

Simple Memory Machine Models for GPUs

KOJI NAKANO^{†1}

The main contribution of this paper is to introduce two parallel memory machines, the Discrete Memory Machine (DMM) and the Unified Memory Machine (UMM). Unlike well studied theoretical parallel computational models such as PRAMs, these parallel memory machines are practical and capture the essential feature of memory access of NVIDIA GPUs. As a first step of the development of algorithmic techniques on the DMM and the UMM, we first evaluated the computing time for the contiguous access and the stride access to the memory on these models. We also present parallel algorithms to transpose a two dimensional array on these models. Since the computing time of our transposing algorithms on the DMM and the UMM is equal to the sum of the lower bounds obtained from the memory bandwidth limitation and the latency overhead, they are optimal from the theoretical point of view.

1. Introduction

1.1 Background

The research of parallel algorithms has a long history of more than 40 years. Sequential algorithms have been developed mostly on the RAM (Random Access Machine)³⁾. In contrast, since there are a variety of connection methods and patterns between processors and memories, many parallel computing models have been presented and many parallel algorithmic techniques have been shown on them. The most well-studied parallel computing model is the PRAM (Parallel Random Access Machine)^{4),6),12)}, which consists of processors and a shared memory. Each processor on the PRAM can access any address of the shared memory in a time unit. The PRAM is a good parallel computing model in the sense that parallelism of each problem can be revealed by the performance of parallel algorithms on the PRAM. However, since the PRAM requires a shared memory that can be accessed by all processors in the same time, it is not feasible.

The GPU (Graphical Processing Unit), is a specialized circuit designed to accelerate computation for building and manipulating images⁷⁾⁻⁹⁾. Latest GPUs are designed for general purpose computing and can perform computation in applications traditionally handled by the CPU. Hence, GPUs have recently attracted the attention of many application developers^{7),11)}. NVIDIA provides a parallel computing architecture called *CUDA* (Compute Unified Device Architecture)²⁾, the computing engine for NVIDIA GPUs. CUDA gives developers access to the virtual instruction set and memory of the parallel computational elements in NVIDIA GPUs. In many cases, GPUs are more efficient than multicore processors¹⁰⁾, since they have hundreds of processor cores.

CUDA uses two types of memories in the NVIDIA GPUs: *the global memory* and *the shared memory*²⁾. The global memory is implemented as a off-chip DRAM, and has large capacity, say, 1.5-4 Gbytes, but its access latency is very long. The shared memory is an extremely fast on-chip memory with lower capacity, say, 16-64 Kbytes. The efficient usage of the global memory and the shared memory is a key for CUDA developers to accelerate applications using GPUs. In particular, we need to consider *the coalescing* of the global memory access and *the bank conflict* of the shared memory access^{1),9),10)}. To maximize the bandwidth between the GPU and the DRAM chips, the consecutive addresses of the global memory must be accessed in the same time. Thus, threads should perform coalesced access when they access to the global memory.

The address space of the shared memory is mapped into several physical memory banks. If two or more processor cores access to the same memory banks in the same time, the access requests are processed sequentially. Hence to maximize the memory access performance, processor cores should access to distinct memory banks to avoid the bank conflicts of the memory access.

1.2 Our Contribution: Introduction to the Discrete Memory Machine and the Unified Memory Machine

The first contribution of this paper is to introduce simple parallel memory bank machine models that capture the essential features of the coalescing of the global memory access and the bank conflict of the shared memory access. More specifically, we present two models, *the Discrete Memory Machine (DMM)* and *the Unified Memory Machine (UMM)*, which reflect the essential features of the

^{†1} Hiroshima University

shared memory and the global memory of NVIDIA GPUs.

The architectures of the DMM and the UMM are illustrated in Figure 1. In both architectures, the processing elements (PEs) are connected to the memory banks (MBs) through the memory management unit (MMU). A single address space of the memory is mapped to the MBs in an interleaved way such that the data of address i is stored in the $(i \bmod w)$ -th bank, where w is the number of MBs. The main difference of the two architectures is the connection of the address line between the MMU and the MBs, which can transfer an address value. In the DMM, the address lines connect the MBs and the MMU separately, while a single address line from the MMU is connected to the MBs in the UMM. Hence, in the UMM, the same address value is broadcast to every MB, and the same address of the MBs can be accessed in each time unit. On the other hand, different addresses of the MBs can be accessed in the DMM. Since the memory access of the UMM is more restricted than that of the DMM, the UMM is less powerful than the DMM.

