

References

- Slonczewski, J. C. *Phys. Rev. B* **39**, 6995–7002 (1989).
- Slonczewski, J. C. *J. Magn. Magn. Mater.* **159**, L1–L7 (1996).
- Berger, L. *Phys. Rev. B* **54**, 9353–9358 (1996).
- Stiles, M. D. & Zangwill, A. *Phys. Rev. B* **66**, 014407 (2002).
- Jullière, M. *Phys. Lett.* **54A**, 225–226 (1975).
- Mooodera, J. S., Kinder, L. S., Wong, T. M. & Meservey, R. *Phys. Rev. Lett.* **74**, 3273–3276 (1995).
- Miyazaki, T. & Tezuka, N. *J. Magn. Magn. Mater.* **139**, L231–L234 (1995).
- Wang, D., Nordman, C., Daughton, J. M., Qian, Z. & Fink, J. *IEEE Trans. Magn.* **40**, 2269–2271 (2004).
- Butler, W. H., Zhang, X.-G., Schulthess, T. C. & MacLaren, J. M. *Phys. Rev. B* **63**, 054416 (2001).
- Yuasa, S., Nagahama, T., Fukushima, A., Suzuki, Y. & Ando, K. *Nature Mater.* **3**, 868–871 (2004).
- Parkin, S. S. P. *Nature Mater.* **3**, 862–867 (2004).
- Min, T. *et al. IEEE Trans. Magn.* **46**, 2322–2327 (2010).
- Liu, J., Jaiyen, B., Veras, R. & Mutlu, O. in *Proc. 39th Annu. Int. Symp. Computer Architecture (ISCA '12)* 1–12 (IEEE Computer Society, 2012).
- Lee, K., Kan, J. J. & Kang, S. H. in *Proc. 2014 Int. Symp. Low Power Electronics and Design (ISLPED '14)* 131–136 (ACM, 2014).
- Thomas, L. *et al. J. Appl. Phys.* **115**, 172615 (2014).
- Ikeda, S. *et al. Nature Mater.* **9**, 721–724 (2010).
- Worledge, D. C. *et al. Appl. Phys. Lett.* **98**, 022501 (2011).
- Sun, J. Z. *Phys. Rev. B* **62**, 570–578 (2000).
- Bedau, D. *et al. Appl. Phys. Lett.* **97**, 262502 (2010).
- Liu, H. *et al. J. Magn. Magn. Mater.* **358–359**, 233–258 (2014).
- Kent, A. D., Özyilmaz, B. & del Barco, E. *Appl. Phys. Lett.* **84**, 3897–3899 (2004).
- Huai, Y., Albert, F., Nguyen, P., Pakala, M. & Valet, T. *Appl. Phys. Lett.* **84**, 3118–3120 (2004).
- Hosomi, M. *et al. in Proc. IEDM Tech. Dig.* 459–462 (2005).
- Kishi, T. *et al. in Proc. IEEE Int. Electron Devices Meeting (IEDM)* 1–4 (2008).
- Worledge, D. C. *et al. Proc. IEEE Int. Electron Devices Meeting (IEDM)* 296–299 (2010).
- Nowak, J. J. *et al. IEEE Magn. Lett.* **2**, 3000204 (2011).
- Kim, W. *et al. IEEE Int. Electron Devices Meeting (IEDM)* 24.1.1–24.1.4 (2011).
- <http://www.everspin.com/>
- Miron, I. M. *et al. Nature Mater.* **9**, 230–234 (2010).
- Liu, L. *et al. Science* **336**, 555–558 (2012).
- Mellnik, A. R. *et al. Nature* **511**, 449–451 (2014).
- Fan, Y. *et al. Nature Mater.* **13**, 699–704 (2014).

Acknowledgements

A.D.K. thanks G. Wolf for comments on the manuscript and for preparing Fig. 2. He acknowledges support from the National Science Foundation, grant number NSF-DMR-1309,202. D.C.W. thanks J. DeBrosse for comments on the manuscript.

Competing financial interests

A.D.K. is the founder of Spin Transfer Technologies.

