

Surrounding Theorem: Developing Parallel Programs for Matrix-Convolutions

Kento Emoto, Kiminori Matsuzaki, Zhenjiang Hu, and Masato Takeichi

Department of Mathematical Informatics,
University of Tokyo

{emoto, kmatsu, hu, takeichi}@ipl.t.u-tokyo.ac.jp

Abstract. Computations on two-dimensional arrays such as matrices and images are one of the most fundamental and ubiquitous things in computational science and its vast application areas, but development of efficient parallel programs on two-dimensional arrays is known to be hard. To solve this problem, we have proposed a skeletal framework on two-dimensional arrays based on the theory of constructive algorithmics. It supports users, even with little knowledge about parallel machines, to develop systematically both correct and efficient parallel programs on two-dimensional arrays. In this paper, we apply our framework to the matrix-convolutions often used in image filters and difference methods. We show the efficacy of the framework by giving a general parallel program for the matrix-convolutions described with the skeletons, and a theorem that optimizes the general program into an application-specific one.

1 Introduction

Computations on two-dimensional arrays, such as matrix computations, image processing, and difference methods, are both fundamental and ubiquitous in scientific computations and other application areas [7, 15, 11]. However, development of efficient parallel programs on two-dimensional arrays is known to be a hard task due to the necessity of considering data allocation, synchronization and communication between processors. *Skeletal parallel programming* is one promising solution to the situation [5, 16]. In this model, users build parallel programs by composing ready-made components (called *skeletons*) implemented efficiently in parallel for various parallel architectures. Since low-level parallelism is concealed in the skeletons, users can obtain a comparatively efficient parallel program without needing technical details of parallel computers or being conscious of parallelism explicitly.

We have proposed a skeletal framework on two-dimensional arrays [9], based on the theory of constructive algorithmics (also known as *Bird-Meertens Formalism*) [2, 4]. Our framework provides users, even with little knowledge about parallel machines, with a concise way to describe safe and efficient parallel computations over two-dimensional arrays, and theorems for deriving and optimizing parallel programs. The main features of our framework are: (1) *a novel use of the abide-tree representation* [2] in developing parallel programs for manipulating two-dimensional arrays; (2) *a strong support* for systematic development of both efficient and correct parallel programs in a highly abstract way; (3) *an efficient implementation* of basic skeletons in C++ and MPI on PC

clusters, guaranteeing that programs composed with these parallel skeletons can run efficiently in parallel. To develop parallel programs in our framework, users construct a simple and general program that covers a class of problems, derive its efficient version using general techniques such as fusion, tupling and generalization, and then instantiate the general program to solve concrete problems. Usually, this derivation is summarized as a theorem (tool).

In this paper, we give a domain-specific tool and show the efficacy of the framework. We focus on computations known as matrix-convolutions [12], in which each element in the resulting array depends on its surrounding elements. This set of computations includes important and fundamental problems such as image filters, difference methods and the N -body problem (although this last problem seems more difficult than the others, it merely refers to not only the nearest neighbors but all the surrounding elements). The most general form *mconv* is described with three components:

$$mconv\ f\ shrink = \text{map}\ f \circ \text{map}\ shrink \circ \text{surrounds}.$$

Here, *surrounds* gathers all the surrounding elements for each element, *shrink* picks the necessary parts up from those gathered elements, and *f* calculates the resulting element from them. This general form is parameterized by the two functions *shrink* and *f*, and users can solve many problems by specifying suitable ones. For example, users can develop a sharpen-filter by choosing the function *shrink* that reduces the surroundings into a 3×3 matrix, and the function *f* that calculates the weighted sum of the nine values. We can further optimize instances of the general program to application-specific ones with the *surrounding theorem*. The main contributions of this paper are as follows.

- We show the general parallel program for the matrix-convolutions described with parallel skeletons. Users can solve their problems as its instance.
- We give the *surrounding theorem* that enables users to get an efficient program easily. The experimental results show that the derived program can be executed efficiently in parallel.

Technical details of this paper are available in the master's thesis [8].

2 Notations

Notation in this paper follows that of Haskell [3], a pure functional language that can describe both algorithms and algorithmic transformation concisely.

