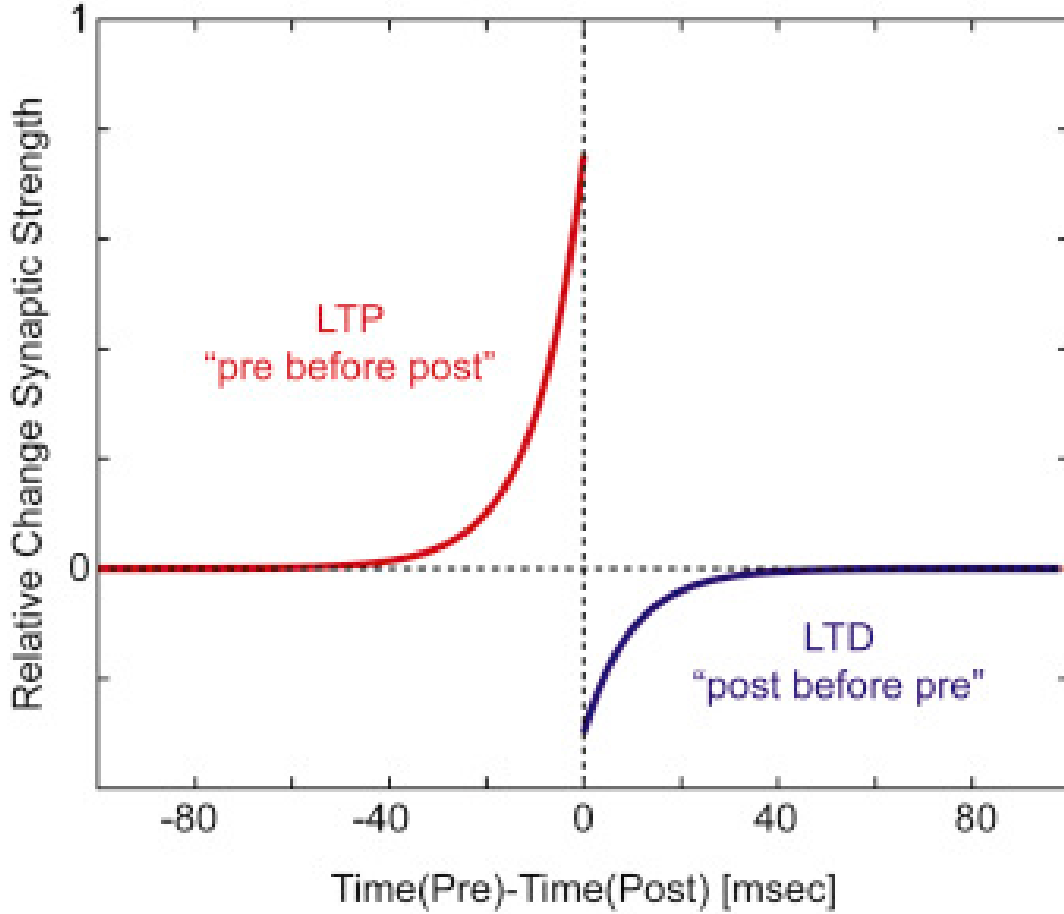


Figure 13: LTP and LTD

3.2.4 Spike-Time-Dependent-Plasticity

Spike-timing-dependent plasticity (STDP) represents a training approach for SNNs that draws even greater inspiration from the learning rules observed in biological neurons. STDP is a biological learning rule that adjusts the strength of synaptic connections based on the precise timing of pre-synaptic and post-synaptic spikes. If a pre-synaptic neuron fires just before a post-synaptic neuron, the connection between them is strengthened. Conversely, if the order is reversed, the connection is weakened. This temporally sensitive form of Hebbian learning, often summarized by the maxim “cells that fire together, wire together,” operates locally at each synapse and does not require global error signals. This makes it inherently biologically plausible, asynchronous, and capable of unsupervised learning.

Reward-modulated STDP (R-STDP) extends the basic STDP rule by incorporating a global reward or punishment signal that modulates the synaptic weight changes. This allows the network to learn task-specific features by reinforcing connections that contribute to correct outputs and weakening those associated with errors. R-STDP bridges the gap between unsupervised STDP and supervised learning, enabling SNNs to tackle more complex, goal-oriented tasks like speech recognition. The reward signal in R-STDP can potentially be implemented using neuromodulators, further enhancing the biological plausibility of SNN training.

