

Tiling on Systems with Communication/Computation Overlap*

Pierre-Yves Calland¹, Jack Dongarra^{2,3} and Yves Robert²

¹ LIP, Ecole Normale Supérieure de Lyon, 69364 Lyon Cedex 07, France

² Department of Computer Science, University of Tennessee, Knoxville, TN 37996-1301, USA

³ Mathematical Sciences Section, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

e-mail: pycallan@lip.ens-lyon.fr

e-mail: [dongarra, yrobert]@cs.utk.edu

Revised – October 1997

Abstract

In the framework of fully permutable loops, tiling is a compiler technique (also known as “loop blocking”) that has been extensively studied as a source-to-source program transformation. Little work has been devoted to the mapping and scheduling of the tiles on to physical parallel processors. We present several new results in the context of limited computational resources, and assuming communication-computation overlap. In particular, under some reasonable assumptions, we derive the optimal mapping and scheduling of tiles to physical processors.

Key-words: tiling, loop blocking, communication-computation overlap, scheduling and mapping, limited resources, parallel processors.

*This work was supported in part by the National Science Foundation Grant No. ASC-9005933; by the Defense Advanced Research Projects Agency under contract DAAH04-95-1-0077, administered by the Army Research Office; by the Office of Scientific Computing, U.S. Department of Energy, under Contract DE-AC05-84OR21400; by the National Science Foundation Science and Technology Center Cooperative Agreement No. CCR-8809615; by the CNRS-ENS Lyon-INRIA project *ReMaP*; and by the Eureka Project *EuroTOPS*. Yves Robert is on leave from Ecole Normale Supérieure de Lyon and is partly supported by DRET/DGA under contract ERE 96-1104/A000/DRET/DS/SR. Pierre-Yves Calland is supported by a grant from *Région Rhône-Alpes*. A short version of this paper has appeared in the proceedings of the IEEE International Conference *ASAP'97*.

1 Introduction

Tiling is a widely used compiler technique to increase the granularity of computations and the locality of data references. Indeed, as pointed out by Carter et al [8], “Good parallel algorithms are not enough; computer features such as the memory hierarchy and processor architecture need to be exploited to achieve high-performance on parallel machines”. The basic idea of tiling (also known as “loop blocking”) is to group elemental computation points into tiles that will be viewed as computational units. The larger the tiles, the more efficient the computations performed using state-of-the-art processors with pipelined arithmetic units and a multilevel memory hierarchy. This is best illustrated by the recasting of numerical linear algebra algorithms in terms of blocked Level 3 BLAS kernels [12, 10]. Another advantage of tiling is the decrease in communication time (which is proportional to the surface of the tile) relative to the computation time (which is proportional to the volume of the tile). The price to pay for tiling may be an increased latency (if there are data dependencies, for example, we need to wait for the first processor to complete the whole execution of the first tile before another processor can start the execution of the second one, and so on), as well as some load-imbalance problems (the larger the tile, the more difficult to distribute computations equally among the processors).

The tiling technique was originally restricted to perfect loop nests with uniform dependencies, as defined by Banerjee [4], but has been extended to sets of fully permutable loops [24, 16, 11]. Tiling has been studied by several researchers and in different contexts [15, 21, 23, 20, 22, 5, 6, 18, 1, 9, 17, 7, 14, 3]¹. Most of the work amounts to partitioning the iteration space of a uniform loop nest into tiles whose shape and size are optimized according to some criteria (such as the communication-to-computation ratio): see Section 2 for an example. Once the tile shape and size are defined, there remains to distribute the tiles to physical processors and to compute the final scheduling. Little attention has been paid to this last problem. For example, if each physical processor is assigned several tiles, what should be the computation ordering of these tiles? An in-depth study has been presented by Ohta et al [18], who have extended results of Hiranandani et al. [13] on fine grain pipelining for DOACROSS loops. We survey their work in Section 3.

In this paper, we build upon the work of Ohta et al [18]. We reformulate the problem of tiling with limited resources using more realistic assump-

¹This small list is far from being exhaustive.

tions on data dependences and communication-computation overlap than those used in [18]. Our most important result is the derivation of an optimal mapping to assign tiles to physical processors. This result is important because it has clear practical implications: indeed, it turns out that in most situations, a columnwise or rowwise mapping is optimal, which greatly simplifies the task of code generation. All our results are presented in Sections 4 and 5. Finally, we state some conclusions in Section 6.

2 Tiling as a loop transformation technique

When targeting a data-parallel or SPMD style of programming, classical constraints in the literature to define tiles are the following:

Tiles are bounded For scalability reasons, we want the number of points inside a tile to be bounded by a constant independent of the domain size.

Tiles are identical by translation This constraint is imposed to allow for automatic code generation: a tile must be the image by a translation of any other one unless it crosses the computation domain boundaries.

Tiles are “atomic” Each tile is a unit of computation: all synchronization points are beginnings and ends of tiles. The order on tiles must be compatible with the order on nodes: one must avoid that two distinct tiles depend upon each other.