Memory leads the way to better computing

H.-S. Philip Wong and Sayeef Salahuddin

New non-volatile memory devices store information using different physical mechanisms from those employed in today's memories and could achieve substantial improvements in computing performance and energy efficiency.

Current memory devices store information in the charge state of a capacitor; the presence or absence of charges represents logic 1's or 0's. Several technologies are emerging to build memory devices in which other mechanisms are used for information storage. They may allow the monolithic integration of memories and computation units in three-dimensional chips for future computing systems¹. Among those promising candidates are spin-transfer-torque magnetic random access memory (STT-MRAM) devices, which store information in the magnetization of a nanoscale magnet. Other candidates that are approaching commercialization include phase change memory (PCM), metal oxide resistive random access memory (RRAM) and conductive bridge random access memory (CBRAM).

Today's computing systems use a hierarchy of volatile and non-volatile data storage devices to achieve an optimal trade-off between cost and performance². The portion of the memory that is the closest to the processor core is accessed frequently, and therefore it requires the fastest operation speed possible; it is also

the most expensive memory because of the large chip area required. Other levels in the memory hierarchy are optimized for storage capacity and speed (Fig. 1). The main memory is often located in a separate chip because it is fabricated with a different technology from that of the microprocessor.

For over 30 years, static random access memory (SRAM)³ and dynamic random access memory (DRAM)³ have been the workhorses of this memory hierarchy⁴. Both SRAM and DRAM are volatile memories — that is, they lose the stored information once the power is cut off. For non-volatile data storage, magnetic hard disk drives (HDDs) have been in use for over five decades^{5–7}. Since the advent of portable electronic devices such as music players and mobile phones, however, solid-state non-volatile memory known as Flash memory⁸ has been introduced into the information storage hierarchy between the DRAM and the HDD. Flash has become the dominant data storage device for mobile electronics; increasingly, even enterprise-scale computing systems and cloud data storage systems are using Flash to complement the storage capabilities of HDD.

Resistive switching memory technologies

The design specifications for memory (volatile data storage, fast, expensive) and for storage (non-volatile data storage, slow, inexpensive) are different, and they often have different data access standards and protocols. Around 15 years ago, researchers started exploring the possibility of blurring the design boundary between memory and storage^{9,10}, and coming up with new data access modes and protocols that are neither 'memory' nor 'storage'. Indeed, the adoption of Flash in the memory hierarchy (albeit on a separate chip from the processor) inspired the exploration of computing architectures that capitalize on the salient features of Flash: non-volatility and high density¹¹. At the same time, new types of non-volatile memory have emerged that can easily be integrated on-chip with the microprocessor cores because they use a different set of materials and require different device fabrication technologies from Flash¹². Some of them can be programmed and read quickly; others can have very high data storage density. Importantly, all of these memories are free from the limitations of Flash — that is, low endurance, need for high voltage supply, slow write speed

