



Determining the Idle Time of a Tiling: New Results.

Frédéric Desprez, Jack Dongarra, Fabrice Rastello, Yves Robert

► To cite this version:

Frédéric Desprez, Jack Dongarra, Fabrice Rastello, Yves Robert. Determining the Idle Time of a Tiling: New Results.. [Research Report] Laboratoire de l'informatique du parallélisme. 1997, 2+17p. hal-02102014

HAL Id: hal-02102014

<https://hal-lara.archives-ouvertes.fr/hal-02102014>

Submitted on 17 Apr 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Laboratoire de l'Informatique du Parallélisme

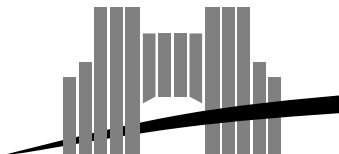
Ecole Normale Supérieure de Lyon
Unité de recherche associée au CNRS n°1398

Determining the Idle Time of a Tiling: New Results

Frédéric Desprez, Jack Dongarra,
Fabrice Rastello and Yves
Robert

October 1997

Research Report N° 97-35



Ecole Normale Supérieure de Lyon

46 Allée d'Italie, 69364 Lyon Cedex 07, France

Téléphone : (+33) (0)4.72.72.80.00 Télécopieur : (+33) (0)4.72.72.80.80

Adresse électronique : lip@lip.ens-lyon.fr

Determining the Idle Time of a Tiling: New Results

Frédéric Desprez, Jack Dongarra, Fabrice Rastello and Yves Robert

October 1997

Abstract

In the framework of fully permutable loops, tiling has been studied extensively as a source-to-source program transformation. We build upon recent results by Högsted, Carter, and Ferrante [12], who aim at determining the cumulated idle time spent by all processors while executing the partitioned (tiled) computation domain. We propose new, much shorter proofs of all their results and extend these in several important directions. More precisely, we provide an accurate solution for all values of the *rise* parameter that relates the shape of the iteration space to that of the tiles, and for all possible distributions of the tiles to processors. In contrast, the authors in [12] deal only with a limited number of cases and provide upper bounds rather than exact formulas.

Keywords: Tiling, fully permutable loops, idle time.

Résumé

Dans le cadre des boucles complètement permutable, le pavage a été beaucoup étudié comme une transformation source-à-source. Nous nous basons sur des travaux récents de Högsted, Carter, et Ferrante [12], dont le but est de déterminer le temps d'attente cumulé passé par tous les processeurs pendant l'exécution le domaine de calcul partitionné (pavé). Nous proposons des nouvelles preuves, plus courtes, de tous leurs résultats et nous les étendons dans plusieurs directions importantes. Nous donnons une solution plus précise pour toutes les valeurs du paramètre *rise* qui relie la forme de l'espace d'itération à celle des tuiles, et pour toutes les distributions possibles des tuiles sur les processeurs. Les auteurs dans [12] ne traitent qu'un nombre limité de cas et fournissent des bornes supérieures plutôt que des formules exactes.

Mots-clés: Pavage, boucles complètement permutable, temps d'attente.

Determining the Idle Time of a Tiling: New Results*

J. Dongarra and Y. Robert[†]
Department of Computer Science
University of Tennessee
Knoxville, TN 37996-1301, USA
(dongarra,yrobert)`@cs.utk.edu`

F. Desprez and F. Rastello
LIP ENS Lyon
46, Allée d'Italie,
Lyon Cedex 07, France
(desprez,frastell)`@lip.ens-lyon.fr`

October 7, 1997

*Jack Dongarra and Yves Robert are with the Department of Computer Science, University of Tennessee, Knoxville, TN 37996-1301, USA. Jack Dongarra is also with the Mathematical Sciences Section, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA. 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. 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 Department of Energy Office of Computational and Technology Research, Mathematical, Information, and Computational Sciences Division 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*. The authors acknowledge the use of the Intel Paragon XP/S 5 computer, located in the Oak Ridge National Laboratory Center for Computational Sciences, funded by the Department of Energy's Mathematical, Information, and Computational Sciences Division subprogram of the Office of Computational and Technology Research. Corresponding author: Yves Robert, `Yves.Robert@inria.fr`.

1 Introduction

Tiling is a widely used technique to increase the granularity of computations and the locality of data references. This technique was originally restricted to perfect loop nests with uniform dependencies, as defined by Banerjee [3], but has been extended to sets of fully permutable loops [22, 14, 10]. Tiling is a widely used technique to increase the granularity of computations and the locality of data references. The basic idea 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 (illustrated by recasting numerical linear algebra algorithms in terms of blocked Level 3 BLAS kernels [11, 9]). 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).

Tiling has been studied by several researchers and in different contexts [13, 19, 21, 17, 20, 4, 5, 16, 1, 8, 15, 6, 12, 2]¹. Rather than providing a detailed motivation for tiling, we refer the reader to the papers by Calland, Dongarra, and Robert [6] and by Högsted, Carter and Ferrante [12], which provide a review of the existing literature. 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). Once the tile shape and size are defined, there remains to distribute the tiles to physical processors and to compute the final scheduling.