Consider the following simple example:

```

for  $i = 0$  to  $N_1$  do
  for  $j = 0$  to  $N_2$  do
     $a(i, j) = a(i - 3, j) + a(i, j - 2)$ 
     $b(i, j) = a(i - 2, j - 3) + b(i - 2, j - 1)$ 
  enddo
enddo

```

This loop nest has depth 2. The dependences are uniform (intuitively, dependence vectors are translations), and they can be encapsulated into the dependence matrix

$$D = \begin{pmatrix} 3 & 0 & 2 & 2 \\ 0 & 2 & 3 & 1 \end{pmatrix}.$$

The atomicity constraint can be expressed by the analytical condition $HD \geq 0$, where H is the matrix of vectors normal to the faces (or the edges in two-dimensional problems) of the tile [15]. In Figure 1, we sketch a valid tiling for our example. The matrix H is the one derived using the scalable communication-to-computation criteria of Boulet et al. [5]:

$$H = \frac{1}{16} \begin{pmatrix} 0 & \frac{1}{3} \\ \frac{1}{2} & 0 \end{pmatrix}.$$

We check that $HD \geq 0$. Note that the volume of the tile, which represents the number of computations per tile, is given by the determinant of H^{-1} : $V_{comp} = \det(H^{-1}) = 96$. The number of communications is the following: each tile sends

- 24 data items to its right neighbor,
- 34 data items to its lower neighbor,
- and 6 data items to its lower-right neighbor.

Note that the third message (the diagonal communication) may be routed horizontally and then vertically, or the other way round, and even may be combined with any of the first two messages.

Please insert Figure 1 here

It is important to point out that the dependences between tiles are summarized by the vector pair

$$\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}.$$

In other words, the computation of a tile cannot be started before both its left and upper neighbor tiles have been executed.

As stated above, the atomicity constraint implies that inter-processor communications only take place at the end of the processing of each tile. Note that current architectures do allow for communications and computations to overlap, and it is important to point out that the atomicity constraint does not prevent a given processor from simultaneously communicating boundary data of one tile (whose execution it just completed) and starting the computation of another tile. Also, minimizing communication

start-up overheads is a “sine-qua-non” condition towards achieving good performance. Even though sophisticated routing strategies are designed and implemented in hardware, communication start-up costs remain very expensive as opposed to the elemental time for communicating one data item (and even worse for performing a computation). Frequent exchanges of short messages should therefore be replaced by fewer sends and receives of longer messages. To summarize, in the context of distributed memory architectures, tiling is a technique that permits to optimize communications while increasing the granularity of computations.

3 Tiling with resource constraints

Ohta et al. [18] aim at determining the best tile size under the following hypotheses:

- (H1) There are P available processors interconnected as a ring.
- (H2) The computation domain is a two-dimensional rectangle of size $N_1 \times N_2$.
- (H3) Tiles are rectangular and their edges are parallel to the axes (see Figure 2). The size of a tile is $n_1 \times n_2$, where n_1 and n_2 are unknowns.
- (H4) Dependences between tiles are summarized by the vector pair $\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$ (as in the example of Section 2).
- (H5) Tiles are assigned to processor using a one-dimensional cyclic distribution: in other words, tile (i, j) is allocated to processor $j \bmod P$.
- (H6) The scheduling of the tiles is column-wise: at step 0, processor P_0 executes tile $(0, 0)$ and the computed value is communicated to the adjacent processor P_1 (more precisely, a rectangular slice of width w and height n_2 is sent). At step 1, processors P_0 and P_1 execute tiles $(0, 1)$ and $(1, 0)$ simultaneously. After having executed a whole column of tiles, a processor moves on to its next column.

Please insert Figure 2 here

A step is the time required to compute a tile and to communicate data. Ohta et al. [18] use the following expression:

$$T_{tile} = T_{comp} + T_{comm} = n_1 n_2 t + (a + b n_2)$$

where t is the elemental computation time, a is a communication start-up and b is the inverse of the communication bandwidth times the width w of the slice being communicated (the communication cost is a linear expression in the message size).

To compute the total execution time, Ohta et al. [18] use the formula $(M_l + M_p)T_{tile}$, where $M_l = P - 1$ is the latency (the step at which the last processor begins to work) and $M_p = \frac{N_1 \times N_2}{P \times n_1 \times n_2}$ is the number of tiles per processor (assumed to be an integer). Using the approximation $M_l = P$, they derive the total execution time T as

$$T = (P + \frac{N_1 N_2}{P n_1 n_2})(n_1 n_2 t + a + b n_2).$$

The execution time is found to be minimal when choosing $n_1 = \frac{N_1}{P}$ and $n_2 = \sqrt{\frac{N_2 a}{N_1 t}}$.

The objective of this paper is to discuss the hypotheses (H1) to (H6) of Ohta et al., and to reformulate their results using a more accurate modeling of current architectures. Indeed, their study is conducted while assuming that processors cannot simultaneously communicate bordering data items of the last tile and perform computations for the next tile. However, overlapping computations and communications is a facility provided by all distributed memory computers, so we relax this restriction. This simple modification has a tremendous effect on the determination of the best tile size.

4 Allowing for communication-computation overlap

4.1 On the model

Hypotheses (H2), (H3) and (H4) may appear very restricting. However, we point out the following justifications:

Tile shape We assume that the tiles are rectangular. This is to be understood as the outcome of previous program transformations. The

first step in tiling amounts to determining the best shape and size of the tiles, assuming an infinite grid of virtual processors. Because this step will lead to tiles whose edges are parallel to extremal dependence vectors in the convex hull of the dependence cone, we can perform a unimodular transformation and rewrite the original loop nest along the edge axes. The resulting domain may not be a rectangular, but we can approximate it using the smallest bounding box (however, this approximation may impact the accuracy of our results).

