

Hierarchical Multi-Chip Architecture for High Capacity Scalability of Fully Parallel Hamming-Distance Associative Memories

Yusuke OIKE^{†a)}, Student Member, Makoto IKEDA[†], and Kunihiro ASADA[†], Members

SUMMARY In this paper, we present a hierarchical multi-chip architecture which employs fully digital and word-parallel associative memories based on Hamming distance. High capacity scalability is critically important for associative memories since the required database capacity depends on the various applications. A multi-chip structure is most efficient for the capacity scalability as well as the standard memories, however, it is difficult for the conventional nearest-match associative memories. The present digital implementation is capable of detecting all the template data in order of the exact Hamming distance. Therefore, a hierarchical multi-chip structure is simply realized by using extra register buffers and an inter-chip pipelined priority decision circuit hierarchically embedded in multiple chips. It achieves fully chip- and word-parallel Hamming distance search with no throughput decrease, additional clock latency of $O(\log P)$, and inter-chip wires of $O(P)$ in a P -chip structure. The feasibility of the architecture and circuit implementation has been demonstrated by post-layout simulations. The performance has been also estimated based on measurement results of a single-chip implementation.

key words: associative memory, content addressable memory, CAM, Hamming distance, capacity scalability, multi-chip structure

1. Introduction

Associative processing has a wide variety of application fields such as pattern recognition, code-book-based data compression, multi-media, intelligent processing and learning systems. It generally requires considerable memory access and data processing time. Therefore some high-speed associative memories [1]–[9] have been developed for Hamming- or Manhattan-distance estimation, which is essentially required for the associative processing. Some conventional associative memories [1]–[6] employ analog/digital circuit techniques for quick nearest-match detection. However, there are difficulties in operating them with faultless precision in a deep sub-micron (DSM) and future process with a low-voltage supply. Moreover the feasible database capacity is limited by the analog operation. Therefore the digital implementations [7]–[9] have been proposed to achieve the large database capacity and the applicability to a system-on-a-chip VLSI in the latest process technologies.

High capacity scalability is important for the associative memories since the required database capacity depends on the various applications. A multi-chip structure is most efficient for the capacity scalability as well as the standard

memories. In the complete-match detection such as [10]–[12], all the detected data are the correct results because they are exactly the same as the input. Therefore, the complete-match data can be compiled without additional comparison among the detected data even in a multi-chip structure. On the other hand, in the conventional nearest-match associative memories [1]–[6], each module provides just the local nearest data since the search operation is executed independently of each other module. Thus, the global nearest detection requires additional memory access and distance calculation because the exact Hamming distance is not provided in the local nearest-match detection. Furthermore, it requires an inter-chip distance comparison among all the local nearest data. These features make it difficult for [1]–[6] to attain high capacity scalability by a multi-chip structure. The digital implementations have a potential capacity scalability by a multi-chip structure. [7] reports an 8-chip structure with extra winner-take-all (WTA) processors. It requires extra 4th, 5th and more pipelined WTA processors on each chip in order to build up a larger database capacity. On the other hand, a fully word-parallel architecture such as [8]–[9] is more efficient for high-speed associative processing than [7].

In this paper, we propose a hierarchical multi-chip architecture which employs our proposed fully word-parallel associative memory [8] to achieve the high capacity scalability. The associative memory is capable of detecting all the template data in order of the exact Hamming distance. It enables high-speed data sorting in addition to nearest-match detection for the conventional use. The hierarchical multi-chip structure is simply realized with the original functions by extra register buffers and an inter-chip pipelined priority decision (PPD) circuit. All the chips are composed of the same circuit configuration, and hierarchically connected via a PPD node embedded in a chip. The present architecture and circuit implementation achieve fully chip- and word-parallel Hamming distance search with no throughput decrease, additional clock latency of $O(\log P)$, and inter-chip wires of $O(P)$ in case of a P -chip structure.

Section 2 reviews a basic architecture and circuit implementation of the fully word-parallel associative memory [8]. A hierarchical multi-chip structure is described in Sect. 3. Section 4 presents the circuit implementation and operation. Section 5 introduces a module generator of the associative memories with various database capacities. Section 6 shows the performance evaluation of the multi-chip structure based on post-layout simulations with mea-

Manuscript received April 1, 2004.

Manuscript revised July 15, 2004.

[†]The authors are with the Faculty of Engineering, and VLSI Design and Education Center (VDEC), the University of Tokyo, Tokyo, 113-8656 Japan.

a) E-mail: y-oike@silicon.u-tokyo.ac.jp

surement results of a single-chip implementation. Finally, Sect. 7 concludes this paper.

2. Fully Parallel Associative Processing

2.1 Basic Operations

Figure 1 shows an operation diagram of the fully digital and word-parallel associative memory [8]. First, the input (D_{in}) is compared with all template data (D_0, D_1, \dots, D_M) by using an XOR/XNOR circuit embedded in a memory cell. Next, the number of mismatch bits are counted by a search signal propagation via hierarchically chained search circuits in word parallel. The search circuit is also embedded in a memory cell and controls the search signal propagation based on the comparison results ($D_{in} \oplus D_M$). A mismatch bit is masked in every word, and then the next mismatch bit is detected by a search signal propagation during a search clock period. The mask and search operations are carried out during a search clock period regardless of where a mismatch bit exists. Therefore the nearest-match data are detected faster than the others, and the 2nd- and 3rd-nearest data are also detected in order of the distance. The associative processing architecture is capable of exact Hamming-distance search for all the template data in the distance order. Finally, the detected address is provided by a priority address encoder.