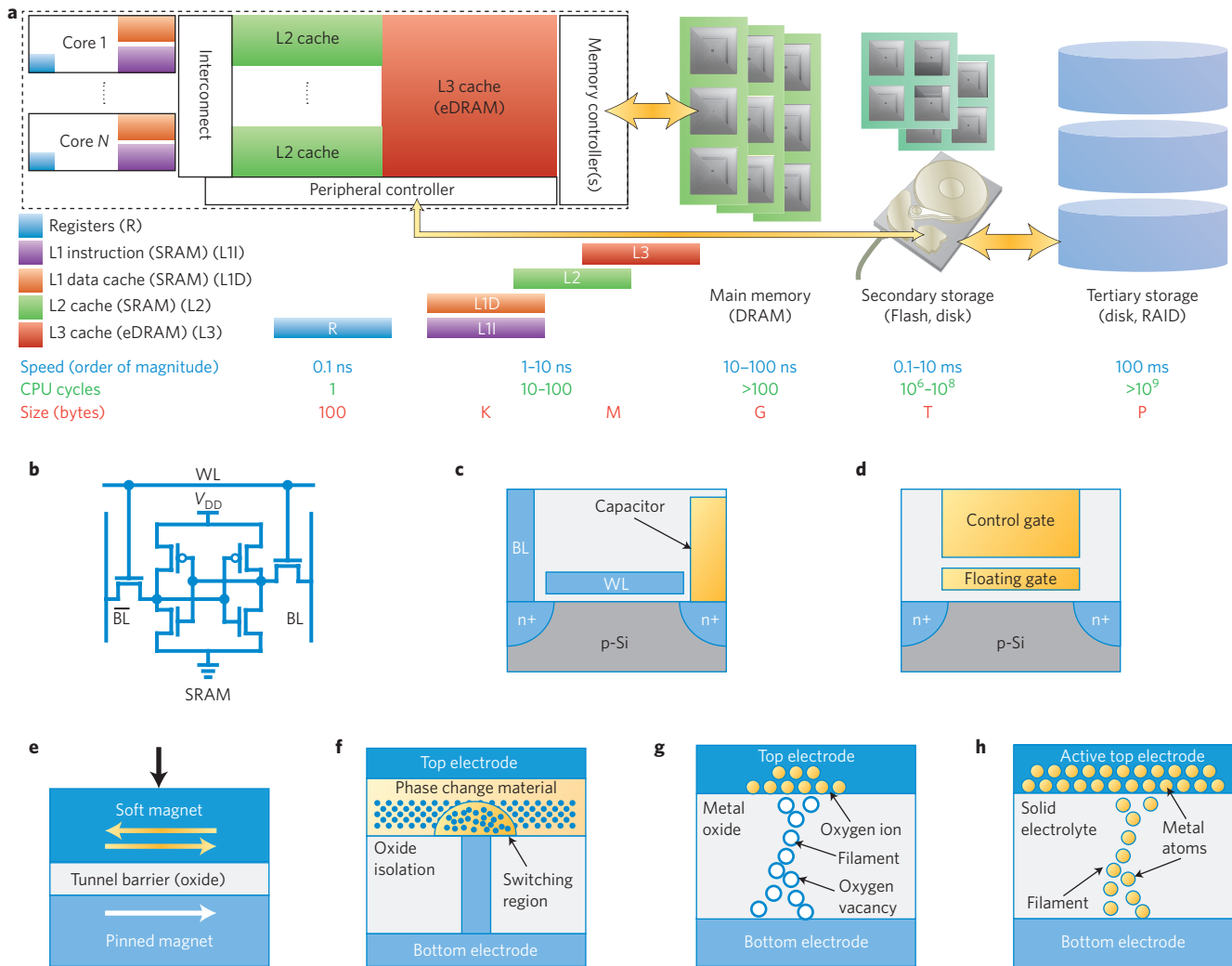


Figure 1 | Memory hierarchy and various memory types. **a**, Physical distribution and hierarchy of memory in a computer. The lower panel shows speed (as order of magnitude), number of processor (CPU) cycles needed to access the memory, and size of the different memories. RAID, redundant array of independent disks. **b**, Schematic of a circuit of a SRAM cell consisting of six transistors. BL, bit line; \overline{BL} , logic complement of BL; WL, word line; V_{DD} , supply voltage. **c-h**, Schematics of DRAM (**c**), Flash (**d**), STT-MRAM (**e**), PCM (**f**), RRAM (**g**) and CBRAM (**h**). The black downward arrow in **e** indicates the direction of current flow.

and cumbersome erase procedure. Coincidentally, these new memories store information using new types of physics that do not rely on storing charge on a capacitor as is the case for SRAM, DRAM and Flash.

Spin-transfer-torque magnetic random access memory. The basic physics of operation of STT-MRAMs is discussed by Kent and Worledge¹³. Here, we focus on the aspect of scaling and its relation to the data retention time of this technology. In a ferromagnet, spin up and spin down states are separated by a well-defined energy barrier (E_B). A loss of memory happens when a state spontaneously flips to the other by overcoming E_B , for example by absorbing energy from the thermal bath. A performance metric for STT-MRAM is therefore $\Delta = E_B/k_B T$ (k_B , Boltzmann constant; T , temperature) such

that the memory retention time is given by $\tau = \tau_0 \exp(\Delta)$. The retention time τ_0 is typically of the order of 1 ns. E_B depends on the energy stored in the magnet. To switch a ferromagnet, one needs to supply it with an amount of energy that is at least equal to E_B ; thus, within this simple scenario, the switching current — and therefore the write energy — is proportional to E_B and varies in a logarithmic fashion with the required retention time. For a memory array, statistical considerations lead to an increase in the required value of Δ (ref. 14), but the simple relation shown above provides a reasonable insight.