Function application is denoted by a space and the argument may be written without brackets. Thus, $f\ a$ means $f(a)$ in ordinary notation. Functions are curried, i.e. functions take one argument and return a function or a value, and the function application associates to the left. Thus, $f\ a\ b$ means $(f\ a)\ b$. The function application binds more strongly than any other operator, so $f\ a \otimes b$ means $(f\ a) \otimes b$, but not $f\ (a \otimes b)$. Function composition is denoted by \circ , so $(f \circ g)\ x = f\ (g\ x)$ from its definition. Binary operators can be used as functions by sectioning as follows: $a \oplus b = (a \oplus)\ b = (\oplus)\ a = (\oplus)\ a\ b$. Two binary operators \ll and \gg are defined by $a \ll b = a$, $a \gg b = b$. Pairs are Cartesian products of plural data, written like (x, y) . A function that applies functions f and g respectively to the elements of a pair (x, y) is denoted by $(f \times g)$. Thus, $(f \times g)\ (x, y) = (f\ x, g\ y)$.

3 Skeletal Framework on Two-Dimensional Arrays

In this section, we introduce our parallel skeletal framework on two-dimensional arrays [9] based on the theory of constructive algorithmics [2, 4].

3.1 Abide-Trees for Two-Dimensional Arrays

To represent two-dimensional arrays, we define the abide-trees, which are built up by three constructors $|\cdot|$ (singleton), \ominus (above) and \oplus (beside) following the idea in [2].

$$\begin{aligned} \text{data } AbideTree \alpha = & |\cdot| \alpha \\ & | (AbideTree \alpha) \ominus (AbideTree \alpha) \\ & | (AbideTree \alpha) \oplus (AbideTree \alpha) \end{aligned}$$

Here, $|\cdot| a$, or abbreviated as $|a|$, means a singleton array of a , i.e. a two-dimensional array of a single element a . For two-dimensional arrays x and y of the same width, $x \ominus y$ means that x is located above y . Similarly, for arrays x and y of the same height, $x \oplus y$ means that x is located on the left of y . Moreover, \ominus and \oplus are associative operators and satisfy the following *abide* (a coined term from above and beside) property.

Definition 1 (Abide Property). *Two binary operators \oplus and \otimes are said to satisfy the abide property or to be abiding, if the following equation is satisfied:*

$$(x \otimes u) \oplus (y \otimes v) = (x \oplus y) \otimes (u \oplus v).$$

In the rest of the paper, we will assume that x has the same width as y when $x \ominus y$ appears, and that u has the same height as v for $u \oplus v$.

Note that one two-dimensional array may be represented by many abide-trees, but these abide-trees are equivalent because of the abide property of \ominus and \oplus . For example, we can express the following 2×2 two-dimensional array by two equivalent abide-trees.

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \Rightarrow \begin{cases} (|1| \oplus |2|) \ominus (|3| \oplus |4|) \\ (|1| \ominus |3|) \oplus (|2| \ominus |4|) \end{cases}$$

This is in sharp contrast to the quadtree representation of matrices [10], which does not allow such freedom.

From the theory of constructive algorithmics [4], it follows that each constructively built-up data structure (i.e., algebraic data structure) is equipped with a powerful computation pattern called homomorphism.

Definition 2 ((Abide-tree) Homomorphism). *A function h is said to be an abide-tree homomorphism, if it is defined as follows for a function f and binary operators \oplus, \otimes .*

$$\begin{aligned} h |a| &= f a \\ h (x \ominus y) &= h x \oplus h y \\ h (x \oplus y) &= h x \otimes h y \end{aligned}$$

For notational convenience, we write $\langle f, \oplus, \otimes \rangle$ to denote h . When it is clear from the context, we just call $\langle f, \oplus, \otimes \rangle$ homomorphism. Note that \oplus and \otimes in $\langle f, \oplus, \otimes \rangle$ should be associative and satisfy the abide property, inheriting the properties of \ominus and \oplus .