2.2 Circuit Configurations

Figure 2 and Fig. 3 show a timing diagram and a block diagram of the associative memory. The associative memory consists of a memory array with hierarchically chained search circuits, a memory read/write circuit, a row decoder, data latch and mask circuits, and a priority address encoder. All the search paths are refreshed by setting search signals (SS) to a low level before a distance search period as shown in Fig. 2. In a distance search period, the search signals are set to a high level, and they pass through the match bits in word parallel. Then, they stop at the first-encountered mismatch bit in each word. During the next clock period, the

detected mismatch bit is masked and the search propagation starts again from the masked bit. Finally, all the mismatch bits are masked in a word, and the search signal (SS) reaches the end of the word as a search result (SO). A priority encoder with word mask circuits provides the address of the search result. In case that several words of the same Hamming distance are found, a search result is masked with the priority decision signals (PO_m) one after another. The present associative memory sequentially provides all the detected addresses of the same Hamming distance.

The operated clock cycles represent Hamming distance

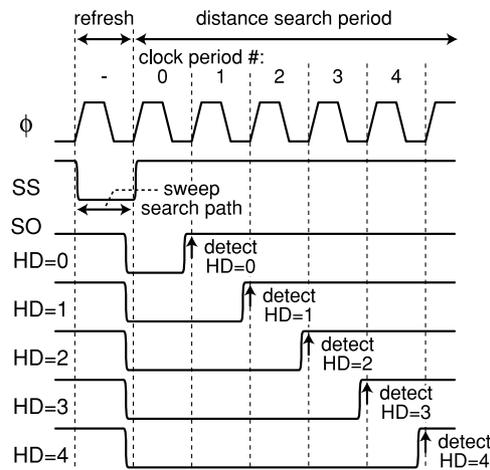
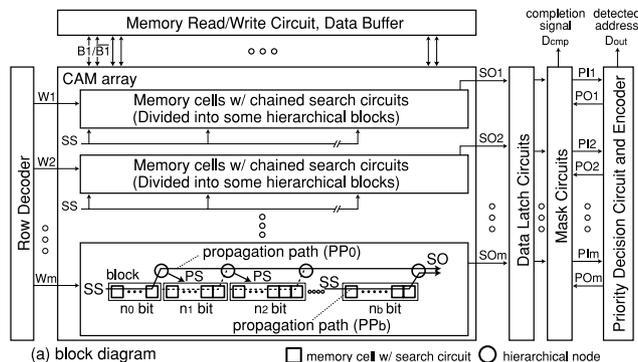


Fig. 2 Timing diagram of Hamming-distance search.



(a) block diagram

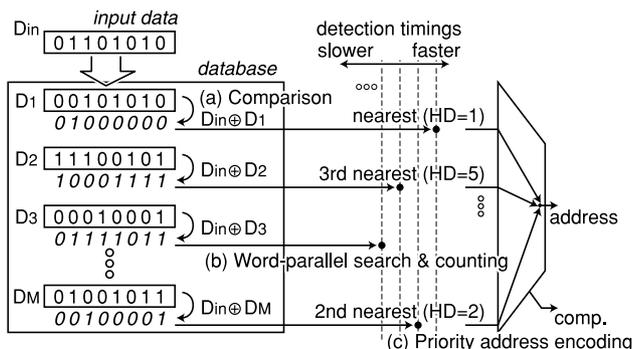


Fig. 1 Operation diagram of a fully digital and word-parallel associative memory.

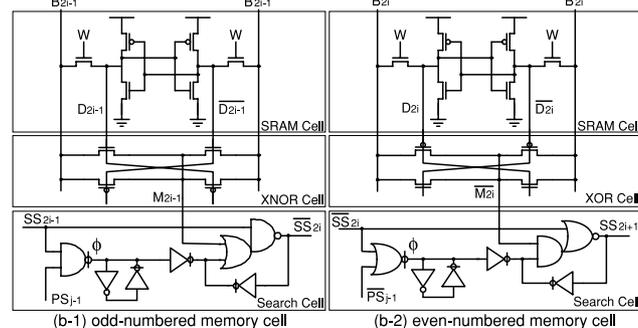


Fig. 3 Associative memory configuration: (a) block diagram; (b-1) schematic of odd-numbered memory cell, (b-2) schematic of even-numbered memory cell.

(HD) of the detected data. For example, data of $HD = 0$ (i.e. complete-match data) are detected in the 0th clock period as shown in Fig. 2 since the search signal passes through all the bits without any interruption. And then, data of $HD = 1$ are detected in the 1st clock period. That is, data of $HD = D$ are detected in the D -th clock period. This feature contributes to save the clock cycles in the bit-serial digital implementation especially for a practical use since the required data are generally close to the input. Furthermore it enables additional associative functions such as high-speed data sorting by the exact distance search.

2.3 Performance Characteristics

The search cycle time is linearly proportional to the bit length in a serial search path structure as shown in Fig. 4(a). It becomes a bottleneck of the associative processing, hence a hierarchical search structure [8] is implemented for the search signal paths as shown in Fig. 4(b). The template data are divided into blocks, which are connected by hierarchical nodes. The hierarchical node receives a search signal from the previous block, and then it provides a permission signal (PS) to the next block. The permission signal makes a mismatch bit maskable in case that the mismatch bit interrupts a search signal propagation. In other words, a mismatch bit that has received both a search signal (SS) and a permission signal (PS) is set to a maskable bit in each word. The search and mask operations are carried out by a search circuit embedded in a memory cell. Even-numbered and odd-numbered search circuits are complementary in order to reduce the propagation delay and the circuit area as shown in

Fig. 3(b). The number of bits in each block needs to be optimized for the critical path minimization since the two-stage hierarchical search structure has several propagation paths (PP_0, PP_1, \dots, PP_b) as shown in Fig. 3(a). The search cycle time is limited by $O(\sqrt{N})$ at an N -bit length database due to the two-stage hierarchical structure. Search results are transferred to a priority address encoder to acquire the address output during the next search operation. The priority address encoder is implemented using a binary-tree structure, hence the address encoding time is limited by $O(\log M)$ at an M -word database as reported in [8]. The search cycle time (T_c), which determines the search throughput, is given by

$$T_c = \max(T_1, T_2), \tag{1}$$

$$T_1 \propto O(\sqrt{N}), \tag{2}$$

$$T_2 \propto O(\log M), \tag{3}$$

where T_1 is a search propagation time and T_2 is a priority address encoding time. N and M are the bit length and the number of words, respectively. The total search time (T_s) is given by

$$T_s = T_c \times (D + 1), \tag{4}$$

where D is Hamming distance between the input and the detected data.