In this paper, we build upon the work of Högsted, Carter, and Ferrante [12]. Given a tiled domain, they aim at determining the cumulated idle time spent by all processors. This cumulated idle time heavily depends upon the tile and domain shapes. A new parameter, called the tile *rise*, is introduced in [12] in order to relate the shape of the iteration domain to that of the tiles, and it is shown to have a significant impact on the idle time. Both parallelogram-shaped and trapezoidal-shaped iteration spaces are considered. We summarize the results of [12] in Section 2. Then we introduce a slightly different model of computation, which enables us to propose new, much shorter proofs of these results, and we extend them in several important directions. More precisely, we provide an accurate solution for all values of the *rise* parameter and for all possible distributions of the tiles to processors, while the authors in [12] deal only with a limited number of cases and provide upper bounds rather than exact formulas. These new results are presented in Section 3. In Section 4, we apply our results to the problem of *hierarchical tiling*, that is, when multiple levels of memory and parallelism hierarchy are involved. In Section 5, we state our conclusions and discuss directions for future research.

2 Determining the Idle Tile of a Tiling

In this section, we summarize the results of Högsted, Carter, and Ferrante [12], who make the following hypotheses:

(H1) There are P available processors interconnected as a ring. Processors are numbered from 0 to $P - 1$.

¹This small list is far from being exhaustive.

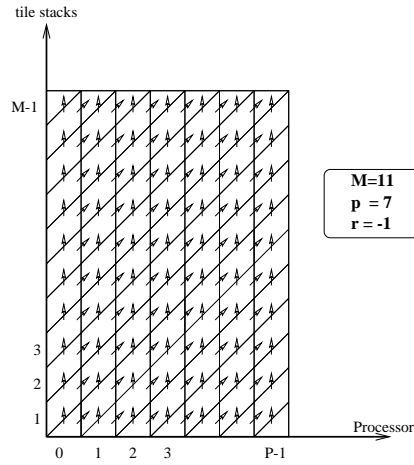


Figure 1: An example of parallelogram-shaped iteration space with parallelogram-shaped tiles. Arrows represent dependences between tiles.

- (H2) Tiles are parallelograms with vertical left and right edges. The size and shape of the tiles are given, so that we deal only with a partitioned (already tiled) iteration space, as in Figure 1.
- (H3) The iteration space is a two-dimensional parallelogram or trapezoid, with vertical left and right boundaries. The first column (and all columns in case of a parallelogram-shaped iteration space) has M tiles.
- (H4) Tiles are assigned to processors using either a one-dimensional full block distribution or a one-dimensional cyclic distribution. In other words,
 - for a block distribution, there are P columns in the iteration space, and all the tiles in column j , $0 \leq j \leq P - 1$, are assigned to processor j ; and
 - for a cyclic distribution, there are bP columns in the iteration space, and all the tiles in column j , $0 \leq j \leq bP - 1$, are assigned to processor $j \bmod P$.
- (H5) The *rise* parameter relates the shape of the iteration space to that of the tiles. It is defined as follows:
 - Let the slope (in reference to the horizontal axis) of the top and bottom edges of the tiles be r_{tile} .
 - If the iteration domain is a parallelogram, let r_{iter} be the slope of the top and bottom boundaries. In this case, Högsted, Carter, and Ferrante [12] define the *rise* r as
$$r = r_{iter} - r_{tile}.$$
 - If the iteration domain is a trapezoid, let r_{iter_top} and r_{iter_bottom} be the slopes of the top and bottom boundaries, respectively. In this case, Högsted, Carter, and Ferrante [12] let $r_t = r_{iter_top} - r_{tile}$ be the rise at the top of the iteration space and $r_b = r_{iter_bottom} - r_{tile}$ be the rise at the bottom of the iteration space.

See Figure 2 for an illustration.

- (H6) Each tile depends upon both its left neighbor and its bottom neighbor (see Figure 1).

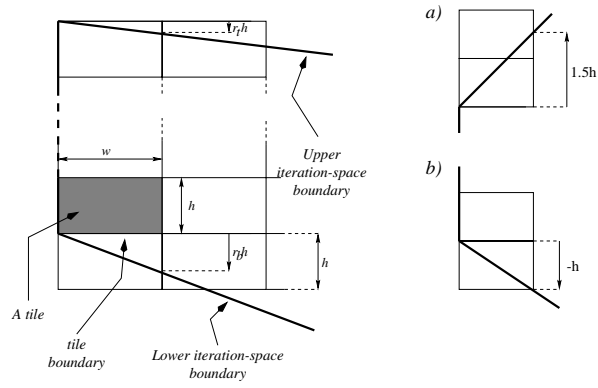


Figure 2: Shape of the iteration space; the rise is positive in (a) and negative in (b).

(H7) Because of hypotheses (H4) and (H6), the scheduling of the tiles is by column. Each processor starts executing its first column of tiles as soon as possible. After having executed a whole column of tiles, a processor moves on to its next column. The time needed to process a tile is T_{comp} (with the notations of [12], $T_{comp} = h \times w$, where h and w are the normalized height and width of a tile). The time needed to communicate data from a tile to its right neighbor is $T_{comm} = c \times T_{comp}$. As stated in [12], c may take any positive value (even though we expect $c < 1$ for large tiles, because the communication volume grows linearly with the tile perimeter, while the computation volume is proportional to the tile volume).

Communications can be overlapped with the computations of other (independent) tiles. Moreover, no communication cost is paid between a tile and its top neighbor, because both are assigned to the same processor².

