

# Unified Data/Instruction Cache with Distributed Crossbar, Hidden Precharge Pipeline and Dynamic CMOS Logic

K. Johguchi<sup>1</sup>, Z. Zhu<sup>1</sup>, K. Aoyama<sup>1</sup>, Y. Mukuda<sup>1</sup>, H. J. Mattausch<sup>1</sup>, T. Koide<sup>1</sup> and T. Hironaka<sup>2</sup>

<sup>1</sup>Hiroshima Univ., Research Center for Nanodevices and Systems, 1-4-2 Kagamiyama, Higashi-Hiroshima, 739-8527, Japan

<sup>2</sup>Hiroshima City Univ., Faculty of Computer Sciences, 3-4-1 OzukaHigashi, AsaMinami-Ku, Hiroshima, 731-3194, Japan

Phone: +81-82-424-6265, FAX: +81-82-424-3499, E-mail: {jouguchi, hjm, koide}@sxsys.hiroshima-u.ac.jp

## 1. Introduction

Modern processors simultaneously fetch, decode and execute many instructions. This results in the demand for a large access bandwidth of the processor's memory components and has already led to register files with many ports. However, for the cache memory the conventional solution of 1-port-data and instruction caches is still in use. Since the demand for parallelism tends to increase at a high rate, the cache system will become the bottleneck of processor performance and has to be innovated.

We propose to improve the access bandwidth of the cache with a 1-port-bank-based Hierarchical Multi-port memory Architecture (HMA), which can simultaneously realize small area, high bandwidth and low power dissipation [1, 2]. Moreover, by combining instruction and data caches into a single multi-port cache, we are able to dynamically schedule the memory amount used for data and instructions, resulting in a more efficient usage of the caches storage capacity.

## 2. Multi-port Cache with HMA Structure

HMA is a 1-port-bank based multi-port memory architecture which further improves area consumption and performance of the conventional crossbar architecture. The crossbar's switch network is distributed into the bank structure, which decreases global wiring and transistor number. A two dimensional bank decoder reduces the overhead for bank selection and allows easy matrix arrangement of the banks [1].

Fig. 1 shows the structure-example of a direct mapped cache which uses HMA [2]. The cache index, consisting of line number (LN) and line offset (LO), is divided into two portions, a bank internal address (BI) and a bank number (BN). BI is used for selecting a cache word or tag within memory banks, and BN is used for selecting the respective banks within data/instruction or tag memory. BN uses the lower rank bits in order to make sure that consecutive lines and words within lines are located in different banks, so that they can be accessed in parallel without access conflict.

## 3. Unification of Instruction and Data Cache

Conventional cache systems adopt split instruction and data caches to enable independent and parallel access to data and instructions, although only 1-port caches are used. A unified data/instruction cache can provide lower miss rate, because dynamical allocation of the effective storage capacity for data and instructions becomes possible. However, the accesses to the instruction cache are normally consecutive. For a 4-way superscalar processor, it is therefore expected, that one instruction port with 4-time larger word length will deliver sufficient instruction-fetch performance. The optimum number of data-access ports is estimated to be 2 or 3. Then, we propose an HMA unified write-through cache which have different word length for data and instruction ports[3].

## 4. Circuit Methodology and Test-Chip Design

### 4.1. Dynamic CMOS Circuits and Hidden Precharge

Increased access time is a general problem for multi-port memories, because larger integration area leads to larger capacitances in the critical access path and port-switching as well as access-conflict management are required. The conventional access method of a synchronous memory divides the clock cycle into two phases: the memory access phase and

the precharge phase, the latter to prepare the memory for the access in the following clock cycle. Therefore, the cycle time becomes the sum of the actual memory-access time and the precharge time. We propose to exploit the hierarchical structure of the multi-bank cache for a 2-stage synchronous access scheme, which hides the precharge by overlapping the access phase of one stage with the precharge phase of the other stage as illustrated in Fig. 2. Bank-conflict management, bank selection, port conversion and bank-internal wordline decoding constitute the 1st stage of the access path and are carried out simultaneously with the bank precharge when the clock is "0". After the clock becomes "1", the second stage of the access path, starting with the wordline-driver activation is executed, while the first stage of the access path is precharged. In this way equal clock-cycle time and actual multi-port cache access time are realized.