Dependence vectors We assume that dependencies are summarized by the vector pair $\mathcal{V} = \{(1, 0)^t, (0, 1)^t\}$. Note that these are dependencies between tiles, not between elementary computations. Hence, having right- and top-neighbor dependencies is a very general situation if the tiles are large enough. In the example of Section 2, we had 4 dependence vectors in the original loop nest, but we ended up with \mathcal{V} after tiling. Technically, since we deal with a set of fully permutable loops, all dependence vectors have nonnegative components only, so that \mathcal{V} permits all other dependence vectors to be generated by transitivity. Note that having a dependence vector $(0, a)^t$ with $a \geq 2$ between tiles, *instead of* having vector $(0, 1)^t$, would mean unusually long dependencies in the original loop nest (in the example of Section 2, $a(i, j)$ would depend upon $a(i, j - 8)$ but not on $a(i, j - x)$ for $x \leq 7$), while having $(0, a)^t$ *in addition to* $(0, 1)^t$ as a dependence vector between tiles is simply redundant. In practical situations, we might have an additional diagonal dependence vector $(1, 1)^t$ between tiles, but the diagonal communication may be routed horizontally and then vertically, or the other way round, and even may be combined with any of the other two messages (induced by dependence vectors $(0, 1)^t$ and $(1, 0)^t$).

On the other hand, hypotheses (H5) and (H6) are unnecessarily restrictive, because the mapping and scheduling of the tiles should be an output decision of the procedure of tiling with limited resources, rather than being given a priori. We overcome this restriction in Section 5.

4.2 Revisiting the results of Ohta et al.

The total execution time is given by the following proposition:

Proposition 1 *Under the hypotheses (H1) to (H6) of Section 3, and allowing for communication-computation overlap, the total computation time T*

is (assuming all fractions to be integer):

$$T = \begin{cases} T_1 = (P - 1)(n_1n_2t + a + bn_2) + \frac{N_1N_2}{P}t & \text{if } N_2n_1t \geq P(n_1n_2t + a + bn_2) \\ T_2 = (\frac{N_1}{n_1} - 1)(n_1n_2t + a + bn_2) + N_2n_1t & \text{otherwise} \end{cases} \quad (1)$$

Proof According to hypothesis (H4), the computation goes column-wise. When a processor has completed the execution of a whole column of tiles, it starts the next column that has been assigned to it. The time to process a whole column of tiles is the number of tiles in the column, namely $\frac{N_2}{n_2}$, times the time to compute a tile, namely $T_{comp} = n_1n_2t$. We obtain the value N_2n_1t for processing a whole tile column.

Now, according to hypothesis (H5), tile columns are distributed cyclically to processors. If a processor starts the execution of the first tile in a given column at time-step t , its right neighbor cannot start the execution of the first tile in the next column before time-step $t + T_{tile}$, where $T_{tile} = T_{comp} + T_{comm} = n_1n_2t + (a + bn_2)$ (this is due to the dependence vector $(1, 0)^t$). Note that T_{tile} is the same as in Section 3, but we pay a communication cost only when the processors owning the tiles are not the same. Two cases can occur:

Please insert Figure 3 here

- Either there are enough tiles in each column so that when a processor has completed the execution of a whole tile column, it does not have to wait for its next tile column to be ready. This will happen if $\frac{N_2}{n_2}T_{comp} = N_2n_1t$ is greater than or equal to the delay imposed by horizontal constraints, i.e. if

$$\frac{N_2}{n_2}T_{comp} \geq P T_{tile}.$$

If this condition holds, all processors remain active throughout the entire computation, once they have started execution. Since the last processor starts at time $(P - 1)T_{tile}$ and has $\frac{N_1N_2}{Pn_1n_2}$ tiles to execute (each in time $T_{comp} = n_1n_2t$), we obtain T_1 , the first expression in Equation (1). See Figure 3 where $T_{comp} = T_{comm} = 1$, and $P = 3$. There are $\frac{N_2}{n_2} = 8$ tiles per column, and $PT_{tile} = 6$, hence the condition is satisfied.

Please insert Figure 4 here

- Or each processor has to wait upon finishing a tile column until the next column is ready. This translates into the condition $\frac{N_2}{n_2}T_{comp} \leq PT_{tile}$. In that case, the total computation time is equal to the time at which the last processor starts the execution of the first tile in the last column, namely $(\frac{N_1}{n_1} - 1)T_{tile}$ plus the time needed to process this column. We obtain the expression $(\frac{N_1}{n_1} - 1)T_{tile} + \frac{N_2}{n_2}T_{comp}$, as stated in the second formula of Equation (1). See Figure 4 where $T_{comp} = 1$, $T_{comm} = 2$, and $P = 3$. There are $\frac{N_2}{n_2} = 8$ tiles per column, and $PT_{tile} = 9$, hence $\frac{N_2}{n_2}T_{comp} \leq PT_{tile}$. Processors remain idle at the end of each tile column, waiting for their next column to be ready. ■

The optimal number of processors that should be used so as to minimize the total execution time is given by the following proposition:

Corollary 1 *Under the hypotheses (H2) to (H6) of Section 3, and allowing for communication-computation overlap, let*

$$P_\alpha = \sqrt{\frac{N_1 N_2 t}{n_1 n_2 t + a + b n_2}} \quad \text{and} \quad P_\beta = \frac{N_2 n_1 t}{n_1 n_2 t + a + b n_2}$$

The number of processors P that minimizes the total execution time is given by:

- *if $P_\beta \leq 1$ or $P_\alpha \leq 1 \leq P_\beta$, then $P = 1$,*
- *if $1 \leq P_\beta \leq P_\alpha$ then $P = P_\beta$,*
- *if $1 \leq P_\alpha \leq P_\beta$ then $P = P_\alpha$.*