We summarize in Table 1 the results obtained in [12]. In this table, I_a denotes the cumulated idle time spent by the P processors while executing the tiled iteration space. As pointed out in [12], idle time can occur for two different reasons: (i) a processor may have to wait for data from another processor; or (ii) a processor may have finished all of the tiles assigned to it, and it is waiting for the last processor to terminate execution. In Table 1, condition (C) is a technical condition ($M \geq (1 + c + r)P$) that states that no processor is kept idle when ending the processing of one column of tiles assigned to it; in other words, it can move on to its next column without waiting for any data to be communicated.

3 New Results

In this section we propose new proofs and extend the work of [12].

3.1 Task Graph Framework

The key to our approach is the following: rather than laboriously computing the idle time of each processor, and then summing up the results to get the total idle time I_a , we compute the parallel

²A precise modeling of the communication and computation costs for state-of-the-art machines can be found in [2].

Parallelogram shaped-Block distribution	
$r \geq -1$	$I_a = P(P-1)(1+r+c)T_{comp}$
$r \leq -2$	$I_a \leq \max(c - \frac{1}{2r}, \frac{-r}{2})PT_{comp}$

Parallelogram shaped-Cyclic distribution	
$r \geq -1$ & condition (C)	$I_a = P(P-1)(1+r+c)T_{comp}$
$r \leq -2$	$I_a \leq \max(c - \frac{1}{2r}, \frac{-r}{2})bPT_{comp}$

Trapezoidal shaped-Block distribution- $r_b < r_t$	
$r_b \geq -1$	$I_a = P(P-1)(1 + \frac{r_t+r_b}{2} + c)T_{comp}$
$r_b \leq -2 < -1 \leq r_t$	$I_a \leq (\max(c - \frac{1-r^2}{2r}, 0) + (P-1)\frac{r_t+r_b}{2})PT_{comp}$
$r_b < r_t \leq -1, r_b \leq -2$	$I_a \leq (\max(c - \frac{1-r^2}{2r}, 0) + (P-1)\frac{r_t+r_b}{2} - \frac{r_t}{2})PT_{comp}$

Trapezoidal shaped-Block distribution- $r_t < r_b$	
$-1 \leq r_t \leq r_b$	$I_a = P(P-1)(1 + \frac{r_t+r_b}{2} + c)T_{comp}$
$-(1+c) \leq r_t < -1$ $-1 \leq r_b$	$I_a \leq ((P-1)(1 + \frac{r_t+r_b}{2} + c) - \frac{r_t}{2}) \times PT_{comp}$
$r_t \leq -(1+c) < -1$ $-1 \leq r_b$	$I_a \leq ((P-1)\frac{r_b-r_t}{2} - \frac{r_t}{2})PT_{comp}$
$r_t < r_b \leq -2, r_b \leq -2$	$I_a \leq (\max(c - \frac{1-r^2}{2r}, 0) - (P-1)\frac{r_t-r_b}{2} - \frac{r_t}{2})PT_{comp}$

Table 1: Summary of the results of Högsted, Carter, and Ferrante

execution time t_P with P processors, and we state that

$$P \times t_P = I_a + T_{seq},$$

where T_{seq} is the sequential time, that is, the sum of all tile weights (see Figure 3).

We describe the tiled iteration space as a task graph $G = (V, E)$, where vertices represent the tiles and edges represent dependencies between tiles. A handy view of the graph is obtained by “rotating” the iteration space so that $r_{tile} = 0$. Dependencies between tiles are now summarized

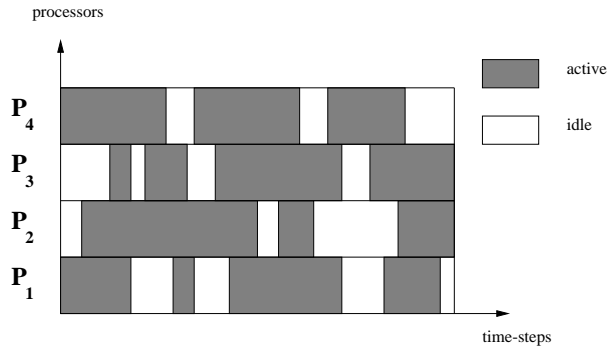


Figure 3: Active and idle processors during execution (illustrating the formula $P \times t_P = I_a + T_{seq}$).

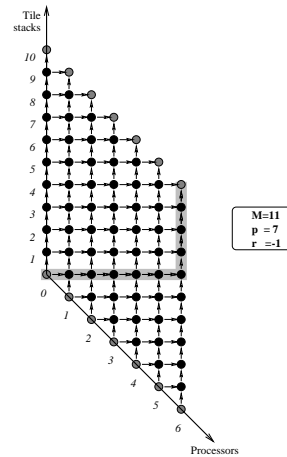


Figure 4: Another view of Figure 1 after rotation.

by the vector pair

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

See Figure 4, where we have rotated the iteration space of Figure 1.

Computing the parallel execution time t_P is a well-known task graph scheduling problem. Since the allocation of tiles to processors is given, the task amounts to computing the longest path in the dependency graph, where the weight of a path is the sum of the weights of its vertices and edges. All vertices have same weight T_{comp} . Horizontal edges have weight T_{comm} (they imply a communication cost), while vertical edges have zero weight (no communication cost due to the allocation). The problem has complexity $O(|V| + |E|)$ (simply traverse the direct acyclic graph G), but we aim at finding a closed-form formula for t_P , specifically, an analytical expression in the problem parameters M , P , r , and c .

