

Improving the Interconnection Network of a Brain Simulator

Jonathan Heathcote

Short Report

1 Introduction

Despite its significance in nature, the brain is one of the least understood entities known to science. Current attempts to understand the brain centre on building models of its behaviour. Though models such as Spaun are able to produce high-level, realistic behaviour, they run slowly with conventional computers taking 2.5 hours to simulate one second of neural activity [1].

Neural models are typically a large graph of ‘neurons’ each being connected to potentially thousands of others. Signals, known as spikes, are produced by each neuron at an average rate of 10 Hz which each spike being destined for around 1,000 other neurons. With billions of neurons in a large neural model this results in very large numbers of very small messages being passed around the simulator [2]. Because neurons are often cheap to simulate, typical super computers featuring powerful processors with comparatively limited interconnection networks perform poorly.

The Blue Brain project has built a model with extremely realistic neuron behaviour which can exploit typical super-computer resources [3]. However, models are severely limited in size to hundreds of thousands of neurons compared with the 85 billion in a human brain. This work instead focuses on the simulation of large networks of simple neurons such as Spaun.

Due to the unsuitability of conventional architectures, special-purpose systems have been built for neural simulation. In this short report, current attempts to overcome these limitations are described followed by an overview of preliminary work carried out on the SpiNNaker brain simulation architecture. The report concludes with the research plan proposed to develop this research eventually to yield an improved architecture for brain simulation with a focus on the topology of the interconnection network.

2 Brain Simulators

Current special purpose brain-simulation architectures can be divided into two categories: those based on conventional digital circuits and those based on analogue electronics inspired by the analogue mechanisms in the brain. Though analogue systems can be very power efficient the maturity of digital design techniques has meant that modern analogue

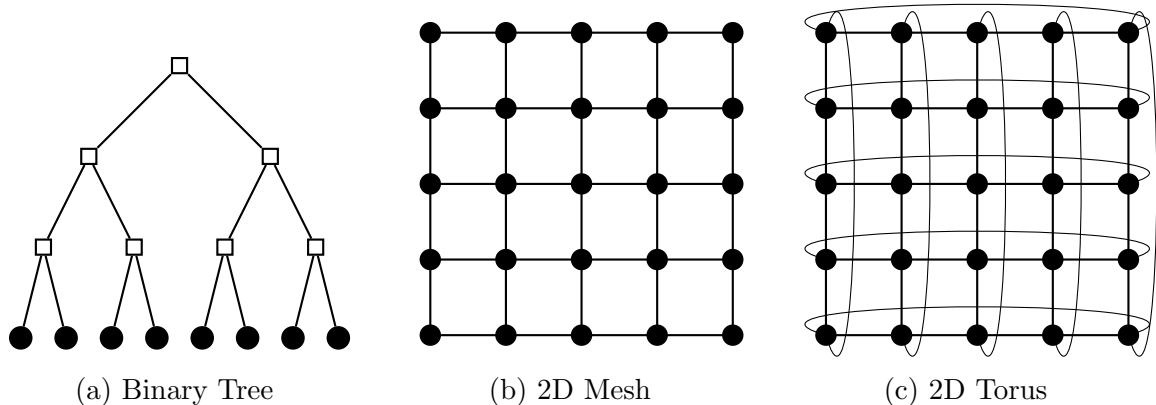


Figure 1: Network topology examples. Dots represent chips, boxes represent switches which forward messages but perform no other computation.

simulators are ‘mixed mode’ implementing only the neurons using analogue electronics and interconnecting them with digital electronics. As a result, the interconnection networks of modern analogue and digital simulators are typically directly comparable.

This section describes the interconnection network of three notable architectures. A wider and more detailed survey of current simulation architectures is available by Misra and Saha [4].

2.1 Neurogrid

The Neurogrid architecture consists of chips with analogue hardware to simulate tens of thousands of neurons. Spikes from these neurons are output by the chip serially and must be routed to each of the target neurons which may reside in other chips. Current prototypes feature 16 chips which form the leaves of a binary tree (figure 1a). This architecture can simulate around a million neurons with 100,000 times less power than a conventional super-computer [5].

The binary tree interconnection network topology is scalable with network latency, in the ideal case, increasing $O(\log N)$ with the number of chips. For neural simulation, however, tree topologies have been shown by Vainbrand and Ginosar to be non-optimal requiring large amounts of hardware to achieve the bandwidth required to transmit large numbers of spikes between nodes [2].

