

Crossbar arrays for Storage Class Memory and non-Von Neumann computing

Geoffrey W. Burr

IBM Research – Almaden

650 Harry Road, San Jose, California USA 95120

Phone: +1-408-927-1512 E-mail: gwburr@us.ibm.com

Abstract

I discuss recent work towards large crossbar arrays of NVM for Storage Class Memory and non-Von Neumann computing, incorporating advancements in nonlinear Access Devices and in the understanding of how device imperfections can adversely affect neural network performance.

1. Introduction

For more than 50 years, the capabilities of Von Neumann-style information processing systems — in which a "memory" delivers operations and then operands to a dedicated "central processing unit" — have improved dramatically. While it may seem that this remarkable history was driven by ever-increasing density (Moore's Law), the actual driver was Dennard's Law: a device-scaling methodology which allowed each generation of smaller transistors to actually perform better, in every way, than the previous generation.

Unfortunately, Dennard's Law terminated some years ago, and as a result, Moore's Law is now slowing considerably. In a search for ways to continue to improve computing systems, the attention of the IT industry has turned to Non-Von Neumann algorithms, and in particular, to computing architectures motivated by the human brain.

2. Storage Class Memory

At the same time, memory technology has been going through a period of rapid change, as new nonvolatile memories (NVM) — such as Phase Change Memory (PCM), Resistance RAM (RRAM), and Spin-Torque-Transfer Magnetic RAM (STT-MRAM) — emerge that complement and augment the traditional triad of SRAM, DRAM, and Flash. While these NVM candidate technologies are still relatively unproven compared to Flash, there is a strong opportunity for one or more of them to find success in applications that do not involve simply "replacing" NAND Flash.

Such memories could enable Storage-Class Memory (SCM, Figure 1) — an emerging memory category that seeks to combine the high performance and robustness of solid-state memory with the long-term retention and low cost of conventional hard-disk magnetic storage. Storage Class Memory creates two entirely new and distinct levels within the memory and storage hierarchy. These levels are differentiated from each other by access time, with both levels located within the more than two orders of magnitude between the latencies of off-chip DRAM (~80ns) and NAND Flash (20μs).

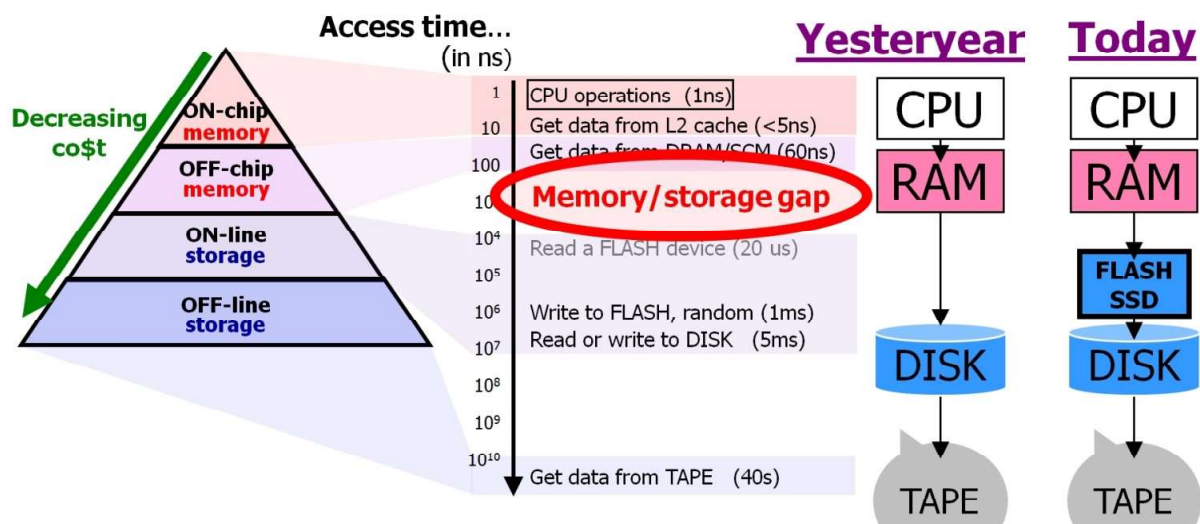


Figure 1 – Storage Class Memory [4-5] can be thought of as the realization that many emerging alternative nonvolatile memory technologies – such as Phase Change Memory (PCM) [1-3], Resistance RAM (RRAM), and Spin-Torque-Transfer Magnetic RAM (STT-MRAM) – can potentially offer significantly more than Flash, in terms of higher endurance, significantly faster performance, and direct-byte access capabilities.