STDP has been employed to train SNNs for speech recognition, often in conjunction with temporal coding schemes such as time-to-first-spike for efficient processing. It has been used for unsupervised feature extraction from speech signals, with subsequent classification using methods like Hidden Markov Models (HMMs) or tempotrons. STDP-based convolutional SNNs have shown promise in achieving accurate and efficient speech recognition, with performance levels approaching those of traditional ANNs in certain scenarios.

Memristors are two-terminal devices whose conductance changes based on the history of voltage or current, closely mimicking how biological synapses adjust their efficacy [25]. In memristor-based STDP, each memristor stores a synaptic weight in its conductance, and weight updates occur directly on-chip whenever spikes arrive, following the device’s own switching dynamics [26]. By collocating memory and computation, this approach avoids the von Neumann bottleneck and enables energy-efficient, on-device learning in neuromorphic hardware [27].

Vlasov et al. (2022) demonstrated this concept on a spoken-digit recognition task by training spiking neural networks with memristor-based STDP using two memristor types—poly-p-xylylene (PPX) and CoFeB–LiNbO₃ nanocomposite [28]. Their networks, deployed entirely on neuromorphic hardware, achieved classification accuracies between 80 % and 94 % depending on network topology and decoding strategy, rivalling more complex off-chip learning algorithms while consuming minimal power and memory [29].

While STDP offers advantages such as biological plausibility, local learning, and energy efficiency, it also has limitations. Training deep networks and achieving state-of-the-art accuracy on complex tasks can be challenging compared to backpropagation-based methods [30]. Directly applying STDP to supervised learning tasks often requires extensions such as R-STDP or classifiers which interpret the output of the model [31]. Furthermore, STDP can be sensitive to the choice of hyperparameters and network architecture [32], and it may tend to extract frequently occurring features that are not necessarily the most discriminative for a specific task [33]. Despite these limitations, ongoing research continues to explore variations and extensions of STDP to enhance its capabilities for diverse learning tasks.

3.2.5 Eventprop

Staying on the theme of implementing the ubiquitous backpropagation for SNNs in a more efficient manner, we come to a novel training algorithm called eventprop. Eventprop utilises precise spike gradients to train the network, as opposed to approximate surrogate gradients typically used in BPTT, and it does so while using less compute and memory resources. It reaches SOTA performance while using 4x less memory and running 3x faster. [34]

More detail and explain why I have chosen it

Why Eventprop is the best

4 Improving Training Using New Loss

5 Methods

o Introduce why loss functions seems funky o Provide evidence for this. I.e. early spikes shouldn’t be weighted more since they may just be random o Implement a new one which is bell curve looking

Focus on hyper-parameter hyper parameter optimisation. include past paper

Choose which one is better and more unique

5.1 Training Process

Walk through how eventprop works mathematically.

The input is provided into the model and it is allowed to develop under neuronal dynamics. The loss is calculated according to Equation (19). Using the adjoint method the loss is propagated backwards to find out how the timing of each event contributed to the loss so that the weights of the synapses can be adjusted accordingly.

$$\mathcal{L}_{sum-exp} = -\frac{1}{N_{batch}} \sum_{m=1}^{N_{batch}} \log\left(\frac{\exp(\int_0^T e^{-t/T} V_{l(m)}^m(t) dt)}{\sum_{k=1}^{N_{out}} \exp(\int_0^T e^{-t/T} V_k^m(t) dt)}\right) \quad (19)$$

5.2 Model Description

First I implemented the same model used in [34] to get SOTA performance. The model uses Leaky Integrate-and-Fire neuron model with exponential synapses. It has an input layer of 700 neurons - corresponding to the 700 channels of the cochlea model used by SHD - and 20 output neurons for digits 0 to 9 in English and German. The model has a single hidden layer which has been tested with a size of 64, 128, 256, 512, and 1024. The hidden layer was tested using feed-forward only connections and fully connected recurrent connections, showing best results with a recurrent architecture.

This loss was chosen as it deals with the problem highlighted by the Gedanken Experiment.

5.3 Bayesian optimisation of hyperparameters

Optimise variables like tau, random shift delay, etc. by implementing Bayesian optimisaion in software.

5.4 Exploring novel loss functions

The current loss function is very simple. The later the spike occurs the less it contributes to the loss - exponentially. Surely spikes happening at the very beginning - when the word has barely begun to be said - shouldn't be so highly weighted? What if the loss is some kind of bell curve. where spikes in the middle contribute to the loss the most.