3.2 Preview of Results

A summary of our results is given in Table 2. A few comments are in order:

- In Table 2 we assume that M is sufficiently large (see Sections 3.3 to 3.6 for a more precise statement). This hypothesis was implicit in the results in Table 1 quoted from [12] (see Remark 2).
- We use a slightly different model for border tiles (the first and the last one in each column). We view the number of tiles in each column as a continuous function of the rise r , which introduces partial tiles. These partial tiles are assigned a weight that is proportional to their area, exactly as in [12]. For instance consider Figure 5 where $r = -1.5$. In the first column, there is no tile below the horizontal axis. In the second column, there are two tiles below the horizontal axis: the first one is a full tile, while the area (hence the weight) of the second one,

Parallelogram shaped–Block distribution	
$1 + r + c \leq 0$	$t_P = MT_{comp}, I_a = 0$
$1 + r + c \geq 0$	$t_P = [M + (P - 1) \times (1 + r + c)]T_{comp}$ $I_a = P(P - 1)(1 + r + c)T_{comp}$

Parallelogram shaped–Cyclic distribution	
$1 + r + c \leq 0$	$t_P = bMT_{comp}, I_a = 0$
$1 + r + c \geq 0$	$t_P = [M + (1 + r + c) \times (bP - 1)]T_{comp}$ $I_a = P[(bP - 1)(1 + r + c) - (b - 1)M]T_{comp}$
$M \leq (1 + r + c)P$	
$1 + r + c \geq 0,$	$t_P = [bM + (1 + r + c) \times (P - 1)]T_{comp}$
$M \geq (1 + r + c)P$	$I_a = P(P - 1)(1 + r + c)T_{comp}$

Trapezoidal shaped–Block distribution	
$1 + r_b + c \leq 0,$ $r_t \leq r_b$	$t_P = MT_{comp}$
$1 + r_b + c \leq 0,$ $r_t \geq r_b$	$t_P = [M + (P - 1)(r_t - r_b)]T_{comp}$
$1 + r_b + c \geq 0,$ $1 + r_t + c \leq 0$	$t_P = MT_{comp}$
$1 + r_b + c \geq 0,$ $1 + r_t + c \geq 0$	$t_P = [M + (P - 1)(1 + r_t + c)]T_{comp}$

Trapezoidal shaped–Cyclic distribution	
$1 + r_b + c \leq 0,$ $r_t \leq r_b$	$t_P = [bM + \frac{b(b-1)}{2}P(r_t - r_b)]T_{comp}$
$1 + r_b + c \leq 0,$ $r_t \geq r_b$	$t_P = [bM + (\frac{b(b-1)}{2}P + b(P - 1))(r_t - r_b)]T_{comp}$
$1 + r_b + c \geq 0$	see Proposition 4

Table 2: Summary of the results of this paper

on the border of the domain, is half that of a full tile. If we had $r = -1.3$, the weight would be $0.3 \times T_{comp}$ for the boundary tile. Since partial tiles may only occur at the bottom and at the top of the iteration space, their weight has a little impact on the total execution time³.

- For trapezoidal shaped iteration spaces, the total idle time I_a is not reported in the table. However, it can be computed straightforwardly from the relation $Pt_P = I_a + T_{seq}$, where

³We have checked this using a comprehensive simulation program, see [18].

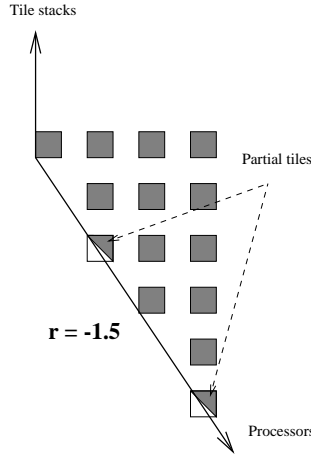


Figure 5: Partial tiles on the boundary.

$$T_{seq} = bP[M + \frac{bP-1}{2}(r_t - r_b)]T_{comp} \text{ (let } b = 1 \text{ for a block distribution).}$$

3.3 Parallelogram shaped–Block distribution

This is the simplest case, and we work it out in full detail. In the formula below, we use the notation a^+ to denote the positive part of a real number a :

$$a^+ = \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{if } a \leq 0 \end{cases}$$

Proposition 1 *Assume a parallelogram-shaped iteration space of size $M \times P$ (block distribution) and rise r . If $M \geq (P - 1)|r|$. Then*

$$t_P = [M + (P - 1)(1 + r + c)^+]T_{comp}.$$

Equivalently,

$$I_a = \begin{cases} P(P - 1)(1 + r + c)T_{comp} & \text{if } 1 + r + c \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Proof All processors have the same workload MT_{comp} . Because of the dependencies, processor $P - 1$ is always the last one to terminate execution. We discuss separately the case $r \leq 0$ and the case $r \geq 0$.

If $r \leq 0$, processor $q = P - 1$ can start processing its first $(P - 1) \cdot (-r)$ tasks at time-step $t = 0$. Then, at time-step $t = -(P - 1)rT_{comp}$, it can continue the processing of its column (i.e., the remaining $M + (P - 1)r$ tiles) *only if* data communicated along the horizontal axis is already available. Otherwise it must wait. To process (and communicate data from) the first $(P - 1)$ tasks of the horizontal axis takes $(P - 1)(1 + c)T_{comp}$. Therefore, the longest path in the dependency graph has length

$$(P - 1) \max(-r, 1 + c)T_{comp} + (M + (P - 1)r)T_{comp}.$$

This longest path is represented in Figure 6.

