



ADVANCED POWER REDUCTION TECHNIQUE FOR DRAM USING LOSSLESS COMPRESSION

ALPHA J.S and ANAND S,
PG student, Assistant professor
alphaece17@gmail.com, anands5@ymail.com,
immanuelarasar JJ college of engineering

ABSTRACT

Dynamic random access memory (DRAM) has served as the major role in computer systems and mobile phones. In DRAM the electrical charges in the form of storage capacitors so the refresh operation increases in proportion to the memory capacity. We propose a new method to reduce the memory capacity by using a Bit mass compression technique. In this the bitstream compression is important in reconfigurable system design since it reduces the bitstream size and the memory requirement. It also improves the communication bandwidth and thereby decreases the reconfiguration time. Existing research in this field has explored two directions: efficient compression with slow decompression or fast decompression at the cost of compression efficiency. This paper proposes a decoding compression technique to improve both compression and decompression efficiencies. The three major contributions of this paper are: 1) smart placement of compressed bitstreams that can significantly decrease the overhead of decompression engine; 2) selection of profitable parameters for bitstream compression; and 3) efficient combination of bitmask-based compression and run length encoding of repetitive patterns. Our proposed technique outperforms the existing compression approaches by 15%, while our decompression hardware for variable-length coding is capable of operating at the speed closest to the best known field-programmable gate array-based decoder for fixed-length coding.

Index Terms—Dynamic random access memory (DRAM), bitstream compression, variable-length coding, field-programmable gate array-based decoder

I INTRODUCTION

Dynamic random access memory (DRAM) has served as the main role of storage in computer systems including high performance systems, personal computers, and mobile phones for more than 30 years. The memory capacity of a DRAM has increased to meet system demands, supported by the development of a semi-conductor process technology that follows Moore's law, whereby the number of elements on a fixed silicon die doubles every 18 months. In fact, needs from the system side have included not only memory capacity, but also data transfer speed, operation current reduction, and standby current reduction. In a DRAM, each bit is stored as an amount of electrical charge in a storage capacitor, and

the increase in memory capacity has directly caused two problems: disturbance and power consumption.

Both of these problems are attributed to the rewrite operation of a memory cell associated with a finite data retention time. In particular, standby power consumption has become one of the most serious problems for using a DRAM in mobile applications. Following the constant-electric-field scaling theory, the voltages should be lowered at the same rate by which the dimensions are reduced. The rewrite voltages are, however, saturated by the difficulty associated with the read operation. This means that the standby power increases with the memory capacity. A new method to reduce the refresh current in a DRAM by extending the retention time effectively when the amount of the data to be stored is small. We call this low-power mode as the partial access mode (PAM). The retention time has been shown to exhibit both tail and

nd main distributions [6]. Most of the cells belong to the main distribution and have retention times significantly higher than the product specification. Only a minor portion suffers from increased leakage.

Although the active power increases by a factor of 2^N , the refresh time increases by more than $2N$ as a consequence of the fact that the majority decision does better than averaging for the tail distribution of retention time. The conversion can be realized very simply from the structure of the DRAM array circuit. This method can reduce the frequency of the disturbance and its power consumption by two orders of magnitude. The proposed DRAM is fully compatible with a conventional DRAM. In its usual operating mode, the full memory capacity is used. In the PAM, the capacity is limited to 2^{-N} of the total capacity; however, memory cells are fully used to share the storage charge to extend the retention time.

This paper is organized as follows: in Section II, the refresh operation and retention time are explained on the basis of the measured data of the fabricated DRAM. Then, conventional methods to reduce the refresh current are examined. In Section III, we propose a PAM and describe the difference between our method and the conventional partial array self refresh method. In Section

II. DRAM OPERATION AND POWER-REDUCTION MODE

A. DRAM Refresh Operation

The DRAM memory capacity has been increasing, even though its die size has almost remained constant, as listed in Table I. In 2011, a 2-Gb DRAM was fabricated by a 30-nm [minimum feature size (F) value] process and was placed on the market. A DRAM stores a single bit in a memory cell as an amount of electrical charge on a storage capacitor. Charge is lost by the leakage current of the p-n junction, sub-threshold current, and gate-induced drain leakage (GIDL) [26]. This means that a DRAM requires a rewrite operation before the memory cell loses its storage charge. This rewrite operation is called refresh. Refresh is performed by issuing an auto-refresh command (AREF). Because refresh is a type of disturbance in the system where sense amplifier (SA) activation, pre-charging, and read or write operations are forbidden, the frequency of the AREF command should be minimized.

