

Hybrid CPU-GPU Co-Processing Scheme for Simulating Spiking Neural Networks

Sreenivasa.N, S. Balaji

Abstract--- Many attempts have been made to study the neural networks and to model them. These attempts have led to the development of neural network simulation software packages such as GENESIS and NEURON which have been the de-facto simulators for some time now. However, further studies have found that one of the major hindrances in using the aforementioned simulators is speed. These simulators use time driven technique which isolates the mimicked time to brief time periods and in every progression the factors of neural states are estimated and reiterated through a numerical examination strategy. This method includes complex calculations which do not foster the development of scalable neural systems. The interest for quick re-encounters of neural systems has offered ascend to the use of alternative reproduction strategy: event driven simulation. The event driven simulation technique just processes and appraises the neural state factors when another event alters the typical advancement of the neuron, that is, the point at which information is created. In the meantime, it is realized that the data communication in neural networks is done by the purported spikes. Less than 1% of the neurons is at the same time dynamic which catalysis the efficiency of event-driven Spiking Neural Networks (SNN) simulation. In this work, we present our study on hybrid CPU-GPU based model for simulating SNNs.

Keywords: SNN, CPU-GPU co-processing, high performance computing, neural networks, numerical analysis

1. INTRODUCTION

The SNNs is a branch of artificial neural networks that thoroughly mimic the natural neural networks [1][6]. The simplified model of an SNN is demonstrated in Figure 1. Neurons are represented by N_j and the presynaptic links X_i of a neurons are labelled as small circle W_i . Then the synaptic weight formula is as follows: Synaptic weight[j] = $\sum_{i=1}^n X_i W_i$.

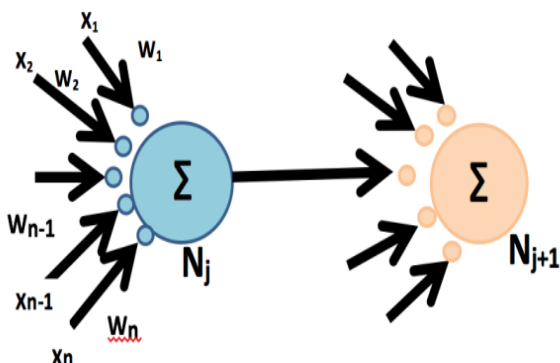


Fig. 1: Simplified spiking neural network

SNNs also comprise the concept of time in addition to the neuronal and the synaptic state in their operating model [13]. Unlike the multi-layer perceptron the neurons in this model fires only when the threshold of the membrane potential is reached [12][14]. A signal is generated as the response to the firing of neurons which travels to other neurons thus altering the potentials depending upon the generated signal [4][5]. SNNs aim to link the gap between neuroscience and machine learning with biologically realistic prototypes of neurons to execute computations. As the name suggests spikes are mathematically discrete events [9]. Membrane potential of a neuron is the most important biological parameter whose differential equations can be used to determine the occurrence of spike in a neuron. When a neuron reaches a certain threshold potential, it spikes which causes the potential to be reset and Leaky Integrate-and-Fire (LIF) is used to model this particular behaviour of neurons[8]. Also, since SNNs are inherently sparse we can leverage the corresponding network topology to model the connections in an SSN.

The overall architecture of the SNN in general is as shown in Figure 2.

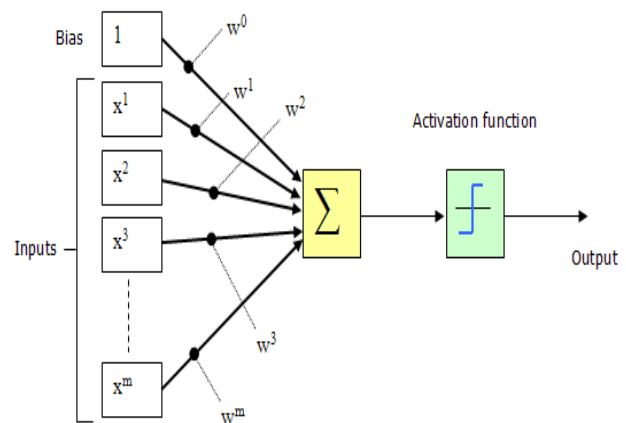


Fig. 2: SNN Architecture

The spike trains offer us enhanced ability to process spatio-temporal data or real-world sensory data. The spatial aspect refers to neurons that are connected to other neurons which are near to them and hence these inherently process chunks of the input independently (similar to a Convolutional Neural Networks would using a filter). The temporal aspect refers to the fact that spike trains occur over time, so what we lose in binary encoding, we gain in the temporal information of the spikes[7]. This permits to naturally process temporal data without the extra complexity

Revised Manuscript Received on December 22, 2018.