2.2 BrainScaleS

The BrainScaleS project has developed an architecture which, unconventionally, uses an entire silicon wafer on which tens of chips have been produced side-by-side. Each chip contains a number of analogue neuron simulators which are interconnected via a 2D mesh network (figure 1b). Spikes can be forwarded from chip to chip until they reach their destination [6]. It is intended that multiple wafers will be combined with conventional Internet Protocol (IP) switches.

Though mesh networks are easily scaled in principle, the BrainScaleS architecture is limited by the maximum size of a silicon wafer. The topology to be used to connect

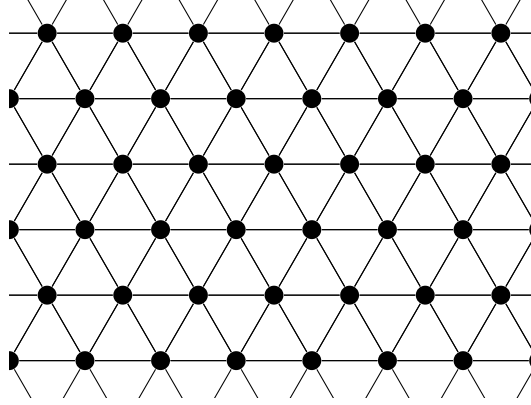


Figure 2: Section of the toroid topology used by SpiNNaker. Wrap-around links omitted.

wafers together has not yet been decided. Since silicon wafers are circular this means that the mesh networks on a wafer don't tessellate. Since gaps will be left in the network this may detract from the system's performance.

2.3 SpiNNaker

The SpiNNaker architecture is a completely digital simulator which uses a large number of small, low-power general purpose processors [7]. Eighteen processors are combined into chips which are connected in a toroid network. A toroid is a generalisation of a torus (figure 1c) with chips connecting to more than four other chips as in figure 2. This type of network was found to be optimal amongst common network topologies by Vainbrand and Ginosar for neural simulations [2]. This topology is easily scaled to arbitrary sizes by adding chips and is only limited by the number of bits used to address neurons in the system.

Preliminary work, however, has found that further improvements can be made over torus-like topologies for use in neural simulation and this is described in the next section.

3 Preliminary Work

Current work has focused on the SpiNNaker architecture and its interconnection network. This section outlines the work completed so far.

3.1 SpiNNaker Simulation

SpiNNaker systems are built by combining groups of 48 chips onto circuit boards which are then connected together via cables with the largest planned system containing 1,200 boards. Connections between chips on the same board use a low-cost technology using 16 wires. Connections between chips on different boards would require cables containing 768 wires which would be prohibitively expensive. Instead these board-to-board connections use more complex high-speed serial signalling which requires only 24 wires provided by commodity cables.

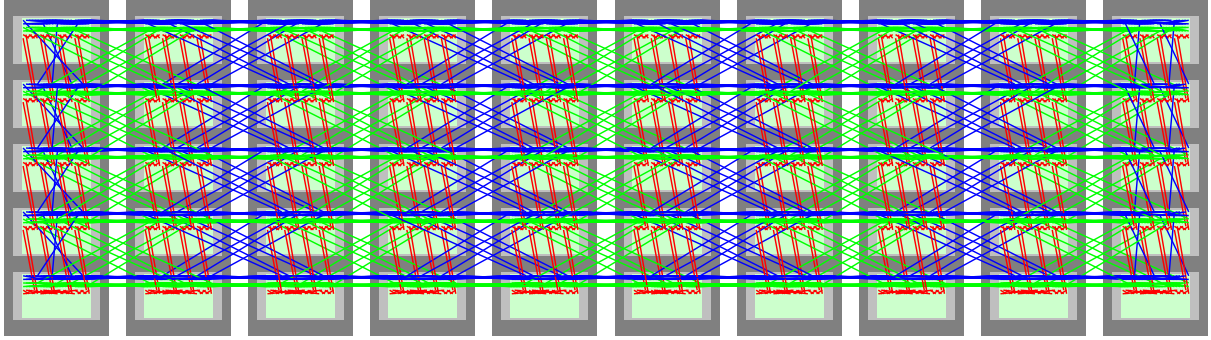


Figure 3: Wiring scheme for a SpiNNaker machine with 1,200 boards arranged in 10 cabinets. Line colour represents which of three sockets the wire is connected to.

A network interconnect simulator was developed which models the behaviour of the different kinds of link in a SpiNNaker system. The model showed that the use of high-speed serial links between boards increased the latency of spikes in the system by 80.4%. High-accuracy neural simulations require that spikes are delivered with low latency and may suffer as a result of this increased latency. Though the models currently planned for use on SpiNNaker are coarser-grained, this latency penalty may pose a problem for future networks.