Dynamic CMOS technology is used for the design of the cache circuits in order to reduce in particular the capacitive loads of the global routing to the banks. Fig. 3 explains the construction of a bank and its 1-to- $N$  port converter. In the shown example of the dynamic data-transfer subcircuit within the 1-to- $N$  port converter, simple NMOS transfer gates and a PMOS transistor for precharging the internal node are used. Read and write path are divided in order to enable the proposed 2-stage synchronous access scheme. Reduced parasitic capacitances and smaller area of the 1-to- $N$  port converter are thus achieved.

### 4.2. Test Chip Design

For the test chip of an HMA cache memory, a configuration with 4 ports was chosen and the design was carried out in a 5 metal, 0.18 $\mu$ m CMOS technology. The chip-layout shown in Fig. 4 contains all needed new functional units. The design data are summarized in Table I. Small area and short delay are achieved with a dynamic CMOS circuit technology and effective floor planning. The area-overhead of the 1:4-port converter for the 1Kbyte bank of Fig. 5 is less than 25%. We also applied a new access method which overlaps bank-conflict management and bank decoding with the precharging phase of the banks. The timing diagram and spice simulation results of the netlist extracted from the layout are shown in Fig. 6 for the critical access cases. In the 1st stage the write access is critical, which takes longer than the read access because write data latching in the 1-to-4 port converter completes after address latching. The 2nd access stage is limited by the read access because the read data must be transmitted from the bank to a cache-output port. Simulated critical delay times for 1st and 2nd stage are both slightly below 1.7 ns. This leads to a duty ratio of 50%, which is beneficial for the processor implementation. Simulated power dissipation is 247 mW at 250 MHz when all ports actively access the cache.

## 5. Conclusions

In this paper, a bank-based unified data/instruction cache with multiple ports has been proposed and its advantages have been verified by simulation. Especially important is our method of providing a different word length for data and instruction ports, which takes advantage of the internal bank structure. To minimize bank conflicts, we use an addressing method, which insures that the words in one cache-line and also consecutive cache-lines are located in different banks, so that access performance close to an ideal multi-port cache can be realized. Adopting an efficient floorplan without routing-only

areas restricts the area-overhead including 2nd level circuits and a conflict manager for the 4 ports to only 25% in comparison to a 1-port cache. A minimum clock cycle time of 3.4 ns could be achieved with a dynamic CMOS circuit technology and by overlapping the external bank access with the bank-internal precharge.

The proposed bank-based multi-port cache is also very attractive for low power dissipation, because the number of activated banks, determining power dissipation, corresponds to the port number and is independent of the total number of banks in the cache.

### Acknowledgements

This research is sponsored by Semiconductor Technology Academic Research Center (STARC).

The VLSI chip in this study has been fabricated in the chip fabrication program of VLSI Design and Education Center (VDEC),

the University of Tokyo in collaboration with Hitachi Ltd., Dai Nippon Printing Corporation and Cadence Design Systems, Inc.

This work was supported in part by the 21st Century COE program and a JSPS research fellowship, No.1605246, Japanese Government.

### References

- [1] S. Fukae, *et al.*, "Optimized Bank-Based Multi-Port Memories through a Hierarchical Multi-Bank Structure," Proc. of SASIMI2003, pp.323-330, 2003.
- [2] H. J. Mattausch, *et al.*, "Area-efficient multi-port SRAMs for on-chip data-storage with high random access bandwidth and large storage capacity," IEICE Trans.Electron., Vol.E84-C, No.37, pp.410-417, 2001.
- [3] K. Johguchi, *et al.*, "Distributed-crossbar architecture for area-efficient combined data/instruction caches with multiple ports," Elec. Lett. 40, pp. 160-162, 2004.

Table I Datasheet of the design.