Intuitively, a homomorphism $\langle f, \oplus, \otimes \rangle$ is a function to replace the constructors $|\cdot|$, \ominus and \oplus in an input abide-tree by f , \oplus and \otimes respectively.

$$\begin{array}{l}
\text{map } f \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mn} \end{pmatrix} = \begin{pmatrix} f \ x_{11} & \cdots & f \ x_{1n} \\ \vdots & \ddots & \vdots \\ f \ x_{m1} & \cdots & f \ x_{mn} \end{pmatrix} \\
\text{reduce}(\oplus, \otimes) \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mn} \end{pmatrix} = \begin{pmatrix} (x_{11} \otimes \cdots \otimes x_{1n}) \oplus \\ \vdots \\ (x_{m1} \otimes \cdots \otimes x_{mn}) \end{pmatrix} \\
\text{zipwith } f \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mn} \end{pmatrix} \begin{pmatrix} y_{11} & \cdots & y_{1n} \\ \vdots & \ddots & \vdots \\ y_{m1} & \cdots & y_{mn} \end{pmatrix} = \begin{pmatrix} f \ x_{11} \ y_{11} & \cdots & f \ x_{1n} \ y_{1n} \\ \vdots & \ddots & \vdots \\ f \ x_{m1} \ y_{m1} & \cdots & f \ x_{mn} \ y_{mn} \end{pmatrix} \\
\text{scan}(\oplus, \otimes) \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mn} \end{pmatrix} = \begin{pmatrix} y_{11} & \cdots & y_{1n} \\ \vdots & \ddots & \vdots \\ y_{m1} & \cdots & y_{mn} \end{pmatrix} \quad \text{where } y_{ij} = \begin{pmatrix} (x_{11} \otimes \cdots \otimes x_{1j}) \oplus \\ \vdots \\ (x_{i1} \otimes \cdots \otimes x_{ij}) \end{pmatrix} \\
\text{scanr}(\oplus, \otimes) \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mn} \end{pmatrix} = \begin{pmatrix} z_{11} & \cdots & z_{1n} \\ \vdots & \ddots & \vdots \\ z_{m1} & \cdots & z_{mn} \end{pmatrix} \quad \text{where } z_{ij} = \begin{pmatrix} (x_{ij} \otimes \cdots \otimes x_{in}) \oplus \\ \vdots \\ (x_{mj} \otimes \cdots \otimes x_{mn}) \end{pmatrix}
\end{array}$$

Fig. 1. Intuitive Definition of Parallel Skeletons on Two-Dimensional Arrays

Table 1. Parallel Complexity of the Skeletons for a Two-Dimensional Array of $n \times n$

| | P processors | n^2 processors |
|--------------|----------------------------------|------------------|
| map, zipwith | $O(n^2/P)$ | $O(1)$ |
| reduce | $O(n^2/P + \log P)$ | $O(\log n)$ |
| scan, scanr | $O(n^2/P + \sqrt{n^2/P} \log P)$ | $O(\log n)$ |

3.2 Parallel Skeletons on Two-Dimensional Arrays

We introduce the parallel skeletons **map**, **reduce**, **zipwith**, **scan** and **scanr** for manipulating two-dimensional arrays. In the theory of constructive algorithmics [2, 4], these functions are known to be the most fundamental computation components for manipulating algebraic data structures and for being glued together to express complicated computations. Intuitive definitions of the skeletons are shown in Fig. 1. All the skeletons are implemented efficiently in parallel and their costs are shown in Table 1.

The skeletons **map** and **reduce** are two special cases of homomorphism. The skeleton **map** applies a function f to each element of a two-dimensional array while keeping the shape of the structure. The skeleton **reduce** collapses a two-dimensional array to a value using two abiding binary operators \oplus and \otimes . They are defined formally as $\text{map } f = (\llbracket \cdot \rrbracket \circ f, \oplus, \otimes)$, and $\text{reduce}(\oplus, \otimes) = (\llbracket id, \oplus, \otimes \rrbracket)$.

