

Effects Of Real-time Synaptic Plasticity Using Spiking Neural Network Architecture

Matthew Frazier Nishant Shukla Worthy Martin*

Abstract—Artificial Neural Networks is a promising approach to study human brain computation in hopes of achieving similar learning by artificial agents. Recent architecture design of a low-power supercomputer by the University of Manchester, the SpiNNaker, has made it easier to design highly parallel brain-inspired algorithms. We used the SpiNNaker machine to implement a neural network capable of rewiring its connection in real-time while trying to minimize information loss and maximize the decrease in statistical dependence. We believe this is the first use of synaptogenesis on a spiking neural network architecture, laying the framework for future efficient brain-like neural networks. Additionally, we explored scalability issues and unintended pitfalls with this approach.

Keywords—Neural Networks, Synaptogenesis, SpiNNaker

I. INTRODUCTION

ARTIFICIAL neural networks often require vast amounts of computation memory and space [1]. Some of the most effective algorithms such as back propagation do not scale well for a network consisting of a billion or more neurons [Cite]. Furthermore, most neural network models do not fully utilize the machine on which they run. The ones that do fully utilize their machine are typically implemented on Graphical Processing Units (GPUs) which disproportionately drain energy.

Finding a solution to revive neural networks to the same computational caliber as other paradigms such as Support Vector Machines (SVM) would better adhere to biologically based design. Studying biologically inspired structures as opposed to statistically based ones allows us to grasp a better understanding of the human brain. By designing a neural network experiment on a new kind of computer architecture, this paper aims to demonstrate the high performing nature of neural networks.

Recently, deep feedforward neural networks have risen in popularity [Cite]. Multiple deep-learning approaches have proven to be practical. However, a major issue with deep-learning is the computation time for large networks. Deep-learning does not scale well due to the innate use of back-propagation in the algorithm.

Our approach to neural network learning is through an unsupervised model on a spiking neural network architecture. By utilizing every core of the 72-core energy-efficient embedded system, we are capable of reproducing complex neural network patterns in less than half the energy-cost. The network creates and removes synapses in real-time, so not only are weights learned and adjusted simultaneously, but so are the actual connections between the neurons. We demonstrate that introducing synaptic modification on a neural network leads

to a biologically accurate model that outperforms previous methods in both performance and energy costs.

II. BACKGROUND

THE subfield of AI known as neural network computation has recently received a large and growing amount of attention. Put simply; neural networks (NN) are computational models inspired by biological brains that are capable of machine learning. A NN usually consists of interconnected "neurons" (in quotation marks for they are artificially simulated to varying degrees of biological accuracy – but that topic is outside the scope of this paper) which compute from their inputs and produce various outputs, depending on the model.

This paper focuses on types of spiking neural network (SNN) models, in which the neuronal communication (and therefore, the computation of the system as a whole) is achieved via message spikes, or action potentials, from one neuron that is synapsed onto another. In this introduction, we first describe how a SNN is a good computational representation of the brain. We then describe various different spike-dependent spike-time learning methods used on a SNN, followed by an exploration of various types of SNN. We discuss why the neuromimetic hardware we use is a good platform for experimentation with SNN, and then describe our approach to the problem of simulating synaptogenesis on such a device and propose a solution. Our design section explores this solution in detail, and our implementation section describes our results. Furthermore, we discuss future work that remains to be conducted, including interesting implications such as computability beyond Von Neuman architecture and Turing machines.

A. The Brain And Other Spiking Neural Networks

The brain is a spiking neural network. That is to say; biological neurons pass information around the brain via action potentials which travel from a neurons' soma (cell body) down the axon and generate synaptic events which are received by the dendrites of all other connected neurons. Figure 1 diagrams the relevant biological prerequisites.

The process has been simplified greatly as neurobiology is not the focus of this paper, nor is the biological process yet completely understood. What's important are the overarching properties exhibited by brains. Namely; that brains are fast, power-efficient, and capable of obtaining, representing, and integrating complex multi-dimensional information sets into useful knowledge (and from noisy sensory inputs, at that). Moreover, the property of being able to fairly reliably solve highly complex problems is something we want computers

to have. Progress is however being made towards the goal of making machines more like minds; indeed, the bus in a computers' hardware may operate at tens of MHz, while an axon may carry only tens to hundreds of action potentials per second, about five orders of magnitude slower than a machine [9]. The power efficiency of the brain is still far ahead of even the best neuromimetic hardware but the gap is closing quickly, considering that biological evolution had several hundred million years of a head start.

