

# The New DRAM Interfaces: SDRAM, DRDRAM and Variants

Brian Davis<sup>1</sup>, Bruce Jacob<sup>2</sup>, Trevor Mudge<sup>1</sup>

<sup>1</sup> Electrical Engineering & Computer Science, University of Michigan,  
Ann Arbor, MI 48109  
{tnm, btdavis}@eecs.umich.edu

<sup>2</sup> Electrical & Computer Engineering, University of Maryland at College Park,  
College Park, MD 20742  
blj@eng.umd.edu

**Abstract.** For the past two decades, developments in DRAM technology, the primary technology for the main memory of computers, have been directed towards increasing density. As a result 256 M-bit memory chips are now commonplace, and we can expect to see systems shipping in volume with 1 G-bit memory chips within the next two years. Although densities of DRAMs have quadrupled every 3 years, access speed has improved much less dramatically. This is in contrast to developments in processor technology where speeds have doubled nearly every two years. The resulting “memory gap” has been widely commented on. The solution to this gap until recently has been to use caches. In the past several years, DRAM manufacturers have explored new DRAM structures that could help reduce this gap and reduce the reliance on complex multilevel caches. The new structures have not changed the basic storage array that forms the core of a DRAM; the key changes are in the interfaces. This paper presents an overview of these new DRAM structures.

## 1 Introduction

For the past two decades developments in DRAM technology, the primary technology for the main memory of computers, have been directed towards increasing density. As a result, 256 M-bit memory chips are now commonplace, and we can expect to see systems shipping with 1 G-bit memory chips in volume within the next two years. Many of the volume applications for DRAM, particularly low cost PCs, will thus require only one or two chips for their primary memory. Ironically, then, the technical success of the commodity DRAM manufacturers is likely to reduce their future profits, because the fall in units shipped is unlikely to be made up by unit price increases.

Although densities of DRAMs have quadrupled every 3 years, access speed has improved much less dramatically. This is in contrast to developments in processor technology where speeds have doubled nearly every two years. The resulting “memory gap” has been widely commented on. The solution to this gap until recently has been to use

caches. In the past several years, DRAM manufacturers have explored new DRAM structures that could help reduce this gap, and reduce the reliance on complex multilevel caches. These structures are primarily new interfaces that allow for much higher bandwidths to and from chips. In some cases the new structures provide lower latencies too. These changes may also solve the problem of diminishing profits by allowing DRAM companies to charge a premium for the new parts. This paper will overview some of the more common new DRAM interfaces.

## **2 DRAM Architectures — Background**

DRAMs are currently the primary memory device of choice because they provide the lowest cost per bit and the greatest density among solid-state memory technologies. When the 8086 was introduced, DRAM were roughly matched in cycle time with microprocessors. However, as we noted, since this time processor speeds have improved at a rate of 80% annually, DRAM speeds have improved at a rate of 7% annually [1]. In order to reduce the performance impact of this gap, multiple levels of caches have been added, and processors have been designed to prefetch and tolerate latency.

The traditional asynchronous DRAM interface underwent some limited changes in response to this growing gap. Examples were fast-page-mode (FPM), extended-data-out (EDO), and burst-EDO (BEDO). Each provided faster cycle times (if accesses were from the same row) and thus more bandwidth than the predecessor.

In recent years more dramatic changes to the interface have been made. These built on the idea, first exploited in fast-page-mode devices, that each read to a DRAM actually reads a complete row of bits or word line from the DRAM core into an array of sense amps. The traditional asynchronous DRAM interface would then select through a multiplexer a small number of these bits (x1, x4, x8 being typical). This is clearly wasteful and, by clocking the interface, it is possible to serially read out the entire row or parts of the row. This describes the now popular synchronous DRAM (SDRAM). The basic two-dimensional array of bit cells that forms the core of every DRAM is unchanged, although its density can continue to undergo improvement. The clock runs much faster than the access time and, thus, the bandwidth of memory accesses are greatly increased, provided they are to the same word line.

There are now a number of improvements and variants on the basic single-data-rate (SDR) SDRAM. These include double-data-rate (DDR) SDRAM, direct Rambus DRAM (DRDRAM), and DDR2, which is under development by a number of manufacturers. DDR signalling is simply a clocking enhancement where data is driven and received on both the rising and the falling edge of the clock. DRDRAM includes high speed DDR signalling over a relatively narrow bus to reduce pin count and a high level of banking on each chip. The core array itself can be subdivided at the expense of some area. This multibanking allows for several outstanding requests to be in flight at the same time, providing an opportunity to increase bandwidth through pipelining of the unique bank requests. Additional core enhancements, which may be applied to many of these new interface specifications include Virtual Channel (VC) caching, Enhanced Memory System (EMS) caching and Fast-Cycle (FC) core pipelining. These additional improvements provide some form of caching of recent sense amp data. Even with these

significant redesigns, the cycle time — as measured by end-to-end access latency — has continued to improve at a rate significantly lower than microprocessor performance. These redesigns have been successful at improving bandwidth, but latency continues to be constrained by the area impact and cost pressures on DRAM core architectures [2].