The skeleton **zipwith**, an extension of **map**, takes two arrays of the same shape, applies a function f to corresponding elements of the arrays and returns a new array of the same shape. The skeletons **scan** and **scanr**, extensions of **reduce**, hold all values generated in reducing an array by **reduce**. The **scan** generates the result of reducing

```

int sharpen_filter(int **b, int **a, int n, int m){
  for(int i = 0; i < m; i++){
    for(int j = 0; j < n; j++){
      b[i][j] = f(a[i][j], a[i-1][j], a[i+1][j], a[i][j+1], a[i][j-1],
                  a[i-1][j+1], a[i-1][j-1], a[i+1][j+1], a[i+1][j-1]);
    }
  }
  int f(int c, int n, int s, int e, int w, int ne, int nw, int se, int sw){
    return 5*c + (-1)*n + (-1)*s + (-1)*e + (-1)*w + 0*ne + 0*nw + 0*se + 0*sw;}
  }

```

Fig. 2. C++ Code of the Sharpen Filter (Sequential Program)

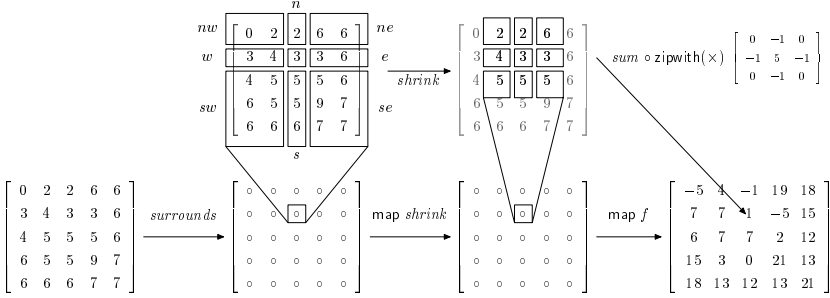


Fig. 3. An Image of the Sharpen Filter in the General Program

the upper-left subarray, while the `scanr` generates that of the lower-right subarray. We omit the formal definition of `zipwith`, `scan` and `scanr` for the space limitation.

4 Developing Parallel Programs for Matrix-Convolutions

In this section, focusing on the matrix-convolutions such as image filters and difference methods, we give the general form described with parallel skeletons, and then give the theorem to get optimized program from the general form.

The matrix-convolution is computation in which each element of the resulting array depends on the surrounding elements. For example, the sharpen-filter that sharpens the input image is one instance of the matrix-convolution. A pixel of the resulting image is the weighted sum of the surrounding pixels of the input image. Similarly, the difference method is another instance of matrix-convolution since it calculates the new value of each point from the old values of the surrounding points. We show a code in C++ for the sharpen-filter in Fig. 2, to give a concrete image of the problems dealt with here.

4.1 A General Form Described with Parallel Skeletons

As argued in the introduction, the most general form of this kind of computation is thought to consist of three components: gathering all the surrounding elements of each element to it, shrinking those to the necessary amount, and applying a function to get a new element from them. Thus, the program is described as follows:

$$mconv \ f \ shrink = \text{map } f \circ \text{map } shrink \circ \text{surrounds} .$$

The idea of our general form is illustrated in Fig. 3 that shows an image of execution of the sharpen-filter: (1) *surrounds* gathers all the surrounding elements for each element, (2) *shrink* picks the necessary parts up from those gathered elements, and (3) *f* calculates the resulting element from them. This general form has clear correspondences to the code in Fig. 2. The function *f* corresponds to \mathbb{f} of the code, *shrink* corresponds to which elements are the arguments passed to *f*, and *surrounds* corresponds to *for*-loops. Thus, users can easily write their programs using the general form.

This general form is parameterized by the two functions *shrink* and *f*, and users can solve many problems by specifying application-specific ones, as shown below. The function *surrounds*, which is commonly used in those problems, has two-phase calculation as follows: (1) calculation of the parts of the northwest (i.e. *c*, *n*, *w* and *nw*) by *scan*, and (2) that of the other parts by *scanr*. Its definition is as follows.