References

- [1] "The state of AI in 2025: Global survey | McKinsey," <https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai>.
- [2] B. Kindig, "AI Power Consumption: Rapidly Becoming Mission-Critical," <https://www.forbes.com/sites/bethkindig/2024/06/20/ai-power-consumption-rapidly-becoming-mission-critical/>.
- [3] W. Maass, "Networks of spiking neurons: The third generation of neural network models," *Neural Networks*, vol. 10, no. 9, pp. 1659–1671, Dec. 1997.
- [4] W. Maass and H. Markram, "On the computational power of circuits of spiking neurons," *Journal of Computer and System Sciences*, vol. 69, no. 4, pp. 593–616, Dec. 2004.
- [5] P. A. Merolla, J. V. Arthur, R. Alvarez-Icaza, A. S. Cassidy, J. Sawada, F. Akopyan, B. L. Jackson, N. Imam, C. Guo, Y. Nakamura, B. Brezzo, I. Vo, S. K. Esser, R. Appuswamy, B. Taba, A. Amir, M. D. Flickner, W. P. Risk, R. Manohar, and D. S. Modha, "A million spiking-neuron integrated circuit with a scalable communication network and interface," *Science*, vol. 345, no. 6197, pp. 668–673, 2014.
- [6] E. Izhikevich, "Simple model of spiking neurons," *IEEE Transactions on Neural Networks*, vol. 14, no. 6, pp. 1569–1572, Nov. 2003.
- [7] —, "Which model to use for cortical spiking neurons?" *IEEE Transactions on Neural Networks*, vol. 15, no. 5, pp. 1063–1070, Sep. 2004.
- [8] L. Deng and X. Li, "Machine Learning Paradigms for Speech Recognition: An Overview," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 5, pp. 1060–1089, May 2013.

- [9] G. Hinton, L. Deng, D. Yu, G. Dahl, A.-r. Mohamed, N. Jaitly, V. Vanhoucke, P. Nguyen, T. Sainath, and B. Kingsbury, “Deep Neural Networks for Acoustic Modeling in Speech Recognition,” 2012.
- [10] J. Wu, E. Yilmaz, M. Zhang, H. Li, and K. C. Tan, “Deep Spiking Neural Networks for Large Vocabulary Automatic Speech Recognition,” *Frontiers in Neuroscience*, vol. 14, p. 199, Mar. 2020.
- [11] A. Bittar and P. N. Garner, “A surrogate gradient spiking baseline for speech command recognition,” *Frontiers in Neuroscience*, vol. 16, p. 865897, Aug. 2022.
- [12] G. Bellec, F. Scherr, E. Hajek, D. Salaj, R. Legenstein, and W. Maass, “Biologically inspired alternatives to backpropagation through time for learning in recurrent neural nets,” *arXiv.org*, Feb. 2019.
- [13] E. O. Neftci, M. Hesham, and Z. Friedemann, “Surrogate Gradient Learning in Spiking Neural Networks: Bringing the Power of Gradient-Based Optimization to Spiking Neural Networks,” *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51–63, Nov. 2019.
- [14] C. Zhou, H. Zhang, L. Yu, Y. Ye, Z. Zhou, L. Huang, Z. Ma, X. Fan, H. Zhou, and Y. Tian, “Direct training high-performance deep spiking neural networks: A review of theories and methods,” *Frontiers in Neuroscience*, vol. 18, Jul. 2024.
- [15] S. Y. Arnaud Yarga and S. U. N. Wood, “Accelerating spiking neural networks with parallelizable leaky integrate-and-fire neurons*,” *Neuromorphic Computing and Engineering*, vol. 5, no. 014012, Mar. 2025.
- [16] S. Y. A. Yarga and S. U. N. Wood, “Accelerating SNN Training with Stochastic Parallelizable Spiking Neurons,” in *2023 International Joint Conference on Neural Networks (IJCNN)*, Jun. 2023, pp. 1–8.
- [17] W. Maas, “Networks of spiking neurons: The third generation of neural network models,” *Trans. Soc. Comput. Simul. Int.*, vol. 14, no. 4, pp. 1659–1671, Dec. 1997.
- [18] R. Koopman, A. Yousefzadeh, M. Shahsavari, G. Tang, and M. Sifalakis, “Overcoming the Limitations of Layer Synchronization in Spiking Neural Networks,” Aug. 2024.
- [19] Y. Zhong, R. Zhao, C. Wang, Q. Guo, J. Zhang, Z. Lu, and L. Leng, “SPiKE-SSM: A Sparse, Precise, and Efficient Spiking State Space Model for Long Sequences Learning,” Oct. 2024.
- [20] G. Bellec, F. Scherr, E. Hajek, D. Salaj, A. Subramoney, R. Legenstein, and W. Maass, “Eligibility traces provide a data-inspired alternative to backpropagation through time,” in *Real Neurons $\{\&\}$ Hidden Units: Future Directions at the Intersection of Neuroscience and Artificial Intelligence @ NeurIPS 2019*, Oct. 2019.
- [21] —, “Eligibility traces provide a data-inspired alternative to backpropagation through time,” in *Real Neurons $\{\&\}$ Hidden Units: Future Directions at the Intersection of Neuroscience and Artificial Intelligence @ NeurIPS 2019*, Oct. 2019.
- [22] A. Rostami, B. Vogginger, Y. Yan, and C. G. Mayr, “E-prop on SpiNNaker 2: Exploring online learning in spiking RNNs on neuromorphic hardware,” *Frontiers in Neuroscience*, vol. 16, Nov. 2022.
- [23] W. van der Veen, “Including STDP to eligibility propagation in multi-layer recurrent spiking neural networks,” Master’s thesis, 2021.
- [24] A. Rostami, B. Vogginger, Y. Yan, and C. G. Mayr, “E-prop on SpiNNaker 2: Exploring online learning in spiking RNNs on neuromorphic hardware,” *Frontiers in Neuroscience*, vol. 16, p. 1018006, Nov. 2022.
- [25] W. Chen, L. Song, S. Wang, Z. Zhang, G. Wang, G. Hu, and S. Gao, “Essential Characteristics of Memristors for Neuromorphic Computing,” *Advanced Electronic Materials*, vol. 9, no. 2, p. 2200833, 2023.
- [26] Y. Li, K. Su, H. Chen, X. Zou, C. Wang, H. Man, K. Liu, X. Xi, and T. Li, “Research Progress of Neural Synapses Based on Memristors,” *Electronics*, vol. 12, no. 15, p. 3298, Jan. 2023.
- [27] C. Weilenmann, A. N. Ziogas, T. Zellweger, K. Portner, M. Mladenović, M. Kaniselman, T. Moraitis, M. Luisier, and A. Emboras, “Single neuromorphic memristor closely emulates multiple synaptic mechanisms for energy efficient neural networks,” *Nature Communications*, vol. 15, no. 1, p. 6898, Aug. 2024.

