

ADVANCED
MICROELECTRONICS

Kiyoo Itoh

VLSI Memory Chip Design



Springer

Springer Series in
ADVANCED MICROELECTRONICS

5

Springer

*Berlin
Heidelberg
New York
Barcelona
Hong Kong
London
Milan
Paris
Singapore
Tokyo*

Physics and Astronomy

ONLINE LIBRARY

<http://www.springer.de/phys/>

Springer Series in
ADVANCED MICROELECTRONICS

Series editors: K. Itoh, T. Sakurai

The Springer Series in Advanced Microelectronics provides systematic information on all the topics relevant for the design, processing, and manufacturing of microelectronic devices. The books, each prepared by leading researchers or engineers in their fields, cover the basic and advanced aspects of topics such as wafer processing, materials, device design, device technologies, circuit design, VLSI implementation, and subsystem technology. The series forms a bridge between physics and engineering and the volumes will appeal to practicing engineers as well as research scientists.

- 1 **Cellular Neural Networks**
Chaos, Complexity and VLSI Processing
By G. Manganaro, P. Arena, and L. Fortuna
- 2 **Technology of Integrated Circuits**
By D. Widmann, H. Mader, and H. Friedrich
- 3 **Ferroelectric Memories**
By J. F. Scott
- 4 **Microwave Resonators and Filters for Wireless Communication**
Theory, Design and Application
By M. Makimoto and S. Yamashita
- 5 **VLSI Memory Chip Design**
By K. Itoh

Series homepage – <http://www.springer.de/phys/books/ssam/>

Kiyoo Itoh

VLSI Memory Chip Design

With 416 Figures and 26 Tables



Springer

Dr. Kiyoo Itoh
Hitachi Ltd., Central Research Laboratory
1-280, Higashi-Koigakubo
Kokubunji-shi
Tokyo 185-8601
Japan
e-mail: k-ito@crl.hitachi.co.jp

SEL
TK7895
M4 I876
2001

Series Editors:

Dr. Kiyoo Itoh
Hitachi Ltd., Central Research Laboratory
1-280 Higashi-Koigakubo
Kokubunji-shi
Tokyo 185-8601
Japan

Professor Takayasu Sakurai
Center for Collaborative Research
University of Tokyo
7-22-1 Roppongi, Minato-ku,
Tokyo 106-8558
Japan

Library of Congress Cataloging-in-Publication Data

Itoh, Kiyoo, 1941-
VLSI memory chip design / Kiyoo Itoh.
p. cm. -- (Springer series in advanced microelectronics ; 5)
Includes bibliographical references and index.
ISBN 3540678204 (alk. paper)
1. Semiconductor storage devices--Design and construction. 2. Integrated
circuits--Very large scale integration--Design and construction. I. Title. II. Series.

TK7895.M4 I876 2001
621.39732--dc21

00-068735

ISSN 1437-0387

ISBN 3-540-67820-4 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

Springer-Verlag Berlin Heidelberg New York
a member of BertelsmannSpringer Science+Business Media GmbH

<http://www.springer.de>

© Springer-Verlag Berlin Heidelberg 2001
Printed in Germany

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Typesetting by the author.

Data conversion and figure processing by LE-TeX Jelonek, Schmidt & Vöckler GbR, 04229 Leipzig.

Cover concept by eStudio Calmar Steinen using a background picture from Photo Studio "SONO". Courtesy of Mr. Yukio Sono, 3-18-4 Uchi-Kanda, Chiyoda-ku, Tokyo.

Cover design: *design & production* GmbH, Heidelberg

Printed on acid-free paper SPIN: 10771899 57/3144/mf - 5 4 3 2 1 0

Preface

The VLSI memory era truly began when the first production of semiconductor memory was announced by IBM and Intel in 1970. The announcement had a profound impact on my research at Hitachi Ltd., and I was forced to change fields: from magnetic thin film to semiconductor memory. This change was so exceptionally sudden and difficult, I felt like a victim of fate. Looking back, however, I realize how fortunate I was. I have witnessed an unprecedented increase in memory capacity (DRAM, for example, has had a 6-order increase in the last three decades – from the 1-Kb level in 1970 to the 1-Gb level today). I have contributed to this progress with full involvement in memory-chip development over my career. Such rapid progress would have been impossible without many of the inventions and innovative technologies, and without the effort of many talented people. Unfortunately, few systematic books on memory-chip design have been written by experts. This is a result of two factors: the difficulty of involving university professors because of rapidly changing technology requiring huge investments and development resources, and a shortage of time on the part of chip designers in industry due to severe competition in the memory-chip business. Therefore, LSI memory-chip design has been isolated from the outside, preventing a deeper understanding of the technology.

This book is based on my 30-year memory-chip (particularly DRAM) design career. In addition to memory circuits and subsystem design issues, I describe boundary issues between processes, devices, and circuits. I also attempt to systematically describe concepts that remain unclear, and discuss state-of-the-art memory-chip design. This book will be beneficial to students and engineers interested in memory-chip design, and also to process and device engineers involved in memory-chip development.

Chapter 1 describes the basics of various VLSI memory chips including DRAM, SRAM, and nonvolatile memory. Particular emphasis is paid to internal organization, operation principles and general trends in chip performance. Chapter 2 deals with the basics of RAM design and technology. The elements constituting a memory chip (MOSFETs, capacitors, and resistors), MOS memory circuits, the scaling law, and other relevant technologies are discussed. The first two chapters lay the groundwork for understanding the rest of the book. Chapter 3 focuses on DRAM chip design. After the catalog

specifications, the determinants of chip performance, the basic technologies for memory-array configuration and each of the peripheral circuits are described. Then, refreshing schemes and redundancy are explained. Chapter 4 discusses the signal-to-noise (S/N) issue in DRAM which strongly influences stable operation in the memory cell and, thus, in the chip. The relationship between memory-cell structure and its driving/sensing is explained in relation to the S/N issue. Chapter 5 describes on-chip voltage generators used for power-supply conversion. These generators are essential for power-supply standardization and stable operation. Chapter 6 discusses subsystem-memory architectures. These are increasingly important in providing wide bandwidth (i.e. throughput) for modern DRAMS. Chapters 7 and 8 describe low-power/low-voltage memory circuits, emphasizing the importance of the partial activation of multi-divided arrays, and of lowering power-supply voltage. Low voltage inevitably needs the subthreshold-current reduction which is the key to future LSI design.