$$\begin{aligned}
\text{surrounds} &= \text{scanr}(\oplus_r, \otimes_r) \circ \text{map } f_r \circ \text{scan}(\oplus_f, \otimes_f) \circ \text{map } f_f \\
\text{where} \\
f_f a &= (a, \text{Nil}, \text{Nil}, \text{Nil}) \\
(c_a, n_a, w_a, nw_a) \oplus_f (c_b, n_b, w_b, nw_b) &= (\underbrace{c_b}_c, \underbrace{n_a \oplus |c_a| \oplus n_b}_n, \underbrace{w_b}_w, \underbrace{nw_a \oplus w_a \oplus nw_b}_{nw}) \\
(c_a, n_a, w_a, nw_a) \otimes_f (c_b, n_b, w_b, nw_b) &= (\underbrace{c_b}_c, \underbrace{n_b}_n, \underbrace{w_a \phi |c_a| \phi w_b}_w, \underbrace{nw_a \phi n_a \phi nw_b}_{nw}) \\
f_r (c, n, w, nw) &= (c, n, \text{Nil}, \text{Nil}, w, \text{Nil}, nw, \text{Nil}, \text{Nil}) \\
(c_a, n_a, s_a, e_a, w_a, ne_a, nw_a, se_a, sw_a) \oplus_r (c_b, n_b, s_b, e_b, w_b, ne_b, nw_b, se_b, sw_b) \\
&= (\underbrace{c_a}_c, \underbrace{n_a}_n, \underbrace{s_a \oplus |c_b| \oplus s_b}_s, \underbrace{e_a}_e, \underbrace{w_a}_w, \underbrace{ne_a}_{ne}, \underbrace{nw_a}_{nw}, \underbrace{se_a \oplus e_b \oplus se_b}_{se}, \underbrace{sw_a \oplus w_b \oplus sw_b}_{sw}) \\
(c_a, n_a, s_a, e_a, w_a, ne_a, nw_a, se_a, sw_a) \otimes_r (c_b, n_b, s_b, e_b, w_b, ne_b, nw_b, se_b, sw_b) \\
&= (\underbrace{c_a}_c, \underbrace{n_a}_n, \underbrace{s_a}_s, \underbrace{e_a \phi |c_b| \phi e_b}_e, \underbrace{w_a}_w, \underbrace{ne_a \phi n_b \phi ne_b}_{ne}, \underbrace{nw_a \phi nw_b \phi nw_b}_{nw}, \underbrace{se_a \phi s_b \phi se_b}_{se}, \underbrace{sw_a \phi sw_b \phi sw_b}_{sw})
\end{aligned}$$

Here, *Nil* is a special value to indicate that there is no value, and we treat it as an identity of \oplus and ϕ for simplification of the notation. Thus, $\text{Nil} \oplus x = x$, $x \oplus \text{Nil} = x$, $\text{Nil} \phi x = x$, and $x \phi \text{Nil} = x$. Each element of the resulting array is a tuple of nine elements. The meaning of each element of the tuple is as follows: *c* is the center element; *s* is an array of the elements on the south of the element; similarly *n*, *e* and *w* are arrays of the elements on the north, east and west respectively; *ne*, *nw*, *se* and *sw* are arrays of the elements on the northeast, northwest, southeast and southwest. Note that this *surrounds* needs $O(n^4)$ memory space for a matrix of $n \times n$.

We show some examples written with the general form.

$$\begin{aligned}
\text{imagefilter } \text{ker} &= \text{mconv} (\text{conv } \text{ker}) \text{shrink}_1 \\
\text{FDM } n \text{ ker} &= \text{iter } n (\text{mconv} (\text{conv } \text{ker}) \text{shrink}_1) \\
\text{where} \\
\text{shrink}_1 &= \text{id} \times B \times T \times L \times R \times BL \times BR \times TL \times TR \\
B &= (\lfloor \cdot \rfloor, \gg, \phi, \rfloor), \quad T = (\lfloor \cdot \rfloor, \ll, \phi, \rfloor), \quad L = (\lfloor \cdot \rfloor, \oplus, \ll, \rfloor), \quad R = (\lfloor \cdot \rfloor, \oplus, \gg, \rfloor), \\
BL &= (\lfloor \cdot \rfloor, \gg, \ll, \rfloor), \quad BR = (\lfloor \cdot \rfloor, \gg, \gg, \rfloor), \quad TL = (\lfloor \cdot \rfloor, \ll, \ll, \rfloor), \quad TR = (\lfloor \cdot \rfloor, \ll, \gg, \rfloor)
\end{aligned}$$

The function *imagefilter ker* is an image filter with the coefficient matrix *ker*, which is used to compute weighted sum of the surrounding pixels. The *shrink*₁ reduces

each part of the gathered surrounding elements to the element closest to the center, and the function *conv ker* calculates the weighted sum of them. The functions *B* and *T* take the bottom row and the top row of the input array respectively. Similarly, each of *L*, *R*, *BL*, *BR*, *TL* and *TR* takes corresponding part of the input array. Figure 3 shows an image of execution of the sharpen-filter by the above general program. The function *FDM n ker* performs the finite difference method, where *iter* is an iteration function and each iteration step is the same as image filters with specific coefficients.