- [28] D. Vlasov, Y. Davydov, A. Serenko, R. Rybka, and A. Sboev, “Spoken Digits Classification Based on Spiking Neural Networks with Memristor-Based STDP,” in *2022 International Conference on Computational Science and Computational Intelligence (CSCI)*, Dec. 2022, pp. 330–335.
- [29] A. Sboev, M. Balykov, D. Kunitsyn, and A. Serenko, “Spoken Digits Classification Using a Spiking Neural Network with Fixed Synaptic Weights,” in *Biologically Inspired Cognitive Architectures 2023*, A. V. Samsonovich and T. Liu, Eds. Cham: Springer Nature Switzerland, 2024, pp. 767–774.
- [30] Y. Guo, X. Huang, and Z. Ma, “Direct learning-based deep spiking neural networks: A review,” *Frontiers in Neuroscience*, vol. 17, p. 1209795, Jun. 2023.
- [31] F. Liu, W. Zhao, Y. Chen, Z. Wang, T. Yang, and L. Jiang, “SSTDP: Supervised Spike Timing Dependent Plasticity for Efficient Spiking Neural Network Training,” *Frontiers in Neuroscience*, vol. 15, p. 756876, Nov. 2021.
- [32] C. Lee, P. Panda, G. Srinivasan, and K. Roy, “Training Deep Spiking Convolutional Neural Networks With STDP-Based Unsupervised Pre-training Followed by Supervised Fine-Tuning,” *Frontiers in Neuroscience*, vol. 12, Aug. 2018.
- [33] A. Bittar and P. N. Garner, “A surrogate gradient spiking baseline for speech command recognition,” *Frontiers in Neuroscience*, vol. 16, Aug. 2022.
- [34] T. Nowotny, J. P. Turner, and J. C. Knight, “Loss shaping enhances exact gradient learning with Eventprop in spiking neural networks,” *Neuromorphic Computing and Engineering*, vol. 5, no. 1, p. 014001, Jan. 2025.