I am indebted to many people including colleagues and the office administration staff members, Ms. Hosoda and Ms. Ohta, at Hitachi Ltd. They offered support, advice, and the material needed to finalize my work. Special thanks go to my wife, Kyoko. Without her continuing support and patience this book would not have been possible.

Stanford, January 2001

Kiyoo Itoh

Contents

1. An Introduction to Memory Chip Design	1
1.1 Introduction	1
1.2 The Internal Organization of Memory Chips	3
1.2.1 The Memory Cell Array	3
1.2.2 The Peripheral Circuit	5
1.2.3 The I/O Interface Circuit	6
1.3 Categories of Memory Chip	6
1.4 General Trends in DRAM Design and Technology	11
1.4.1 The History of Memory-Cell Development	11
1.4.2 The Basic Operation of The 1-T Cell	15
1.4.3 Advances in DRAM Design and Technology	19
1.5 General Trends in SRAM Design and Technology	24
1.5.1 The History of Memory-Cell Development	24
1.5.2 The Basic Operation of a SRAM Cell	26
1.5.3 Advances in SRAM Design and Technology	29
1.6 General Trends in Non-Volatile Memory Design and Technology	31
1.6.1 The History of Memory-Cell Development	31
1.6.2 The Basic Operation of Flash Memory Cells	34
1.6.3 Advances in Flash-Memory Design and Technology	46
2. The Basics of RAM Design and Technology	49
2.1 Introduction	49
2.2 Devices	49
2.2.1 MOSFETs	49
2.2.2 Capacitors	57
2.2.3 Resistors	60
2.2.4 Wiring and Wiring Materials	61
2.2.5 Silicon Substrates and CMOS Latch-Up	65
2.2.6 Other Devices	67
2.3 NMOS Static Circuits	67
2.3.1 The dc Characteristics of an Inverter	68
2.3.2 The ac Characteristics of an Inverter	70
2.3.3 The Improved NMOS Static Inverter	74

VIII Contents

2.4	NMOS Dynamic Circuits	76
2.4.1	The Dynamic Inverter	76
2.4.2	The Bootstrap Driver	77
2.5	CMOS Circuits	79
2.5.1	The dc Characteristics	80
2.5.2	The ac Characteristics	82
2.6	Basic Memory Circuits	83
2.6.1	The Inverter and the Basic Logic Gate	83
2.6.2	The Current Mirror	83
2.6.3	The Differential Amplifier	83
2.6.4	The Voltage Booster	87
2.6.5	The Level Shifter	88
2.6.6	The Ring Oscillator	88
2.6.7	The Counter	89
2.7	The Scaling Law	90
2.7.1	Constant Electric-Field Scaling	90
2.7.2	Constant Operation-Voltage Scaling	92
2.7.3	Combined Scaling	92
2.8	Lithography	93
2.9	Packaging	94
3.	DRAM Circuits	97
3.1	Introduction	97
3.1.1	High-Density Technology	98
3.1.2	High-Performance Circuits	100
3.2	The catalog Specifications of the Standard DRAM	102
3.2.1	Operational Conditions	102
3.2.2	Modes of Operation and Timing Specifications	105
3.3	The Basic Configuration and Operation of the DRAM Chip ..	110
3.3.1	Chip Configuration	110
3.3.2	Address Multiplexing	111
3.4	Fundamental Chip Technologies	113
3.4.1	A Larger Memory Capacity and Scaled-Down Devices ..	113
3.4.2	High S/N Ratio Circuits	116
3.4.3	Low Power Circuits	117
3.4.4	High-Speed Circuits	123
3.4.5	The Multidivision of a Memory Array	128
3.5	The Multidivided Data Line and Word Line	131
3.5.1	The Multidivided Data Line	132
3.5.2	The Multidivided Word Line	139
3.6	Read and Relevant Circuits	141
3.6.1	The Address Buffer	141
3.6.2	The Address Decoder	144
3.6.3	The Word Driver	147
3.6.4	The Sensing Circuit	157

3.6.5	The Common I/O-Line Relevant Circuit	167
3.6.6	The Data-Output Buffer	172
3.7	Write and Relevant Circuits	174
3.8	Refresh-Relevant Circuits	175
3.8.1	Refresh Schemes	175
3.8.2	The Extension of Data-Retention Time in Active Mode	176
3.8.3	Current Reduction Circuits in Data-Retention Mode ..	176
3.9	Redundancy Techniques	178
3.9.1	Issues for Large-Memory-Capacity Chips	184
3.9.2	Intra-Subarray Replacement Redundancy	185
3.9.3	Inter-Subarray Replacement Redundancy	189
3.9.4	The Repair of dc-Characteristics Faults	191
3.10	On-Chip Testing Circuits	192
4.	High Signal-to-Noise Ratio	
	DRAM Design and Technology	195
4.1	Introduction	195
4.2	Trends in High S/N Ratio Design	195
4.2.1	The Signal Charge	197
4.2.2	Leakage Charge	204
4.2.3	The Soft-Error Critical Charge	208
4.2.4	The Data-Line Noise Charge	210
4.3	Data-Line Noise Reduction	210
4.3.1	Noise Sources and Their Reduction	210
4.3.2	Word-Line Drive Noise	213
4.3.3	Data-Line and Sense-Amplifier Imbalances	217
4.3.4	Word-Line to Data-Line Coupling Noise	230
4.3.5	Data-Line Interference Noise	237
4.3.6	Power-Supply Voltage Bounce	240
4.3.7	Variation in the Reference Voltage	241
4.3.8	Other Noises	244
4.4	Summary	247
5.	On-Chip Voltage Generators	249
5.1	Introduction	249
5.2	The Substrate-Bias Voltage (V_{BB}) Generator	251
5.2.1	The Roles of the V_{BB} generator	251
5.2.2	Basic Operation and Design Issues	256
5.2.3	Power-On Characteristics	258
5.2.4	Characteristics in the High- V_{DD} Region	264
5.2.5	The V_{BB} Bump	266
5.2.6	Substrate-Current Generation	269
5.2.7	Triple-Well Structures	272
5.2.8	Low-Power V_{BB} Generators	273