In this paper the refresh operation using a 256-Mb DRAM as an example is used. The refreshing of all memory cells is completed by issuing $2^{13} = 8192$ AREF commands, called an 8-K AREF operation. The data retention time of each memory cell is expressed by t_{ret} . The minimum retention time of the memory cells, $t_{ret,min}$, during which all memory cells maintain their own charges, should be longer than 64 ms, standardized by the cell refresh time t_{ref} [27]. Thus, the maximum time interval of an AREF

IV, we examine the distribution of the retention time of $2N$ cells/bit statistically and show that the tail distribution can be eliminated when $N \geq 2$. For $2N$ cells/bit, the cell signal is determined not by the average, but by the majority rule of 2^N cells.

TABLE I. DRAM COMPARISON

Total capacity	256-Mb	1-Gb	2-Gb
BANK structure	4BANKs	8BANKs	8BANKs
BL length	512	512	512
WL length	512	512	512
AREF	8-K	8-K	8-K
X address	X0-X12	X0-X13	X0-X14
Y address	Y0-Y9	Y0-Y9	Y0-Y9
Data	X8	X8	X8
Page size	8K	8K	8K
VDD	2.5V	1.5V	1.5V

command, t_{REF} , is $64 \text{ ms} / 2^{13} = 7.8 \mu\text{s}$ in an 8-K AREF operation.

The 256-Mb DRAM in Table I consists of four banks, each of which is composed of 16×16 mats. One mat is 256 kb, including 512 word-lines (WLs) and 512 bit-lines (BLs) as shown in Fig. 1(a) [28], [29]. In an 8-K AREF operation, a single AREF command is accomplished by one refresh operation for all four banks. In each bank, one refresh operation is applied to 16 mats located in the WL direction. The corresponding 16 WLs are selected at the same time, and memory node voltages are read to the BLs through transfer NMOSFETs. These signals are amplified by sense amplifiers that perform rewrite operations of an 8-K memory cells. Once an SA is activated, 16 WLs allocated on the same X address of the 16 mats in one bank are selected at once, and data are read to an 8-K BLs. This number, 8-K, is called the page size and shows the maximum data size written or read during one SA activation. If the data size for writing is greater than this page size, the system must issue a precharge command and another activate command. According to Fig. 1(a), each SA located next to the mat has common power supplies CSP for a PMOSFET and CSN for an NMOSFET: these perform pre-charging and SA amplification, as shown in Fig. 2. In the pre-charge mode, CSP, CSN, and the BL (BLT and BLB) voltages are set to $V_{ARY}/2$, which is half of the voltage of V_{ARY} . The high (H) and low (L) levels of the BL voltage are V_{ARY} and V_{SS} (0 V), respectively. Initially, the SA amplification

voltage VDD is applied to BLEQB, which reduces BLT/B from VARY/2.

