# IBM Research Report

## $10^{14}$

**Theodore M. Wong, Robert Preissl, Pallab Datta, Myron Flickner, Raghavendra Singh, Steven K. Esser, Emmett McQuinn, Rathinakumar Appuswamy, William P. Risk, Horst D. Simon\*, Dharmendra S. Modha**

IBM Research Division
Almaden Research Center
650 Harry Road
San Jose, CA  95120-6099
USA

\* Lawrence Berkeley National Lab
Berkeley, CA  94720

# $10^{14}$

Theodore M. Wong, Robert Preissl, Pallab Datta, Myron Flickner,
Raghavendra Singh, Steven K. Esser, Emmett McQuinn, Rathinakumar Appuswamy, William P. Risk,
Horst D. Simon[†], and Dharmendra S. Modha
IBM Research - Almaden, San Jose CA 95120
[†] Lawrence Berkeley National Lab, Berkeley, CA 94720
Contact e-mail: dmodha@us.ibm.com

## I. INTRODUCTION

Since the final submission of our work on the Compass scalable simulator for the IBM TrueNorth Cognitive Computing architecture [1], we have simulated an unprecedented 2.084 billion neurosynaptic cores containing $53 \times 10^{10}$ neurons and $1.37 \times 10^{14}$ synapses running at only $1542\times$ slower than real time. We attained this scale by using the Sequoia 96-rack IBM® Blue Gene®/Q supercomputer at Lawrence Livermore National Labs. By comparison, the ultimate vision of the DARPA SyNAPSE program is to build a cognitive computing architecture with $10^{10}$ neurons and $10^{14}$ synapses, inspired by the following: Shepherd [2] estimates the number of synapses in the human brain as $0.6 \times 10^{14}$, and Koch [3] estimates the number of synapses in the human brain as $2.4 \times 10^{14}$.

It is important to clarify that we have not built a biologically realistic simulation of the complete human brain. Rather, we have simulated a novel modular, scalable, non-von Neumann, ultra-low power, cognitive computing architecture at the DARPA SyNAPSE metric of $10^{14}$ synapses that itself is inspired by the number of synapses in the human brain. We mathematically abstract away computation ("neurons"), memory ("synapses"), and communication ("axons", "dendrites") from biological detail towards engineering goals of maximizing function (utility, applications) and minimizing hardware cost (power, area, delay) and design complexity.

## II. RESULTS

We completed a new experimental evaluation of Compass that simulated up to an unprecedented 2.08 billion neurosynaptic cores and showed near-perfect weak and excellent strong scaling behavior when simulating a CoCoMac macaque network model [4] on the Sequoia Blue Gene/Q. Sequoia has 96 Blue Gene/Q racks, with each rack comprising 1024 compute nodes. Each compute node is composed of 17 processor cores, hosted on a single multi-core CPU paired with 16 GB of dedicated physical memory, and is connected to other nodes in a five-dimensional torus through 10 bidirectional 2-GB/second links. HPC applications run on 16 of the processor cores, and can spawn up to 4 hardware threads on each core for a total of 64 threads; per-node system software runs on the remaining core. The maximum Blue Gene/Q size was therefore 98,304 nodes, equivalent to 1,572,864 Blue Gene/Q application cores addressing 1.5 Petabytes of memory.

We show the strong and weak scaling behavior of Compass on Sequoia in figure 1. We first evaluated weak scaling behavior by simulating a CoCoMac model of 21,206 TrueNorth cores per Blue Gene/Q node while increasing the Blue Gene/Q core count from 16,384 to 1,572,864 (i.e., increasing the rack count from one to 96). We then evaluated strong scaling behavior by taking each of the weak-scaling models that ran on 393,216, 786,432, 1,179,648, and 1,572,864 Blue Gene/Q cores (i.e., 24, 48, 72, and 96 racks), and simulating the model on 393,216, 786,432, and 1,179,648 Blue Gene/Q cores; we had a sufficient Blue Gene/Q time allocation to finish all but one strong scaling experiment. We ran all simulation experiments for 500 simulated one-millisecond TrueNorth ticks.
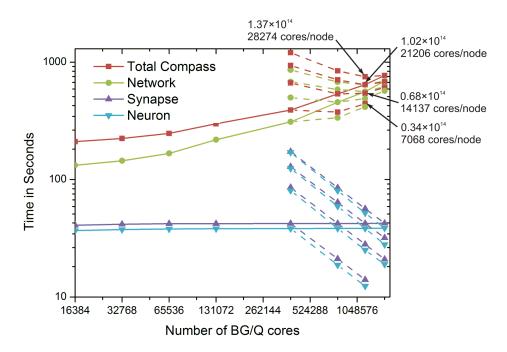
## REFERENCES

[1] R. Preissl, T. M. Wong, R. Appuswamy, P. Datta, M. Flickner, R. Singh, S. K. Esser, E. McQuinn, W. P. Risk, H. D. Simon, and D. S. Modha, "Compass: A scalable simulator for an architecture for cognitive computing," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis (SC 2012)*, Nov. 2012.

[2] G. M. Shepherd, Ed., *The Synaptic Organization of the Brain*, 5th ed. Oxford University Press, USA, Nov. 2003.

[3] C. Koch, *Biophysics of Computation: Information Processing in Single Neurons*, 1st ed. Oxford University Press, USA, Oct. 2004.

[4] R. Kötter, "Online retrieval, processing, and visualization of primate connectivity data from the cocomac database," *Neuroinformatics*, vol. 2, pp. 127–144, 2004.

(a) Average total and per-execution-phase runtime per Compass OpenMP thread when executing 500 simulated one-millisecond TrueNorth ticks. "Network" is the time spent sending and receiving spikes via MPI messaging at a Compass processes. "Synapse" is the time spent propagating spikes through all of the simulated TrueNorth crossbars hosted at a Blue Gene/Q node. "Neuron" is the time spent leaking, integrating, and firing at all of the simulated TrueNorth neurons hosted at a Blue Gene/Q node. For each "Total Compass" strong scaling curve, we include an arrow showing the synapse count and the number of TrueNorth cores per Blue Gene/Q node when simulating using 1,179,648 Blue Gene/Q cores. We discuss the details of the execution phases in the final submission of our work on the Compass simulator [1].



(b) Average counts of MPI messages per Compass MPI process and spikes received across all processes, per simulated one-millisecond TrueNorth tick. We show the counts for the weak-scaling experiments only.

Fig. 1. Strong and weak scaling in Compass when simulating cores on the Sequoia 96-rack Blue Gene/Q, showing (a) the average total and per-execution-phase runtime per Compass OpenMP thread and (b) the average counts of MPI messages and spikes received per Compass MPI process. The solid lines show weak scaling results when simulating 21,206 TrueNorth cores per Blue Gene/Q node. The dashed lines show strong scaling results when simulating the weak-scaling models that ran on 393,216, 786,432, 1,179,648, and 1,572,864 Blue Gene/Q cores (0.52 billion, 1.04 billion, 1.56 billion, and 2.08 billion TrueNorth cores respectively). We ran with one 64-thread Compass MPI process per Blue Gene/Q node.