5.3	The Voltage Up-Converter	276
5.3.1	The Roles of the Voltage Up-Converter	276
5.3.2	Design Approaches and Issues	278
5.3.3	High Boost-Ratio Converters	283
5.3.4	Low-Power, High Supply Current Converters	285
5.4	The Voltage Down-Converter	290
5.4.1	The Roles of the Voltage Down-Converter	290
5.4.2	The Negative-Feedback Converter and Design Issues ..	293
5.4.3	Optimum Design	297
5.4.4	Phase Compensation	301
5.4.5	Reference-Voltage Generators	316
5.4.6	Burn-In Test Circuits	323
5.4.7	Voltage Trimming	327
5.4.8	Low-Power Circuits	329
5.5	The Half- V_{DD} Generator	332
5.6	Examples of Advanced On-Chip Voltage Generators	333
6.	High-Performance Subsystem Memories	339
6.1	Introduction	339
6.2	Hierarchical Memory Systems	341
6.2.1	Memory Hierarchy	341
6.2.2	Improvements in Memory-Subsystem Performance	344
6.2.3	Memory-Chip Performance	349
6.3	Memory-Subsystem Technologies	354
6.3.1	Wide-Bit I/O Chip Configurations	354
6.3.2	Parallel Operation of Multidivided Arrays	354
6.3.3	Multibank Interleaving	357
6.3.4	Synchronous Operation	358
6.3.5	Pipeline/Prefetch Operations	362
6.3.6	High-Speed Clocking Schemes	363
6.3.7	Terminated I/O Interfaces	363
6.3.8	High-Density Packaging	364
6.4	High-Performance Standard DRAMs	365
6.4.1	Trends in Chip Development	365
6.4.2	Synchronous DRAM	368
6.4.3	Rambus DRAM	380
6.5	Embedded Memories	383
7.	Low-Power Memory Circuits	389
7.1	Introduction	389
7.2	Sources and Reduction of Power Dissipation in a RAM Subsystem	392
7.2.1	Wide-Bit I/O Chip Configuration	393
7.2.2	Small Package	394
7.2.3	The Low-Voltage Data-Bus Interface	396

7.3	Sources of Power Dissipation in the RAM Chip.....	402
7.3.1	Active Power Sources	402
7.3.2	Data-Retention Power Sources	405
7.4	Low-Power DRAM Circuits	406
7.4.1	Active Power Reduction	406
7.4.2	Data-Retention Power Reduction	412
7.5	Low-Power SRAM Circuits	413
7.5.1	Active Power Reduction	413
7.5.2	Data-Retention Power Reduction	423
8.	Ultra-Low-Voltage Memory Circuits	425
8.1	Introduction	425
8.2	Design Issues for Ultra-Low-Voltage RAM Circuits	426
8.2.1	Reduction of the Subthreshold Current	426
8.2.2	Stable Memory-Cell Operation	432
8.2.3	Suppression of, or Compensation for, Design Parameter Variations	433
8.2.4	Power-Supply Standardization	435
8.3	Ultra-Low-Voltage DRAM Circuits	437
8.3.1	Gate Boosting Circuit	439
8.3.2	The Multi- V_T Circuit	440
8.3.3	The Gate-Source Back-Biasing Circuit	442
8.3.4	The Well Control Circuit	456
8.3.5	The Source Control Circuit	461
8.3.6	The Well and Source Control Circuit	462
8.4	Ultra-Low-Voltage SRAM Circuits	463
8.5	Ultra-Low-Voltage SOI Circuits	466
	References	473
	Index	489

SRAM Speed. Due to the possession of a gain, the SRAM cell allows the small signal voltage to be directly transmitted to the main amplifier on the I/O line without imposing any restriction on YL activation timing. Moreover, the non-destructive read-out characteristics of SRAM cell allow write data to be inputted to the selected data line at the earliest timing, eliminating the rewrite operation that is necessary for DRAM. Even so, the resultant high voltage swing on the selected data line does not destroy small read voltages on the adjacent data lines, despite capacitive coupling, since the small voltages are static ones and thus are immune to various noises, unlike the floating ones of the DRAM. Thus, the necessary word pulse can be shortened, enabling fast row access and cycle times. If the sum of the data-replacing time involved in a write operation and the equalizing time on the data lines is equal to the sum of the column delay (T_C) and the equalizing time on the I/O lines, as in usual SRAM designs, the row access time is equal to the row cycle time. The column speeds are faster than the row speed, as in the DRAM. If multi-stage pipelining is used, even the cycle time faster than the access time is achievable [6.11].

6.3 Memory-Subsystem Technologies

It is obvious that wide-bit I/O chip configurations combined with parallel operation of multidivided arrays increases the throughput. Moreover, to improve the traditional DRAM performance discussed in Chap. 3 modern DRAMs incorporate memory-subsystem technologies such as multibank interleaving, synchronous operation with a latch function, pipeline/prefetch operations, high-speed clocking schemes, and terminated interfaces combined with high-density packaging. These technologies are supported by command operations, on-chip mode registers, and packet protocols.