Proof The “steady-state” condition $N_2 n_1 t \geq P(n_1 n_2 t + a + b n_2)$ in Equation (1) can be rewritten as

$$P \leq P_\beta.$$

Consider $T_1 = (P - 1)(n_1 n_2 t + a + b n_2) + \frac{N_1 N_2}{P} t$ (see Equation (1)). The minimum of T_1 is obtained for $P = P_\alpha$. The expression of T_1 shows that is a non-increasing function of P when $P \leq P_\alpha$, and then a non-decreasing

function of P when $P \geq P_\alpha$. Also, note that T_2 does not depend on P (except than through the condition $P \geq P_\beta$). Then the result follows according to a simple case analysis. \blacksquare

For large domains, we will have $1 \leq P_\beta \leq P_\alpha$, and it is no surprise that the optimal number of processors is the one required to ensure steady-state execution: Equation (1) in Proposition 1 states that all processors remain active once started if

$$N_2 n_1 t \geq P(n_1 n_2 t + a + b n_2).$$

So far, we have assumed that n_1 and n_2 were input parameters, because the size and shape of the tiles may be imposed by some a priori considerations (such as the cache size). We can try to determine the values of n_1 and n_2 in the range $1 \leq n_1 \leq N_1$, $1 \leq n_2 \leq N_2$ that would minimize the total execution time. We rewrite the steady-state inequality by introducing the following function f :

$$f(n_1) = \frac{N_2 n_1 t - P a}{P(n_1 t + b)} \quad (2)$$

f is defined so that $N_2 n_1 t \geq P(n_1 n_2 t + a + b n_2) \iff n_2 \leq f(n_1)$. We have the following result:

Corollary 2 *Under the hypotheses (H1) to (H6) of Section 3, and allowing for communication-computation overlap, the total execution time is minimum for*

- $n_1 = \sqrt{\frac{N_1(a+b)}{(N_2-1)t}}$ and $n_2 = 1$ if $f(\frac{N_1}{P}) \leq 1$
- $n_1 = \frac{P(a+b)}{t(N_2-P)}$ and $n_2 = 1$ otherwise.

Proof We break down the problem into two sub-cases depending on the values taken by the function f , whose argument n_1 ranges from 1 to $\frac{N_1}{P}$;

- $\forall n_1, f(n_1) \leq 1$. Since f is a nondecreasing function of n_1 , this condition is equivalent to $f(\frac{N_1}{P}) \leq 1$. In this case, Equation (2) is never satisfied ($n_2 \geq 1$). Then the minimum of T is obtained by minimizing T_2 with $n_2 = 1$, namely

$$T = \left(\frac{N_1}{n_1} - 1\right)(n_1 t + a + b) + N_2 n_1 t$$

This easily leads to $n_1 = \sqrt{\frac{N_1(a+b)}{(N_2-1)t}}$, as stated in the theorem

- $\exists n_1, f(n_1) \geq 1$. We take n_1^0 such that $f(n_1^0) = \min f(n_1)$ and $n_1^0 \geq 1$. Note that all values of $n_1 \geq n_1^0$ will lead to admissible values for n_2 , because we always have $f(n_1) \leq \frac{N_2}{P}$ by definition of f . Now consider the expression of T for arbitrary n_1 and n_2 :
 - if $n_2 \leq f(n_1)$, then $T = T_1$, T is a non-increasing function of both n_1 and n_2 decreases, then the minimum is obtained with $n_2 = 1$ and $n_1 = n_1^0$.
 - if $n_2 \geq f(n_1)$ then $T = T_2$ and is a non-increasing function of n_2 . Then the minimum of T is reached if $n_2 = f(n_1)$. In that case $T_2 = T_1$, and again the minimum is reached when $n_2 = 1$ and $n_1 = n_1^0$.

■

This result is disappointing in that we end up with degenerate tiles in most practical situations. For instance if $P \ll N_2$ (which is very likely to happen in practice), $f(1) \geq 1$, and the optimal tile size is $n_1 = n_2 = 1$, not a very coarse-grain tiling indeed! For other values of the problem parameters we would have an optimal tile size that depends upon the domain size, thereby contradicting the assumption that tiles are bounded (Section 2). Note that Ohta et al [18] also have this problem in their original model. The flaw is that the model is not accurate enough to take the impact of data locality and data reuse into account (which are the main objectives of tiling, and the main motivation for designing blocked linear algebra algorithms [12]). A first solution is to model the computation cost of a tile by an affine expression like $T_{comp} = n_1 n_2 t + u$, where u represents some access overhead. It is not difficult to plug this expression into the formula given for the total execution time T , and to derive the optimal tile size. Another solution is to assume a fixed tile size that would be imposed by some a priori considerations (such as the cache size). Again, we can let $n_1 n_2 = S$ in Equation (1), and minimize T for n_1 , say.

4.3 Generalizing the model

Assuming communication-computation overlap seems a reasonable hypothesis for current machines which have communication coprocessors and allow for asynchronous communications (posting instructions ahead, or using active messages). We can think of independent computations going along a

thread while communication is initiated and performed by another thread. As written in Pacheco [19, p. 268], “if we have communication coprocessors (and smart compilers) ... the actual running time [for performing k computations and sending/receiving a message of length m] ... might be $\max\{t_s + mt_c, kt_a\}$ ” (with our notations, $t_a = t$, $t_s = a$ and $t_c = b/w$).