Sreenivasa. N, Research Scholar-Jain University, Dept. of Computer Science & Engineering., Nitte Meenakshi Institute of Technology, Yelahanka Bengaluru-560064, India (meetcna@yahoo.co.in)

S. Balaji, Centre for Incubation, Innovation, Research and Consultancy, Jyothy Institute of Technology, Tataguni, Bengaluru-560082, India (drsbalaji@gmail.com)

that Recurrent Neural Networks improve. In fact, that spiking neuron is fundamentally more influential computational units than conventional artificial neurons.

There has been a lot of on-going research in this area Francisco Naveros et al [3], EDLUT, Hamid Soleimani et al, Jesus A. Garrido et al, Govind Narasimman et al [2]. However, the major challenges faced were speed of simulation and the scalability of the simulation model which still persist. In this work, we have studied the existing models and proposed a hybrid CPU-GPU co-processing model to address the speed and scalability of simulation challenges.

2. THE MODEL

The critical aspect of developing any simulator is its flexibility and accuracy. This is the driving factor for leveraging the parallelism feature of GPU and endeavour towards co-processing (CPU-GPU) model. The parallelism enables us to use the time-driven model and shorter integration time intervals to achieve better accuracy in real-time. In the conventional parallel processing model for SNNs as in NeMo and GeNN, the parallelization in GPU is driven by the CPU to initialize the simulation and, therefore, we cannot leverage the complete potential of the GPU for modelling such a simulation.

On the other hand, when CPU-GPU co-process the events, all the parallelizable tasks can be independently run on GPU, while the CPU handles the sequential events as in Brian [11][13]. This intuitively means that when the CPU is processing the sequential events, the GPU updates the neural-state. The membrane potential (V_{m-c}) determines the neural state which can be expressed as:

$$C_m \frac{dV_{m-c}}{dt} = g_{AMPA}(t)(E_{AMPA} - V_{m-c}) + g_{GABA}(t)(E_{GABA} - V_{m-c}) + G_{rest}(E_{rest} - V_{m-c})$$

Where:

C_m : The capacitance membrane

E_{AMPA} & E_{GABA} : Reverse potential of every synaptic conductance

E_{rest} : Resting potential

g_{AMPA} & g_{GABA} : Conductance (decaying exponential functions)

The work described in provides a detailed explanation about the various parameters of a neural model. The state variables of a neuron i.e. V_{m-c} , g_{AMPA} & g_{GABA} help in defining the state of a neuron.

- 1) V_{m-c} potential membrane, output spike is generated if the membrane potential reaches threshold value.
- 2) g_{AMPA} & g_{GABA} - excitatory and inhibitory conductance. These conductance variables modify the membrane potential V_{m-c} . These parameters are inversely proportional to the integration and directly proportional to the synaptic weight.

A parallel combination of a resistor (whose conductance is g_L) and a capacitor (C) is used to represent a neuron in the LIF model. To charge the capacitor to produce a potential $V(t)$ a current source $I(t)$ is used as synaptic current as an input. The current input $I(t)$ charges the capacitor and

when the potential exceeds the threshold potential V^{th} i.e ($V(t) \geq V_{th}$), the capacitor discharges to the potential E_L which uses a switch controlled by voltage to simulate a biological neuron. This model can be represented using the differential equation as shown below:

$$C \frac{dV}{dt} = -g_L(V(t) - E_L) + I(t)$$

If the $I(t)$ is low, the potential $V(t)$ manages to be within the threshold potential V_{th} and thus spike is not produced as the necessary condition is not satisfied. In order to solve the differential of the LIF model, the fourth-order Runge-Kutta method was implemented. E. Ro et. Al. [10] in their works has processed this differential equation for both event driven (offline) and time-driven (online) techniques. The key performance indices for these methods are size of lookup table and size of integration step, respectively, as the execution time in former case would be directly proportional to the size of lookup table.