One very important property of the brain that has received little attention by Computer Scientists, which is also believed to be quintessential in mankind's ability to adapt to and solve new and ever-more difficult problems, is neural plasticity involving the dynamic creation or removal of synapses. It would be erroneous to say that no degree of plasticity has been implemented in any type of SNN, but to the best of our knowledge there is yet no model which allows networks to add or remove synapses at runtime. Given that we do not yet fully understand the mechanisms governing neural plasticity in our own heads, it suffices to say that there is yet much work to be done in this field of computational neuroscience in which artificial neural networks dynamically alter their own topology.

B. Overview of Spike-time coding and learning methods

Before discussing different types of artificial spiking neural networks, it is necessary to describe various spike-time coding and learning methods that are relevant to understanding how the different types of spiking neural networks function.

Firstly, the distinction between rate- and rank-order coding must be addressed. Rate order coding is, as it sounds, a way to describe neural responses to stimulus in terms of firing rate, paying attention to the timing between spikes of a single input. Rank order, on the other hand, focuses on the relative timing of spikes from all inputs. Rate order has a number of drawbacks, though the primary one is its inability to transmit and process information in a short period of time; with n neurons being able to transmit over the course of 10 ms only $\log_2(nC1)$ bits of information [22]. Rank order, under the same constraints, is capable of achieving $\log_2(n!)$ bits of information [15].

As detailed in [22], rate order can be considered an *analog to frequency* converter, while rank order would be an *analog to delay* converter, as the time an integrate-and-fire neuron takes to reach threshold is dependent upon input strength. During an initial learning phase, in which a new output neuron i is created for each N -dimensional training input pattern, connection weights $w_{j,i}$ ($j = 1; 2; \dots; N$) are calculated based on the rank order learning rule:

$$w_{j,i}(t) = \text{mod}^{\text{order}(j;i)} \quad (1)$$

where $w_{j,i}$ is the connection weight between post-synaptic neuron i and pre-synaptic neuron j , mod is a learning parameter, mod is a modulation factor (determining the importance of the order of the first spikes), with $\text{order}(j;i)$ representing the order of spikes arriving at neuron i from all synapses to the neuron i (has value 0 for the first spike to neuron i , increases incrementally for each subsequent spike from other pre-synaptic neurons) [15]. Some models using rank

order coding will also merge output neurons with similar weight vectors based on the Euclidean distance between them, thus reducing the number of redundant neurons that must be simulated. In short, rank order (RO) coding is assumed to be a better basis for model construction as it allows networks to process more information in a shorter amount of time with fewer redundant neurons [22]. One of the only downsides to rank order over rate order is an increased overhead.

Moving on to learning methods, spike time dependent plasticity (STDP) implements plasticity through the use of long-term potentiation (LTP) and depression (LTD) [4], [5], [12]. Put simply, the connection weight between two neurons increases if the pre-synaptic neuron spikes before the post-synaptic neuron, and the weight decreases in the reverse case. Spike driven synaptic plasticity (SDSP) is a semi-supervised variant of STDP in which a threshold V_{mth} is given to the post-synaptic neurons' membrane potential [15]. If the membrane potential is above V_{mth} when an input spike arrives, potentiation occurs. Otherwise, the synapse experiences depression. These two cases are typically either shortly before or shortly after a post-synaptic spike is emitted, respectively.

STDP is the most commonly discussed paradigm for learning in neural networks. Its benefits include its unsupervised nature, and, significantly, the ability for the LTP and LTD weight modifications to be calculated by the post-synaptic neuron upon receipt of a spike. This allows the updates to be performed at a single point in the simulation: at the spike arrival event, which can further increase the efficiency of a simulation [4], [5].

We will later describe a new learning method referred to as synaptogenetic spike time dependent plasticity (SSTDP), utilizing STDP and governed by biologically inspired synaptogenetic formulas taken from [3], and influenced by Information Theory from Shannon.

C. Types of Spiking Neural Networks

- (a) **A spiking neural network (SNN)** is the most basic form, in which synaptic weights are static and neurons communicate via message spikes. While they have some drawbacks, an SNN is still a fairly powerful biologically-inspired model.
- (b) **An evolving spiking neural network (eSNN)** evolves its functionality and structure based on incoming information using the RO rule, but only during a learning phase [15]. Once a neurons' weight has been set, however, this model does not allow for any further tuning of synaptic weights to reflect on other incoming spikes at the same synapse. That is to say; these synapses can capture some long-term memory during the learning phase, but have little ability to capture any short term memory, which is detrimental to the computational capabilities of the network as a whole [15].
- (c) **A dynamically evolving spiking neural network (deSNN)** represents an improvement over eSNN in that it implements not only RO learning during an initial learning phase, but also a dynamic synaptic plasticity mechanism so that it will continue to learn and improve

performance during a recall phase; the deSNN proposed in [15] uses the SDSP mechanism. This fixes the primary shortcoming of the eSNN model; allowing for short term memory to be stored easily, and to be potentiated as long as it remains useful [15].