3. Hierarchical Multi-Chip Structure

Figure 5 shows possible multi-chip structures of the present associative memory. Figure 5(a) shows a bus structure with a scan controller, which has the high capacity scalability and flexibility. It is, however, difficult to attain a high-speed search operation since the scan controller sequentially searches all the chips for a detected address during a search clock period. Figure 5(b) shows a star structure with a winner-take-all (WTA) processor. The WTA processor simultaneously collects all the detected addresses. It is capable of acquiring a detected address during a search clock period. On the other hand, it requires a special WTA processor according to the number of chips. The address signal wires increase in proportion to $O(P \times \log P)$ in case of a P -chip structure, and all the output signals concentrate on the same WTA processor chip. It becomes a potential problem on the capacity scalability and flexibility.

We propose a hierarchical structure using an inter-chip pipelined priority decision (PPD) circuit as shown in Fig. 5(c). In the present architecture, an associative memory chip interacts with each other using a completion signal ($Dcmp$) via a hierarchical PPD node embedded in a chip. A completion signal ($Dcmp$) represents whether any data are detected in a chip or not, which is provided by intra-chip priority decision results (PO_m). The inter-chip PPD circuit determines whether any chip contains a detected address and which chip is given priority for providing a search result. Therefore, a search result can be autonomously provided from the associative memory chip with priority. A

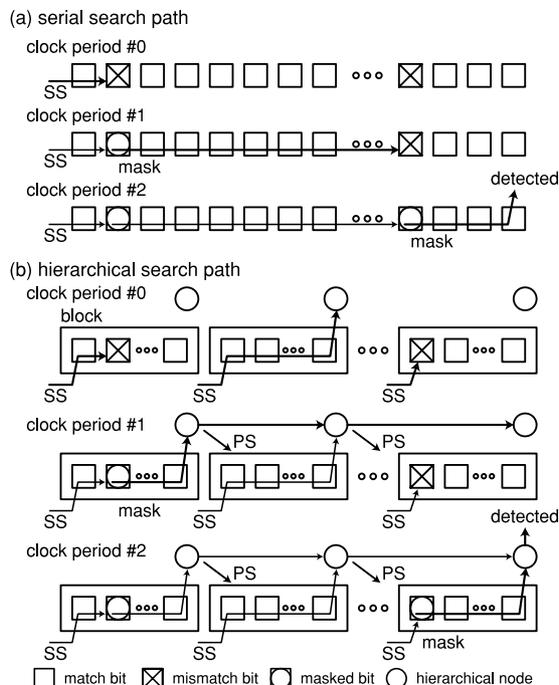


Fig. 4 Search path structure: (a) serial search path, (b) hierarchical search path.

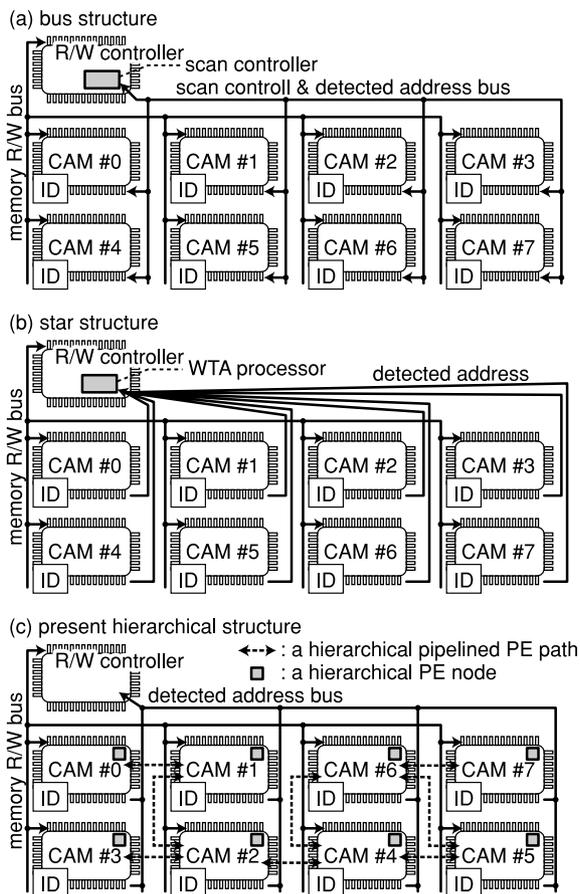


Fig. 5 Possible multi-chip structures: (a) a bus structure with a scan controller, (b) a star structure with a WTA processor, (c) the present hierarchical structure.

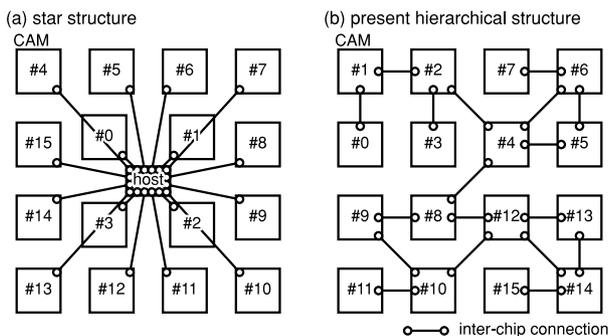


Fig. 6 Examples of inter-chip wiring in a multi-chip structure: (a) a star structure, (b) the present hierarchical structure.

long signal wire between chips limits the search operation speed. The present multi-chip structure, however, realizes a two-dimensional chip array with a tree network by short signal wires as shown in Fig. 6 since a chip is adjacently connected by peer-to-peer interaction with four chips at a maximum. Therefore, it requires short signal wires of $O(P)$ for an inter-chip PPD circuit and output bus wires of $O(\log P)$. The present multi-chip architecture enables fully chip- and word-parallel Hamming distance search with no throughput decrease, additional clock latency of $O(\log P)$, and inter-chip wires of $O(P)$ in a P -chip structure. Table 1 shows comparison among the multi-chip structures in case of a capacity of $256 \text{ bit} \times 256 \text{ word}$ per chip. In the comparison, CAM chips are placed in a two-dimensional array, and they are connected by straight wires as shown in Fig. 6. The wire length is normalized by a pitch of the chip array. In a star structure, we assume that an additional WTA host processor compiles all the detected addresses from CAM chips and searches them in a single search clock.