It is noteworthy that such a one-to-one relation between write current and retention time is not common to the other emerging technologies and provides an important design tool for STT-MRAMs. Although the traditional industry standard

is a retention time τ of around 10 years, for most modern-day RAM applications this standard is simply irrelevant. In particular, for embedded memory applications, a much smaller τ (even in the range of seconds, depending on application) could be sufficient, allowing significant reduction in the energy needed for the write operation by trading-off in reduced retention times. In embedded applications where the computing system is normally in the off state and requirements on speed are relaxed, such design optimizations — combined with an appropriate architecture that exploits the non-volatility — could lead to significant performance improvements. In addition, emerging all-spin schemes¹⁵⁻¹⁷ that combine logic and memory devices in the same structure could lead to significant energy savings and increase in data storage density.

Phase change memory. In PCM, two electrodes sandwich a chalcogenide glass that can change between the crystalline and amorphous phase on heating and cooling. Such phase changes are induced by passing a current through the material to heat it up. Information is stored in the phase of the active material; crystalline or amorphous phases have a different resistance R and correspond to the two logic states. One of the advantages of PCM¹⁸ is that materials that have been extensively studied¹⁹ and that can be mass produced can be used, for example in DVDs. The resistance ratio between the amorphous phase and the polycrystalline phase is over 100 times larger than in STT-MRAM devices. Thus, multi-bit storage might seem achievable; however, in phase change materials the amorphous intermediate resistance states drift with time, t , towards higher resistance, following a $R(t) = R_0(t/t_0)^\nu$ power-law relationship (R_0 is the resistance at t_0 , ν is material- and device-dependent), thus making it difficult to distinguish the programmed states over time. Because PCM is programmed by Joule heating, the programming current scales down with device area. Programming current of the order of microamps is shown to be possible when the device area is scaled to sizes smaller than 10 nm (ref. 20). Advances in the development of materials that are based on superlattices, and switch phase without melting, promise to push the programming current even lower²¹.

Metal oxide resistive random access memory. In devices of this type, a metal oxide is sandwiched between two metal electrodes. An applied electric field induces the creation and motion of oxygen vacancies, resulting in the formation of conductive filaments in the oxide. This changes the device resistance, which varies between high and low states. Industry has high hopes for RRAM^{22–24} because it uses materials that are common in semiconductor manufacturing. Typical metal oxides include HfO_x , TaO_x , TiO_x and AlO_x , all of which can be deposited using atomic layer deposition. Although the device concept is simple, the physics is anything but. There are controversies surrounding the shape of the conductive filament and the role of the top and the bottom electrodes. The mobility, energy and stability of the oxygen vacancies remain topics of intense study²². As a result of these open issues, projection of device reliability becomes difficult. Furthermore, RRAMs have issues of reproducibility of their electrical characteristics; there are large resistance variations not just between devices, but also between cycles of programming of the