6.3.1 Wide-Bit I/O Chip Configurations

These configurations [6.4, 6.17] offer high throughput as well as ease of use, which is realized by reducing the chip count needed by the system and by adding flexible add-on memory capability. Despite the possibility of at least 256 b organization, compared to the 32 b organization for current experimental 1 Gb chips [6.4, 6.17], the number of I/O pins is eventually restricted by the following drawbacks: the chip power increases rapidly with an increase in the pin count, because the number of simultaneously charged and discharged DLs (i.e. m in the logical array) increases. The chip area also increases due to an increase in the I/O relevant circuits. The details are given in Chap. 7.

6.3.2 Parallel Operation of Multidivided Arrays

The concept of the multidivided array [6.4, 6.17], a combination of multidivided data lines (DLs) and multidivided word lines (WLs), shown in Fig. 6.11,

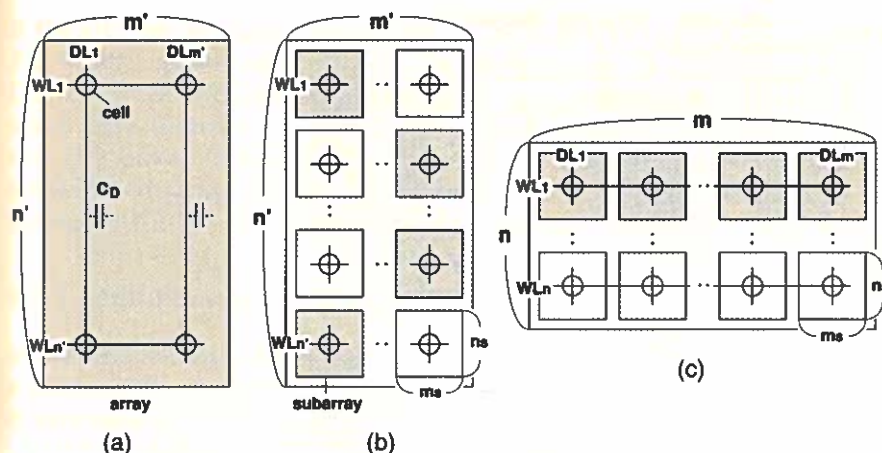


Fig. 6.11. The concept of a multidivided RAM array [6.14, 6.17]: the shading denotes the activated area. $M = n'm' = nm$. (a) A non-divided array; (b) a multidivided array; (c) a logical array

is the key to designing a high-performance RAM. The division of a DL and its partial activation dramatically increases the inherently small signal voltage and reduces the DL power dissipation that dominates the total chip power. The division of a WL is also essential to improve the ever-increasing WL delay with increasing memory capacity. A multidivided array realizes the high performance of a resulting subarray, if low-resistivity multilevel metal wiring and high-speed subarray-selection circuits are adopted. Multilevel metal wiring also minimizes the additional increases in area at the divisions. Any combination of a number of subarrays could be simultaneously activated, since each subarray could be randomly accessed. The parallel operation capability of subarrays enables multibank interleaving, if each bank in a memory system is asserted by each subarray. However, for DRAMs the number of simultaneously activated subarrays is restricted by the DL power dissipation and the maximum refresh time (the data-retention time of the cell), t_{REFmax} , which is specified in the catalog for the chip. The activation for the complicated physical array is simplified by using the logical array comprising n virtual word lines shown in Fig. 6.11c. Here, n is the number of refresh cycles in the catalog specification, which are usually distributed within t_{REFmax} ; and m is the number of simultaneously activated DLs, taking m_s to M/n_s , where M , n_s , and m_s are the memory capacity of the chip, the number of sub-WLs, and the number of sub-DLs in a subarray, respectively. Here, n_s is less than 1k in the megabit era, which is determined by the minimum signal voltage for a successful sensing, while m_s is 256 or 512 in terms of the WL delay. Figure 6.12a shows a more detailed multidivided array, widely used in multimegabit DRAM products, featuring the DL orthogonally aligned to the internal I/O line. In an actual design, in addition to the shared Y decoders

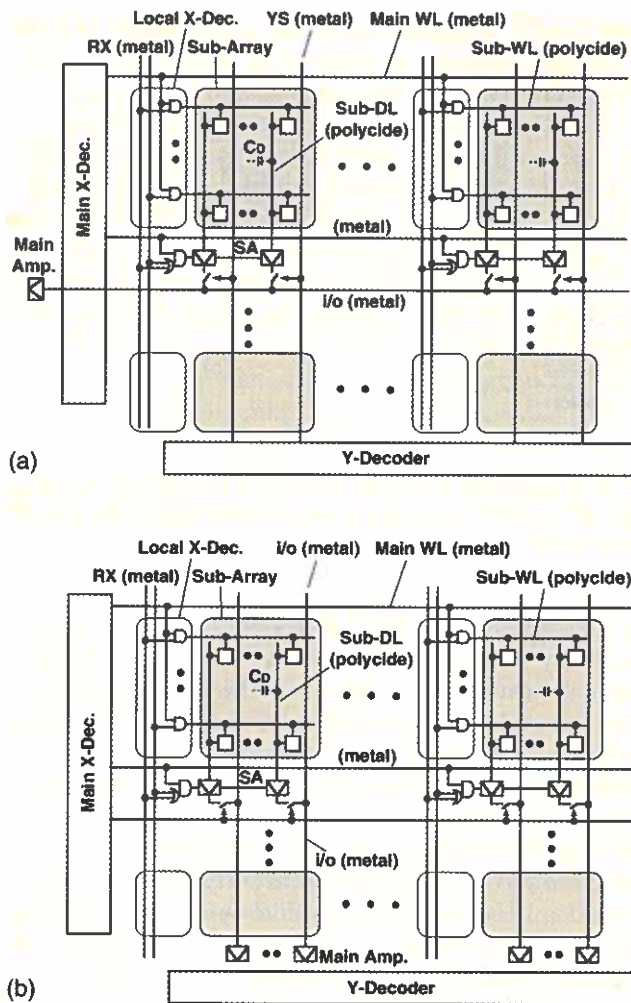


Fig. 6.12. An actual multidivided DRAM array [6.4, 6.17]: the paired-line arrangement is actually applied for each DL and each i/o line. (a) DL orthogonal to i/o line; (b) DL parallel to i/o line

shown in the figure, a combination of a shared sense amplifier (SA) and a shared I/O is used to further reduce the DL charging capacitance, as well as the chip area. The details are discussed in Chap. 3.