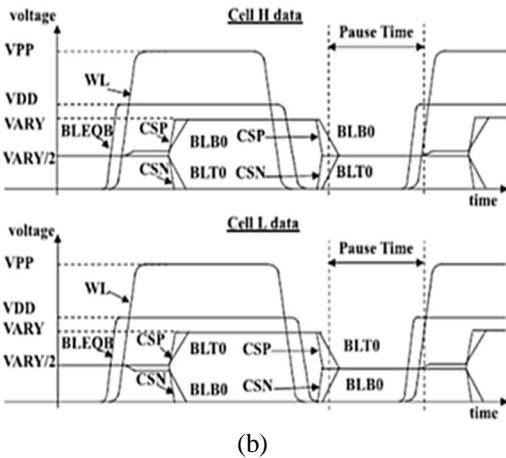
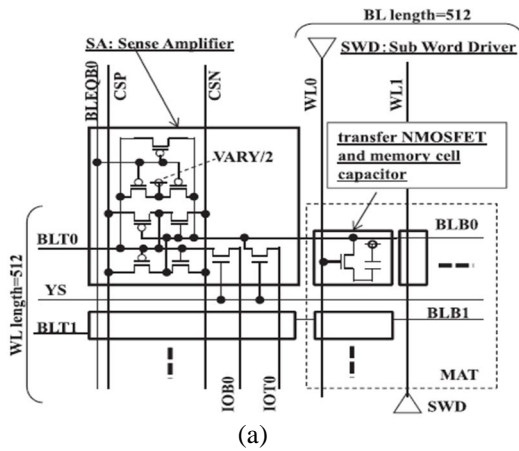
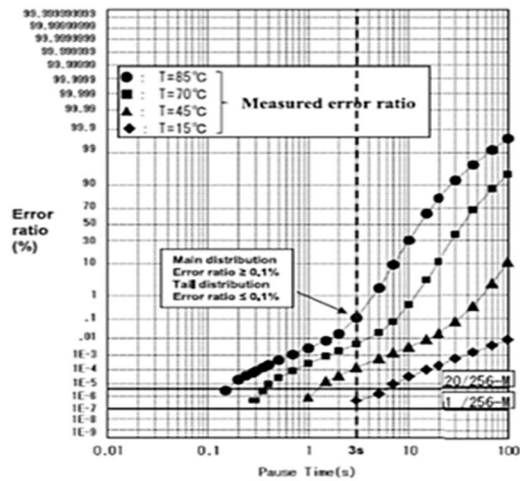


Fig. 1. (a) Array circuit and (b) its waveforms.

Then, the selection of the WL connects the memory node to the BL. The charge in the storage capacitor changes the BL voltage to a higher or lower value than VARY/2. The difference from VARY/2 is called the signal amount. After the voltage of the BL has stabilized, CSP is changed to VARY and CSN is changed to VSS at the same time. The SA increases the voltage difference between BLT and BLB to VARY. In the pre-charge mode, the WL voltage is changed to VSS to reduce the memory node from the BL voltage, and then, CSP and CSN are reduced from VARY and VSS, respectively. Finally, BLEQB is lowered to set CSP, CSN, and the BL voltage to VARY/2. Memory cell structure and leakage current path. Measured error ratio including temperature dependency. The results show zero fail bits at 100 ms and eight fail bits at 150 ms for T = 85 °C, which indicates 100 ms < t_{ret}, min < 150 ms.

	ACH, Pause	ACH, Disturb
cross section		
leakage current path	body	source-drain body

(a)



(b)

Fig. 2. (a) Memory cell structure and leakage current path. (b) Measured error ratio including temperature dependency. The results show zero fail bits at 100 ms and eight fail bits at 150 ms for T = 85 °C, which indicates 100 ms < t_{ret}, min < 150 ms.

B. Charge Retention Time of the Storage Capacitor

The AREF interval time, t_{REF}, should be as long as possible to meet the system requirements. The user usually sets this time to 7.8 μs, which is the maximum value defined from the specification, t_{ref} of 64 ms. The retention time t_{ret} depends on the characteristics of a memory cell that has a leakage current that reduces the charge on the storage capacitor. This leakage current is caused by the diffusion and generation of electrons and holes at the p-n junction in the silicon substrate, the sub-threshold current, and GIDL [1], [3], and [4]. Fig. 2(a) shows the structure of the storage capacitor and transfer NMOSFET in a memory cell. The leakage current is influenced by the voltages of the WL, BL, and body (p-well). This causes a variety of t_{ret} values among several conditions [4], [14], and [15]. There are two states that hold the storage charge in the memory cell, described as follows.