3.2 SpiNNaker Wiring

Large SpiNNaker machines will be constructed by placing boards into standard computer cabinets and with cables connecting them together. There are two key constraints which must be considered when deciding how to connect the boards together. Firstly, the high-speed signalling technology requires that the cables used be short, disallowing connections between cabinets at opposite ends of the system. Secondly, since the cables must be manually connected, the system must be simple to assemble.

These constraints are typically universal to any large architecture and are an important consideration when designing network topologies. A tool was developed to aid the design of practical wiring schemes and used to develop a scheme for SpiNNaker which meets the constraints specified above.

A wiring scheme, shown for illustrative purposes in figure 3, was created using the tool for the largest planned SpiNNaker machine containing 3,600 wires. Though complex, only 53 repeating patterns are required for a technician to assemble the entire machine.

3.3 Small-World Super-Computers

The connections in the brain, along with many graph-like networks observed in nature exhibit the ‘small world’ property, first described by Frigyes Karinthy in 1929 [8]. It is best known as the theory of six degrees of separation which states that any two randomly-selected people in the world are connected via a chain of no more than six acquaintances. That is, for a very large graph with the small-world property, only a small number of ‘hops’ between nodes are required between any two places.

Watts and Strogatz proposed an algorithm for generating random networks with the

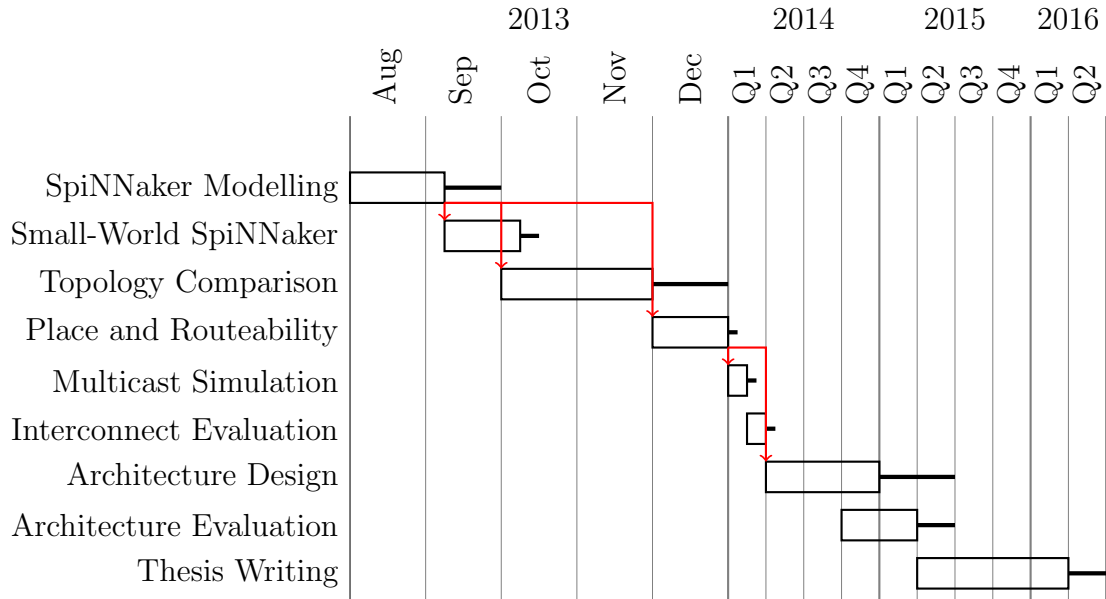


Figure 4: Gantt chart of proposed research plan. Arrows indicate dependencies, thick lines indicate slack. Note non-linear scale.

small-world property in which a regular network topology (such as a torus) is augmented with a small number of random connections [9]. Watts-Strogatz style topologies have been shown to increase the bandwidth available over simple torus networks at a low additional wiring cost [10].

Preliminary work has shown that latency is also improved in small-world networks based on tori. As a result, the latency of spikes in the system can be reduced potentially allowing increased simulation accuracy or speed.

The work was extended to account for the wiring constraints outlined in the previous subsection. Even when random links requiring long cables are disallowed, the improvement in latency is only slightly reduced.

4 Research Plan

Initial work will focus on the development of a more detailed simulator for SpiNNaker's interconnect topology. This will contribute to work comparing actual prototype SpiNNaker hardware against various software models of its behaviour in collaboration with other researchers. This work is hoped for journal publication later this year. This work may extend the existing INSEE [11] network simulator to allow greater flexibility in simulated designs. Alternatively, the work may extend the simulator built as part of the preliminary work on SpiNNaker.