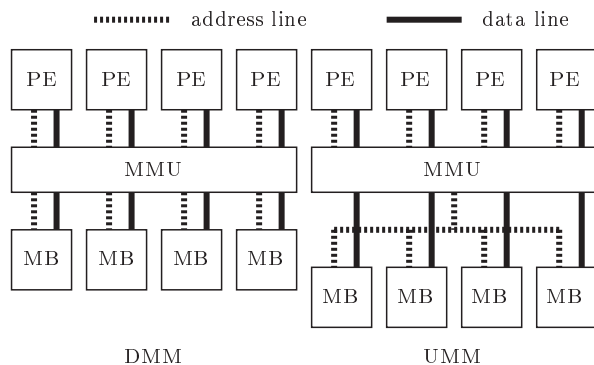


Fig. 1 The architectures of the DMM and the UMM

The performance of algorithms of the PRAM is usually evaluated using two parameters: the size n of the problem and the number p of processors. For example, it is well known that the sum of n numbers can be computed in $O(\frac{n}{p} + \log p)$ time on the PRAM⁴⁾. We will use additional two parameters, the width w

and the latency l of the memory when we evaluate the performance of algorithms on the DMM and on the UMM. The width w is the number of memory banks and the latency l is the number of time units to complete the memory access. Hence the performance of algorithms on the DMM and the UMM is evaluated as a function of n (the size of a problem), p (the number of processors), w (the width of a memory), and l (the latency of a memory). In typical NVIDIA GPUs, the width w of global and shared memory is 16 or 32. Also the latency l of the global memory is 400-800 clock cycles.

We also introduce *the bandwidth limited PRAM* (BPRAM). In the BPRAM of width w , any w processors out of the p processors can access to the memory in a time unit. Clearly, the BPRAM is less powerful than the PRAM and is more powerful than DMM and the UMM. Unlike the DMM and the UMM, the BPRAM has no restriction of the addresses of memory access and 1 memory access latency. We use the BPRAM to show the goodness of the performance of algorithms on the DMM and the UMM. If the computing time of an algorithm to solve some problem on the DMM or the UMM is almost the same as that on the BPRAM, we can say that the algorithm on the DMM or the UMM is close to optimal.

Please note that the DMM and the UMM are theoretical models of parallel computation, that capture the essential feature of the global memory and the shared memory of NVIDIA GPUs. NVIDIA GPUs have other features such as hierarchical architecture grid/block/thread, and the cache of the global memory. However, if these aspects are incorporated in our theoretical parallel models, they will be complicated and need more parameters. The development of algorithms on such complicated model may have too much non-essential and tedious optimizations. Thus, we have introduces two simple parallel models, the DMM and the UMM, which focuses on the memory accesses of the global memory and the shared memory of NVIDIA GPUs.

In⁵⁾, the authors have presented a GPU memory model and presented a cache-efficient FFT. However, their model focuses on the cache mechanism and ignores the coalescing and the bank conflict. Also, in¹³⁾, acceleration techniques for GPU have been discussed. Although they are taking care of the limited bandwidth of the global memory, the details of the memory architecture are not considered.

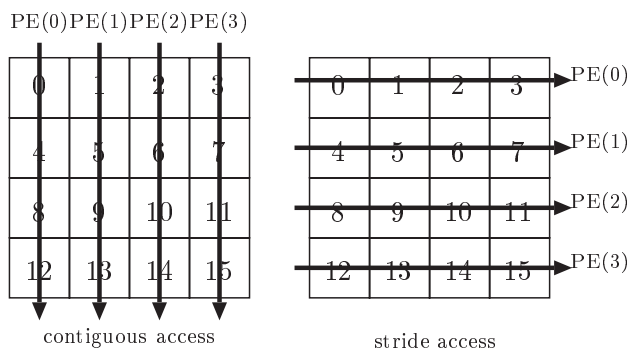


Fig. 2 The contiguous access and the stride access for $p = 4$ and $n = 16$.

As far as we know, this paper is the first work that introduces simple theoretical parallel computing models for GPUs. We believe that the development of algorithms on these models are useful to investigate algorithmic techniques for the GPUs.