1). The observed memory cell is not selected, and all memory cells in the same mat are not selected. WL voltage = VSS, and BL voltage = VARY/2. If the memory node voltage is H in the observed memory cell, leakage current flows from the memory node to the body [2], [3], [7], and [8]. The left panel in Fig. 2(a) shows this state. The memory node voltage H is changed to L in t_{ret} , which is the specific finite time of the cell. We call this destruction mode as the all cell high (ACH) pause. All in ACH means that the memory node voltages in all DRAM memory cells are H at the last restore time. If the memory node voltage is L, there is no destruction to H as the voltage of the body is the same voltage, L.

2). The observed memory cell is not selected, and one of the other WLs in the same mat is selected. WL voltage = VSS, and BL voltage = H or L. If the memory node voltage is H and BL voltage is L, leakage current flows not only into the body, but also into the BL. The electric field induces a current into the BL. This destruction mode is called ACH disturb. If the memory node voltage is L and BL voltage is H, only leakage current into the BL appears and there is no current into the body. This is called all cell low disturb. The case of ACH disturb is the worst because there are two current paths [2], [3], and [7]. The right panel in Fig. 2(a) shows this state.

The measured t_{ret} data in state 1 is obtained from the time when the fail-bit judgment of the read command appears after a pause time in the precharge mode from the write operation of memory node voltage H. Fig. 1(b) shows the pause time. For a 256-Mb DRAM, the fail-bit count of ACH pause denotes the number of cells of signal Lin 256 Mb after the write operation of H. In Fig. 1(b), 1 cell/bit measurement data show the error ratio of the 256-Mb DRAM from 100 ms to 100 s, where the error ratio is defined by the number of fail bits divided by 256 Mb. The error ratio is converted with an inverse cumulative distribution function to check whether the distribution coincides with a standard normal distribution.

If the distribution appears as a standard normal distribution, the line becomes straight [5], [6]. Fig. 2(b) shows that all lines have a kink indicating the existence of two kinds of standard normal distributions, tail and main distributions. The variation of t_{ret} is greater than three orders of magnitude at a temperature of 85°C, and t_{ret} has one order of variation in the area whose error ratio is less than 0.1%. This means that, although t_{ret} of 99.9% for 256 Mb is longer than 3s, $t_{ret, min}$ of the 256-Mb DRAM is 100 ms. This difference is caused by the tail distribution, whose error ratio is less than 0.1%. Furthermore, the variation among temperatures is very large. According to Fig. 2(b), $t_{ret, min}$ of 100 ms at a temperature of 85°C, which is the worst value, is improved to 1 s at 45°C. In the DRAM

specification, t_{ref} ensures a $t_{ret, min}$ of 64 ms at temperatures between 0°C and 85°C [9], [10], [14], [24], and [25]. A DRAM supplier can replace the worst 20 or 30 memory cells with other good ones prepared in advance. Fig. 2(b) shows the effect of 20-bit replacement.

C. DRAM Standby Current

A DRAM consumes standby power, mainly by refresh, even when there are no read and write accesses to memory cells. The storage capacitor needs around 20 fF, independent of F, as t_{ref} is always 64 ms through all generations. If the minimum feature size is cut in half by process improvement, the memory capacity of the DRAM for the same die size increases by a factor of four. The voltage to rewrite data should be lowered to maintain the same power consumption. The BL rewrite voltage VARY, however, saturates at around 1 V caused by real operation difficulties. Table I shows the external voltage saturation that is caused by SA amplification where VARY is the source voltage. Therefore, the amount of refresh current increases by a factor of four as well. In fact, the total refresh current is less than this because the other currents used to drive signals for the refresh operation become less than half owing to the miniaturization. Considering this factor, the standby power increases in every generation.

If there is no access to the DRAM, t_{ret} is determined only by ACH pause 1 without disturb 2 discussed in the previous section. The DRAM has a SELF mode in the specification, which means that the DRAM performs a refresh operation by itself in the longer interval $t_{ref, cl}$. The oscillator circuit in the DRAM can adjust the period for t_{REF} in the SELF mode. This adjustment based on the measurement of the manufactured DRAM enables a greater reduction in power consumption in the SELF mode. This t_{REF} adjustment is set on the basis of the t_{ret} -measured DRAM data that were compiled during the manufacturing process of the silicon wafer. If it is known in advance that the system has no access to the DRAM in a certain interval, a user issues a SELF ENTRY command, which is the longest time interval during which storage charge is not lost. Once the DRAM receives a SELF ENTRY command, it remains in the SELF mode until receiving a SELF EXIT command, and it accepts no commands except SELF EXIT [11]-[13], [16].