A very interesting approach has been proposed by Andonov and Rajopadhye [3]: they introduce the *tile period* P_t as the time elapsed between corresponding instructions of two successive tiles that are mapped to the *same* processor, while they define the *tile latency* L_t to be the time between corresponding instructions of two successive tiles that are mapped to *different* processors. With these notations, the parallel execution time becomes [3]

$$T = \begin{cases} T_1 = (P - 1)L_t + \frac{N_1}{n_1} \frac{N_2}{n_2} \frac{1}{P} P_t & \text{if } \frac{N_2}{n_2} P_t \geq P L_t \\ T_2 = (\frac{N_1}{n_1} - 1)L_t + \frac{N_2}{n_2} P_t & \text{otherwise} \end{cases} \quad (3)$$

The power of this approach is that the expressions for L_t and P_t can be modified to take into account several architectural models, while Equation (3) still remains valid. A very detailed architectural model is presented in [3], and several other models are explored in [2].

With our notations, $P_t = T_{comp}$ and $L_t = T_{comp} + T_{comm}$. We can rewrite Equation (1) as

$$T = \begin{cases} T_1 = (P - 1)(T_{comp} + T_{comm}) + \frac{N_1}{n_1} \frac{N_2}{n_2} \frac{1}{P} T_{comp} & \text{if } \frac{N_2}{n_2} T_{comp} \geq P(T_{comp} + T_{comm}) \\ T_2 = (\frac{N_1}{n_1} - 1)(T_{comp} + T_{comm}) + \frac{N_2}{n_2} T_{comp} & \text{otherwise} \end{cases} \quad (4)$$

Equation (3), or its variant Equation (4), is the key to our tiling problem, because it expresses the parallel execution time as a function of the domain size, of the number of processors, and of the tile parameters P_t and L_t , or equivalently T_{comp} and T_{comm} .

5 Optimal mapping and scheduling

Hypotheses (H5) and (H6) are very restrictive in that they impose the mapping of tiles to processors as well as their scheduling. The intuitive motivation for (H5) is that a cyclic distribution of tiles is quite natural to load-balance computations. Once the distribution of tiles to processors is fixed,

there are several possible schedulings (any wavefront execution that goes along a left-to-right diagonal is valid). Specifying a column-wise execution may lead to the simplest code generation.

It turns out that (H5) and (H6) provide the best solution among all possible distributions of tiles to processors, which is a very strong result. This result holds true under the assumption that the communication cost for a tile is not larger than its computation cost. Since the communication cost for a tile grows linearly with its size, while the computation costs grows quadratically, this hypothesis will be satisfied if the tile is large enough². This result is formally stated in the theorem below. Beforehand, we need to refine the communication cost as follows:

- $T_{comm_horiz} = a + bn_2$ is the cost of communicating data from (the processor owning) tile (i, j) to (the processor owning) its right neighbor tile $(i + 1, j)$,
- $T_{comm_vert} = a' + b'n_1$ is the cost of communicating data from (the processor owning) tile (i, j) to (the processor owning) its bottom neighbor tile $(i, j + 1)$.

We pay a communication cost only when the two processors that own the neighboring tiles are not the same. So far we never paid any cost for vertical communications, and we always did for horizontal communications, because of hypothesis (H5). We had to refine the communication cost because in this section, we do not make any assumption on the mapping of tiles to processors. Depending upon the values of T_{comm_horiz} and T_{comm_vert} , the best mapping will be column-wise or row-wise:

Theorem 1 *Under the hypotheses (H2) to (H4) of Section 3, and allowing for communication-computation overlap, let n_1 and n_2 be chosen so that*

$$\max\{T_{comm_horiz}, T_{comm_vert}\} = \max\{a + bn_2, a' + b'n_1\} \leq T_{comp} = n_1n_2t.$$

1. *If $T_{comm_horiz} \leq T_{comm_vert}$, assume that the steady state equation holds: $N_2n_1t \geq P(n_1n_2t + a + bn_2)$. Then the absolute minimum for the total execution time is*

$$T_1 = (P - 1)(T_{comp} + T_{comm_horiz}) + \frac{N_1N_2}{P}t$$

and it is achieved by mapping and scheduling tiles according to hypotheses (H5) and (H6),

²Of course, we can imagine theoretical situations where the communication cost is so large that a sequential execution would lead to the best result.

2. If $T_{comm_vert} \leq T_{comm_horiz}$, assume that the steady state equation holds: $N_1 n_2 t \geq P(n_1 n_2 t + a' + b' n_1)$. Then the absolute minimum for the total execution time is

$$T'_1 = (P - 1)(T_{comp} + T_{comm_vert}) + \frac{N_1 N_2}{P} t$$

and it is achieved by mapping rows of tiles using a one-dimensional cyclic distribution (tile (i, j) is allocated to processor $i \bmod P$), and by scheduling the tiles with a row-wise ordering.