With the resulting simulator, more detailed experiments will be carried out on the small-world topologies examined in the preliminary work. In particular, simulation of realistic network traffic and a more detailed model of the interconnect will allow better judgement of this unconventional topology.

Experiments on other less common topologies, such as express cubes [12] which offer latency advantages over standard torus networks, will be examined and compared.

These experiments will follow the process developed during the small-world topology experiments.

A further consideration when designing network topologies is the difficulty of assigning resources in the system. For example, when simulating a pair of connected neurons, they should be placed such that spikes can be quickly transmitted between them. In a torus-like network this typically means attempting to place connected neurons on chips which are near to each other. Though these challenges are well understood for common networks such as tori [13] work will need to be done to extend this to less conventional networks.

The penultimate stage of the project will be to carry out detailed research into the technologies available for implementing current interconnection networks. This will likely have an impact on the topology chosen as it may place constraints on cabling and performance as found in the preliminary studies of SpiNNaker.

The final stage of the work will be the development of a new architecture for neural simulation intended specifically to improve on the SpiNNaker architecture. The architecture will be tested by adding further detail to the models developed earlier in the project. In addition, comparisons with the actual performance of the mature SpiNNaker hardware should be possible.

5 Conclusion

The simulation of models of the brain has yielded promising results with high levels of realism being achieved. Conventional super-computer architectures are poorly suited to the task due to their focus on computational power rather than communication. This has led to the development of numerous special-purpose architectures for neural simulation.

The SpiNNaker architecture presents an interconnection network which allows the system to scale-up to fit large neural models. Though well suited to neural simulation, preliminary work has shown that improvements may be possible through the use of alternative network topologies. In addition, tools have been developed which will enable the performance and practicality of new interconnection topologies to be evaluated.

The project will eventually develop a new architecture for neural simulation with a focus on the topology of the interconnection network. Work will develop from the preliminary studies of SpiNNaker with the construction of a flexible network simulator which will allow new topologies to be tested and evaluated.

References

- [1] Chris Eliasmith, Terrence C Stewart, Xuan Choo, Trevor Bekolay, Travis DeWolf, Charlie Tang, and Daniel Rasmussen. A large-scale model of the functioning brain. *Science*, 338(6111):1202–1205, 2012.
- [2] Dmitri Vainbrand and Ran Ginosar. Scalable network-on-chip architecture for configurable neural networks. *Microprocessors and Microsystems*, 35(2):152–166, 2011.
- [3] Henry Markram. The blue brain project. *Nature Reviews Neuroscience*, 7(2):153–160, 2006.

- [4] Janardan Misra and Indranil Saha. Artificial neural networks in hardware: A survey of two decades of progress. *Neurocomputing*, 74(13):239 – 255, 2010.
- [5] Swadesh Choudhary, Steven Sloan, Sam Fok, Alexander Neckar, Eric Trautmann, Peiran Gao, Terry Stewart, Chris Eliasmith, and Kwabena Boahen. Silicon neurons that compute. In *Artificial Neural Networks and Machine Learning–ICANN 2012*, pages 121–128. Springer, 2012.
- [6] Johannes Schemmel, Johannes Fierens, and Karlheinz Meier. Wafer-scale integration of analog neural networks. In *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on*, pages 431–438. IEEE, 2008.
- [7] SB Furber, Steve Temple, and AD Brown. High-performance computing for systems of spiking neurons. In *AISB06 workshop on GC5: Architecture of Brain and Mind*, volume 2, pages 29–36, 2006.
- [8] Frigyes Karinthy. *Chain-Links*. Public Domain, 1929. Translated from Hungarian and annotated by Adam Makkai and Enikő Jankó.
- [9] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393(6684):440–442, 1998.
- [10] Ji-Yong Shin, Bernard Wong, and Emin Gün Sirer. Small-world datacenters. In *Proceedings of the 2nd ACM Symposium on Cloud Computing*, pages 2:1–2:13. ACM, 2011.
- [11] Javier Navaridas, Jose Miguel-Alonso, Jose A Pascual, and Francisco J Ridruejo. Simulating and evaluating interconnection networks with INSEE. *Simulation Modelling Practice and Theory*, 19(1):494–515, 2011.
- [12] William J. Dally. Express cubes: improving the performance of k -ary n -cube interconnection networks. *Computers, IEEE Transactions on*, 40(9):1016–1023, 1991.
- [13] William James Dally and Brian Towles. *Principles and Practices of Interconnection Networks*. Morgan Kaufmann Publishers, 2004.