1.3 Our Contribution: Fundamental Algorithms on the DMM and the UMM

The second contribution of this paper is to evaluate the performance of two memory access methods, *the contiguous access* and *the stride access* on the DMM and the UMM. The reader should refer to Figure 2 for illustrating these two access methods. We will show that the contiguous access of an array of size n can be done in $O(\frac{n}{w} + \frac{nl}{p})$ time units on the DMM and the UMM. Also, the contiguous access and the stride access can be done in $O(\frac{n}{w})$ time units on the BPRAM. Thus, the contiguous access of the DMM and the UMM is optimal when $wl \leq p$, because $\frac{n}{w} \geq \frac{nl}{p}$ if this is the case. Further, we will show that the stride access of the DMM can be done in $O(\frac{n}{w} \cdot \text{GCD}(\frac{n}{p}, w) + \frac{nl}{p})$ time units on the DMM, where $\text{GCD}(\frac{n}{p}, w)$ is the greatest common divisor of $\frac{n}{p}$ and w . Hence, the stride access of the DMM is optimal if $\frac{n}{p}$ and w are co-prime and $wl \leq p$. The stride access of the UMM can be done in $O(\min(n, \frac{n}{w} \cdot \frac{n}{p} + \frac{nl}{p}))$ time units. Hence, the stride access of the UMM needs an overhead of a factor of $\frac{n}{p}$.

From the results of the second and the third contributions, we have one im-

portant observation as follows. The factor $\frac{n}{w}$ in the computing time comes from *the bandwidth limitation* of the memory. It takes at least $\frac{n}{w}$ time units to access whole data in an array of size n from the memory bandwidth w . Also, the factor $\frac{nl}{p}$ comes from *the latency overhead*. From the memory access latency l , each processor cannot send next access request in l time units. It follows that, each processor can access to the memory once in l time units and each of the l time units can have expected $\frac{p}{l}$ access requests by processors. Hence, $\frac{nl}{p}$ time units are necessary to access all of the elements in an array of size n . Further, to hide the latency overhead factor $\frac{nl}{p}$ from the bandwidth limitation factor $\frac{n}{w}$, the number p of the processors must be no less than wl . We can confirm this fact from a different aspect. We can think that the memory access request are stored in a pipeline buffer of size l for each memory bank. Since we have w memory banks, we have wl pipeline registers to store memory access requests at all. Since at most one memory request per processor are stored in the wl pipeline registers, $wl \leq p$ must be satisfied to fill the pipeline registers full of memory access requests.

2. Parallel Memory Machines: DMM and UMM

Let us start with defining *PRAM* (Parallel Random Access Machine), the most popular shared memory parallel machine model. The PRAM consists of p processors and a shared memory. The shared memory is an array of memory cells, each of which can store a word of data. Each of the processors can select a memory cell in the array independently, and can perform read/write operation in a time unit. Please see⁴⁾ for the details of the PRAM.

We introduce a memory bandwidth limited PRAM. We assume that w memory cells can be read/written in a time unit. If more than w memory cells are accessed, the w operations are automatically serialized. More specifically, if p memory cells are accessed, w read/write operations performed in each time unit. and it takes $\lceil \frac{p}{w} \rceil$ time to complete the p read/write operations. We call such PRAM *the Bandwidth-limited PRAM (BPRAM)*.

We next introduce *the Discrete Memory Machine (DMM)* of width w and latency l . Let $m[i]$ ($i \geq 0$) denote a memory cell of address i in the memory. Let $B[j] = \{m[j], m[j+w], m[j+2w], m[j+3w], \dots\}$ ($0 \leq j \leq w-1$) denote *the j -th bank*. Clearly, a memory cell $m[i]$ is in the $(i \bmod w)$ -th memory bank.

We assume that memory cells in different banks can be accessed in a time unit, but no two memory cells in the same bank can be accessed in a time unit. Also, we assume that l time units are necessary to complete an access request and continuous requests are processed in a pipeline fashion through the memory management unit. Thus, it takes $k + l - 1$ time units to complete k continuous access requests to a particular bank.