4. Circuit Realization and Operation

4.1 Inter-Chip Connections with PPD Circuits

Figure 7 shows a hierarchical multi-chip structure using a binary-tree pipelined priority decision (PPD) circuit. All

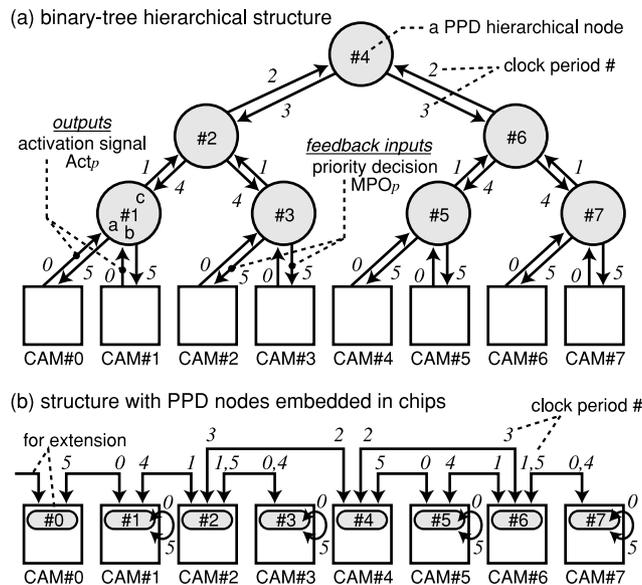


Fig. 7 Hierarchical multi-chip structure using embedded binary-tree pipelined priority decision circuits.

Table 1 Comparison among multi-chip structures.

	bus structure		star structure		hierarchical structure	
	16 chips	64 chips	16 chips	64 chips	16 chips	64 chips
Num. of wires	12	16	208	1088	$60 + 12^*$	$252 + 16^*$
Total wire length	180.0	1008.0	311.5	3311.3	$65.0 + 180.0^*$	$272.3 + 1008.0^*$
Search clock latency	16	64	1	1	7	11
Throughput	1/16	1/64	1 (lossless)	1 (lossless)	1 (lossless)	1 (lossless)

* A hierarchical tree network and address output buses, respectively.

CAM chips are hierarchically connected via PPD nodes as shown in Fig. 7(a). A CAM chip that detects data of $HD = D$ during the D -th clock period provides an activation signal (Act_p) to a PPD node. The activation signal is generated by an intra-chip completion signal ($Dcmp$). The hierarchical PPD nodes transfer the activation signals to the next stage while it determines which one is a priority result. Finally, they return the priority decision results (MPO_p) to the CAM chips. The priority decision is carried out in the pipeline. Therefore, it requires additional latency of L_c clock cycles, which is given by

$$L_c = 2 \times \log_2 P - 1, \quad (5)$$

where P is the number of chips in the multi-chip structure. For example, the pipelined priority decision with eight CAM chips is completed in five clocks as shown by clock period numbers in Fig. 7(a).

The number of hierarchical PPD nodes (N_{ppd}) is given by

$$N_{ppd} = P - 1, \quad (6)$$

due to a binary-tree structure. Therefore each PPD node can be efficiently embedded in a CAM chip as shown in Fig. 7(b). All CAM chips are implemented by the same circuit configuration. This feature enables a multi-chip structure without any additional processor chip. In the multi-chip structure, one PPD node always remains as shown by CAM#0 in Fig. 7(b). The remaining PPD node is used for extension of the capacity, hence it attains the high capacity scalability by the flexible number of chips.

4.2 Extended Associative Memory Configuration

Figure 8 shows a block diagram of an associative memory chip extended for the multi-chip structure. It requires

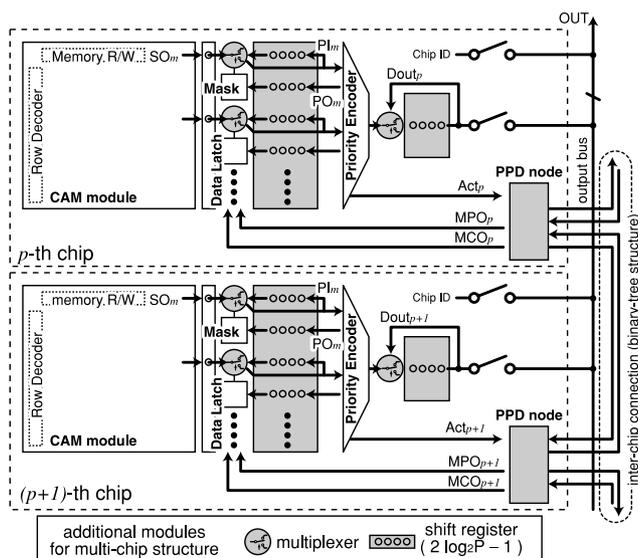


Fig. 8 Block diagram of associative memory for multi-chip configuration.

two-input multiplexers and shift registers in addition to the single-chip circuit configuration presented in Fig. 3(a). An associative memory chip provides an activation signal (Act_p) to a PPD node in case that it detects a search output (SO_m). In the chip- and word-parallel Hamming distance search, some data of the same Hamming distance can be simultaneously detected ranging over all chips. Therefore, the inter-chip PPD circuit determines which chip is given priority over the other activated chips. An activated chip that receives the priority from the inter-chip PPD circuit provides the detected address and the chip ID as a search result. After the priority word is masked, the other detected words are evaluated again by the intra- and inter-chip priority decision circuits. In this case, all the search signal propagations are interrupted. And then, the search results (SO_m), which are temporarily buffered by shift registers, are provided to the intra-chip priority encoder again. The search signal propagations start again after all the detected addresses are processed since the priority decision circuit becomes available for the next search results. MCO_p is a completion signal of the inter-chip PPD circuit. The number of shift registers (N_{reg}) is a logarithmic order of the number of chips as follows:

$$N_{reg} = 2 \times \log P - 1, \quad (7)$$

since it is determined by the additional clock latency resulting from a hierarchical PPD circuit.