If $r \geq 0$, processor $q = P - 1$ must wait $(P - 1)rT_{comp}$ time-steps for processor $q = 0$ to complete its first $(P - 1)r$ tasks. Then processor $q = P - 1$ must wait another $(P - 1)(1 + c)T_{comp}$ time-steps

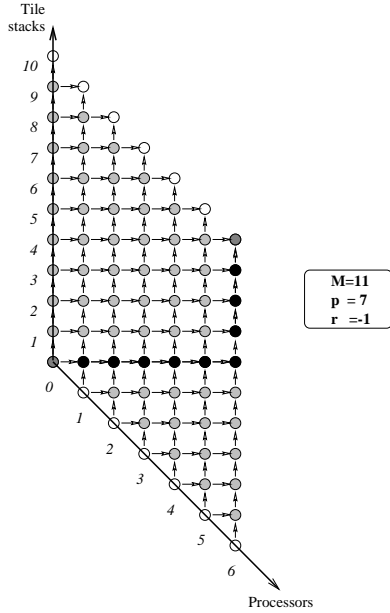


Figure 6: Longest path when $r \leq 0$ (and $M \geq P|r|$)

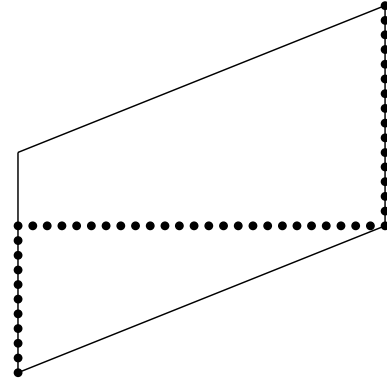


Figure 7: Longest path when $r \geq 0$ (and $M \geq P|r|$)

for executing tiles and communicating data along the horizontal dependence path that leads to its first task. Only then, at time-step $t = (P - 1)(1 + r + c)T_{comp}$, can processor $q = P - 1$ start the execution of its M tasks, and it will not be further delayed during this processing. The longest path in the dependency graph is represented in Figure 7.

We summarize both cases with the single formula

$$t_P = [M + (P - 1)(1 + r + c)^+]T_{comp}.$$

The formula for I_a is derived from the equation $Pt_P = I_a + T_{seq}$, with $T_{seq} = MPT_{comp}$. ■

Remark 1 We see that the results in [12], as reported in Table 1, are inaccurate. A small rise does not prevent from a quadratic idle time; the precise condition is $1 + c + r \leq 0$, which makes good sense because the communication-to-computation ratio of the target architecture has to play a role. In a word, when $1 + r + c \leq 0$, the rise is so small ($r \leq -(1 + c)$) that all tile columns can be processed independently. On the other hand, when $1 + r + c \geq 0$ (which is always true when $r \geq -1$), the total idle time grows quadratically with the number of processors.

Remark 2 The assumption $M \geq (P - 1)|r|$ is needed to ensure the validity of the formulas. See Figure 8: (one of) the longest path is given by the processor $q = Q$, where $M = (Q - 1)|r|$, and the parallel time is

$$t_P = [M + (Q - 1)(1 + r + c)^+]T_{comp}.$$

Note that Q may be much smaller than P in this case. The idle time becomes

$$I_a = \begin{cases} P(\frac{M}{|r|} - 1)(1 + r + c)T_{comp} & \text{if } 1 + r + c \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

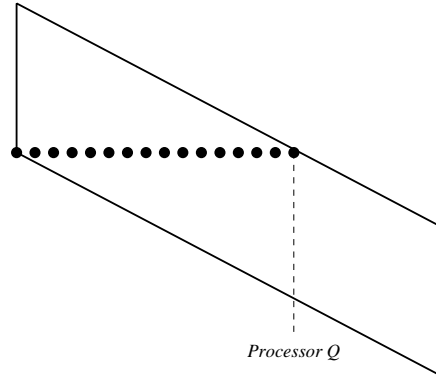


Figure 8: $r \leq 0$ and $M < P|r|$

3.4 Trapezoidal Shaped-Block Distribution

Proposition 2 *Assume a trapezoidal-shaped iteration space of size $M \times P$ (block distribution) and rises r_b (bottom) and r_t (top). If $M \geq (P - 1)(|r_t| + |r_b|)$, then*

$$t_P = [M + (P - 1)((1 + r_b + c)^+ + r_t - r_b)^+]T_{comp}.$$

Proof Let $C(j) = M + j(r_t - r_b)$ for $0 \leq j \leq P - 1$ be the workload of processor j , that is, the total weight of column j . If $1 + r_b + c \leq 0$, then all columns can be processed independently. The total time is given by the largest processor workload: $t_P = C(0) = M$ if $r_t \leq r_b$, and $t_P = C(P - 1) = M + (P - 1)(r_t - r_b)$ otherwise.