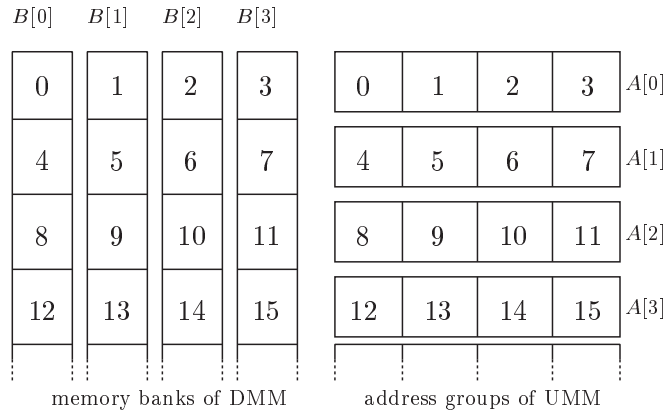


Fig. 3 Banks and address groups for $w = 4$

We assume that p processors are partitioned into $\frac{p}{w}$ groups of w processors called *warps*. More specifically, p processors are partitioned into $\frac{p}{w}$ warps $W(0), W(1), \dots, W(\frac{p}{w} - 1)$ such that $W(i) = \{\text{PE}(i \cdot w), \text{PE}(i \cdot w + 1), \dots, \text{PE}((i + 1) \cdot w - 1)\}$ ($0 \leq i \leq \frac{p}{w} - 1$). Warps are activated for memory access in turn, and w processors in a warp try to access the memory in the same time. In other words, $W(0), W(1), \dots, W(\frac{p}{w} - 1)$ are activated in a round-robin manner if at least one processor in a warp requests memory access. If no processor in a warp needs memory access, such warp is not activated and is skipped. When $W(i)$ is activated, w processor in $W(i)$ sends memory access requests, one request per processor, to the memory. We also assume that a processor cannot send a new memory access request until the previous memory access request is completed. Hence, if a processor send a memory access request, it must wait for l time units

to send a next memory access request.

Let us evaluate the time for memory access using Figure 4 on the DMM for $p = 8, w = 4$, and $l = 3$. Suppose that processors in $W(0)$ try to access $m[0], m[1], m[5]$, and $m[10]$, and those in $W(1)$ try to access $m[8], m[9], m[14]$, and $m[15]$. First, memory access requests to $m[0], m[1]$, and $m[10]$ are sent to the banks $B[0], B[1], B[2]$ first, and then access requests to $m[5]$ are sent to the bank $B[1]$. After that, memory access requests to $m[8], m[9], m[14], m[15]$ are sent to the bank $B[0], B[1], B[2], B[3]$. Since we have latency $l = 3$, all of these memory requests are completed in $3 + l - 1 = 5$ time units on the DMM.

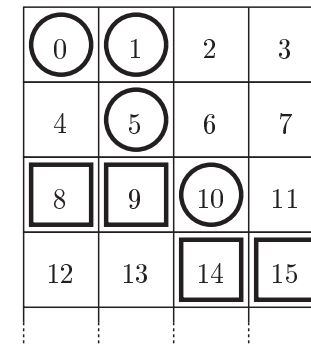


Fig. 4 An example of memory access

We next define the *Unified Memory Bank Machine (UMM)* for short) of width w as follows. Let $A[j] = \{m[j \cdot w], m[j \cdot w + 1], \dots, m[(j + 1) \cdot w - 1]\}$ denote the j -th address group. We assume that memory cells in the same address group are processed in the same time. However, if they are in the different groups, one time unit is necessary for each of the groups. Also, similarly to the DMM, p processors are partitioned into warps and each warp access to the memory in turn.

Again, let us evaluate the time for memory access using Figure 4 on the UMM for $p = 8, w = 4$, and $l = 3$. To complete the memory access requests by $W(0)$, those to $A[0]$ is performed for $m[0]$ and $m[1]$. After that, access requests to $A[1]$ is

performed for $m[5]$ and then that to $A[2]$ is performed for $m[10]$. Next, memory access requests by $W(1)$ are processed. First, memory access to $A[2]$ is requested for $m[8]$ and $m[9]$, and then that to $A[3]$ is requested for $m[14]$ and $m[15]$. Since we have latency $l = 3$, all of these memory request are processed in $5 + l - 1 = 7$ time units on the UMM.

3. Sequential memory access operations

We begin with simple operations to evaluate the potentiality of the DMM and the UMM. Let p and w be the number of processors and the width of the machines. We assume that an array m of size n is arranged in the memory. Let $m[i]$ ($0 \leq i \leq n - 1$) denote the i -th word of the memory. We assume that $n \gg p > w$. We consider two access operations to the memory such that each of the p processors reads $s = \frac{n}{p}$ memory cells out of the n memory cells. Again, the readers should refer to Figure 2. In the *contiguous access*, the first p memory cells are accessed by the p processors. Next, the second p memory cells are accessed. In this way, all n memory cells are accessed in turn. In the *stride access*, the memory is partitioned into p groups of $\frac{n}{p}$ consecutive memory cells each. A processor is assigned to each group and the assigned processor accesses memory cells in the group sequentially. The contiguous access and the stride access are written as follows.

[Contiguous Access]

for $t \leftarrow 0$ to $s - 1$

 for $i \leftarrow 0$ to $p - 1$ do in parallel

 PE(i) accesses to $m[i + t \cdot p]$

[Stride Access]

for $t \leftarrow 0$ to $s - 1$

 for $i \leftarrow 0$ to $p - 1$ do in parallel

 PE(i) accesses to $m[i \cdot s + t]$

In the contiguous access, w processors in each warp access memory cells in different memory banks. Hence, the memory access by a warp takes l time unit. Also, the memory access by a warp is processed in every 1 time unit. Since we have $\frac{n}{w}$ warps, p memory cells can be accessed by p processors in $\frac{n}{w} + l - 1$ time

units. Consequently, the contiguous access takes $(\frac{n}{w} + l - 1) \cdot s = O(\frac{n}{w} + \frac{nl}{p})$ time units on the DMM. In the contiguous access on the UMM, each warp access to the memory cells in the same address group. Thus, the memory access by a warp takes l time unit and the whole contiguous access runs in $O(\frac{n}{w} + \frac{nl}{p})$.