4.3 Pipelined Priority Decision Circuit Configuration

The intra-chip priority decision is carried out by a binary-tree priority address encoder as reported in [8]. It consists of a priority decision circuit and an address encoder as shown in Fig. 9(a). An inter-chip PPD circuit is designed based on the binary-tree priority decision circuit. A PPD node consists of a priority decision cell, an ID decoder, and register buffers. A priority decision cell has three inputs (Pin) and three outputs ($Pout$) in a similar configuration to the intra-chip priority decision circuit as shown in Figs. 9(a) and (b). In the intra-chip priority decision circuit, an input of Pin_a is also used for a return path from the upper hierarchical level. On the other hand, the inter-chip priority decision circuit loses the original input of Pin_a since the operations are pipelined. Therefore, an input of Pin_a is buffered by shift registers in each PPD node. The shift registers are prepared according to the maximum number of chips. The number of buffer stages is set by the chip ID since the return path length is different for the hierarchical levels. Figure 10 shows a timing diagram of the inter-chip PPD circuit. The number of buffer stages can be determined by the least true bit of a chip ID because of a binary-tree structure. An inter-chip completion signal MCO_p is acquired by $Pout_c$ of the top node, for example, $Pout_c$ of CAM#4 in a multi-chip structure with eight chips. The completion signal is provided to each chip along a return path.

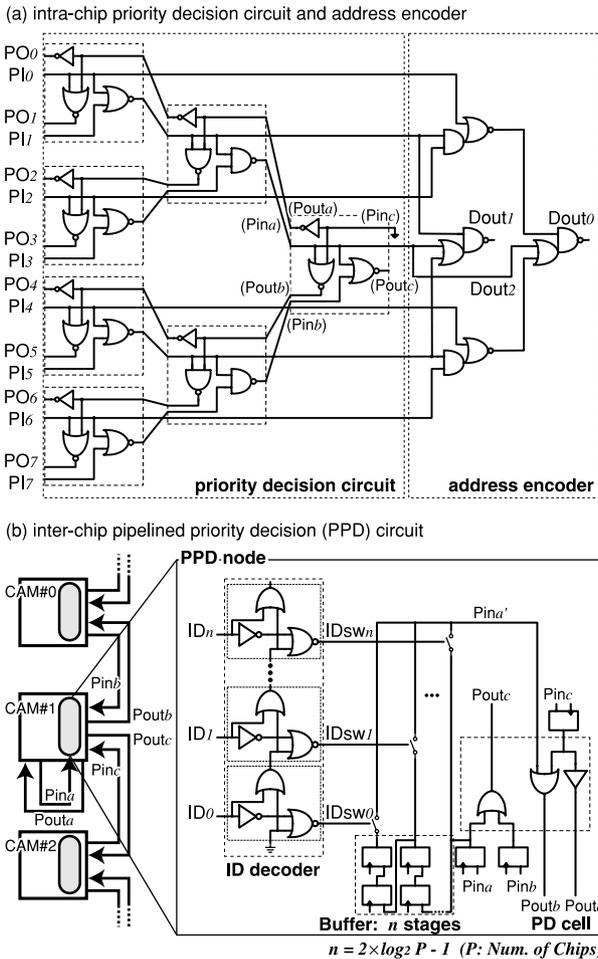


Fig. 9 Simplified binary-tree priority decision circuit: (a) intra-chip priority decision circuit and address encoder, (b) inter-chip pipelined priority decision circuit.

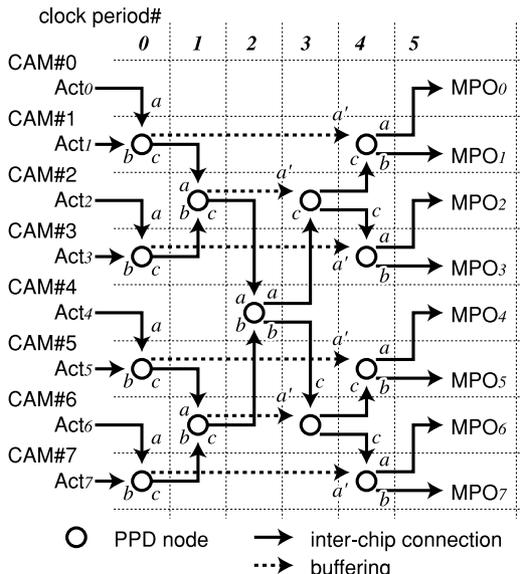


Fig. 10 Timing diagram of PPD circuit in case of 8 chips.