The following example calculates the array of which element at (i, j) is the maximum in the i -th row and the j -th column, i.e. the maximum in the cross. The *shrink_{max}* reduces each part of the gathered surrounding elements to the biggest element in the part, where the binary operator \uparrow takes the bigger element. The function *max₅* takes the maximum of the column and the row including the center element.

$$\begin{aligned} \text{crossmax} &= \text{mconv max}_5 \text{ shrink}_{\text{max}} \\ \text{where } \text{shrink}_{\text{max}} &= \text{max} \times \cdots \times \text{max} \\ \text{max} &= (\uparrow \text{id}, \uparrow, \uparrow) \\ \text{max}_5 (c, n, s, e, w, -, -, -, -) &= c \uparrow n \uparrow s \uparrow e \uparrow w \end{aligned}$$

As shown in this example, *shrink* is allowed not only to shrink the shape of the surroundings but to perform some calculation.

4.2 Surrounding Theorem

In this section, we give the theorem to optimize the general form by fusing *shrink* to *surrounds*.

Image filters and difference methods usually have the *shrink* of the fixed size window that takes the fixed-size rectangle region (window) of the surrounding elements. The function that takes a fixed number of columns (rows) can be written as a homomorphism. For example, the function *right* = $(\uparrow \mid, \oplus, \gg)$ takes the right-most column, which is used in the examples in the previous section. Thus, we here consider the general *shrink* that consists of homomorphisms. It is defined as follows.

$$\begin{aligned} \text{shrink} &= g_c \times h_n \times h_s \times h_e \times h_w \times h_{ne} \times h_{nw} \times h_{se} \times h_{sw} \\ \text{where} \\ h_n &= (g_n, \oplus_n, \otimes_n), \quad h_s = (g_s, \oplus_s, \otimes_s), \quad h_e = (g_e, \oplus_e, \otimes_e) \\ h_w &= (g_w, \oplus_w, \otimes_w), \quad h_{ne} = (g_{ne}, \oplus_{ne}, \otimes_{ne}), \quad h_{nw} = (g_{nw}, \oplus_{nw}, \otimes_{nw}) \\ h_{se} &= (g_{se}, \oplus_{se}, \otimes_{se}), \quad h_{sw} = (g_{sw}, \oplus_{sw}, \otimes_{sw}) \end{aligned}$$

Here, \oplus_X and \otimes_X are extended to satisfy the following equations: $\text{Nil} \oplus_X x = x$, $x \oplus_X \text{Nil} = x$, $\text{Nil} \otimes_X x = x$, and $x \otimes_X \text{Nil} = x$. The general form using this *shrink* uses $O(n^4)$ operations for a two-dimensional array of $n \times n$.

Then, we give the result of the optimization by fusing *shrink* to *surrounds*.

Theorem 1 (Surrounding). *Let the function shrink be defined by homomorphisms as above. Then, there exist a projection function proj and operators \oplus'_f , \otimes'_f , \oplus'_r and \otimes'_r , whose complexity is bounded by the largest of \oplus_X and \otimes_X , and the program*

$$\text{mconv } f \text{ shrink}$$

is optimized to the following program.

$$\text{map } (f \circ \text{proj}) \circ \text{scanr}(\oplus'_r, \otimes'_r) \circ \text{map } f'_r \circ \text{scan}(\oplus'_f, \otimes'_f) \circ \text{map } f'_f$$

Proof. The theorem is proved by the promotion of `map shrink` with extending the tuples. See the master's thesis [8] for details .

The resulting program uses $O(n^2)$ operations for a two-dimensional array of $n \times n$, while the original general form uses $O(n^4)$ operations. The parallel complexity of the resulting program is $O((n^2/P + \sqrt{n^2/P} \log P)T_{(\oplus_X, \otimes_X)})$ for P processors, provided that the calculational complexity of \oplus_X and \otimes_X in the homomorphisms are $T_{(\oplus_X, \otimes_X)}$.