The performance analysis of the stride access on the DMM is a bit complicated. Let us start with a simple case: $s = w$. In this case, the p processors access to p memory cells $m[t], m[w+t], m[2w+t], \dots, m[(p-1)w+t]$ for each t ($0 \leq t \leq w-1$). Unfortunately, these memory cells are in the same memory bank $B[t]$. Hence, memory access by a warp takes $w + l - 1$ time units and the memory access to these p memory cells takes $w \cdot \frac{p}{w} + l - 1 = p + l - 1$ time units. Thus, the stride access when $s = w$ takes at least $(p + l - 1) \cdot \frac{n}{p} = O(n + \frac{nl}{p})$ time units.

Next, let us consider general case. The w processors in the first warp access to $m[t], m[s+t], m[2s+t], \dots, m[(w-1)s+t]$ for each t ($0 \leq t \leq w-1$). These w memory cells are allocated in the banks $B[t \bmod w], B[(s+t) \bmod w], B[(2s+t) \bmod w], \dots, B[((w-1)s+t) \bmod w]$. Let $L = \text{LCM}(s, w)$ and $G = \text{GCD}(s, w)$ be the Least Common Multiple and the Greatest Common Divisor of s and w , respectively. From the basic number theory, it should be clear that $t \bmod w = (\frac{L}{s} \cdot s + t) \bmod w$, and the values of $t \bmod w, (s+t) \bmod w, \dots, ((\frac{L}{s} - 1) \cdot s + t) \bmod w$ are distinct. Thus, the w memory cells are in the $\frac{L}{s} = \frac{w}{G}$ banks $B[t \bmod w], B[(s+t) \bmod w], B[(2s+t) \bmod w], \dots, B[((\frac{w}{G} - 1)s + t) \bmod w]$ equally, and each bank has G memory cells of the w memory cells. Hence, the w processor in a warp takes $G + l - 1$ time units for each t , and the p processors takes $G \cdot \frac{p}{w} + l - 1$ time units for each t . Therefore, the DMM takes $(G \cdot \frac{p}{w} + l - 1) \cdot s = O(\frac{nG}{w} + \frac{nl}{p})$ time to complete the stride access. If $s = w$ then $G = w$ and the time for the stride access is $O(n + \frac{nl}{p})$. If s and w are co-prime, $G = 1$ and the stride access takes $O(\frac{n}{w} + \frac{nl}{p})$ time units.

Finally, we will evaluate the computing time of the stride access on the UMM. If $s \geq w$ (i.e. $n \geq pw$), then the w memory cells are accessed by w processors in a warp are in the different address group. Thus, w processors access to w memory cells in $w + l - 1$ time units, and the stride access takes $(w \cdot \frac{p}{w} + l - 1) \cdot \frac{n}{p} = O(n + \frac{nl}{p})$ time. When $s < w$ (i.e. $n < pw$), the w memory cells accessed by w processors in a warp are in at most $\lceil \frac{(w-1)s+1}{w} \rceil \leq s$ address groups. Hence, the stride access

Table 1 The running time for the contiguous access and the stride access

Operation	PRAM	BPRAM	DMM	UMM
Contiguous Access	$O(\frac{n}{p})$	$O(\frac{n}{w})$	$O(\frac{n}{w} + \frac{nl}{p})$	$O(\frac{n}{w} + \frac{nl}{p})$
Stride Access	$O(\frac{n}{p})$	$O(\frac{n}{w})$	$O(\frac{n}{w} \cdot G + \frac{nl}{p})$	$O(\min(n, \frac{n}{w} \cdot \frac{n}{p}) + \frac{nl}{p})$

$n = \#data, p = \#processors, w = \text{memory bandwidth}, G = \text{GCD}(\frac{n}{p}, w)$

by p processors for each t takes at most $s \cdot \frac{p}{w} + l - 1$ time, and thus, the whole stride access takes $(\frac{sp}{w} + l - 1) \cdot \frac{n}{p} = O(\frac{ns}{w} + \frac{nl}{p})$ time. Consequently, the stride access can be completed in $O(\min(n, \frac{ns}{w}) + \frac{nl}{p})$ for all $s = \frac{n}{p}$. Finally, we have,

Theorem 1 The contiguous access and the stride access on the PRAM, the BPRAM, the DMM, and the UMM can be done in time units shown in Table 1.