5. Module Generator for Various Capacities

We have developed a module generator for various capacities of the present associative memories. A required capacity of associative memories is different for various applications. Therefore a module generator which automatically provides an optimized structure with any database capacity is also important for the high capacity scalability. The present architecture of fully chip- and word-parallel Hamming-distance estimation has the simplicity, regularity, and flexibility in structure. Therefore an associative memory module with variable capacities can be designed using a common macro cell library which includes a memory cell with a search circuit, a part of an address decoder, a sense amplifier, a word mask circuit, a shift register, and so on. Figure 11 shows the module generator functions. The module generator partially employs Synopsys HSPICE, Cadence Dracula LPE and Virtuoso. The inputs are hard macro cells and a specification file including their cell sizes and pin locations. First, the library cells are extracted to SPICE netlists by using Dracula LPE, and then the cell performances are characterized by using HSPICE. The characterization can be skipped in case that the module generator has already characterized the library cells. The delay of a hierarchical search node is especially estimated with various fan-outs since the fan-out increases in proportion to bit length of the next block. Then, the module generator divides the database into hierarchical blocks based on the capacity requirements and the characterization results. A hierarchical structure that provides the minimum search path is generated by simulated annealing. Finally, the library cells are arranged, and the module generator provides a layout script file for Virtuoso. An inter-chip PPD node and additional shift registers are automatically added to the associative memory module according to the specified number of chips. Figure 12 shows an execution example of the module generator. Figures 13(a) and (b) are the module generation examples. Fig-

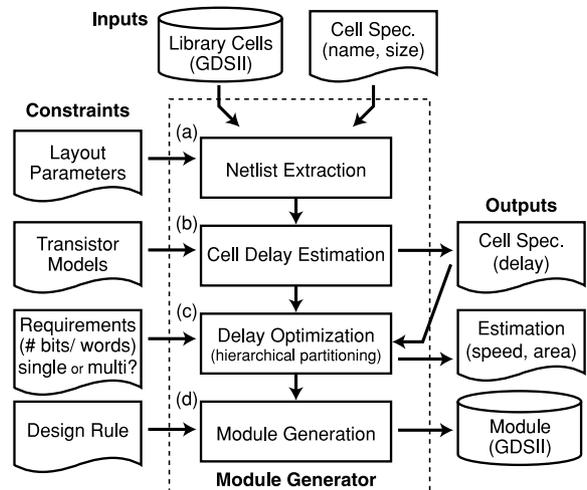


Fig. 11 Module generator functions.

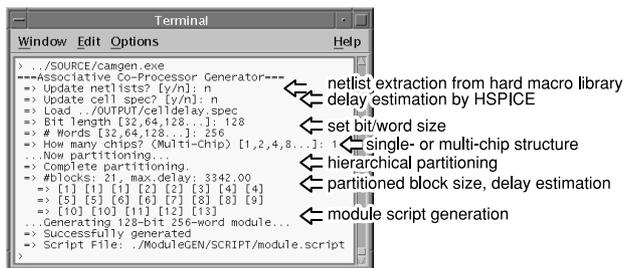


Fig. 12 Module generator execution example.

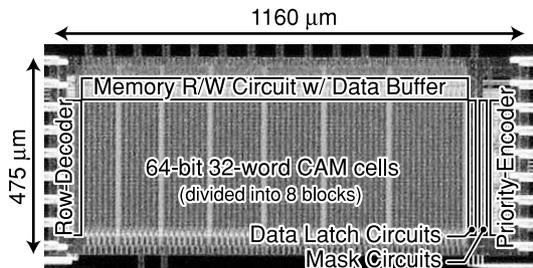


Fig. 14 Chip microphotograph of 64-bit 32-word associative memory for a single-chip structure.

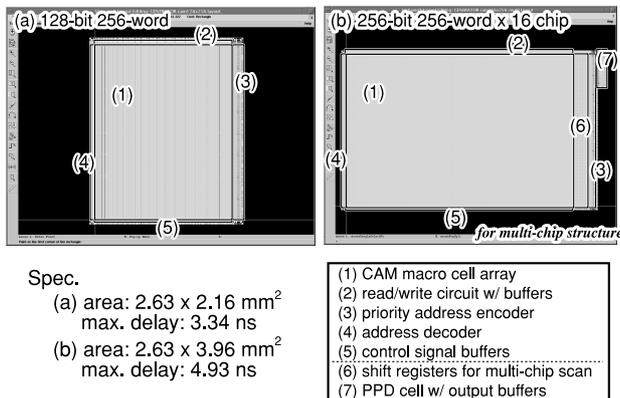


Fig. 13 Examples of module generation: (a) 128-bit 256-word module for a single chip, (b) 256-bit 256-word module for 16-chip structure.

ure 13(a) is a 128-bit 256-word module for a single-chip structure. Figure 13(b) is a 256-bit 256-word module for a 16-chip structure. The module generator also reports the maximum delay.

6. Performance Evaluation

6.1 Chip Implementation

A 64-bit 32-word associative memory has been fabricated using a 1P5M 0.18 μm CMOS process, which has been reported in [9]. In this paper, the associative memory is presented just for the feasibility study and the performance estimation of the present architecture since it has been designed for a single-chip structure. Figure 14 shows a chip microphotograph of the 64-bit 32-word associative memory. It operates at a supply voltage from 0.75 V to 1.8 V. The search operation attains a speed of 411.5 MHz at 1.8 V. The worst-case search time is 158.0 ns since it requires 65 search clocks for the complete-mismatch data (i.e. HD = 64). The chip specifications are summarized in Table 2.

Table 3 shows the estimated areas of an associative memory module with various database capacities. The number of transistors in the present associative memory cell is larger than that applying the conventional analog approaches. However, the analog approaches make it difficult for device scaling to keep the performance and marginal capacity. The present approach can achieve device scaling and operate at a low supply voltage because of the synchronous

Table 2 Specifications of the designed associative memory.

Process	1P5M 0.18 μm CMOS process
Power supply voltage	0.75–1.8 V
Organization	64 bits × 32 words memory cells
Functions	Nearest-match detection Distance ordering
Module size	475 μm × 1160 μm (0.55 mm ²)
Num. of transistors	88.5 k
Memory cell size	9.6 μm × 13.6 μm (130.56 μm ²)
Operation speed	411.5 MHz (@ 1.8 V, measured) 454.5 MHz (@ 1.8 V, simulated)
Worst-case search time	158.0 ns (0-bit to 64-bit distance)
Power dissipation	51.3 mW (@ 1.8 V, 400 MHz)

Table 3 Area of associative memory module.