If $1 + r_b + c \geq 0$, all processors spend some idle time due to horizontal communications. The discussion is similar to that in the proof of Proposition 1. If $r \leq 0$, processor q , $0 \leq q \leq P - 1$, can start processing its first $j \cdot (-r_b)$ tiles at time-step $t = 0$, but then needs to wait until time-step $t = (1 + c)jT_{comp}$ before processing the rest of its column, that is, the remaining $C(j) + r_b j$ tiles. If $r \geq 0$, processor j has to wait until time-step $t = (1 + c + r)jT_{comp}$ before starting to work. In both cases, processor j terminates the execution of its column at time-step

$$t = [(1 + c + r_b)j + M + j(r_t - r_b)]T_{comp}.$$

Depending upon the sign of $(1 + c + r_b) + (r_t - r_b) = 1 + c + r_t$, this quantity is maximum either for $j = 0$ or for $j = P - 1$.

Altogether, we assemble the results of our case analysis in the above formula. ■

Again, we point out that the condition $1 + r_b + c \leq 0$ is the key to minimizing idle time: If this condition holds, the only idle time that remains is due to the unbalanced workload (with a trapezoidal iteration space, processors have different workloads), but no overhead is due to data dependencies (and to the communications they incur).

3.5 Parallelogram Shaped–Cyclic Distribution

Proposition 3 Assume a parallelogram-shaped iteration space of size $M \times Pb$ (cyclic distribution) and rise r . If $M \geq (P - 1)br$, then

$$t_P = \begin{cases} bMT_{comp} & \text{if } 1 + r + c \leq 0 \\ [(1 + r + c)(bP - 1) + M]T_{comp} & \text{if } 1 + r + c \geq 0, M \leq (1 + r + c)P \\ [(1 + r + c)(P - 1) + bM]T_{comp} & \text{if } 1 + r + c \geq 0, M \geq (1 + r + c)P \end{cases}$$

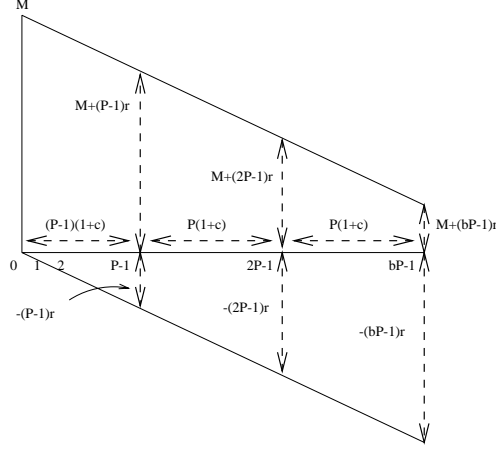


Figure 9: Sketch of the proof with $r_b \leq 0$.

Proof All processors have the same workload bMT_{comp} . If $1 + r + c \leq 0$, all tile columns can be processed independently, and the first part of the result follows.

If $1 + r + c \geq 0$, processor $P - 1$ is always the latest one to terminate execution. We discuss the two cases $r \leq 0$ and $r \geq 0$ separately. If $r \leq 0$, the work of processor $P - 1$ (which is assigned columns $P - 1, 2P - 1, \dots, bP - 1$), can be decomposed as follows (see Figure 9):

$$\begin{aligned} t_P &= \max(1 + c, -r)(P - 1) \\ &\quad \text{max of propagating data along horizontal axis} \\ &\quad \text{and of computing tiles below axis in column } P - 1 \\ &+ \sum_{k=1}^{b-1} \max((1 + c)P, M - Pr) \\ &\quad \text{remaining tiles in column } kP - 1 \\ &\quad \text{and tiles below axis in column } (k + 1)P - 1 \\ &+ M + (bP - 1)r \\ &\quad \text{remaining tiles in column } bP - 1 \end{aligned}$$

For $r \geq 0$, the same decomposition leads to (see Figure 10):

$$\begin{aligned}
t_P = & (1+r+c)(P-1) \\
& \textit{start-up time} \\
& + \sum_{k=0}^{b-2} \max((1+c)P, M) \\
& \textit{tiles in column } j+kP \\
& + M \\
& \textit{tiles in column } j+(b-1)P
\end{aligned}$$

It turns out that, because $1+r+c \geq 0$, the two expressions for t_P coincide: in other words, the last expression is valid for both $r \leq 0$ and $r \geq 0$. This directly leads to the result. \blacksquare

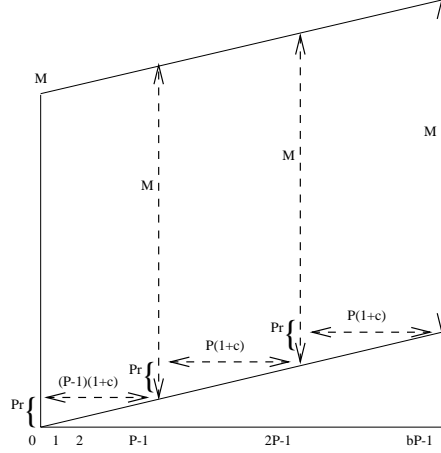


Figure 10: Sketch of the proof with $r_b \geq 0$.

3.6 Trapezoidal Shaped–Cyclic Distribution

The trapezoidal shaped–cyclic distribution is the most difficult case. We have the following result:

Proposition 4 *Assume a trapezoidal-shaped iteration space of size $M \times bP$ (cyclic distribution) and rises r_b (bottom) and r_t (top). If $M \geq (P-1)b(|r_t| + |r_b|)$, then*

$$t_P = \begin{cases} [bM + \frac{b(b-1)}{2}P(r_t - r_b)]T_{comp} \\ \textit{if } 1+r+c \leq 0, r_t \leq r_b \\ [bM + (\frac{b(b-1)}{2}P + b(P-1))(r_t - r_b)]T_{comp} \\ \textit{if } 1+r+c \leq 0, r_t \geq r_b \\ \max(t(j), 0 \leq j \leq P-1) \\ \textit{if } 1+r+c \geq 0 \end{cases},$$

where $t(j) = [(1+c)j + \sum_{k=0}^{b-2} \max((1+c)P, M - Pr_b + (j+kP)(r_t - r_b)) + M + (j+(b-1)P)r_t]T_{comp}$.