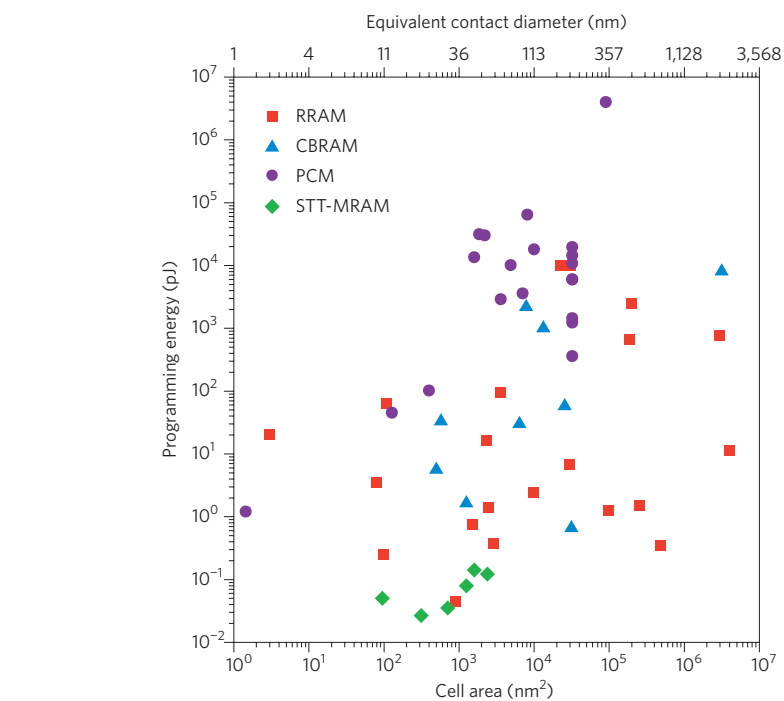


Figure 2 | A comparison of a key attribute (write energy versus device size) of emerging non-volatile memories. Data from ref. 35.

same device. This problem has been holding RRAM back from commercialization despite its many attractive features. However, the ease with which RRAM can be built in three dimensions is yet another strong incentive to develop this technology further²⁵.

Conductive bridge random access memory. This type of memory²⁶ is a close cousin of RRAM, in which the metal oxide is replaced by a solid electrolyte and one of the electrodes is a metal that can easily oxidize into metal ions. On application of an electric field, mobile metal atoms from the top (active) electrode migrate into the solid electrolyte to form conductive filaments that bridge the bottom electrode. Because of the stochastic migration process, the configuration of the conductive filament changes every time it is re-formed and results in large resistance variations similar to those in RRAM. The switching voltages of early CBRAM were too low (<0.5 V) (ref. 22), which led to poor retention. A recent work²⁷ shows significant improvement in both the switching voltage and resistance variation.

Energy, performance, scalability
Phase change memory has much better endurance ($\sim 10^9$ cycles) than Flash and can easily achieve multi-bit data storage. The writing speed of PCM in the tens of nanoseconds and its relatively large

programming current at today's feature size of tens of nanometres make it less attractive than STT-MRAM, RRAM and CBRAM. STT-MRAM excels in endurance cycling and speed, but its low resistance ratio requires a memory cell architecture that limits its device density. RRAM and CBRAM have endurance on a par with or better than PCM, and higher speed (of the order of nanoseconds), but suffer from resistance variation far worse than that of PCM and STT-MRAM. All these memories promise to scale further than Flash and DRAM (Fig. 2).

When these emerging memories were initially proposed, there was hope that one of them would be able to serve the entire memory hierarchy, meeting the need for power, energy, retention, endurance and speed required at each level; it was also envisaged that these memories could achieve high device density at low cost and be scalable for many technology generations — essentially that they could become a 'universal memory'^{28–30}. Now it is generally agreed that the vision of a 'universal memory' is not realistic. Application-driven design requires the optimization of performance at each level of the memory hierarchy, and this requires trade-offs in device characteristics that span many orders of magnitude; this is fundamentally hard to achieve by any individual device technology. For example, the need for low energy consumption during writing operations is