Database capacity	Area (Module size)
4 K (64 b × 64)	0.98 mm ² (0.79 × 1.24)
16 K (128 b × 128)	3.02 mm ² (1.40 × 2.16)
64 K (256 b × 256)	11.05 mm ² (2.63 × 4.20)
256 K (512 b × 512)	38.34 mm ² (5.08 × 7.55)
1 M (1024 b × 1024)	146.40 mm ² (10.00 × 14.64)

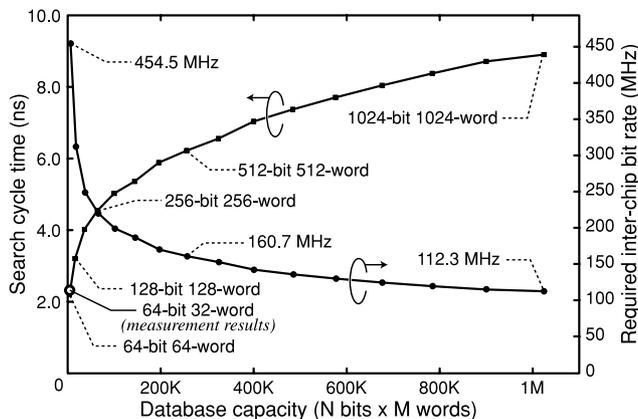


Fig. 15 Search cycle time and inter-chip bit rate.

digital search logics embedded in the memories. Therefore, in comparison with the conventional designs, the associative memory has greater potential for practical use and a larger capacity.

6.2 Search Cycle Time and Inter-Chip Bit Rate

Figure 15 shows a search cycle time of various database

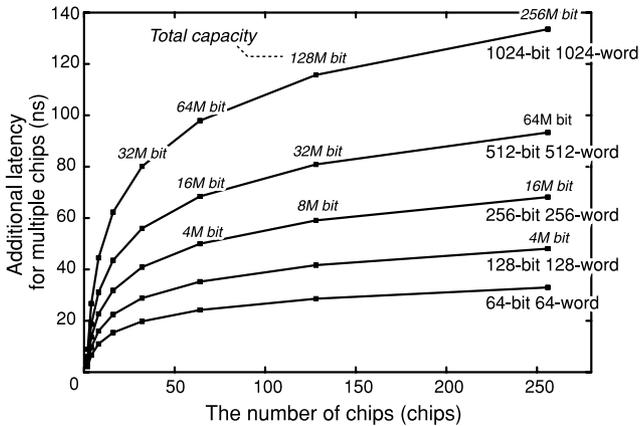


Fig. 16 Additional latency for the multi-chip structure.

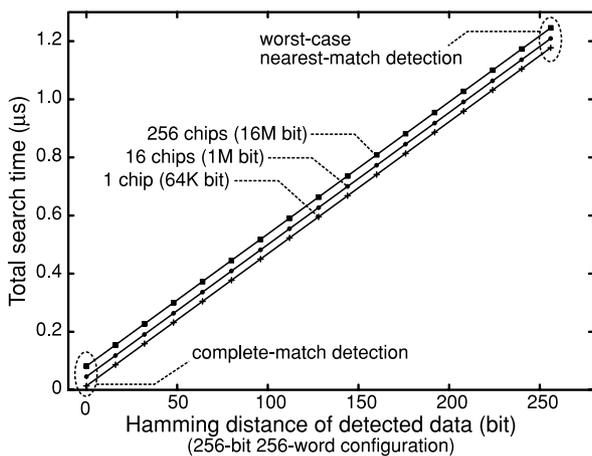


Fig. 17 Total search time as a function of Hamming distance of the detected data.

capacities assuming the bit length (N) and the number of words (M) are the same. The search cycle time is limited by the search signal propagation of $O(\sqrt{N})$ or the priority address encoding of $O(\log M)$ as shown by Eq. (1). Therefore the hierarchical search structure attains a high-speed search operation in a large database. It achieves a search cycle time of 8.90 ns at a 1024-bit 1024-word database (i.e. 1 Mb capacity). The required inter-chip bit rate is determined by the search cycle time. 454.5 MHz and 112.3 MHz inter-chip signalings are required for the associative memories of 4K b/chip and 1M b/chip, respectively. These inter-chip transmission speeds are feasible in the latest chip-to-chip interconnect technologies.

6.3 Hamming-Distance Search Time

Figure 16 shows additional latency for the multi-chip structure. The binary-tree PPD circuit reduces the additional latency to $O(\log P)$ as shown by Eq. (5). Therefore the additional latency is just 133.5 ns even for a 256 Mb database which consists of 256 associative memory chips with a 1024-bit 1024-word capacity. Furthermore the multi-chip

architecture maintains a continuous search operation with no throughput decrease, which enables the detection of data after the 2nd-nearest data. The total search time depends on the Hamming distance between the input and the detected data as shown by Eq. (4). Figure 17 shows the total search time in 1-, 16-, and 256-chip structures of 256-bit 256-word associative memories as a function of Hamming distance of the detected data. In these configurations, the search time for the complete-match data is 13.6 ns, 45.5 ns, and 81.8 ns at 16 Mb, 1 Mb, and 64 Kb capacities, respectively. Furthermore the search time for the nearest-match data is 1.18 μ s, 1.21 μ s, and 1.25 μ s in the worst case, respectively. The hierarchical multi-chip architecture and circuit implementation achieve the capacity scalability with small performance degradation.

7. Conclusions

We have proposed a hierarchical multi-chip architecture using fully digital and word-parallel associative memories based on Hamming distance. The multi-chip structure efficiently realizes the high capacity scalability by using an inter-chip pipelined priority decision (PPD) circuit. The inter-chip PPD circuit enables fully chip- and word-parallel associative processing by taking advantage of the feature of the digital associative processing architecture, which attains no throughput decrease, additional clock latency of $O(\log P)$, and inter-chip wires of $O(P)$ in the P -chip structure. The developed module generator automatically optimizes the hierarchical search structure and provides the associative memory module for various capacity requirements. The feasibility of the architecture and circuit implementation has been demonstrated by post-layout simulations with measurement results of a single-chip implementation. The performance evaluation shows that the hierarchical multi-chip architecture is capable of the high-speed and continuous associative processing based on Hamming distance with a megabit database capacity.