In our study, we have used Hodgkin-Huxley (H&H) Model [12][13] to simulate the SNN. The H&H model focuses on 3 channels: (i) sodium channel denoted by Na (ii) potassium channel denoted by K and (iii) I denotes the leakage channel along with the resistance R . Let us denote the input current as I , which is the sum of the impulses, namely, excitatory impulse which is denoted by I_E , I_I being the inhibitory impulse and I_{offset} being the current offset which is constant. The current which is externally injected to the system is denoted by I_{inj} and the model can be defined by the state equations:

$$\frac{dV}{dt} = \frac{1}{C} [-g_{Na}m^3h(V - E_{Na}) - g_Kn^4(V - E_K) - g_L(V - E_L) - g_e(V - E_e) - g_i(V - E_i) + I_{offset} + I_{inj}]$$

$$\frac{dm}{dt} = \alpha_m(V)(1 - m) - \beta_m(V)m$$

$$\frac{dn}{dt} = \alpha_n(V)(1 - n) - \beta_n(V)n$$

$$\frac{dh}{dt} = \alpha_h(V)(1 - h) - \beta_h(V)h$$

$$\frac{dg_e}{dt} = -\frac{g_e}{\tau_{synE}}$$

$$\frac{dg_i}{dt} = -\frac{g_i}{\tau_{synI}}$$

Where g_e and g_i denote excitatory and inhibitory synapses conductivity and g_{Na} , g_K , g_L ion channels conductance, E_{Na} , E_K , E_L , E_e , E_i ion channels reverse potentials. The definitions of the functions $\alpha_m, \beta_m, \alpha_n, \beta_n, \alpha_h, \beta_h$ are provided in [14]. The necessary condition for an action potential to be produced is that the potential membrane should quickly cross the threshold ($dV / dt \geq V_{th}$).

3. RESULTS AND DISCUSSION

In our study and experiments the key performance parameters like accuracy and scalability of the techniques and models have been evaluated which leverages the hybrid co-processing (CPU-GPU) model where time-driven events are processed in CPU and the event driven events are processed in GPU (to leverage the parallelism offered by GPU). As discussed earlier, the time driven simulation on CPU achieve high performance at a high precision. However, the shortcoming of such a simulation is scaling it to a large-scale neural network. Figure 3 shows the firing rate using Hodgkin-Huxley model.

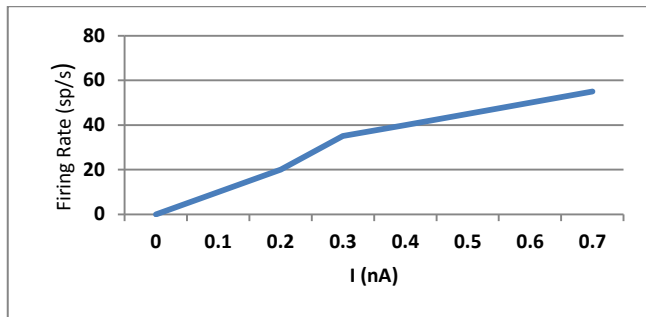


Fig. 3: Firing rate on GPU

The results shows that the firing rate on 100, 1000 and 10,000 neurons and we can see the constant behaviour when it fares with I(nA).

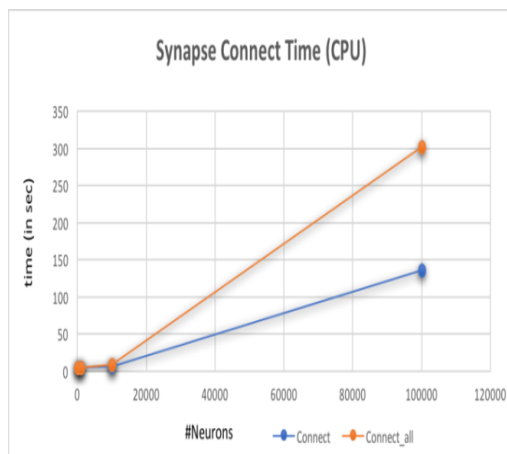


Fig. 4: Synapse connection time (CPU)

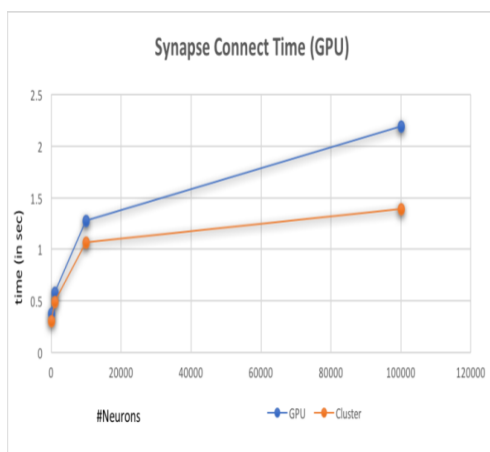


Fig. 5: Synapse connection time (GPU)

Figure 4 shows the Synapse connection time for 100, 1000, 10000 and 100,000 neurons on an Intel Core i7-5960X CPU @ 3.00GHz. We can see that the time grows exponentially. The blue line is an indicator of the conditional connection of synapse i.e. the synapses are connected based on a condition whereas the orange line indicated unconditional connection of all synapses.