4. Transpose of a two dimensional array

Suppose that two dimensional array a of size $\sqrt{n} \times \sqrt{n}$ is arranged in the memory. We assume that $a[i][j]$ ($0 \leq i, j \leq \sqrt{n} - 1$) is located in the $(i \cdot \sqrt{n} + j)$ th memory cell of the memory. The transpose of a two dimensional array is a task to move a word of data stored in $a[i][j]$ to $a[j][i]$ for all $(0 \leq i, j \leq \sqrt{n} - 1)$.

Let us start with a straightforward transpose algorithm using the contiguous access and the stride access. The following algorithm transposes a two dimensional array a size $\sqrt{n} \times \sqrt{n}$.

[Straightforward transpose algorithm]

for $t \leftarrow 0$ to $\frac{n}{p} - 1$

for $i \leftarrow 0$ to $p - 1$ do in parallel

$j \leftarrow (t \cdot p + i) / \sqrt{n}$

$k \leftarrow (t \cdot p + i) \bmod \sqrt{n}$

PE(i) performs exchange $a[j][k] \leftrightarrow a[k][j]$

For each t , the processors performs the contiguous access and the stride access, which takes $O(1)$ time on the PRAM and $O(\frac{p}{w})$ time on the BPRAM. Hence, the PRAM and the BPRAM can transpose a two dimensional array in $O(\frac{n}{p})$ time units and $O(\frac{n}{w})$ time units, respectively.

Since the straightforward algorithm contains the stride access, it is not difficult to see that the DMM and the UMM takes $O(\frac{n}{w} \cdot \text{GCD}(\sqrt{n}, w) + \frac{nl}{p})$ time units

and $O(\min(n, \frac{n}{w} \cdot \frac{n}{p}) + \frac{nl}{p})$ time units for transposing a two dimensional array, respectively. On the DMM, $\text{GCD}(\sqrt{n}, w) = w$ if \sqrt{n} is divisible by w . If this is the case, the transpose takes $O(n)$ time units the DMM. We will show that, regardless of the value of n , the transpose can be done in $O(\frac{n}{w} + \frac{nl}{p})$ time both on the DMM and on the UMM.

We first show an efficient transposing algorithm on the DMM. The key idea is to access the array in diagonal fashion. We use a two dimensional array b of size $n \times n$ as a work space.

[Transpose by the diagonal access on the DMM]

for $t \leftarrow 0$ to $\frac{n}{p} - 1$

for $i \leftarrow 0$ to $p - 1$ do in parallel

$j \leftarrow (t \cdot p + i) / \sqrt{n}$

$k \leftarrow (t \cdot p + i) \bmod \sqrt{n}$

PE(i) performs $b[j][k] \leftarrow a[j][k]$

for $t \leftarrow 0$ to $\frac{n}{p} - 1$

for $i \leftarrow 0$ to $p - 1$ do in parallel

$j \leftarrow (t \cdot p + i) / \sqrt{n}$

$k \leftarrow (t \cdot p + i) \bmod \sqrt{n}$

PE(i) performs

$a[(j + k) \bmod \sqrt{n}][k] \leftarrow b[k][(j + k) \bmod \sqrt{n}]$

The readers should refer to Figure 5 for illustrating the indexes of processors reading from memory cells in b and writing from memory cells in a for $n = p = 16$ and $w = 4$. From the figure, we can confirm that processors PE($j \cdot 4 + 0$), PE($j \cdot 4 + 1$), PE($j \cdot 4 + 2$), PE($j \cdot 4 + 3$) read from memory cells in diagonal location of b and write to memory cells in diagonal location of a for every j ($0 \leq j \leq 3$). Thus, reading and writing memory banks by w processors in a warp are different. Hence p processors can copy p memory cells in $\frac{p}{w} + l$ time units and thus the total computing time is $(\frac{p}{w} + l) \cdot \frac{n}{p} = O(\frac{n}{w} + \frac{nl}{p})$ time. Therefore, we have,

Lemma 2 The transpose of a two dimensional array of size $\sqrt{n} \times \sqrt{n}$ can be done in $O(\frac{n}{w} + \frac{nl}{p})$ time using p processors on the DMM with memory width w and latency l .