Technology	Logic CMOS w/o SRAM rules
Minimum Gate Length	200 nm
Routing Layer	1-Poly, 5-Al layers
Si-Area	6.2 mm <sup>2</sup>
Total Storage Capacity	20.5 KByte
Port Number	4 ports
Minimum Access Cycle Time	3.4 nsec
Power Dissipation	247 mW at 250 MHz
Instruction & Service Port	
Port Number	2
Wordlength	64 bit
Data Ports	
Port Number	2
Wordlength	16 bit
Tag Memory	
Storage Capacity	4.5 KByte
Bank Number	16
Bank Capacity	2304 bit
Data/Instruction Memory	
Storage Capacity	16 KByte
Bank Number	64
Bank Capacity	2Kbit

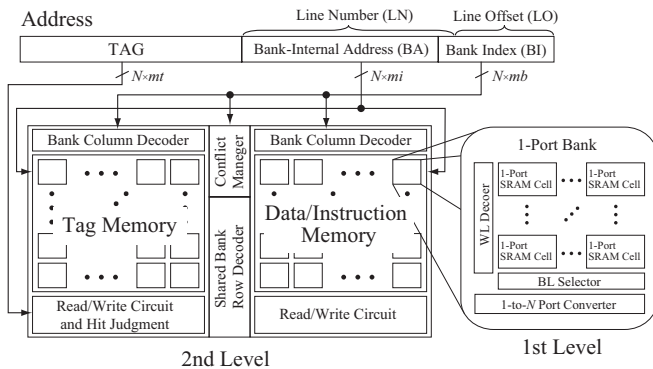


Fig. 1 Block diagram of direct mapped data cache with  $N$  ports in bank-based HMA structure.

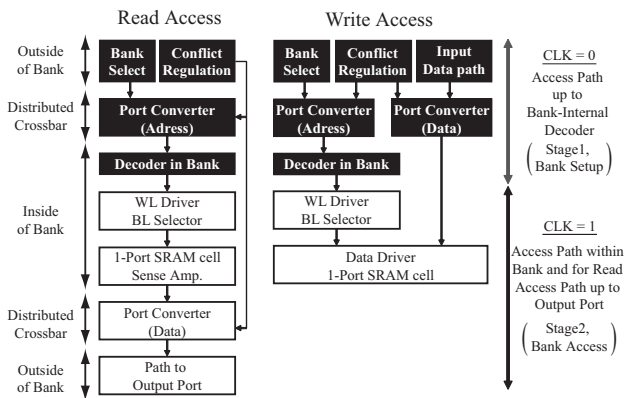


Fig. 2 Synchronous 2-stage access mode for hiding precharge phases.

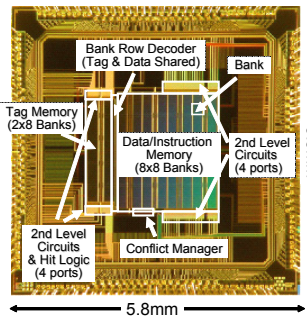


Fig. 4 Microphotograph of the test chip for 4-port cache.

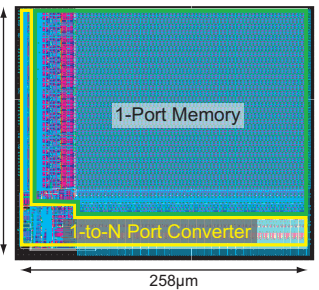


Fig. 5 Layout of a bank with 1-to-4-port converter.

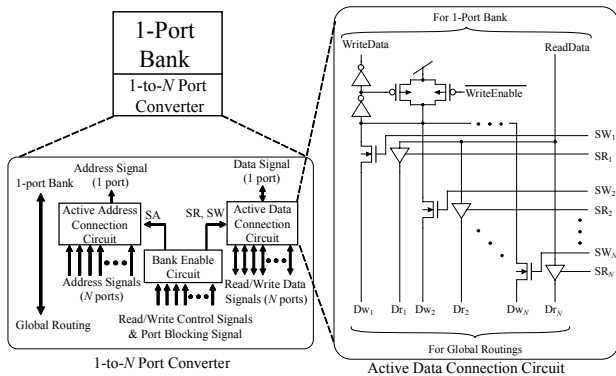


Fig. 3 Block diagram and partial schematic of 1-to- $N$  port converter.