In the next section we will give an introduction to some of the most popular new DRAMs.

### 3 The New DRAM Interfaces

The following subsections will briefly discuss a variety of DRAM architectures. The first four are shown in Table 1. The min/max access latencies given in this table

**Table 1:** Synchronous DRAM Interface Characteristics

|                     | <b>PC100</b>            | <b>DDR266<br/>(PC2100)</b> | <b>DDR2</b>             | <b>DRDRAM</b>           |
|---------------------|-------------------------|----------------------------|-------------------------|-------------------------|
| Potential Bandwidth | 0.8 GB/s                | 2.133 GB/s                 | 3.2 GB/s                | 1.6 GB/s                |
| Interface Signals   | 64(72) data<br>168 pins | 64(72) data<br>168 pins    | 64(72) data<br>184 pins | 16(18) data<br>184 pins |
| Interface Frequency | 100 MHz                 | 133 MHz                    | 200 MHz                 | 400MHz                  |
| Latency Range       | 30-90 nS                | 18.8-64 nS                 | 17.5-42.6 nS            | 35-80 nS                |

assume that the access being scheduled has no bus utilization conflicts with other accesses.

The last two subsections (ESDRAM and FCDRAM) discuss core enhancements which can be mated to almost any interface.

#### 3.1 SDR SDRAM (PC100)

The majority of desktop PC's shipped in 1Q2000 use SDR PC100 SDRAM. SDR DRAM devices are currently available at 133 Mhz (PC133). The frequency endpoint for this line of SDRAMs is in question, though PC150 and PC166 are almost certain to be developed. As we noted, SDRAM is a synchronous adaptation of the prior asynchronous FPM and EDO DRAM architectures that streams or bursts data out synchronously with a clock, provided the data is all from the same row of the core. The length of the burst is programmable up to the maximum size of the row. The clock (133 MHz in this case) typically runs nearly an order of magnitude faster than the access time of the core. As such, SDRAM is the first DRAM architecture with support for access concurrency on a single shared bus. Earlier non-synchronous DRAM had to support access concurrency via externally controlled interleaving.

### **3.2 DDR SDRAM (DDR266)**

Earlier we noted that DDR SDRAM differs from SDR in that unique data is driven and sampled at both the rising and falling edges of the clock signal. This effectively doubles the data bandwidth of the bus compared to an SDR SDRAM running at the same clock frequency. DDR266 devices are very similar to SDR SDRAM in all other characteristics. They use the same signalling technology, the same interface specification, and the same pinouts on the DIMM carriers. The JEDEC specification for DDR266 devices provides for a number of “CAS-latency” speed grades. Chipsets are currently under development for DDR266 SDRAM and are expected to reach the market in 4Q2000.

### **3.3 DDR2**

The DDR2 specification under development by the JEDEC 42.3 Future DRAM Task Group is intended to be the follow-on device specification to DDR SDRAM. While DDR2 will have a new pin-interface and signalling method (SSTL), it will leverage much of the existing engineering behind current SDRAM. The initial speed for DDR2 parts will be 200 MHz in a bussed environment, and 400Mhz in a point-to-point application, with data transitioning on both edges of the clock [3]. Beyond strictly the advancement of clock speed, DDR2 has a number of interface changes intended to enable faster clock speeds or higher bus utilization. The lower interface voltage (1.8 V), differential clocking and micro-BGA packaging are all intended to support a higher clock rate on the bus. Specifying a write latency equal to the read latency minus one ( $WL = RL - 1$ ) provides a time profile for both read and write transactions that enables easier pipelining of the two transaction types and, thus, higher bus utilization. Similarly, the addition of a programmable additive latency (AL) postpones the transmission of a CAS from the interface to the core. This is typically referred to as a “posted-CAS” transaction. This enables the RAS and CAS of a transaction to be transmitted by the controller on adjacent bus cycles. Non-zero usage of the AL parameter is best paired with a closed-page-autoprecharge controller policy, otherwise open-page-hits incur an unnecessary latency penalty. The burst length on DDR2 has been fixed at 4 data bus cycles. This is seen as a method to simplify the driver/receiver logic at the expense of heavier loading on the address signals of the DRAM bus [4]. The DDR2 specification is not finalized, but the information contained here is based upon the most recent drafts for DDR2 devices and conversations with JEDEC members.