III. PARTIAL ACCESS MODE

A. Conventional PAM Method

Mismatches between the AREF operation in the normal mode and the preserved bank data in the PASR have been explained. Our proposed PAM eliminates these mismatches and uses all memory cells efficiently. One characteristic of the PAM is that it holds data using 2N

cells/bit to extend t_{ref} . Another characteristic is the control method between 1 and $2N$ cells/bit that is located higher in the hierarchy than the control of the normal and SELF modes. The PAM reduces the refresh operation frequency for both the normal and SELF modes. Therefore, the AREF command frequency in the normal mode and power consumption in the SELF mode are reduced at once. In contrast, the PASR reduces the power consumption in the SELF mode only.

A PAM ENTRY operation indicates a conversion from 1 to $2N$ cells/bit. This is simply a copy operation from the memory cell connected with one WL to $2N-1$ memory cells connected with $2N-1$ WLs in the same mat. Table II shows how this operation is simply achieved by the DRAM array architecture in Fig. 1(a). This operation is performed by only a delayed WL selection. A PAM ENTRY operation is completed through $8-K/2N$ copy operations applied to all memory cells in the 256-Mb DRAM.

B. Bitmask Selection

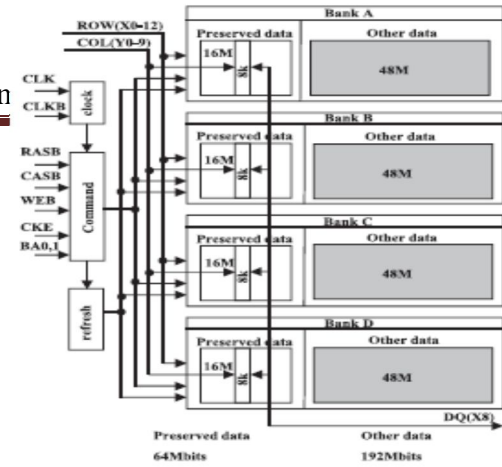
The generic encoding formats of bitmask-based compression technique for various number of bitmasks. A compressed data stores information regarding the bitmask type, bitmask location, and the mask pattern itself. The bitmask can be applied on different places on a vector and the number of bits required for indicating the position varies depending on the bitmask type.

Fig 3.(a) Block diagram for the PAM

PAM	Bits2-0	Access Array	t_{ref}	times	t_{REF}
	000	256M(All Banks)	64ms	8Kref	7.8 μ s
	001	128M(X0=0)	128ms	4Kref	31.3 μ s
	010	64M(X0=X1=0)	256ms	2Kref	125 μ s
	011	restriction	--	--	--
	100	restriction	--	--	--
	101	32M(X0=X1=X2=0)	512ms	1Kref	500 μ s
	110	16M(X0=X1=X2=X3=0)	1024ms	512ref	2000 μ s
111	restriction	--	--	--	

Fig 3(b) Specification of the PAM.

For instance, if we consider a 32-bit vector, an 8-bit mask applied on only byte boundaries requires 2-bits, since it can be applied on four locations. If we do not restrict the placement of the bitmask, it will require 5 bits to indicate any starting position on a 32-bit vector. Bitmasks may be sliding or fixed. Fixed bitmasks are referred with the letter \cdot while sliding bitmasks are referred with the letter \cdot . For example, $\cdot\cdot$ refers to a sliding bitmask of length 2 while $\cdot\cdot$ refers to a fixed bitmask of length 2. A fixed bitmask is one which can be applied to fixed locations, such as byte boundaries. However, sliding bitmasks can be applied anywhere in the test vector. Since the fixed bitmasks can be applied only to fixed locations, the number of positions where they can be applied is significantly less compared



to sliding bitmasks. Hence, the number of bits needed to represent them are less than sliding bitmasks.