Next, we will show that the transpose of a two dimensional array can be also

b				a			
PE(0)	PE(4)	PE(8)	PE(12)	PE(0)	PE(13)	PE(10)	PE(7)
PE(13)	PE(1)	PE(5)	PE(9)	PE(4)	PE(1)	PE(14)	PE(11)
PE(10)	PE(14)	PE(2)	PE(6)	PE(8)	PE(5)	PE(2)	PE(15)
PE(7)	PE(11)	PE(15)	PE(3)	PE(12)	PE(9)	PE(6)	PE(3)

The indices of processors reading from memory cells of b The indices of processors writing in memory cells of a

Fig. 5 Transposing on the DMM

done in $O(\frac{n}{w} + \frac{nl}{p})$ on the UMM if every processor has a local memory that can store w words. As a preliminary step, we will show that the UMM can transpose a two dimensional array of size $w \times w$ in $O(w + l)$ time units using w processors with each processor having a local storage of size w . We assume that each processor has local memory that can store w words. Let l_i denote the local memory of PE(i), and $l_i[0], l_i[1], \dots, l_i[w - 1]$ can store a word of data.

[Transpose by the rotating technique on the UMM]

for $t \leftarrow 0$ to $w - 1$

for $i \leftarrow 0$ to $w - 1$ do in parallel

PE(i) performs $l_i[t] \leftarrow a[t][(t + i) \bmod w]$

for $t \leftarrow 0$ to $w - 1$

for $i \leftarrow 0$ to $w - 1$ do in parallel

PE(i) performs $a[t][(t - i) \bmod w] \leftarrow l_i[(t - i) \bmod w]$

Let (i, j) denote the value stored in $a[i][j]$ initially. The readers should refer to Figure 5 for illustrating how these values are transposed.

Let us confirm that the algorithm above correctly transpose two dimensional array a . In other words, we will show that, when the algorithm terminates, $a[i][j]$ stores (j, i) . It should be clear that, the value stored in $l_i[t]$ is $(t, (t + i) \bmod w)$. Since $((t - i) \bmod w, t)$ is stored in $l_i[(t - i) \bmod w]$, it is also stored

in $a[t][(t - i) \bmod w]$ when the algorithm terminates. Thus, every $a[i][j]$ ($0 \leq i, j \leq w - 1$) stores (j, i) . This completes the proof of the correctness of our transpose algorithm on the UMM.

Let us evaluate the computing time. In the reading operation $l_i[t] \leftarrow a[t][(t + i) \bmod w]$, w memory cells $a[t][(t + 0 \bmod w)], a[t][(t + 1 \bmod w)], \dots, a[t][(t + w - 1 \bmod w)]$ are in the different memory banks. Also, in the writing operation $a[t][(t - i) \bmod w] \leftarrow l_i[(t - i) \bmod w]$, w memory cells $a[t][(t - 0 \bmod w)], a[t][(t - 1 \bmod w)], \dots, a[t][(t - (w - 1) \bmod w)]$ are in the different memory banks. Thus, each reading and writing operation can be done in l time units and this algorithm runs in $O(w + l)$ time units.

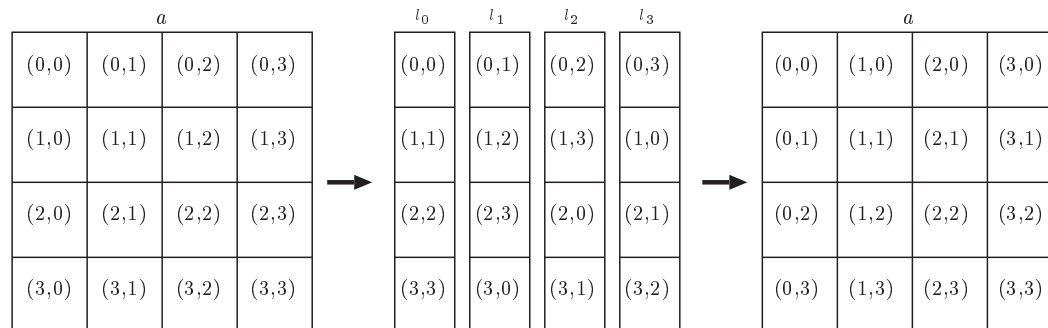
The transpose of a larger two dimensional array of size $\sqrt{n} \times \sqrt{n}$ can be done by repeating the transpose of two dimensional array of size $w \times w$. More specifically, the two dimensional array is partitioned into $\frac{\sqrt{n}}{w} \times \frac{\sqrt{n}}{w}$ subarrays of size $w \times w$. Let $A[i][j]$ ($0 \leq i, j \leq \frac{n}{w} - 1$) denote the subarray of size $w \times w$. First, each subarray $A[i][j]$ are transposed independently using w processors. After that, the corresponding words of $A[i][j]$ and $A[j][i]$ are swapped for all i and j in an obvious way.