The DRAM array is a built-in structure capable of a massively parallel operation at the expense of large DL power. This stems from the refresh operation requirement of the 1T cell, that needs simultaneous activation of all the cells along the selected WL. A multidivided array further increases the available number of data, with at least one unit of data from each subarray.

If the I/O line is arranged in parallel to the corresponding DL, as shown in Fig. 6.12b, the available number is maximized by one unit of data from each DL. Thus, the DRAM array inherently favors embedded DRAM designs, in which high throughput is the first priority, although increases in the DL and I/O power become serious concerns. To achieve high throughput of multidivided arrays while reducing the area, a large logic-gate block in an embedded DRAM chip will require additional layers of metal wiring.

6.3.3 Multibank Interleaving

Multibank interleaving has been widely used to increase the throughput with a substantially parallel operation of the multibank. In this scheme, a memory system consists of N banks, which are sequentially addressed from bank 1 to bank N , as shown in Fig. 6.13. Each bank is composed of a memory module using many DRAM chips, so that the memory bus has an N_{MB} -bit data-bus (i.e. I/O) width. When N words, each of which comprises N_{MB} bits, are

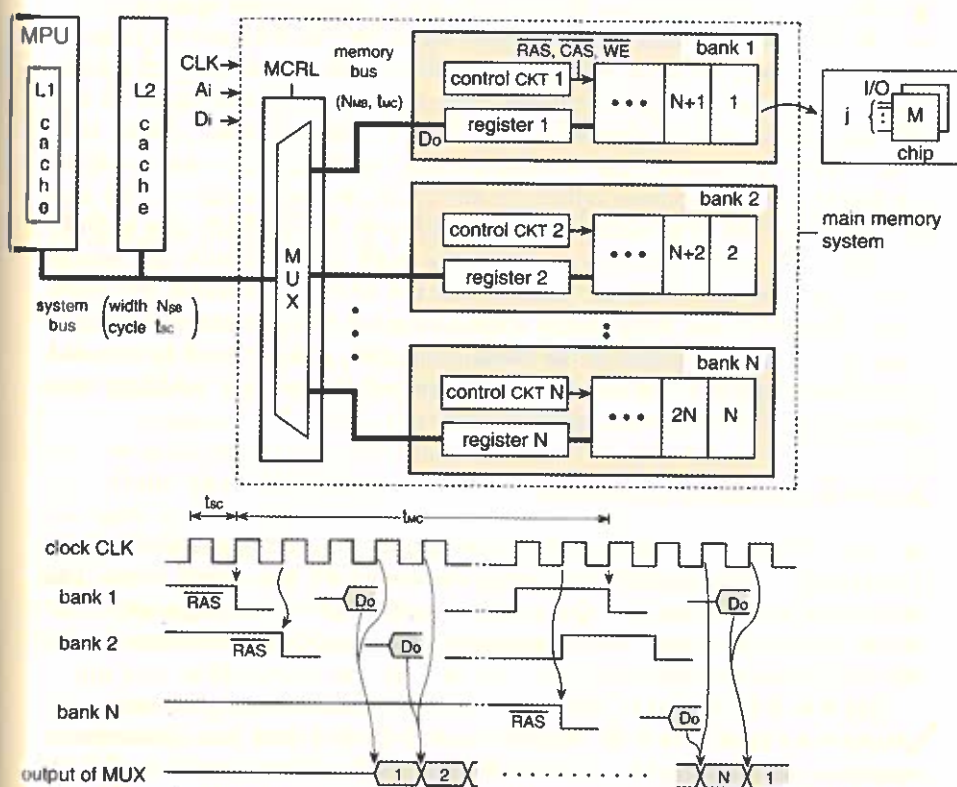


Fig. 6.13. Conventional multibank interleaving [6.8]. MUX, multiplexer; MCRL, memory controller

successively read, the resulting N_{MB} -bit data are outputted on the memory bus of each bank, and stored in each register so as to be able to wait for the next accessing. As a result of multiplexing of data from each bank, N_{MB} -bit data are available every system-clock cycle t_{sc} on the system bus. Since it continues for the memory cycle time (t_{MC}) of bank, the throughput is increased from N_{MB}/t_{MC} to $N_{MB}N/t_{MC}$, with an N -fold increase. However, this approach causes an increase in the minimum add-on memory capacity, which is expressed as $MN_{MB}N/j$, where M and j are the memory capacity of the chip and the I/O pin count of chip, respectively. Moreover, it enables an increase in the number of bus lines and relevant devices with $N_{MB}N$, preventing flexible design, miniaturization, and a low cost for the whole system. Increasing j is beneficial to increasing the throughput for a fixed add-on memory capacity, or to reducing the minimum add-on memory capacity for a fixed throughput. However, an excessively large j causes increased power dissipation, chip area, and package size, as discussed in Chap. 7. It would also degrade the chip speed, with increased noise at the I/O pins.