IV. PROPOSED METHOD

A. Compressed Technique.

The compression technique used is the Bitstream compression which is important in reconfigurable system design since it reduces the bitstream size and the memory requirement. It also improves the communication bandwidth and thereby decreases the reconfiguration time. This technique gives the good compression ratio due to complex and variable-length coding.

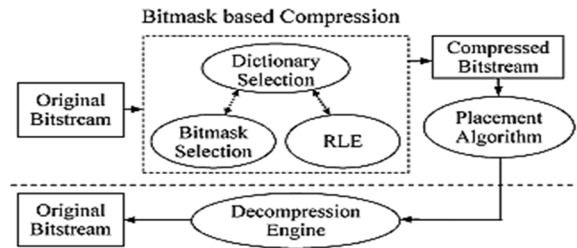
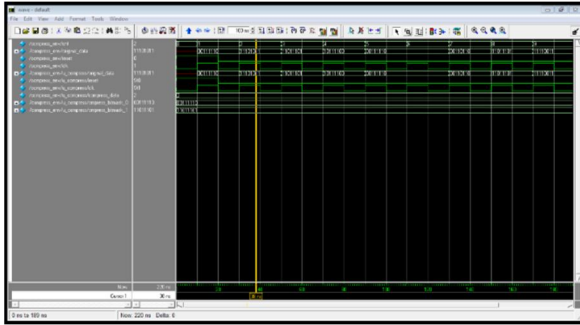


Fig. 4(a) Block diagram for Bitmask compression.

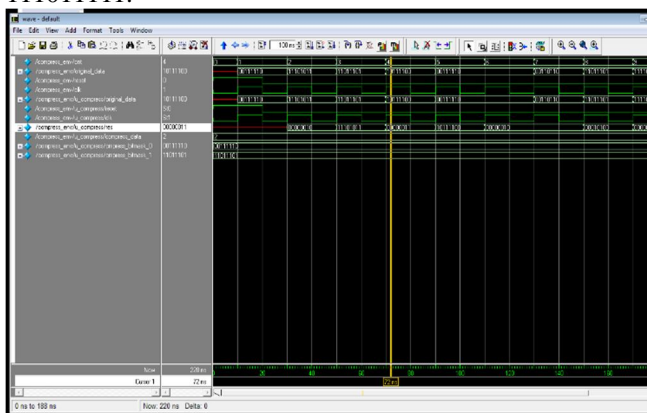
This approach accelerates decompression using simple or fixed-length coding (FLC) where more LUTs are needed because of the dictionary selection method. Fig 5(a) performs smart placement of compressed bitstreams to enable fast decompression of variable-length coding which selects bitmask-based compression parameters suitable for bitstream compression. Finally, it efficiently combines run length encoding and bitmask-based compression to obtain better compression and faster decompression.

B. Result

The given input test vector 00111110, 11101011, 11011101, 10111100, 00111110, 00111110, 00110110, 11011101, 11110011, 11011111.



The Compressed Data's are 010, 111101011, 011, 110111100, 010, 010, 0010100, 011, 111110011, 111011111.



V. CONCLUSION

The specification of the PAM shows the reduction of refresh operation frequency, which contributes to the current reduction. In the normal mode, the PAM is useful if t_{ACT} is longer than 260 ns, even in situations such as continuous addressing. We examined the t_{ret} extension using the fail-bit data of the fabricated DRAM. An important result is that the amount of t_{ret} extension is not proportional to the composed cell number. Partial access conversion does not cause a parallel shift in the t_{ret} distribution, but does reduce the t_{ret} variance. This means that the width of the t_{ret} variation approaches zero with an increasing number of composed cells, which is one example of the law of large numbers.

This paper presented the advantages of bitmask-based compression. This paper developed efficient bitmask for test data compression in order to create maximum matching patterns. Our test compression technique used the bitmask selection methods to significantly reduce the testing time and memory requirements.

Our proposed technique outperforms the existing compression approaches by 15%, while our decompression hardware for variable-length coding is

capable of operating at the speed closest to the best known field-programmable gate array-based decoder for fixed-length coding.