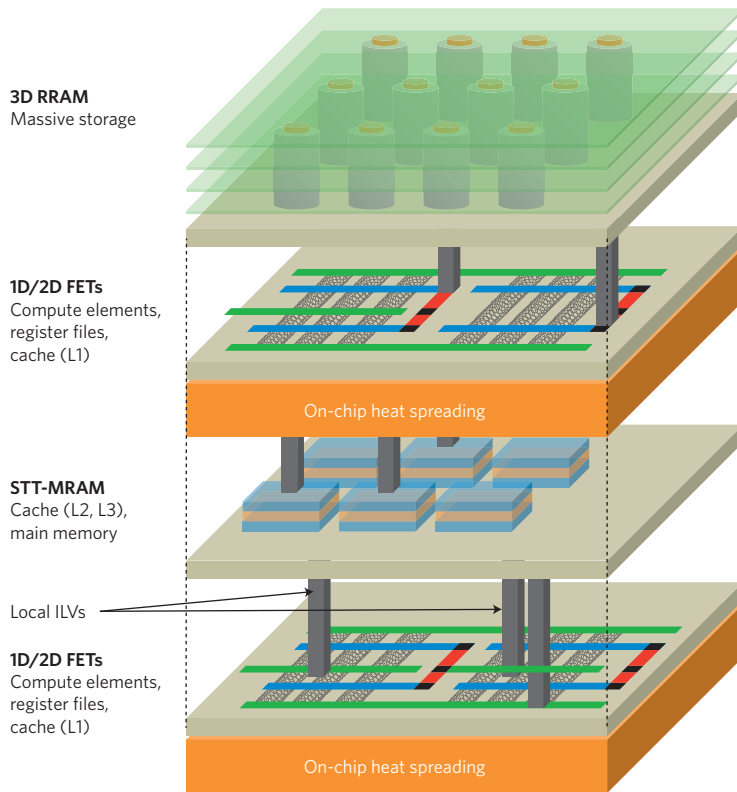


Figure 3 | Monolithic 3D integration of memory interleaved with logic computation layers. ILV, interlevel via; FET, field-effect transistor. Figure adapted with permission from ref. 36, IEEE.

accomplished by lowering the energy barrier to be overcome when erasing or writing data, which limits the non-volatile retention time of data.

Bring memory closer to computation

Shuttling data between the various levels of memory hierarchy incurs substantial latency and energy costs^{4,31}. The effectiveness of the memory hierarchy depends on data locality — data in cache memories are accessed and reused frequently. For twenty-first century applications such as big-data analytics, processing of large graphs in real time and other emerging machine-learning tasks, the memory capacity required is enormous, and the memory locations are accessed in an unpredictable order. The processor often cannot find the data needed in cache memory. Because it takes many processor (central processing unit, CPU) cycles to fetch data from the main memory (Fig. 1), the CPU may stall through unavailability of data, resulting in loss of performance and power efficiency. The lack of data locality associated with important applications as mentioned above demands a different approach from the conventional one.

New memory technologies offer a unique opportunity to bring large amounts

of memory closer to the computing elements, resulting in high-bandwidth, low-latency access. This is because all the emerging memory technologies described above can be fabricated at low temperatures and integrated monolithically in a three-dimensional (3D) chip; the relatively small thickness of these devices means that very short, high-density interlevel vias can be used to connect the memory with the computing units and the memory controller circuits (Fig. 3). Memory can also be woven into logic layers in a fine-grain fashion³². Significant improvements in energy-delay product (a measure of energy efficiency) can be achieved using this approach¹. When memory layers are interleaved and woven with the logic computation layers, long-distance data movement is eliminated. Logic layers that support this monolithic 3D vision are beginning to emerge, as 1D carbon nanotube transistors³³ and 2D layered transition metal dichalcogenide transistors³⁴ advance from prototype devices to circuit and system demonstrations.

With the ability to tightly integrate massive amounts of memory with logic, it is conceivable that future computing chips may be much more compact and energy efficient than they are today. □

H.-S. Philip Wong is in the Department of Electrical Engineering and the Stanford SystemX Alliance, Stanford University, Stanford, California 94305, USA. Sayeef Salahuddin is in the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, California 94720, USA. e-mail: hspwong@stanford.edu; sayeef@EECS.Berkeley.EDU