The multibank interleaving that memory-system designers have taken for granted can be implemented on one chip, if the multidivided array structure is utilized. This is because the subarray shown in Fig. 6.11 can be regarded as a bank, if each subarray equips its own peripheral circuit and each address buffer of the bank, and thus it can latch the input address signal so that different banks can be successively selected at the minimum system clock cycle. When a certain bank is selected, the corresponding address signals are latched at the address buffers of the bank, so that address input lines are ready for the next addressing for a different bank. While the memory operation of the succeeding circuits in the bank proceeds with the latched address signals, the different bank is selected with the different addresses. Thus, the MPU can successively access different banks without having to wait by the memory cycle of the bank. Moreover, while a bank is accessed, other banks could perform precharge or refresh operations, enabling these inherent DRAM operations to be hidden.

6.3.4 Synchronous Operation

In this scheme, all of the input/output signals of chip are latched at the I/O interface circuits of the chip, synchronously with the system clock. This scheme not only allows chip designers to incorporate various high-speed functions, but it also allows system designers to improve the system speed, with an easier timing on the board.

Figures 6.14 and 6.15 show the concepts behind asynchronous- and synchronous operation [6.8], assuming a non-divided data line. Synchronous operation has been used for modern DRAMs, such as the synchronous DRAM (SDRAM), The double-data-rate (DDR) DRAM, and the Rambus DRAM, while asynchronous operation has been used for traditional DRAMs, such as the page mode, nibble mode, static column mode, and extended-data-out

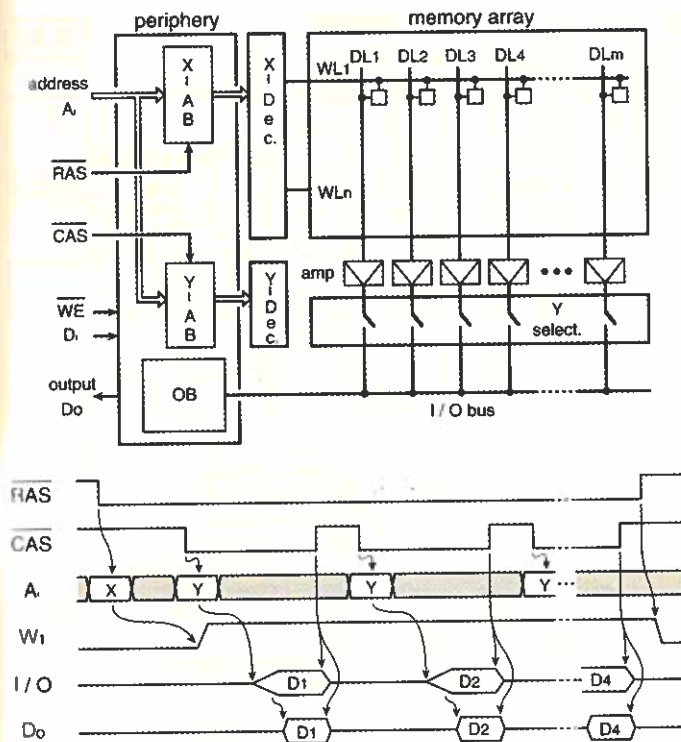


Fig. 6.14. The asynchronous operation of a DRAM chip [6.8]. AB, address buffers; OB, data-output buffer; X, row; Y, column

(EDO) DRAMs. In synchronous operation, a row address strobe signal $RS1$ is generated from the chip-select signal CS at the rising edge of the system clock CLK , so that row addresses are strobed and the corresponding word line (for example, WL_1) is activated. The resulting cell signals on m data lines are amplified in the usual manner. On the other hand, a column address strobe signal $CS1$ is generated by the next CLK , so that the column addresses are strobed and the amplified signal on the corresponding data line (for example, DL_1) is outputted on the common I/O bus line. Then, the data output \overline{D}_o is available from the chip synchronously with the succeeding CLK . The synchronous operation widens the timing margins between the internal control signals. Easier timing designs due to the use of simple latch circuits that work synchronously with the system clock are responsible to the wider margins. On the other hand, in asynchronous operation the internal timing designs are complicated, because they are closely related to many sophisticated set-up/hold timing specifications between external input/output signals. Moreover, the latch function achieves high throughput with the resultant pipeline operation (discussed later) and multibank interleaving. The

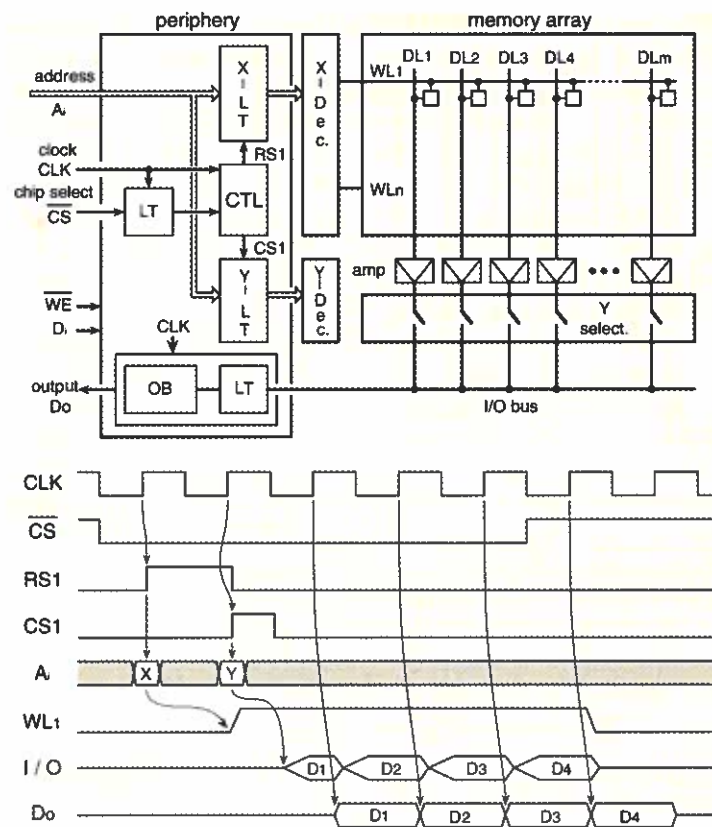


Fig. 6.15. The synchronous operation of a DRAM chip [6.8]. LT, latch; CTL, control circuit; OB, data-output buffer