REFERENCES

- [1] G. A. M. Hurkx, D. B. M. Klaassen, and M. P. G. Knuvers, "A new recombination model for device simulation including tunneling," *IEEE Trans. Electron Devices*, vol. 39, no. 2, pp. 331–338, Feb. 1992.
- [2] S. Amakawa and K. Nakazato, "A new approach to failure analysis and yield enhancement of very large-scale integrated systems," in *Proc. ESSDERC*, Sep. 2002, pp. 147–150.
- [3] O. K. B. Lui and P. Migliorato, "A new generation-recombination model for device simulation including the Poole-Frenkel effect and phonon-assisted tunneling," *Solid-State Electron.*, vol. 41, no. 4, pp. 575–583, 1997.
- [4] W. Shockley and W. T. Read, "Statistics of the recombinations of holes and electrons," *Phys. Rev.*, vol. 87, no. 5, pp. 835–842, 1952.
- [5] A. Hiraiwa, M. Ogasawara, N. Natsuaki, Y. Itoh, and H. Iwai, "Field-effect trap-level-distribution model of dynamic random access memory data retention characteristics," *J. Appl. Phys.*, vol. 81, no. 10, pp. 7053–7060, 1997.
- [6] S. Ueno, Y. Inoue, M. Inuishi, and N. Tsubouchi, "Leakage mechanism of local junctions forming the main or tail mode of retention characteristics for dynamic random access memories," *Jpn. J. Appl. Phys.*, vol. 39, no. 4B, pp. 1963–1968, 2000.
- [7] M. Chang, J. Lin, S. N. Shih, T. C. Wu, B. Huang, J. Yang, and P. I. Lee, "Impact of gate-induced drain leakage on retention time distribution of 256 Mbit DRAM with negative wordline bias," *IEEE Trans. Electron Devices*, vol. 50, no. 4, pp. 1036–1040, Apr. 2003.
- [8] K. Saino, S. Horiba, S. Uchiyama, Y. Takaishi, M. Takenaka, T. Uchida, Y. Takada, K. Koyama, H. Miyake, and C. Hu, "Impact of gate-induced drain leakage current on the tail distribution of DRAM data retention time," in *IEDM Tech. Dig.*, Dec. 2000, pp. 837–840.
- [9] K. Yamaguchi, "Temperature dependence of anomalous currents in worst-bit cells in dynamic random-access memories," *J. Appl. Phys.*, vol. 87, no. 11, pp. 8064–8069, 2000.
- [10] H. Kim, B. Oh, Y. Son, K. Kim, S. Y. Cha, J. G. Jeong, S. J. Hong, and H. Shin, "Study of trap models related to the variable retention time phenomenon in DRAM," *IEEE Trans. Electron Devices*, vol. 58, no. 6, pp. 1643–1648, Jun. 2011.
- [11] J. P. Kim, W. Yang, and H. Y. Tan, "A low-power 256-Mb SDRAM with an on-chip thermometer and biased reference line sensing scheme," *IEEE J. Solid-State Circuits*, vol. 38, no. 2, pp. 329–337, Feb. 2003.