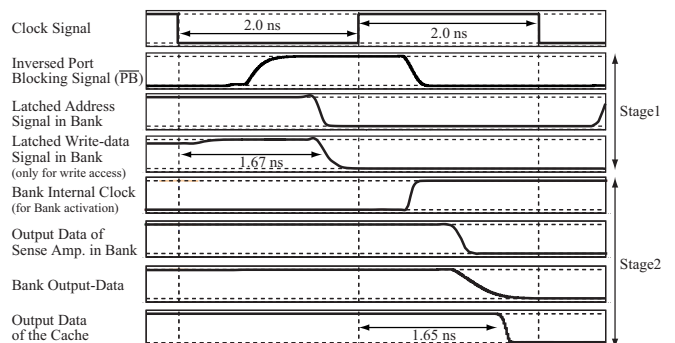


Fig. 6 Simulated access waveforms from extracted layout.



# Unified Data/Instruction Cache with Distributed Crossbar, Hidden Precharge Pipeline and Dynamic CMOS Logic

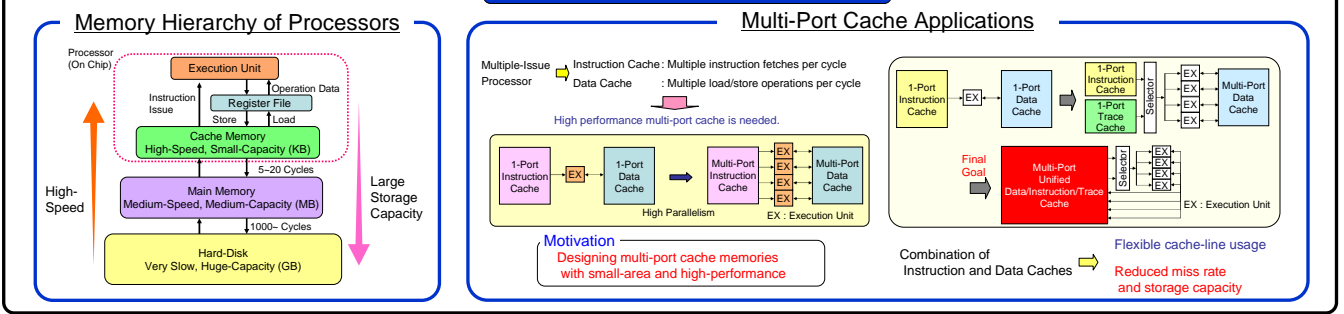
○Koh Johguchi<sup>1</sup>, Zhaomin Zhu<sup>1</sup>, Ken-ichi Aoyama<sup>1</sup>, Yuya Mukuda<sup>1</sup>, Hans Jürgen Mattausch<sup>1</sup>, Tetsushi Koide<sup>1</sup>, Tetsuo Hironaka<sup>2</sup>

<sup>1</sup> Research Center for Nanodevices and Systems (RCNS), Hiroshima University

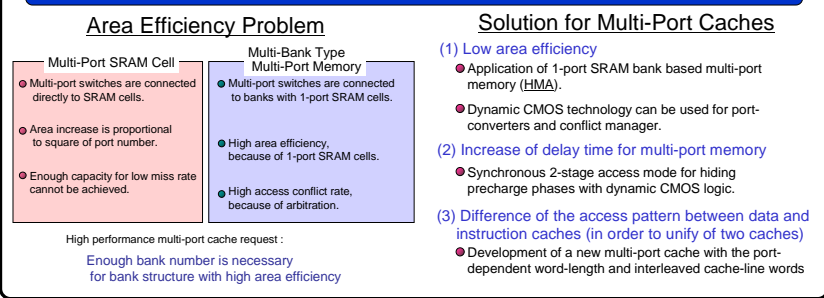
<sup>2</sup> Department of Computer Engineering, Hiroshima City University