Proof Without loss of generality, assume that $T_{comm_vert} \leq T_{comm_horiz}$ (the result is symmetric in the rows and columns), and let $T_{comm} = T_{comm_vert}$. We begin the proof with the following preliminary result, where σ denotes any valid scheduling of the tiles ($\sigma(I)$ is the time-step at which the execution of I begins):

Lemma 1 Let $I = (i, j)$ be a tile index, and let $I' = (i + 1, j)$ and $I'' = (i, j + 1)$ be its successor tiles. We have

$$\max\{\sigma(I'') - \sigma(I), \sigma(I') - \sigma(I)\} \geq T_{comm} + T_{comp}.$$

Proof Let $proc(I)$ be the processor that executes tile I . We have three cases to consider, depending upon whether $proc(I)$ also executes both successors I' and I'' , or exactly one of them, or none of them:

both successors: $proc(I) = proc(I') = proc(I'')$

The same processor executes both successors. They are executed sequentially and the last one being executed cannot begin execution before time-step $\sigma(I) + 2T_{comp}$. As $T_{comm} \leq T_{comp}$ the result is proven.

one successor: $proc(I) = proc(I')$ and $proc(I) \neq proc(I'')$

(respectively $proc(I) = proc(I'')$ and $proc(I) \neq proc(I')$). A communication is needed between I and I'' (respectively I and I'), hence $\sigma(I'') - \sigma(I) \geq T_{comm} + T_{comp}$ (respectively $\sigma(I') - \sigma(I) \geq T_{comm} + T_{comp}$)

no successor: $proc(I) \neq proc(I')$ and $proc(I) \neq proc(I'')$

This case is similar to the previous one. ■

Back to the proof of the theorem, let $T_{//}$ the total execution time using P processors. Let $Idle$ be the cumulated idle time of all processors during execution. Finally, let $T_{seq} = N_1 N_2 t$ be the sequential execution time. Clearly,

$$PT_{//} = Idle + T_{seq}.$$

Hence, to show that $T_{//} \geq T_1 = (P - 1)(T_{comp} + T_{comm}) + \frac{T_{seq}}{P}$, we need to show that

$$Idle \geq P(P - 1)(T_{comp} + T_{comm}).$$

The structure of the dependence graph does impose that some processors are idle at the beginning of the computation, which will lead to a lower bound for $Idle$. For instance, during the execution of tile $(0, 0)$, there are necessarily $P - 1$ idle processors. To go on, we recursively define $pivot_tile(k)$ as follows (see Figure 5):

Please insert Figure 5 here

- $pivot_tile(0) = (0, 0)$, and
- for $k \geq 1$, $pivot_tile(k)$ is the one of the two successors of $pivot_tile(k - 1)$ which is executed last: if $pivot_tile(k - 1) = I = (i, j)$, let $I' = (i + 1, j)$ and $I'' = (i, j + 1)$ be the successors of tile I :
 - If $\sigma(I') \geq \sigma(I'')$, then $pivot_tile(k) = I'$, and we define $S(k)$ as the remaining tiles in column j : $S(k) = \{(i, j + l), l \geq 1\}$,
 - If $\sigma(I'') \geq \sigma(I')$, then $pivot_tile(k) = I''$, and we define $S(k)$ as the remaining tiles in row i : $S(k) = \{(i + l, j), l \geq 1\}$,

We know from Lemma 1 that for all $k \geq 1$,

$$\sigma(pivot_tile(k)) - \sigma(pivot_tile(k - 1)) \geq T_{comm} + T_{comp}.$$

We prove by induction that for $1 \leq k \leq P - 1$, at least $P - k$ processors are kept idle between the beginning of the execution of $pivot_tile(k - 1)$ and that of $pivot_tile(k)$. This will lead to:

$$Idle \geq ((P - 1) + (P - 2) + \dots + 1)(T_{comm} + T_{comp}) = \frac{P(P - 1)}{2}(T_{comm} + T_{comp}).$$

This will prove the desired result, because the same amount of idleness, so to speak, will be spent at the end of the computation (by symmetry of the dependence graph). Now, for the induction:

- Let $k = 1$: $pivot_tile(1)$ is either $(0, 1)$ or $(1, 0)$. See Figure 5 where $pivot_tile(1) = (1, 0)$ and $S(1) = \{(0, 0 + l), l \geq 1\}$. Between the the beginning of the execution of $pivot_tile(0)$ and that of $pivot_tile(1)$, the only successors of $pivot_tile(0)$ that can be executed are in $S(1)$. But all tasks in $S(1)$ must be executed sequentially, hence between the beginning of the execution of $pivot_tile(0)$ and that of $pivot_tile(1)$, at least $(P - 1)$ processors are kept idle.
- Assume that the hypothesis is true until step k . Between the beginning of the execution of $pivot_tile(k)$ and that of $pivot_tile(k + 1)$, at most one processor can be active in $S(1)$, at most another one in $S(2)$, \dots , and at most one processor in $S(k + 1)$, so that at most $k + 1$ processors can be active, or equivalently, at least $P - (k + 1)$ processors remain idle.

■

It is worth to point out that Theorem 1 holds true in a large framework. Whatever the model used for estimating the communication time T_{comm} and the computation time T_{comp} , the parallel execution time for a columnwise allocation of tiles to processors is given by Equation (4). Theorem 1 basically says that such a columnwise or rowwise allocation will be optimal as soon as

1. $T_{comm} \leq T_{comp}$
2. Steady-state condition: the weight of a tile column (or tile row) is greater that the tile latency

$$L_t = P(T_{comm} + T_{comp})$$

The first hypothesis will be fulfilled if the tile is large enough (because the communication cost grows linearly while the computation cost grows quadratically). The second hypothesis will be fulfilled as soon as the domain is large enough in front of the number of processors, a situation very likely to happen in practice.

Finally, note that when the steady-state condition is not satisfied, we can still derive similar results. For instance assume a square $N \times N$ tiled iteration space (N tiles per row and per column). Let T_{comp} be the computation time for a tile, and let T_{comm} be the communication time (either horizontal or vertical). With P processors, if $NT_{comp} \leq P(T_{comm} + T_{comp})$, a columnwise

allocation of tiles to processors leads to the parallel execution time $T = (N - 1)(T_{comp} + T_{comm}) + NT_{comp}$. If $T_{comm} \leq T_{comp}$, this is optimal: use Lemma 1 to show that the execution of diagonal tile (i, i) , $0 \leq i < N$, cannot start before time-step $(i - 1)(2T_{comp} + T_{comm})$.