- [12] C. K. Kim, J. G. Lee, Y. H. Jun, C. G. Lee, and B. S. Kong, "CMOS temperature sensor with ring oscillator for mobile DRAM self-refresh control," *Microelectron. J.*, vol. 38, pp. 1042–1049, Jan. 2007.
- [13] S. S. Pyo, C. H. Lee, G. H. Kim, K. M. Choi, Y. H. Jun, and B. S. Kong, "45 nm low-power embedded pseudo-SRAM with ECC-based autoadjusted self-refresh scheme," in *Proc. IEEE Int. Symp. Circuits Syst.*, May 2009, pp. 2517–2520.
- [14] Y. Ito and H. Iwai, "Data storing method of dynamic RAM and semiconductor memory device," U.S. Patent 6 697 992, Feb. 24, 2004.
- [15] S. H. Kim, W. O. Lee, J. H. Kim, S. S. Lee, S. Y. Hwang, C. I. Kim, T. W. Kwon, B. S. Han, S. K. Cho, D. H. Kim, J. K. Hong, M. Y. Lee, S. W. Yin, H. G. Kim, J. H. Ahn, Y. T. Kim, Y. H. Koh, and J. S. Kih, "A low power and highly reliable 400 Mbps mobile DDR SDRAM with on-chip distributed ECC," in *Proc. Asian Solid-State Circuits Conf.*, 2007, pp. 34–37.
- [16] W. Gao and S. Simmons, "A study on the VLSI implementation of ECC for embedded DRAM," in *Proc. IEEE Can. Conf. Electr. Comput. Eng.*, vol. 1, May 2003, pp. 203–206.
- [17] N. Shinozaki and Y. Matsuzaki, "Semiconductor memory," U.S. Patent 6 829 192, Dec. 7, 2004.
- [18] T. Vogelsang, H. Lorenz, and W. Hokenmaier, "Memory with selectable single cell or twin cell configuration," U.S. Patent 7 254 089, Aug. 7, 2007.
- [19] Y. Mori, S. Kamohara, M. Moniwa, K. Ohyu, T. Yamanaka, and R. Yamada, "Direct observation of worst-bit leakage currents of DRAM," *IEEE Trans. Electron Devices*, vol. 53, no. 2, pp. 398–400, Feb. 2006.
- [20] R. Takemura, K. Itoh, T. Sekiguchi, S. Akiyama, S. Hanzawa, K. Kajigaya, and T. Kawahara, "Long-retention-time, high-speed DRAM array with 12-F2 twin cell for sub 1-V operation," *IEICE Trans. Electron.*, vol. E90-C, no. 4, pp. 758–764, 2007.
- [21] R. Takemura, K. Itoh, and T. Sekiguchi, "A 0.5-V FD-SOI twin-cell DRAM with offset-free dynamic-VT sense amplifiers," in *Proc. ISLPED*, 2006, pp. 123–126.
- [22] T. Iwai, M. Kaku, T. Miyazaki, H. Iwai, H. Takenaka, A. Suzuki, S. Miyano, and M. Hamada, "Low power embedded DRAM using 0.6 V super retention mode with word line data mirroring," in *Proc. Asian Solid-State Circuits Conf.*, 2009, pp. 209–212.
- [23] H. Shimano, F. Morishita, K. Dosaka, and K. Arimoto, "A voltage-scalable advanced DFM RAM with accelerated screening for low power SoC platform," *IEICE Trans. Electron.*, vol. E90-C, no. 10, pp. 1927–1935, 2007.
- [24] R. R. Shaller, "Moore's law: Past, present and future," *IEEE Spectrum*, vol. 34, no. 6, pp. 52–59, Jun. 1997.
- [25] R. H. Dennard, F. H. Gaensslen, V. L. Rideout, E. Bassous, and A. R. LeBlanc, "Design of ion-implanted MOSFET's with very small physical dimensions," *IEEE J. Solid-State Circuits*, vol. 9, no. 5, pp. 256–268, Oct. 1974.
- [26] V. Nathan and N. C. Das, "Gate-induced drain leakage current in MOS devices," *IEEE Trans. Electron Devices*, vol. 40, no. 10, pp. 1888–1890, Oct. 1993.
- [27] V. Cuppu, B. Jacob, B. Davis, and T. Mudge, "A performance comparison of contemporary DRAM architectures," *ACM SIGARCH Comput. Archit. News*, vol. 27, no. 2, pp. 222–233, May 1999.
- [28] H. Kawamoto, T. Shinoda, Y. Yamaguchi, S. Shimizu, K. Ohishi, N. Tanimura, and T. Yasui, "A 288 K CMOS Pseudo-static RAM," *IEEE J. Solid-State Circuits*, vol. 19, no. 5, pp. 619–623, Oct. 1984.
- [29] Y. Oowaki, K. Tsuchida, Y. Watanabe, D. Takashima, M. Ohta, and H. Nakano, S. Watanabe, A. Nitayama, F. Horiguchi, K. Ohuchi, and F. Masuoka, "A 33-ns 64-Mb DRAM," *IEEE J. Solid-State Circuits*, vol. 26, no. 11, pp. 1498–1505, Nov. 1991.