Acknowledgment

This work is supported by VLSI Design and Education Center (VDEC), University of Tokyo in collaboration with Cadence Design Systems Inc. and Synopsys Inc. The VLSI chip in this study has been fabricated through VLSI Design and Education Center (VDEC), University of Tokyo in collaboration with Hitachi Ltd. and Dai Nippon Printing Co.

References

- [1] T. Yamashita, T. Shibata, and T. Ohmi, "Neuron MOS winner-take-all circuit and its application to associative memory," IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers, pp.236-237, Feb. 1993.
- [2] M. Nagata, T. Yoneda, D. Nomasaki, M. Sato, and A. Iwata, "A minimum-distance search circuit using dual-line PWM signal processing and charge-packet counting techniques," IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers, pp.42-43, Feb. 1997.

- [3] M. Ikeda and K. Asada, "Time-domain minimum-distance detector and its application to low-power coding scheme on chip-interface," Proc. Eur. Solid-State Circuit Conf. (ESSCIRC), pp.464–467, Sept. 1998.
- [4] T. Hanyu and M. Kameyama, "Multiple-valued logic-in-memory VLSI architecture based on floating-gate-MOS pass-transistor logic," IEICE Trans. Electron., vol.E82-C, no.9, pp.1662–1668, Sept. 1999.
- [5] H.J. Mattausch, T. Gyohten, Y. Soda, and T. Koide, "An architecture for compact associative memories with deca-ns nearest-match capability up to large distances," IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers, pp.170–171, Feb. 2001.
- [6] H.J. Mattausch, N. Omori, S. Fukae, T. Koide, and T. Gyohten, "Fully-parallel pattern-matching engine with dynamic adaptability to Hamming or Manhattan distance," Symp. on VLSI Circuits Dig. Tech. Papers, pp.252–255, June 2002.
- [7] A. Nakada, T. Shibata, M. Konda, T. Morimoto, and T. Ohmi, "A fully parallel vector-quantization processor for real-time motion-picture compression," IEEE J. Solid-State Circuits, vol.34, no.6, pp.822–830, June 1999.
- [8] Y. Oike, M. Ikeda, and K. Asada, "A high-speed and low-voltage associative co-processor with Hamming distance ordering using word-parallel and hierarchical search architecture," Proc. IEEE Custom Integrated Circuits Conference (CICC), pp.643–646, Sept. 2003.
- [9] S. Nakahara and T. Kawata, "A digital circuit for a minimum distance search using an asynchronous bubble shift memory," IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers, pp.504–505, Feb. 2004.
- [10] T. Ogura, J. Yamada, S. Yamada, and M. Tanno, "A 20-kbit associative memory LSI for artificial intelligence machines," IEEE J. Solid-State Circuits, vol.24, no.4, pp.1014–1020, Aug. 1989.
- [11] T. Ikenaga and T. Ogura, "A fully parallel 1-Mb CAM LSI for real-time pixel-parallel image processing," IEEE J. Solid-State Circuits, vol.35, no.4, pp.536–544, April 2000.
- [12] H. Miyatake, M. Tanaka, and Y. Mori, "A design for high-speed low-power CMOS fully parallel content-addressable memory macros," IEEE J. Solid-State Circuits, vol.36, no.6, pp.956–968, June 2001.



Yusuke Oike was born in Tokyo, Japan on July 4, 1977. He received the B.S. and M.S. degrees in electronic engineering from University of Tokyo, Tokyo, in 2000 and 2002, respectively. He is currently pursuing the Ph.D. degree at the Department of Electronic Engineering, University of Tokyo. His current research interests include architecture and design of smart image sensors, mixed-signal circuits, and functional memories. He received the Best Design Awards from IEEE International Conference on

VLSI Design and IEEE/ACM ASP-DAC in 2002 and 2004, respectively. He is a student member of the Institute of Electrical and Electronics Engineers (IEEE), and the Institute of Image Information and Television Engineers of Japan (ITEJ).



Makoto Ikeda received the B.S., M.S., and Ph.D. in electronics engineering from the University of Tokyo, Tokyo, Japan, in 1991, 1993, and 1996, respectively. He joined the Department of Electronic Engineering, University of Tokyo as a Faculty Member in 1996, and is currently an Associate Professor at VLSI Design and Education Center, University of Tokyo. His research interests include the reliability of VLSI design. He is a member of Institute of Electrical and Electronics Engineers (IEEE), and Information Processing Society of Japan (IPJS).



Kunihiro Asada was born in Fukui, Japan, on June 16, 1952. He received the B.S., M.S., and Ph.D. in electronic engineering from University of Tokyo in 1975, 1977, and 1980, respectively. In 1980 he joined the Faculty of Engineering, University of Tokyo, and became a lecturer, an associate professor and a professor in 1981, 1985 and 1995, respectively. From 1985 to 1986 he stayed in Edinburgh University as a visiting scholar supported by the British Council. From 1990 to 1992 he served as the

first Editor of English version of IEICE (Institute of Electronics, Information and Communication Engineers of Japan) Transactions on Electronics. In 1996 he established VDEC (VLSI Design and Education Center) with his colleagues in University of Tokyo. It is a center supported by the Government to promote education and research of VLSI design in all the universities and colleges in Japan. He is currently in charge of the head of VDEC. His research interest is design and evaluation of integrated systems and component devices. He has published more than 390 technical papers in journals and conference proceedings. He has received Best Paper Awards from IEEJ (Institute of Electrical Engineers of Japan), IEICE and ICMTS1998/IEEE and so on. He is a member of IEEE and IEEJ.