6 Conclusion

In this paper, we have studied tiling techniques aimed at adapting the granularity of permutable loop algorithms towards execution on distributed-memory machines. We view tiling as a two-step process: the first step amounts to determining the best shape and size of the tiles (assuming an infinite grid of virtual processors), while the second step consists in mapping and scheduling the tiles to physical processors. We have concentrated on the second step, assuming a realistic model where (independent) communication and computation may overlap. We have obtained several new results, including a strong result on the optimal mapping and scheduling. However, much remains to be done to extend these results to arbitrary dimensions and domain shapes.

More generally, the relationship between tiling, scheduling and mapping is not yet well understood, and the two-step approach may prove too complicated for practical problems. Yet, such a two-step approach is typical in the field of parallelizing compilers (other examples are general task graph scheduling, software pipelining and loop parallelization algorithms).

Finally, the recent development of heterogeneous computing platforms may well lead to using tiles whose size and shape will depend upon the characteristics of the processors they are assigned to ... a truly challenging problem!

Acknowledgment We are deeply indebted to Sanjay Rajopadhye for his useful comments on a first version of this paper. We also thank the referees for suggesting several improvements.

References

- [1] A. Agarwal, D.A. Kranz, and V. Natarajan. Automatic partitioning of parallel loops and data arrays for distributed shared-memory multiprocessors. *IEEE Trans. Parallel Distributed Systems*, 6(9):943–962, 1995.

- [2] Rumen Andonov, Hafid Bourzoufi, and Sanjay Rajopadhye. Two-dimensional orthogonal tiling: from theory to practice. In *International Conference on High Performance Computing (HiPC)*, pages 225–231, Trivandrum, India, 1996. IEEE Computer Society Press.
- [3] Rumen Andonov and Sanjay Rajopadhye. Optimal tiling of two-dimensional uniform recurrences. *Journal of Parallel and Distributed Computing*, to appear. Available as Technical Report LIMAV-RR 97-1, at <http://www.univ-valenciennes.fr/limav>.
- [4] Utpal Banerjee. An introduction to a formal theory of dependence analysis. *The Journal of Supercomputing*, 2:133–149, 1988.
- [5] Pierre Boulet, Alain Darté, Tanguy Risset, and Yves Robert. (pen)-ultimate tiling? *Integration, the VLSI Journal*, 17:33–51, 1994.
- [6] Pierre-Yves Calland and Tanguy Risset. Precise tiling for uniform loop nests. In P. Cappello et al., editors, *Application Specific Array Processors ASAP 95*, pages 330–337. IEEE Computer Society Press, 1995.
- [7] P.Y. Calland, J. Dongarra, and Y. Robert. Tiling with limited resources. In L. Thiele, J. Fortes, K. Vissers, V. Taylor, T. Noll, and J. Teich, editors, *Application Specific Systems, Architectures, and Processors, ASAP'97*, pages 229–238. IEEE Computer Society Press, 1997. Extended version available on the WEB at <http://www.ens-lyon.fr/~yrobert>.
- [8] L. Carter, J. Ferrante, S. F. Hummel, B. Alpern, and K.S. Gatlín. Hierarchical tiling: a methodology for high performance. Technical Report CS-96-508, University of California at San Diego, San Diego, CA, 1996. Available at <http://www.cse.ucsd.edu/~carter>.
- [9] Y-S. Chen, S-D. Wang, and C-M. Wang. Tiling nested loops into maximal rectangular blocks. *Journal of Parallel and Distributed Computing*, 35(2):123–32, 1996.
- [10] J. Choi, J. Demmel, I. Dhillon, J. Dongarra, S. Ostrouchov, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. ScaLAPACK: A portable linear algebra library for distributed memory computers - design issues and performance. *Computer Physics Communications*, 97:1–15, 1996. (also LAPACK Working Note #95).

- [11] Alain Darté, Georges-André Silber, and Frédéric Vivien. Combining retiming and scheduling techniques for loop parallelization and loop tiling. *Parallel Processing Letters*, 1997. Special issue, to appear. Also available as Tech. Rep. LIP, ENS-Lyon, RR96-34, and on the WEB at <http://www.ens-lyon.fr/LIP>.
- [12] J. J. Dongarra and D. W. Walker. Software libraries for linear algebra computations on high performance computers. *SIAM Review*, 37(2):151–180, 1995.
- [13] Seema Hiranandani, Ken Kennedy, and Chau-Wen Tseng. Evaluating compiler optimizations for Fortran D. *Journal of Parallel and Distributed Computing*, 21(1):27–45, 1992.
- [14] K. Högstedt, L. Carter, and J. Ferrante. Determining the idle time of a tiling. In *Principles of Programming Languages*, pages 160–173. ACM Press, 1997. Extended version available as Technical Report UCSD-CS96-489, and on the WEB at <http://www.cse.ucsd.edu/~carter>.
- [15] François Irigoien and Rmy Triolet. Supernode partitioning. In *Proc. 15th Annual ACM Symp. Principles of Programming Languages*, pages 319–329, San Diego, CA, January 1988.
- [16] Amy W. Lim and Monica S. Lam. Maximizing parallelism and minimizing synchronization with affine transforms. In *Proceedings of the 24th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM Press, January 1997.
- [17] Naraig Manjikian and Tarek S. Abdelrahman. Scheduling of wavefront parallelism on scalable shared memory multiprocessor. In *Proceedings of the International Conference on Parallel Processing ICPP 96*. CRC Press, 1996.
- [18] H. Ohta, Y. Saito, M. Kainaga, and H. Ono. Optimal tile size adjustment in compiling general DOACROSS loop nests. In *1995 International Conference on Supercomputing*, pages 270–279. ACM Press, 1995.
- [19] Peter Pacheco. *Parallel programming with MPI*. Morgan Kaufmann, 1997.
- [20] J. Ramanujam and P. Sadayappan. Tiling multidimensional iteration spaces for multicomputers. *Journal of Parallel and Distributed Computing*, 16(2):108–120, 1992.