D. SpiNNaker

The Spiking Neural Network Architecture (SpiNNaker) board is a neuromimetic hardware system designed by the University of Manchester with the eventual intent of simulating the human brain. Each SpiNNaker processing node is a multi-core system-on-chip (SoC) with 1GB SDRAM, containing 18 ARM968 processor cores, connected in a toroidal triangular mesh with passed messages being managed by an on-chip router for an asynchronous packet-switched network-on-chip (NoC) [18]. It is the brainchild of Steve Furber, and has a number of properties that make it useful for neural computation:

Fast

The SpiNNaker system has a bisection bandwidth of over 5 billion packets/s [7]. It is, by design, meant to model spiking neural networks in real time. However, it can be said to outperform biology in terms of raw speed of neural activity, and must be slowed down to run real-time simulations [7], [19].

Power efficient

A simulation run on the SpiNNaker system has been documented taking 100 nJ per neuron per millisecond and 43 nJ per postsynaptic potential, a smaller power consumption than any other recorded digital simulation [20].

Massively parallelized

The brain is a massively parallel system populated with many low-performance asynchronous components known as neurons [8]. It follows that a digital system meant to effectively and efficiently emulate neural activity would also be massively parallelized, though each ARM core is substantially more complex than a single neuron, and can efficiently model many neurons in real time [7]. Each NoC is extended seamlessly to surrounding chips, creating a power-efficient system-wide interprocessor communication network [18].

Globally asynchronous, locally synchronous (GALS)

In a GALS system, individual processing cores have their own clock signals, potentially with different frequencies, and require no interprocessor phase-alignment [18]. Each SpiNNaker chip is its own GALS system, with an independently-clocked router managing intercore and interchip communications, leading to the SpiNNaker system as a whole also being a GALS system [8]. Major benefits of a GALS approach include that it eliminates any top-level system constraints (i.e. it is decentralized), and any timing closure issues [18]. It also offers increased flexibility regarding process variability [8].

Fault Tolerant

In any large-scale systems fault tolerance is a major concern, and is increases in importance with the number of fallible subsystems [7]. Fault tolerance was paid a great deal of attention by the SpiNNaker development team, and listing all the built-in mechanisms at various levels of abstraction is outside the scope of this paper. To summarize, fault tolerance mechanisms exist at nearly every level of the system, including those related to processors, the interrupt controller, timers, and packet communications, allowing subsystems to generally fail gracefully without greatly affecting the performance of the system as a whole [7].

Scalable

Due to the nature inherent to its GALS design, it is arbitrarily scalable. The only cost associated with scaling up the system is power consumption. The full million-core machine has an expected power budget of around 90 kW [7].

To summarize, the SpiNNaker specifically disregards three significant axioms of conventional supercomputing (memory coherence, synchronization, determinism) in a way that makes it well-suited for a wide range of biological and non-biological applications [7].

E. Synaptic Plasticity

A difficult problem in designing neural networks is configuring the initial parameters to optimize the network. For example, defining how the neurons are connected together greatly affects the efficiency and performance in learning [3]. Perceptrons are often designed as bipartite graphs, where every neuron in the pre-synaptic layer is connected to every neuron in the post-synaptic layer. This brute-force approach to constructing a network topology may achieve satisfactory results if the weights converge to the optimal values within an acceptable time. However, a bipartite graph scales exponentially to the number of nodes in each layer. In our solution, we let the neural network modify itself, providing a “hands-off” approach to designing a network.

Other than updating synaptic weights to achieve learning, we introduce two additional forms of synaptic plasticity. In our study, we allow new synapses to be formed (synaptogenesis), or existing synapses to be entirely disconnected. The three types of plasticity undergo the following principles [3](Cite: Levy 2004):

- (a) **Synaptogenesis** - Unlike most networks, we allow ours to form new connections between neurons while learning occurs. The creation of a new synapse depends on whether the post-synaptic neuron is at an optimal activation. If not, the network considers forming a new synapse. Additionally, synaptogenesis only occurs between neighboring neurons, ensuring a constraint on locality. Synapses are formed by taking into account both receptivity and avidity [3](Cite: Levy 2004).