Proof If $1+r_b+c \leq 0$, all tile columns can be processed independently. In this case, the processor that has the largest workload is processor 0 if $r_t \leq r_b$ and processor $P-1$ otherwise. The workload of processor j is $\sum_{k=0}^{b-1} C(j+kP)$, where $C(j+kP) = M + (j+kP)(r_t - r_b)$ is the weight of column

$j+kP, 0 \leq j \leq P-1, 0 \leq k \leq b-1$. This leads to the first part of the result: if $r_t \leq r_b$, the maximum is achieved for processor 0 (and $\sum_{k=0}^{b-1} C(kP) = bM + \frac{b(b-1)}{2}P(r_t - r_b)$), while if $r_t \geq r_b$, the maximum is achieved for processor $P-1$ (and $\sum_{k=0}^{b-1} C(P-1+kP) = bM + (\frac{b(b-1)}{2}P + b(P-1))(r_t - r_b)$).

If $1+r_b+c \geq 0$, it is more difficult to determine the longest path. If $r_b \leq 0$, we use the same decomposition as in the proof of Proposition 3 to decompose the work of processor $j, 0 \leq j \leq P-1$:

$$\begin{aligned}
& t(j) = \max(1+c, -r_b)j \\
& \quad \text{max of propagating data along horizontal axis} \\
& \quad \text{and of computing tiles below axis in column } j \\
& + \sum_{k=0}^{b-2} \max((1+c)P, M - Pr_b + (j+kP)(r_t - r_b)) \\
& \quad \text{remaining tiles in column } j+kP \\
& \quad \text{and tiles below axis in column } j+(k+1)P \\
& + M + (j+(b-1)P)r_t \\
& \quad \text{remaining tiles in column } j+(b-1)P
\end{aligned}$$

We have to take the maximum value of these quantities to obtain the parallel execution time:

$$t_P = \max(t(j), 0 \leq j \leq P-1).$$

Now if $r_b \geq 0$, the same decomposition leads to the expression (for processor j):

$$\begin{aligned}
& t(j) = (1+c+r_b)j \\
& \quad \text{start-up time} \\
& + \sum_{k=0}^{b-2} \max((1+c+r_b)P, M + (j+kP)(r_t - r_b)) \\
& \quad \text{tiles in column } j+kP \\
& + M + (j+(b-1)P)(r_t - r_b) \\
& \quad \text{tiles in column } j+(b-1)P
\end{aligned}$$

Again, this last expression for $t(j)$ coincides with the one when $r_b \leq 0$, hence the result. ■

Remark 3 It is not difficult to analytically compute the value of $j, 0 \leq j \leq P-1$, that maximizes $t(j)$ in Proposition 4. This is a simple but tedious case analysis depending upon the problem parameters P, M, c, r_b and r_t .

4 Hierarchical Tiling

As pointed out by Högsted, Carter, and Ferrante [12], tiling may be used for multiple levels of memory and parallelism hierarchy. One important motivation for determining the idle tile of a timing in [12] was, in fact, to demonstrate that such an idle time can have a significant impact on real performance for a large application.

We reuse the example in [12] to illustrate this point. A large rectangular iteration space with horizontal and vertical dependencies is partitioned into supertiles. In turn, each supertile is partitioned into second-level tiles that are assigned to processors. See Figure 11, where supertiles and second-level tiles are rectangular, as opposed to the situation in Figure 12, where supertiles and second-level tiles have a parallelogram shape. To motivate this example, think of a large out-of-core problem, where data is stored on disk. Supertiles are brought in from disk and distributed among the processor main memories (there is an implicit synchronization between two consecutive supertiles). Which is the best strategy, rectangular tiles as in Figure 11, or parallelogram-shaped tiles as in Figure 12? It is stated in [12] that rectangular tiles incur a substantial idle time penalty, whereas parallelogram-shaped tiles do not (at least in steady state—partial tiles do incur a penalty, too).

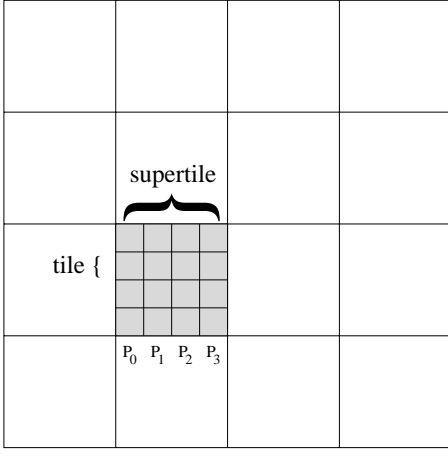


Figure 11: Partitioning the iteration space into rectangular supertiles and second-level tiles (the rise is $r = 0$).