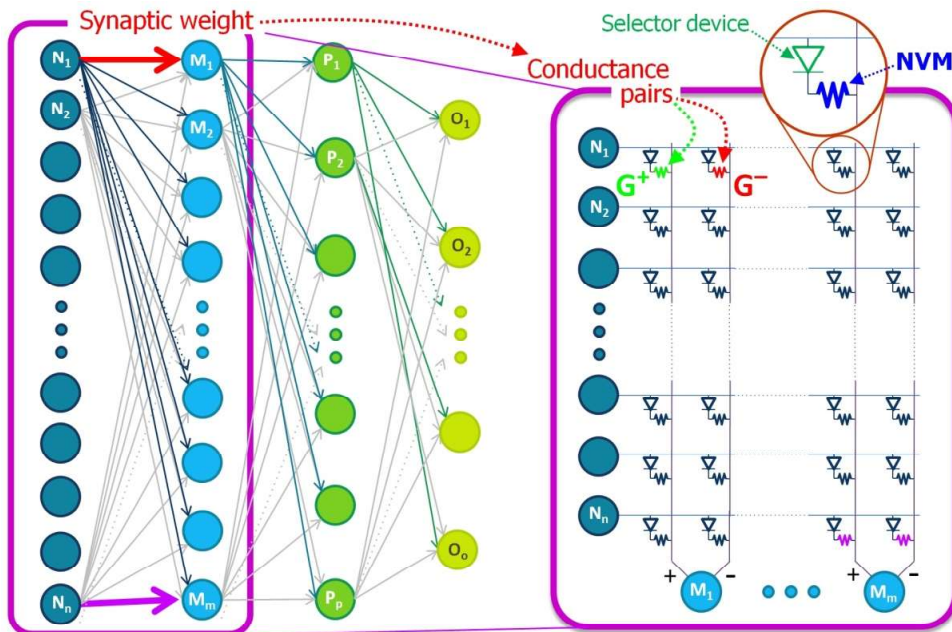


Figure 2 – Neuro-inspired non-Von Neumann computing [9-11], in which neurons activate each other through dense networks of programmable synaptic weights, can be implemented using dense crossbar arrays of nonvolatile memory (NVM)[1-4] and selector [6-8] device-pairs.

3. Neuromorphic Computing

Such large arrays of NVM can also be used in non-Von Neumann neuromorphic computational schemes, with device conductance serving as the plastic (modifiable) “weight” of each “native” synaptic device (Figure 2). This is an attractive application for these devices, because while many synaptic weights are required, requirements on yield and variability can be more relaxed. However, work in this field has remained highly qualitative in nature, and slow to scale in size.

I will discuss our recent work towards large crossbar arrays of NVM for both of these applications. After briefly reviewing earlier work on PCM [1-3], SCM [4-5], and access devices [6] based on copper-containing Mixed-Ionic-Electronic-Conduction (MIEC) [7-8], I will discuss our recent work on quantitatively assessing the engineering tradeoffs inherent in NVM-based neuromorphic systems [9-11].

Acknowledgements

This presentation describes work done together with many co-authors and collaborators, both from IBM and elsewhere, over the past 10 years. While I cannot possibly list everyone, I would particularly like to thank:

- Simone Raoux, Charles Rettner, Yi-Chou Chen, Bong-Sub Lee, Michelle Cheng, Hsiang-Lan Lung, Chung Lam, Matt BrightSky, Sangbum Kim, Martin Salinga and Daniel Krebs, in the context of Phase-Change Memory;
- Rich Freitas and Winfried Wilcke, in the context of Storage Class Memory;
- Kailash Gopalakrishnan, Rohit Shenoy, Kumar Virwani,

Bulent Kurdi and Alvaro Padilla, in the context of access devices based on Mixed-Ionic-Electronic Conduction (MIEC);

- Carmelo di Nolfo, Irem Boybat, Severin Sidler, Alessandro Fumarola, Lucas Sanches, Junwoo Jang, Kibong Moon, Hyunsang Hwang, and Yusuf Leblebici, in the context of Neuromorphic Devices & Architectures;
- Management support from Bulent Kurdi, Winfried Wilcke, Spike Narayan, Chung Lam, T. C. Chen, Sudhir Gowda, Dario Gil, Wilfried Haensch and Heike Riel; and
- Most particularly Bob Shelby (PCM + Neuromorphic) and Pritish Narayanan (MIEC + Neuromorphic), who have been willing to join me on more than one of these adventures.

References

- [1] S. Raoux et al., IBM J. R&D, 52(4/5), 465 (2008).
- [2] G. W. Burr et al., IEEE JETCAS, 6(2), 146 (2016).
- [3] G. W. Burr et al., J. Vac. Sci. Tech. B, 28(2), 223 (2010).
- [4] G. W. Burr et al., IBM J. R&D, 52(4/5), 449 (2008).
- [5] ITRS 2013, ERD chapter (see www.itrs2.net).
- [6] G. W. Burr et al., JVST B, 32(4), 040802 (2014).
- [7] R. S. Shenoy et al., SST, 29(10), 104005 (2014).
- [8] P. Narayanan et al., J. EDS, 3(5), 423 (2016).
- [9] G. W. Burr et al., IEDM Tech. Digest, T29.5 (2014).
- [10] G. W. Burr et al., IEEE Trans. Electr. Dev., 62(11) 3498 (2015).
- [11] G. W. Burr et al., IEDM Tech. Digest, T4.4 (2015).

More references available at
researcher.watson.ibm.com/researcher/view_group.php?id=3631