All the examples shown in the previous section have the *shrink* functions described with homomorphisms. Thus, we can apply this theorem to all of them, and they are executed in $O(n^2/P + \sqrt{n^2/P} \log P)$ complexity using the skeletons.

As mentioned above, the function that takes a fixed number of columns (rows) can be written as a homomorphism. Thus, this theorem holds for the *shrink* of the fixed size window that shrinks the surrounding elements to a fixed size, which is often seen in image filters and difference methods.

Corollary 1 (Fixed Size Window). *Let the function shrink be the fixed size window. Then, the program `mconv f shrink` is optimized to that of $O(n^2)$ operations.*

Note that the homomorphism taking $h \times w$ subarray of a two-dimensional array has the operators of $O(wh)$ complexity. Thus, the total complexity of the program of fixed size window is $O(n^2wh)$.

Finally, we note that we may perform more optimizations by using the shifting of the edges instead of butterfly computations for the global computations of `scan` and `scanr`, provided that the operators influence only a fixed number of elements [8]. This leads to the parallel complexity of $O((n^2/P + \sqrt{n^2/P})T_{(\oplus_X, \otimes_X)})$ for P processors.

5 Experimental Results

We implemented the program¹ using our parallel skeleton library [14] and did our experiment on a cluster (distributed memory). Each of the nodes connected with Gigabit Ethernet has a CPU of Intel® Xeon®2.80GHz and 2GB memory, with Linux 2.4.21 for the OS, gcc 2.96 for the compiler, and mpich 1.2.7 for the MPI.

Figures 4 and 5 show the speedups and the calculation times of the sharpen-filter. The program is an optimized one from the general form (an equivalent of the program in Fig. 2). The inputs are images of 1000×1000 and 2000×2000 . The computation times of the program on one processor are 0.70s and 3.85s respectively.

The result shows programs described with skeletons can be executed efficiently in parallel, and proves the success of our framework. The program achieves almost linear speedups, and the total computational complexity of the optimized program is $O(n^2)$ (thus, its parallel complexity is $O(n^2/P)$ for small P). However, the serial performance

¹ The source code of the test program as well as the skeleton library are available at the web page <http://www.ipl.t.u-tokyo.ac.jp/sketo/>.

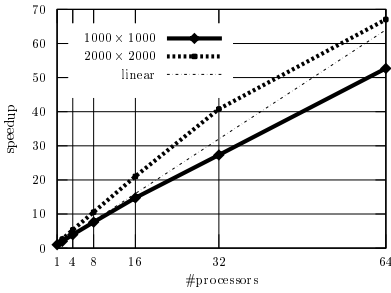


Fig. 4. Speedup of Image Filter

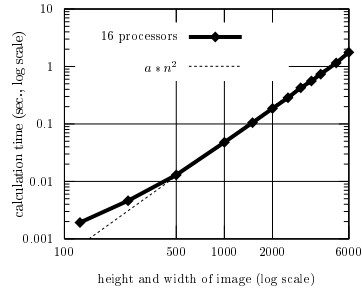


Fig. 5. Calculation Time vs. Size of Image

is rather poor due to the overhead of using general skeletons (i.e. `scan` and `scanr`). We think this problem can be solved by replacing the general skeletons with those specialized for this domain, and it can be automatically done by compilers (future work).

6 Related Work

SKiPPER [17] is a skeleton-based parallel programming environment for real-time image processing. It has skeletons specialized for image processing, while we use general skeletons on two-dimensional arrays. Thus, a program developed with SKiPPER may be faster than that written with our skeletons, but, the latter program can be easily composed with other programs and be optimized by fusion due to generality and solid foundation of our skeletons.

There are several other skeletal parallel approaches (libraries), such as eSkel [1], Muesli [13] and P3L [6]. Their formalizations of skeletons on two-dimensional arrays are not enough (e.g. they have no `scan` skeletons, and the reduction takes only one operator) to deal with matrix-convolutions suitably. Our skeletons have a solid foundation, so that we can easily deal with matrix-convolutions and perform optimizations.