Similarly, Figure 5 shows the synapse connection time on a GPU. However, blue line indicates a standalone GPU and orange indicate a GPU cluster and we can see that the performance stabilizes as we scale the number of neurons. As stated, the time-driven (TD) method faces the challenge of scaling and shows an exponential behaviour whereas the proposed co-processing (CPU-GPU) does address this challenge fairly well in terms of scaling and timing performance.

4. CONCLUSION

From the description of the aforementioned work, we can infer that it is imperative to integrate heterogeneous simulation techniques to develop a scalable, high-performing simulator. The next step in the study will be to perform hyper-parameter tuning to study the effect on the spike propagation and the queue management mechanism.

REFERENCES

1. Liwan, Yuling lu, Shuxiang song, Jim harkin, Junxiu liu "Efficient neuron architecture for FPGA-based spiking neural networks" Signals and systems conference (ISSC), IEEE, 2016
2. Govind Narsimhan, Subhrajit Roy, Xuanyou Fong, Kaushik Roy, Chip-Hong Chang, Arindam Basu "A low-voltage, low power STDP synapse implementation using domain-wall magnets for spiking neural networks" ISCAS, IEEE, 2016
3. Francisco Naveros, Niceto R. Luque, Jesús A. Garrido, Richard R. Carrillo, Mancia Anguita, and Eduardo Ro: "A Spiking Neural Simulator Integrating Event-Driven and Time-Driven Computation Schemes Using Parallel CPU-GPU Co-Processing: A Case Study" IEEE transactions on neural networks and learning systems, vol. 26, no. 7, July 2015.
4. Reza Haghighi and Chien Chern Cheah, "Optical Micromanipulation of Multiple Groups of Cells", IEEE International Conference on Robotics and Automation (ICRA) Washington State Convention Center Seattle, Washington, 2015.
5. Youssef Chahibi, Sasitharan Balasubramaniam et.al., "Molecular Communication Modeling of Antibody-mediated Drug Delivery Systems", IEEE Transactions On Biomedical Engineering, Vol. 62, No. 7, July 2015.
6. Jesus A. Garrido¹, Richard R. Carrillo², Niceto R. Luque¹, and Eduardo Ros¹ J. Cabestany, I. Rojas, and G. Joya (Eds.): "Event and Time Driven Hybrid Simulation of Spiking Neural Networks" IWANN 2011, Part I, LNCS 6691, pp. 554–561, 2011.
7. N. R. Luque, J. A. Garrido, R. R. Carrillo, O. J. D. Coenen and E. Ros, "Cerebellar input configuration toward object model abstraction in manipulation tasks," IEEE Trans. Neural Networks, vol. 22, no. 8, pp. 1321–1328, Aug. 2011.
8. D. F. Goodman and R. Brette, "The Brian simulator," Frontiers Neuroscience., vol. 3, no. 2, pp. 192–197, 2009.
9. Gibson Hu, Ying Guo, Rongxin Li, Adaptive Systems Team, Autonomous Systems Laboratory ICT Centre, CSIRO, A

- Self-Organizing Nano-Particle Simulator and Its Applications”, 978-0-7695-3166-3/08 \$25.00 © 2008 IEEE.
10. R. Brette et al., “Simulation of networks of spiking neurons: A review of tools and strategies,” J. Computer. Neuroscience., vol. 23, no. 3, pp. 349–398, 2007.
 11. E. Kandel, J. Schwartz, and T. Jessell, Principles of Neural Science, 4th ed. Amsterdam, Netherlands: Elsevier, 2000.
 12. W. Gerstner and W. Kistler, Spiking Neuron Models. Cambridge, U.K.: Cambridge Univ. Press, 2002.
 13. [R.Brette et al., “Simulation of networks of spiking neurons: A review of tools and strategies,” J. Comput. Neurosci., vol. 23, no. 3, pp. 349–398,2007.