References

- Shulaker, M. et al. Monolithic 3D integration: a path from concept to reality <http://www.date-conference.com/conference/session/9.8> (2015).
- Kogge, P. (ed.) *ExaScale Computing Study: Technology Challenges in Achieving Exascale Systems* (DARPA Information Processing Techniques Office, 2008); <http://www.cse.nd.edu/Reports/2008/TR-2008-13.pdf>
- Itoh, K. *VLSI Memory Chip Design* (Springer, 2001).
- Borkar, S. & Chien, A. W. *Commun. Assoc. Comput. Machin.* **54**, 67–77 (2011).
- <http://www-03.ibm.com/press/us/en/pressrelease/20209.wss>
- <http://www.hgst.com/science-of-storage/about-hgst-research/innovation-timeline>
- http://www-03.ibm.com/ibm/history/exhibits/storage/storage_350.html
- Brewer, J. & Gill, M. (eds) *Nonvolatile Memory Technologies with Emphasis on Flash: A Comprehensive Guide to Understanding and Using Flash Memory Devices* (Wiley/IEEE, 2011).
- Qureshi, M. K., Srinivasan, V. & Rivers, J. A. *SIGARCH Comput. Archit. News* **37**, 24–33 (2009).
- Parkin, S. S. P. *Proc. Int. Electron Devices Meeting (IEDM)* 903–906 (2004).
- Freitas, R. F. & Wilcke, W. W. *IBM J. Res. Dev.* **52**, 439–447 (2008).
- Burr, G. W. et al. *IBM J. Res. Dev.* **52**, 449–464 (2008).
- Kent, A. D. & Worledge, D. *Nature Nanotech.* **10**, 187–191 (2015).
- Khvalkovskiy, A. V. et al. *J. Phys. D* **46**, 74001 (2013).
- Morris, D., Bromberg, D., Zhu, G.-J. & Pileggi, L. *Design and Automation Conf. (DAC), 49th ACM/EDAC/IEEE* 486–491 (2012).
- Datta, S., Salahuddin, S. & Behin-Aein, B. *Appl. Phys. Lett.* **101**, 252411 (2012).
- Bhowmik, D., You, L. & Salahuddin, S. *Nature Nanotech.* **9**, 59–63 (2014).
- Wong, H.-S. P. et al. *Proc. IEEE* **98**, 2201–2227 (2010).
- Yamada, N. et al. *Jpn. J. Appl. Phys.* **26**, 61–66 (1987).
- Liang, J., Jeyasingh, R. G. D., Chen, H.-Y. & Wong, H.-S. P. *IEEE Trans. Electron Devices* **59**, 1155–1163 (2012).
- Takaura, N. et al. *VLSI Technol. Symp. T130–T131* (2013).
- Waser, R. & Aono, M. *Nature Mater.* **6**, 833–840 (2007).
- Wong, H.-S. P. et al. *Proc. IEEE* **100**, 1951–1970 (2012).
- Kamiya, K. et al. *Phys. Rev. B* **87**, 155201 (2013).
- Prince, B. *Vertical 3D Memory Technologies* (Wiley, 2014).
- Waser, R., Dittmann, R., Staikov, G. & Szot, K. *Adv. Mater.* **21**, 2632–2663 (2009).
- Zahurak, J. et al. *Int. Electron Devices Meeting (IEDM)* 140–144, Paper 6.2 (2014).
- Lai, S. & Lowrey, T. *Int. Electron Devices Meeting (IEDM)* Paper 36.5 (2001).
- Bette, A. et al. *Digest of Technical Papers: 2003 Symp. VLSI Circuits* 217–220 (IEEE, 2003).
- Akerman, J. *Science* **308**, 508–510 (2005).
- <http://www.cccblog.org/2012/05/29/21st-century-computer-architecture>
- Zhu, Q. et al. *IEEE 23rd Int. Conf. Application-Specific Systems, Architectures Processors (ASAP)* 125–132 (2012).
- Shulaker, M. et al. *Nature* **501**, 256–530 (2013).
- Wang, H. et al. *Int. Electron Devices Meeting (IEDM)* 88–91 (2012).
- <https://nano.stanford.edu/stanford-memory-trends>
- Ebrahimi, M. S. et al. *SOI-3D-Subthreshold Microelectronics Technology Unified Conf. (S3S)* 1–2 <http://dx.doi.org/10.1109/S3S.2014.7028198> (IEEE, 2014).

Acknowledgements

The authors acknowledge support from the National Science Foundation Center for Energy Efficient Electronics Science, STARnet FAME, LEAST, and SONIC Centers, IARPA, and member companies of the Stanford Non-Volatile Memory Technology Initiative (NMTRI) and the Stanford SystemX Alliance. Discussions with S. Mitra, M. Sabry, C. Kozyrakis, K. Olukotun, L. Pileggi, F. Franchetti, J. Rabaey and J. Bokor, as well as technical assistance from our students are gratefully acknowledged.