7 Conclusion

In this paper, we proposed a general theorem, called *surrounding theorem*, for optimization of a general skeleton program into an efficient application-specific program. It can deal with a wide class of matrix-convolution problems including image filters and difference methods. The experimental results show that the optimized program can be executed efficiently in parallel. We are now working on making an automatic mechanism for translating the sequential code to our general form with skeletons, and further an optimization mechanism for the application-specific program with respect to its global communication and sequential performance.

Acknowledgment

This work is partially supported by the Grant-in-Aid for Scientific Research (B), No. 17300005, Japan Society for the Promotion of Science. We are grateful to the referees for their detailed and helpful comments.

References

1. A. Benoit, M. Cole, J. Hillston, and S. Gilmore. Flexible skeletal programming with eskel. In *Proceedings of 11th International Euro-Par Conference (Euro-Par'05)*, volume 3648 of *Lecture Notes in Computer Science*, pages 761–770. Springer-Verlag, 2005.
2. R. S. Bird. *Lectures on Constructive Functional Programming*. Technical Report Technical Monograph PRG-69, Oxford University Computing Laboratory, 1988.
3. R. S. Bird. *Introduction to Functional Programming using Haskell*. Prentice Hall, 1998.
4. R. S. Bird and O. de Moor. *Algebras of Programming*. Prentice Hall, 1996.
5. M. Cole. *Algorithmic Skeletons : Structured Management of Parallel Computation*. Research Monographs in Parallel and Distributed Computing, Pitman, London, 1989.
6. M. Danelutto, F. Pasqualetti, and S. Pelagatti. Skeletons for data parallelism in P3L. In *Proceedings of 3rd International Euro-Par Conference (Euro-Par'97)*, volume 1300 of *Lecture Notes in Computer Science*, pages 619–628. Springer-Verlag, 1997.
7. E. Elmroth, F. Gustavson, I. Jonsson, and B. Kagstrom. Recursive Blocked Algorithms and Hybrid Data Structures for Dense Matrix Library Software. *SIAM Review*, 46(1):3–45, 2004.
8. K. Emoto. A Compositional Framework for Parallel Programming on Two-Dimensional Arrays. Master's thesis, Graduate School of Information Science and Technology, the University of Tokyo, 2006. Available at http://www.ipl.t.u-tokyo.ac.jp/~emoto/master_thesis.pdf.
9. K. Emoto, Z. Hu, K. Kakehi, and M. Takeichi. A Compositional Framework for Developing Parallel Programs on Two Dimensional Arrays. Technical Report METR2005-09, Department of Mathematical Informatics, University of Tokyo, 2005.
10. J. D. Frens and D. S. Wise. QR Factorization with Morton-Ordered Quadtree Matrices for Memory Re-use and Parallelism. In *Proceedings of 4th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP'03)*, pages 144–154, 2003.
11. G. Hains. Programming with Array Structures. In A. Kent and J. G. Williams, editors, *Encyclopedia of Computer Science and Technology*, volume 14, pages 105–119. M. Dekker inc, New-York, 1994. Appears also in *Encyclopedia of Microcomputers*.
12. A. K. Jain. *Fundamentals of Digital Image Processing*. Prentice Hall, 1989.
13. H. Kuchen. A Skeleton Library. In *Proceedings of 8th International Euro-Par Conference (Euro-Par'02)*, volume 2400 of *Lecture Notes in Computer Science*, pages 620–629. Springer-Verlag, 2002.
14. K. Matsuzaki, K. Emoto, H. Iwasaki, and Z. Hu. A library of constructive skeletons for sequential style of parallel programming (invited paper). In *Proceedings of the First International Conference on Scalable Information Systems (INFOSCALE 2006)*. IEEE Press, 2006. To appear.
15. L. Mullin, editor. *Arrays, Functional Languages, and Parallel Systems*. Kluwer Academic Publishers, 1991.
16. F. A. Rabhi and S. Gorlatch, editors. *Patterns and Skeletons for Parallel and Distributed Computing*. Springer-Verlag, 2002.
17. J. Serot and D. Ginjac. Skeletons for Parallel Image Processing: an Overview of the SKIPPER Project. *Parallel Computing*, 28(12):1685–1708, 2002.