synchronous operation can eliminate the $\overline{\text{RAS}}$ and $\overline{\text{CAS}}$ functions that are familiar in asynchronous operation, since the same functions are carried out by the internal signals RS1 and CS1. The operation also eliminates long hold-signals such as $\overline{\text{CS}}$ in Fig. 6.15 if the command instruction scheme is adopted. Various operation modes are set by the combination of command signals whose pulse widths are almost equal to that of the system clock, as seen in SDRAMs (see Fig. 6.32).

The synchronous operation of DRAM chips also increases the system speed. The asynchronous operation in Fig. 6.16, in which many control signals are generated at a timing generator (similar to control circuits in Fig. 6.13) by using the system clock, makes it complicated to synchronize address signals with the control signals and/or the system clock. For example, although various signal skews exist as a result of running on the memory board of the bank, the minimum timing specifications of the DRAM chip regarding the row-address set-up time t_{ASR} , the row-address hold time t_{RAH} , the column-

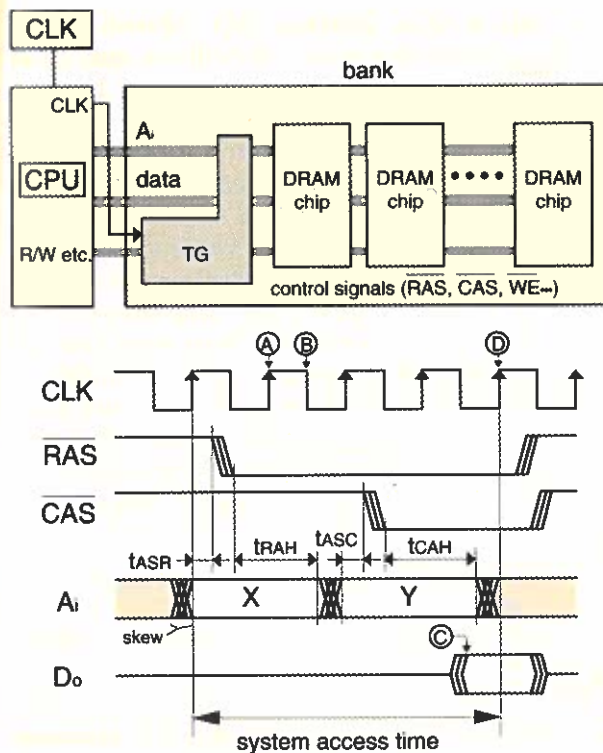


Fig. 6.16. The asynchronous operation of the memory system [6.8]. TG, timing generator; X, row; Y, column

address set-up time t_{ASC} , and the column-address hold time t_{CAH} must be ensured. Obviously, four kinds of skews are added to the minimum cycle time achieved by the chip itself, degrading the memory-bank cycle time. A timing mismatch between the control signals and the system clock further degrades the cycle time. For example, even if t_{RAH} exists between timings A and B of the clock and thus addresses could be switched to the column ones, the address switching must take place at the falling edge B. In addition, even if the chip can output the data at timing C, the bank actually outputs it at timing D. Such redundancies in timing fail to fully bring out the high speed of the chip. On the contrary, the synchronous operation shown in Fig. 6.17 allows the memory access time and cycle time to be expressed as an integer multiple of the system clock cycle, and all the relevant signals to be synchronized with only one system clock. Thus, it is enough for system designers to pay attention only to the set-up and hold time for one clock, without paying any attention to the phase relations between various signals. System designers are thus released from a complicated timing design caused

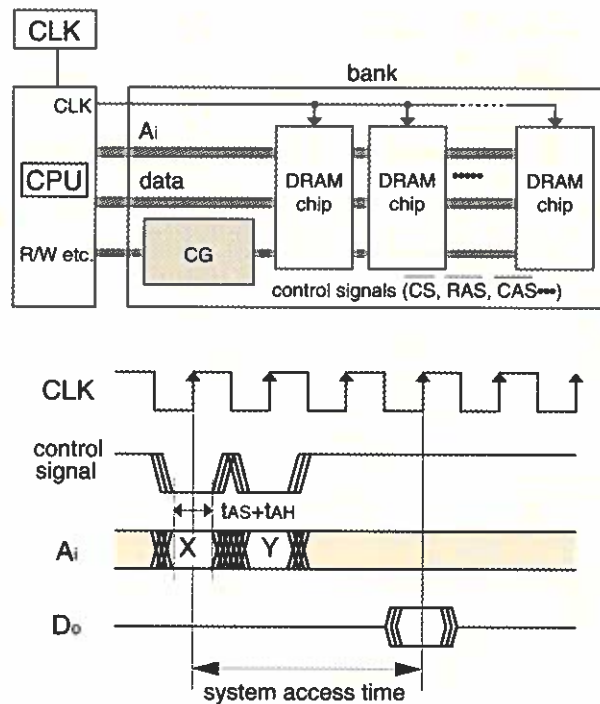


Fig. 6.17. The synchronous operation of the memory system [6.8]. CG, command generator

by asynchronous memory operation, and they can design the memory system just as they design a conventional logic system.