Let us evaluate the computing time to complete the transpose of a $\sqrt{n} \times \sqrt{n}$ two dimensional array. Suppose that we have p ($\leq \frac{n}{w}$) processors and partition the p processors into $\frac{p}{w}$ groups with w processors each. We assign $\frac{n}{w^2} / \frac{p}{w} = \frac{n}{pw}$ subarrays to each group of w processors. The p processors can transpose $\frac{p}{w}$ subarrays in parallel in $O(w \cdot (\frac{p}{w} + l)) = O(p + wl)$ time units. This transpose of $\frac{p}{w}$ subarrays are repeated $\frac{n}{pw}$ times, the total computing time for the subarray transpose is $\frac{n}{pw} \cdot O(p + wl) = O(\frac{n}{w} + \frac{nl}{p})$ time units. Clearly, the swap operation of subarrays can be also done in $O(\frac{n}{w} + \frac{nl}{p})$ time units. Thus we have,

Lemma 3 The transpose of a two dimensional array of size $\sqrt{n} \times \sqrt{n}$ can be done in $O(\frac{n}{w} + \frac{nl}{p})$ time on the UMM with each processor having local memory of w words.

5. Conclusion

In this paper, we have introduced two parallel memory machines, the Discrete Memory Machine (DMM) and the Unified Memory Machine (UMM). We first

**Fig. 6** Transposing on the UMM

evaluated the computing time of the contiguous access and the stride access of the memory on the DMM and the UMM. We also presented an algorithm to transpose a two dimensional array on the DMM and the UMM.

Although the DMM and the UMM are simple, they capture the characteristic of the shared memory and the global memory of NVIDIA GPUs. Thus, these two parallel computing models are promising for developing algorithmic techniques for NVIDIA GPUs. As a future work, we plan to implement various parallel algorithms developed for the PRAM so far on the DMM and on the UMM. Also, NVIDIA GPUs have small shared memory and large global memory. Thus, it is also interesting to consider a hybrid memory machine such that processors are connected to a small memory of DMM and a large UMM.

References

- 1) *NVIDIA CUDA C Best Practice Guide Version 3.1*, 2010.
- 2) *NVIDIA CUDA C Programming Guide Version 4.0*, 2011.
- 3) Alfred V. Aho, Jeffrey D. Ullman, and John E. Hopcroft. *Data Structures and Algorithms*. Addison Wesley, 1983.
- 4) A.Gibbons and W.Rytter. *Efficient Parallel Algorithms*. Cambridge University Press, 1988.
- 5) NagaK Govindaraju, Scott Larsen, Jim Gray, and Dinesh Manocha. A memory model for scientific algorithms on graphics processors. In *Proc. of the ACM/IEEE Conference on Supercomputing*, pages 6–6, 2006.
- 6) Ananth Grama, George Karypis, Vipin Kumar, and Anshul Gupta. *Introduction to Parallel Computing*. Addison Wesley, 2003.
- 7) Wen-mei W. Hwu. *GPU Computing Gems Emerald Edition*. Morgan Kaufmann, 2011.
- 8) Yasuaki Ito, Kouhei Ogawa, and Koji Nakano. Fast ellipse detection algorithm using hough transform on the GPU. In *Proc. of International Conference on Networking and Computing*, pages 313–319, Dec. 2011.
- 9) Duhu Man, Kenji Uda, Yasuaki Ito, and Koji Nakano. A GPU implementation of computing euclidean distance map with efficient memory access. In *Proc. of International Conference on Networking and Computing*, pages 68–76, Dec. 2011.
- 10) Duhu Man, Kenji Uda, Hironobu Ueyama, Yasuaki Ito, and Koji Nakano. Implementations of a parallel algorithm for computing euclidean distance map in multicore processors and GPUs. *International Journal of Networking and Computing*, 1:260–276, July 2011.
- 11) Kazufumi Nishida, Yasuaki Ito, and Koji Nakano. Accelerating the dynamic programming for the matrix chain product on the GPU. In *Proc. of International Conference on Networking and Computing*, pages 320–326, Dec. 2011.
- 12) Michael J. Quinn. *Parallel Computing: Theory and Practice*. McGraw-Hill, 1994.
- 13) Shane Ryoo, Christopher I. Rodrigues, Sara S. Baghsorkhi, Sam S. Stone, David B. Kirk, and Wen mei W. Hwu. Optimization principles and application performance evaluation of a multithreaded gpu using cuda. In *Proc. of the 13th ACM SIGPLAN Symposium on Principles and practice of parallel programming*, pages 73–82, 2008.