- [21] Robert Schreiber and Jack J. Dongarra. Automatic blocking of nested loops. Technical Report 90-38, The University of Tennessee, Knoxville, TN, August 1990.
- [22] S. Sharma, C.-H. Huang, and P. Sadayappan. On data dependence analysis for compiling programs on distributed-memory machines. *ACM Sigplan Notices*, 28(1), January 1993. Extended Abstract.
- [23] M. E. Wolf and M. S. Lam. A data locality optimizing algorithm. In *SIGPLAN Conference on Programming Language Design and Implementation*, pages 30–44. ACM Press, 1991.
- [24] Michael E. Wolf and Monica S. Lam. A loop transformation theory and an algorithm to maximize parallelism. *IEEE Trans. Parallel Distributed Systems*, 2(4):452–471, October 1991.

Figures and Tables

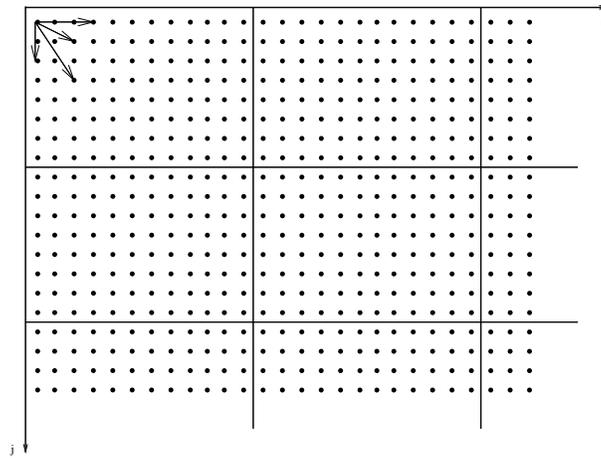


Figure 1: Optimal tiling for a computation volume $V_{comp} = 96$.

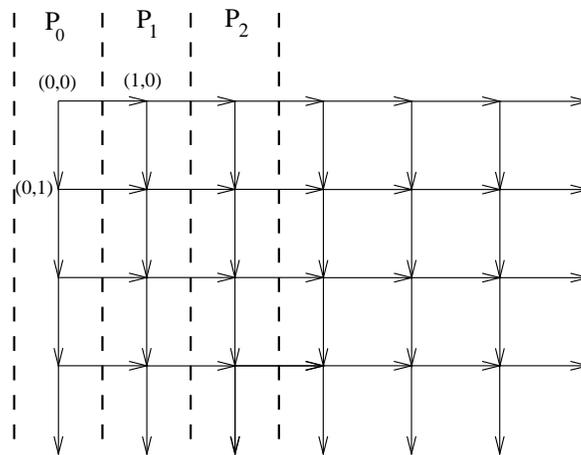


Figure 2: Mapping rectangular tiles onto a ring of processors.

	P_0	P_1	P_2	P_0
0	0	2	4	8
1	1	3	5	9
2	2	4	6	10
3	3	5	7	11
4	4	6	8	12
5	5	7	9	13
6	6	8	10	14
7	7	9	11	15

Figure 3: Scheduling tiles with $T_{comp} = 1$, $T_{comm} = 1$ and $P = 3$.

	P_0	P_1	P_2	P_0
0	0	3	6	9
1	1	4	7	10
2	2	5	8	11
3	3	6	9	12
4	4	7	10	13
5	5	8	11	14
6	6	9	12	15
7	7	10	13	16

Figure 4: Scheduling tiles with $T_{comp} = 1$, $T_{comm} = 2$ and $P = 3$.

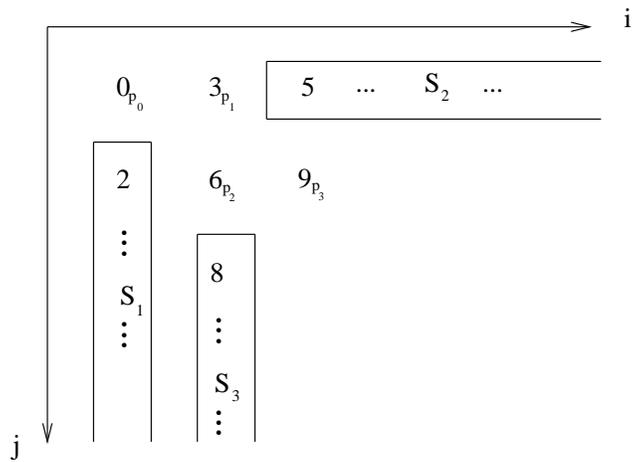


Figure 5: A schedule when $T_{comm} = 1$ and $T_{comp} = 2$. Pivot tiles are labeled, and sets $S(k)$ are framed.