**Hiroshima University**  
Nanoelectronics for Tera-BIT  
Information Processing

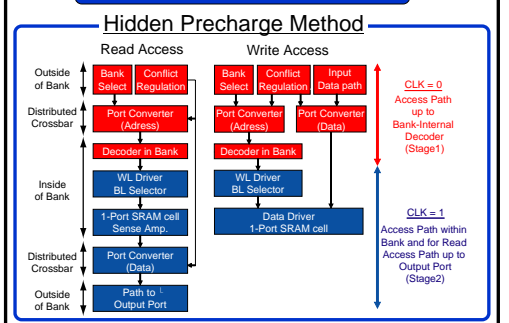
## Background & Motivation



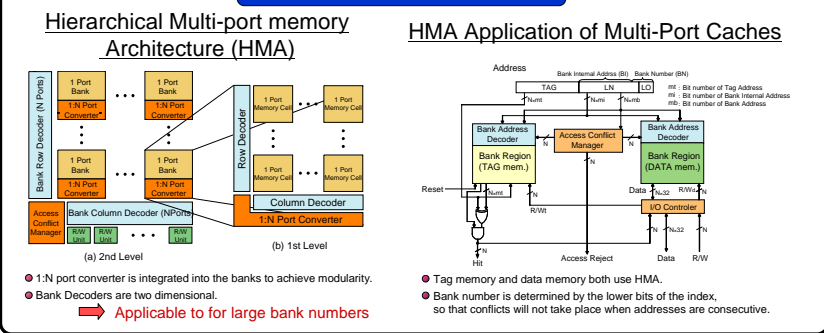
## Special Issues of Conventional Multi-Port Cache & Proposed Solution



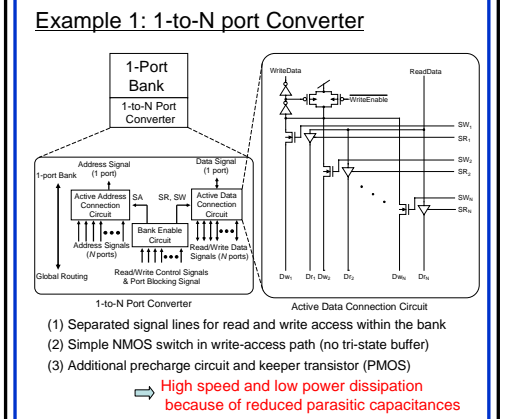
## Solution for Delay Time



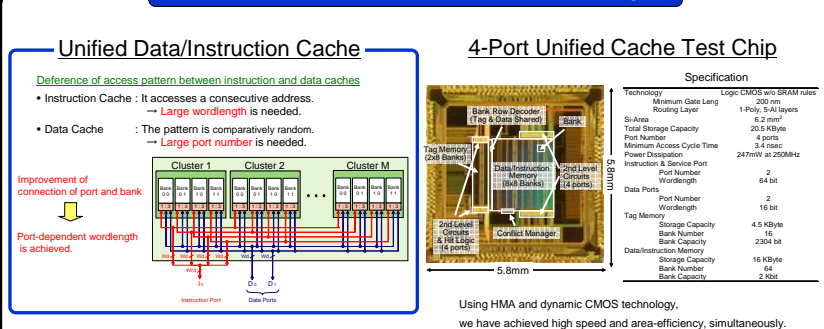
## HMA Cache Architecture



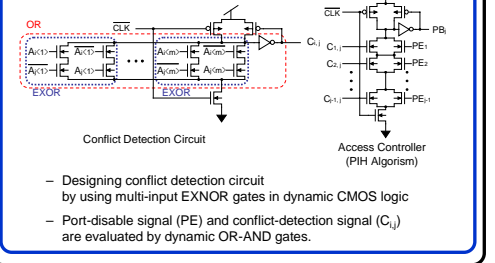
## Dynamic CMOS Logic



## Unified Data/Instruction HMA Cache Design



### Example 2: Conflict Manager



## Conclusions