6.3.5 Pipeline/Prefetch Operations

A pipeline operation, which features the parallel operation of a number of circuit blocks combined with latch circuits, has been proposed to increase the throughput of an asynchronous operation DRAM [6.12]. The combination of a pipeline operation and synchronous operation, however, further enhances the throughput. An application of pipeline operation with a latch circuit to the I/O line, as shown in Fig. 6.15, is especially important. It shortens the cycle time of the I/O relevant circuit, which is usually slow in speed because of its inherently heavy capacitance. If the latch circuit (LT) on the common I/O bus line latches the read data from a memory array, and then isolates the latched data from the I/O line, the I/O line can be used for other operations, such as precharge and equalizing operations, while the latched data is processed by the succeeding circuit. If the column-selection circuits (Y select. in the figure) have shift-register and ring-counter functions, a string of m data is sequentially output on the I/O line by inputting only the first addresses to the

column decoder. This operating mode is called the burst mode. In practice, m -bit data are divided into a number of units, each of which consists of a few bits and is independently operable. Thus, any bit within any unit can be accessed, followed by the sequential data of the remaining bits of the unit. The throughput of the burst mode is higher than that of the page mode or the nibble mode discussed in Chap. 3. This is because the page mode, for example, needs not only timing margins between the external clock $\overline{\text{CAS}}$ and addresses (as shown in Fig. 6.14), but also I/O-line precharging. Note that the maximum number of data bits available with one burst mode (i.e. one word-line selection) is m . When a different word line is successively selected, the next burst mode is interrupted by data-line precharging.

Pipeline and/or parallel operations, combined with synchronous operation using the external system clock, offer a faster column cycle. For example, a prefetch scheme doubles the throughput, if two bits from memory array are simultaneously fetched to the data output buffer, and are outputted through a high-speed, parallel-to-serial operation at the buffer.

6.3.6 High-Speed Clocking Schemes

The high throughput of modern DRAMs is supported by high-speed clocking schemes such as command operations, on-chip mode registers, and packet protocols. Command operations allow all of the operating modes to be designated by combinations of two logic levels of command signals, whose pulse widths are almost equal to that of the system clock, as will be explained with an SDRAM. The resulting simplified timing design enables the high speed. On-chip mode registers defines how the DRAM chip operates. For instance, they define column latency and burst length, as explained later. This information is programmed into the registers by users, so that the chip operates based on this information. Packet protocols feature a request packet of addresses, data, and controls transmitted on the system bus, as will be discussed below for a Rambus DRAM.

6.3.7 Terminated I/O Interfaces

Small-swing impedance-matched (terminated) interfaces provide high-speed performance with the reduction of the high switching I/O power, especially in the high-speed region, of the TTL or LVTTTL interfaces. Rambus, SSTL (Stub Series Terminated Logic), and HSTL (High-Speed Transceiver Logic) interfaces are good examples. Protocol architectures such as Rambus and SyncLink, which have an extremely high clocking scheme supported by high-density packaging, partly solve the problems caused by the simply widened bus architecture. The interfaces will be explained in more detail below.

6.3.8 High-Density Packaging

The rapid and remarkable progress in high-density chip packaging can be summarized into two categories: achievements involving better physical characteristics (thinner, smaller, and lighter) and those involving a higher pin count. As a result, the LOC (Lead-On-Chip) package (Fig. 2.46) can accommodate ever-enlarged chips in a small package. The ball grid array (BGA) and the chip-size package maximize the pin count while reducing the lead-frame inductance and the I/O loading capacitance, which are essential to limit the ever-increasing area, power, and noise of the chip. Multichip module packaging (MCM) further improves the system performance. The data transfer rate of standard DRAM chips dramatically increases with low power, if high-density packaging that keeps up with high-density chip technology, and low-voltage operation, are applied, as shown in Fig. 6.18. The I/O wiring on the module is short enough to eliminate the termination with reduced capacitance. MCM packaging will make standard DRAMs competitive with embedded DRAMs in the gigabit era, while keeping the low-cost advantage of standard DRAMs.

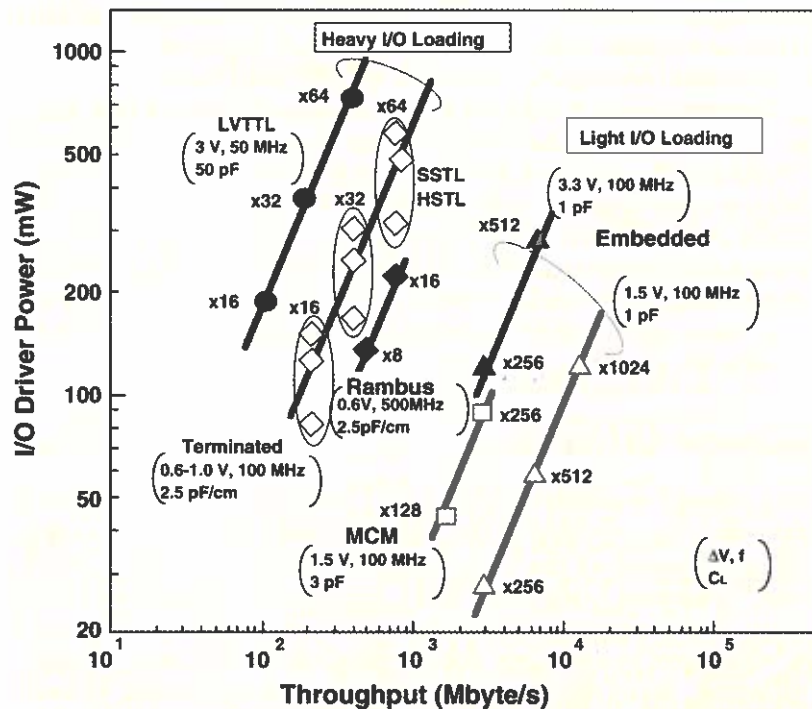


Fig. 6.18. The data transfer rate of a DRAM chip [6.4, 6.17]. x_j , The I/O pin count of a DRAM chip; C_L , the loading capacitance of each I/O line, without termination power