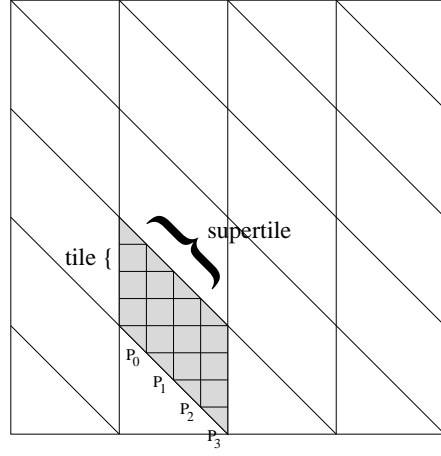


Figure 12: Partitioning the iteration space into parallelogram-shaped supertiles and second-level tiles (the rise is $r = -1$).

The results of the preceding sections enable us to answer the problem: we analytically compute the best partition shape as a function of the iteration space parameters and of the target machine characteristics.

Let h and w be the normalized height and width of second-level tiles whose processing requires $T_{comp} = hw$ time-steps. Assume a block distribution of tiles to processors so that each supertile is of size $Mh \times Pw$: in other words, in a supertile there are P columns of M tiles each. We have $P = 4$ and $M = 5$ in Figures 11 and 12. Let the size of the whole iteration space be $D_1h \times D_2w$, where $D_1 = d_1M$ and $D_2 = d_2P$. With the rectangular partitioning, there are $d_1 \times d_2$ supertiles. With the parallelogram-shaped partitioning, there are $(d_1 + 1) \times d_2$ supertiles, and the first and last supertiles in each column are partial. The following lemma is a direct consequence of the results in Table 2:

Lemma 1 *With the previous notations, assume that $M \geq (P - 1)|r|$. The total execution time to process the iteration space is*

$$\begin{cases} T_{rect} = [M + (P - 1)(1 + c)]d_1d_2T_{comp} \\ \text{for rectangular tiles} \\ T_{rise}(r) = [M + (P - 1)(1 + c + r)^+](d_1 + 1)d_2T_{comp} \\ \text{for parallelogram tiles with rise } r \end{cases}$$

For rectangular tiles, we rewrite T_{rect} as

$$T_{rect} = [M + (P - 1)(1 + c)]\frac{D_1}{M}\frac{D_2}{P}T_{comp},$$

to show that it is a decreasing function of M . In other words, M should be chosen as large as possible, namely, $M = M_{max}$, where M_{max} is such that M_{max} tiles (i.e., $M_{max}hw$ computational points of the iteration space) fit in the (cache) memory of a single processor. For the same value of

M , we choose for r the smallest value such that $(1+r+c)^+ = 0$ (i.e., $r = -(1+c)$), and we derive that

$$T_{rise}(-(1+c)) = M\left(\frac{D_1}{M} + 1\right)\frac{D_2}{P}T_{comp}.$$

We formulate the following proposition:

Proposition 5 *If $M_{max} \leq (P-1)(1+c)\frac{D_1}{M_{max}}$, then $T_{rise}(-(1+c)) \leq T_{rect}$.*

The condition in Proposition 5 will always be true for large enough domains. In other words, parallelogram-shaped supertiles will lead to the best performance.

5 Conclusion

In this paper, we have extended results by Högsted, Carter, and Ferrante [12], and we have been able to accurately determine the idle time of a tiling for both parallelogram-shaped and trapezoidal-shaped iteration spaces. We have provided a closed-form expression of the idle time for all values of the rise parameter, for a block distribution as well as for a cyclic distribution.

Furthermore, we have used our new results in the context of hierarchical tiling. Although we have dealt only with a particular instance of the multilevel tiling problem, we believe our approach is general enough to be applied in several situations (such as those described in [7]).

Finally, we point out that 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. An interesting research direction would be to extend our approach so as to incorporate processor speed as a new parameter of the tiling problem.

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] Rumén 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, <http://www.univ-valenciennes.fr/limav/andonov>.
- [3] Utpal Banerjee. An introduction to a formal theory of dependence analysis. *The Journal of Supercomputing*, 2:133–149, 1988.
- [4] Pierre Boulet, Alain Darte, Tanguy Risset, and Yves Robert. (pen)-ultimate tiling? *Integration, the VLSI Journal*, 17:33–51, 1994.
- [5] 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.
- [6] 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.
- [7] L. Carter, J. Ferrante, S. F. Hummel, B. Alpern, and K.S. Gatlin. 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>.
- [8] 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):108–120, 1996.
- [9] 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).
- [10] Alain Darte, 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.
- [11] J. J. Dongarra and D. W. Walker. Software libraries for linear algebra computations on high performance computers. *SIAM Review*, 37(2):151–180, 1995.
- [12] 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.
- [13] François Irigoin and Rémy Triolet. Supernode partitioning. In *Proc. 15th Annual ACM Symp. Principles of Programming Languages*, pages 319–329, San Diego, CA, January 1988.
- [14] 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.

- [15] 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.
- [16] 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.
- [17] J. Ramanujam and P. Sadayappan. Tiling multidimensional iteration spaces for multicomputers. *Journal of Parallel and Distributed Computing*, 16(2):108–120, 1992.
- [18] Fabrice Rastello. Techniques de partitionnement. Master’s thesis, Ecole Normale Supérieure de Lyon, June 1997.
- [19] Robert Schreiber and Jack J. Dongarra. Automatic blocking of nested loops. Technical Report 90-38, The University of Tennessee, Knoxville, TN, August 1990.
- [20] 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.
- [21] 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.
- [22] 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.