### **3.4 Direct Rambus (DRDRAM)**

Direct Rambus DRAM (DRDRAM) devices use a 400 Mhz 3-byte-wide channel (2 for data, 1 for addresses/commands). DRDRAM devices use DDR signalling, implying a maximum bandwidth of 1.6 G-bytes/s, and these devices have many banks in relation to SDRAM devices of the same size. Each sense-amp, and thus row buffer, is shared between adjacent banks. This implies that adjacent banks cannot simultaneously maintain an open-page or maintain an open-page while a neighboring bank performs an access. The increased number of banks for a fixed address space has the result of increasing ability to pipeline accesses due to the reduced probability of sequential accesses mapping into the same bank. The sharing of sense-amps increases the row-buffer miss

rate as compared to having one open row per bank, but it reduces the cost by reducing the die area occupied by the row buffer [5].

### **3.5 ESDRAM**

A number of proposals have been made for adding a small amount of SRAM cache onto the DRAM device. Perhaps the most straightforward approach is advocated by Enhanced Memory Systems (EMS). The caching structure proposed by EMS is a single direct-mapped SRAM cache line, the same size as the DRAM row, associated with each bank. This allows the device to service accesses to the most recently accessed row, regardless of whether refresh has occurred, and enables the precharge of the DRAM array to be done in the background without affecting the contents of this row-cache. This architecture also supports a no-write-transfer mode within a series of interspersed read and write accesses. The no-write-transfer mode allows writes to occur through the sense-amps, without affecting the data currently being held in the cache-line associated with that bank [6]. This approach may be applied to any DRAM interface, PC100 interface parts are currently available and DDR2 parts have been proposed.

### **3.6 FCDRAM**

Fast Cycle DRAM (FCRAM) developed by Fujitsu is an enhancement to SDRAM which allows for faster repetitive access to a single bank. This is accomplished by dividing the array not only into multiple banks but also into small blocks within a bank. This decreases each block's access time due to reduced capacitance, and enables pipelining of requests to the same bank. Multistage pipelining of the core array hides precharge, allowing it to occur simultaneously with input-signal latching and data transfer to the output latch. FCDRAM is currently sampling in 64M-bit quantities, utilizing the JEDEC standard DDR SDRAM interface, but is hampered by a significant price premium based upon the die area overhead of this technique [7]. Fujitsu is currently sampling FCDRAM devices which utilize both SDR and DDR SDRAM interfaces. Additionally, low-power devices targeted at the notebook design space are available.

## **4 Conclusions**

DRAM are widely used because they provide a highly cost effective storage solution. While there are a number of proposals for technology to replace DRAM, such as SRAM, magnetic RAM (MRAM) [8] or optical storage [9], the new DRAM technologies remain the volatile memory of choice for the foreseeable future. Increasing the performance of DRAM by employing onboard cache, interfaces with higher utilization, or smaller banks may impact the cost of the devices, but it could also significantly increase the performance that system designers are able to extract from the primary memory system.

The DRAM industry has been very conservative about changing the structure of DRAMs in even the most minor fashion. There has been a dramatic change in this attitude in the past few years, and we are now seeing a wide variety of new organizations being offered. Whether one will prevail and create a new standard commodity part remains to be seen.

## References

- [1] W. Wulf and S. McKee. 1995 “Hitting the Memory Wall: Implications of the Obvious.” ACM Computer Architecture News. Vol 23, No. 1. March 1995.
- [2] V. Cuppu, B. Jacob, B. Davis, and T. Mudge. 1999. “A performance comparison of contemporary DRAM architectures.” In Proc. 26th Annual International Symposium on Computer Architecture (ISCA’99), pages 222–233, Atlanta GA.
- [3] JEDEC 2000. Joint Electronic Device Engineering Council; 42.3 Future DRAM Task Group. <http://www.jedec.org/>
- [4] B. Davis, T. Mudge, B. Jacob, V. Cuppu. 2000. “DDR2 and low-latency Variants.” In Proc. Solving the Memory Wall Workshop, held in conjunction with the 27th International Symposium on Computer Architecture (ISCA-2000). Vancouver BC, Canada, June 2000.
- [5] IBM. 1999. “128Mb Direct RDRAM”. International Business Machines, <http://www.chips.ibm.com/products/memory/19L3262/19L3262.pdf>
- [6] EMS. 2000. “64Mbit - Enhanced SDRAM”. Enhanced Memory Systems, [http://www.edram.com/Library/datasheets/SM2603,2604pb\\_r1.8.pdf](http://www.edram.com/Library/datasheets/SM2603,2604pb_r1.8.pdf).
- [7] Fujitsu. 2000. “8 x 256K x 32 BIT Double Data Rate (DDR) FCRAM.” Fujitsu Semiconductor, <http://www.fujitsumicro.com/memory/fcram.htm>.
- [8] “Motorola believes MRAM could replace nearly all memory technologies.” Semiconductor Business News, May 10, 2000. <http://www.semibiznews.com/story/OEG20000510S0031>
- [9] C-3D. 2000. “C3D Data Storage Technology.” Constellation 3D, <http://www.c-3d.net/tech.htm>