Avidity is the measure of ability for each neuron to participate in a new synaptic connection, defined as

$$A_i(t) = \frac{P_j^a}{(a + \sum_k W_{ikj}(t))}$$

$$a = \begin{cases} < 1.0 & 10^{33} & \text{for unlimited avidity} \\ 1.0 & & \text{for moderate avidity} \\ : & 1.0 \cdot 10^{-3} & \text{for limiting avidity} \end{cases}$$

The receptivity of new synaptic connections is inversely proportional to the running average of the neuron's activation. We use the following equation for receptivity,

$$R_j(t) = \frac{r_1}{r_1 + y_j(t)^{r_2}}$$

$$y(t) = 0.99y_j(t-1) + 0.01Y_j(t)$$

(r_1 and r_2 are experimental constants)

Finally, the probability of forming a new synapse between two neurons i and j depends on both avidity and receptivity.

$$\text{Prob}(\text{of new synapse } ij) \propto A_i(t) R_j(t)$$

- (b) **Weight Modification** - Similar to most unsupervised weight modification rules, we use a deterministic process depending on the pre- and post-synaptic activities. This type of weight adjustment is based off the Hebbian rule, where links between nodes that fire together strengthens. Specifically, weights are adjusted using the following formula,

$$\Delta w(t+1) \propto f(\text{post}_j(t)) - g(\text{pre}_j(t); w_{ij}(t))$$

- (c) **Synapse Removal** - Lastly, we allow our network to undergo synaptic removal, a stochastic process where the probability of removal becomes non-zero if excitatory synapses fall below a specified threshold. Given some constant $\alpha > 0$, probability of synaptic removal is defined as

$$\text{Prob}(\text{removal } ij) = \begin{cases} 0 & \text{if } w_{ij}(t) + \Delta w_{ij}(t+1) > 0 \\ > 0 & \text{otherwise} \end{cases}$$

Synaptic connections are slowly modified through additions or removals until optimal output firing levels are obtained.

Levy asserts that “information theory has gained popularity in recent years as a tool for understanding brain recordings.”

III. DESIGN

NEURAL networks consist of biologically inspired relationships between individual nodes (neurons). Learning in such a network occurs by intelligently adjusting weights on the links between the nodes. We define the relationship between nodes into the following three categories:

- (a) **Static** - where the network topology is fixed, and can only change by manually adjusting the relationship between nodes. Most neural networks fall into this category

because regardless of the number of nodes n , only one fixed relationship forms between them:

$$\text{NumberOfTopologies}(n) = 1$$

- (b) **Semi-Dynamic** - in which the network is not static, and there is some freedom for nodes to rewire with other nodes in real time. The complexity grows exponentially. Each node has non-zero probability to be rewired with a constant c number of other nodes. Given n such nodes, the number of possible topologies can be up to

$$\text{NumberOfTopologies}(n) = c^n$$

- (c) **Dynamic** - where each node has non-zero probability to be wired with any other neuron in the entire network. Given n nodes, the number of possible topologies becomes

$$\text{NumberOfTopologies}(n) = n^{n-1}$$

A biological neural network such as the human brain does not follow the static network topology. Connections between neurons in the brain are regularly formed and removed over time [2], [3]. Moreover, such a neural network is not fully dynamic either, since locality is a physical constraint. For example, a neuron on the far end of the left hemisphere of a brain might never directly connect with a neuron on the right hemisphere.

We designed a programming framework on SpiNNaker's interface to enable a semi-dynamic network topology. Each neuron has the flexibility to rewire with none, all, or some of the fixed number of other neurons, following our synaptic plasticity rules. To test performance, the network analyzed the MNIST database of handwritten digits. We compared information loss (see appendix A) and statistical dependency (see appendix B) from a network following our synaptic plasticity model, to that without synaptogenesis and synapse removal.

In order to implement a spiking neural network, we choose an existing neural model. There exists multiple neural models for implementing a spiking neural network. Some of the most popular models are briefly covered below.

- (a) **Integrate-and-fire (IF)** is one such model which increments the voltage while a current is present, and finally resets the voltage once a threshold V_{th} is reached. Every voltage reset produces a spike for that neuron. While computationally efficient, the IF model does not exhibit properties of the cortical spiking neurons.
- (b) **Leaky integrate-and-fire (LIF)** models are a biologically revised version of the previous model. This model allows for “forgetting” isolated input currents. However, it also is insufficient for modeling cortical spiking neurons.
- (c)

