



Ecole polytechnique
Promotion X2012
NICOLET Victor

RAPPORT DE STAGE DE RECHERCHE

A task-based approach to stencil computations

Département d'informatique
Champ : INF 591
Directeur de stage : BOURNEZ Olivier
Maitre de stage : COHEN Albert
23 mars 2015 - 8 juillet 2015
Département d'informatique, Ecole Normale Supérieure
45 rue d'Ulm 75005 Paris
France

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 4 |
| 1.1 | Motivation | 4 |
| 1.2 | Contribution and outline | 4 |
| 2 | Background | 4 |
| 2.1 | Stencil applications | 4 |
| 2.1.1 | Characterization | 4 |
| 2.2 | Program transformations and polyhedral framework | 5 |
| 2.2.1 | Loop transformations | 6 |
| 2.2.2 | Polyhedral model | 6 |
| 2.3 | Tiling | 7 |
| 2.3.1 | Tiling techniques in previous research work | 8 |
| 2.4 | Scheduling | 10 |
| 2.4.1 | Loop scheduling in OpenMP | 10 |
| 2.4.2 | Tasks and libkpn | 10 |
| 2.4.3 | Dynamic task scheduling | 12 |
| 3 | Harris Corner Detection | 13 |
| 3.1 | Harris Corner Detection's pipeline and parallelization opportunities | 13 |
| 3.2 | Algorithm with tasks | 14 |
| 3.3 | Results | 14 |
| 3.4 | 1-d Jacobi implementations | 15 |
| 3.5 | Half-diamonds tiling | 15 |
| 4 | A task-based approach | 16 |
| 4.1 | First results | 17 |
| 4.2 | Tile shape determination at runtime | 17 |
| 5 | Conclusion and future work | 19 |
| 6 | References | 20 |
| | References | 20 |

Acknowledgements

I would like to express my gratitude for my internship advisers, professor Albert Cohen, for giving me the opportunity to work with them, and his advice and discussions. I also want to thank Nhat Minh-Lê for his help and availability with Libkpn and general advice, and also all those who worked on the project, Adrien Guatto and Robin Morisset. Finally, I'd like to thank the whole PARKAS team for their welcome during my stay.

1 Introduction

1.1 Motivation

Multi-core processors with shared memory are now common in devices ranging from mobile phones to supercomputers. There has been a growing interest in providing frameworks that would help the common programmer to build efficient programs that would execute on multicore architectures. Though there have been different approaches, they all try to adress the problem of :

- providing the tools to program parts that can be solved concurrently.
- implementing an overall control and coordination mechanism.

Here, we will adresss computational problems in which breaking the program into parts for concurrent execution is not straightforward. During the computation, time is often spent in nested loops, often with iterations depending on each other, and splitting or “tiling” the iteration space for work distribution has often proved problematic. Previous research has provided many solutions to this problems, aiming to improve both parallelism and memory efficiency. Even compilers can provide automatic parallelization today, thanks to models like the polyhedral model. But these solutions are not always optimal. We will provide a simple tiling for time iterated stencil applications that yields good performance on jacobi 1-d, a simple time-iterated stencil.

The scheduling of the different parts of the program is also an interesting field of research. We will compare the performance static scheduling of for loops used in OpenMP to the dynamic scheduling of a task network in Libkpn, the library developped in PARKAS.

1.2 Contribution and outline

The main purpose of my internship was is to show that tasks and dynamic scheduling can perform as well as the static scheduling used in widely used shared-memory models such as OpenMP. I initially focused on learning about the state of the art in parallel programming on shared-memory architectures. This included research on program transformations and dynamic scheduling.

The first goal I had was to use inner parallelism in an image processing application, Harris Corner Detection, but this didn’t lead to convincing results. The following report will mostly explain what I learned and how the second subject was treated : showing a new tiling for a stencil application, unidimensionnal jacobi (1-d jacobi), and compare two different implementations using OpenMP4 and Libkpn.

2 Background

2.1 Stencil applications

Stencil applications are a broad family of compute intensive applications, often used in the graphic or scientific domain. We can characterize them by an update of a point in a grid, using neighboring points. They have properties that make them a good subject for optimization. We will take as example the jacobi in one spatial dimension as example (see 1a).

2.1.1 Characterization

A stencil computation is a time-iterated nearest-neighbor computation that operates on each point in a grid and updates each grid point using data in its neighboring points, in recent time steps. We use the stencil characterization presented in [1] :

- Grid points are updated using the values of neighboring points in previous and recent time-steps. Nested loops have always a time-step iteration as outermost loop, and we assume there is no horizontal dependence between statements (the values of a points doesn't depend directly on the value of other points in the same time-step).
- If we have d spatial dimensions, the computations uses a $d + 1$ dimensional grid, the outermost dimension being the time dimension. In the implementation, the grid will be reduced to the spatial grid only, for memory efficiency.
- The stencil can always be written in a single-statement of the following form, with a stencil function f :

$$value_p[t + 1] = f(value_p[t], value_{neighbors_of_p}[t])$$

A stencil computation computes the stencil for each point of the grid over many time steps. Theses applications are quite simple to implement using nested loops, but they are subject to many problems of optimization. These are mainly loop transformations, since the computation is contained in the nested loops.

We will consider the following stencil function f in a one-dimensional grid \mathbb{G} , with a data grid $A_{\mathbb{G}}$:

$$\mathbb{G} = [0, N - 1] \subset \mathbb{N} \quad , \quad A_{\mathbb{G}} \subset \mathbb{R} \quad , \quad f(A) = (A[t, i] + [A[t, i - 1] + A[t, i + 1])/3$$

The stencil computation can be written simply with two nested loops and one statement, using a time-space grid 1a or two data grids figure 1b. The use of a two dimensional array with a time dimension permits perfectly nested loops and would allow simple transformations, with only one statement. But this is very memory-inefficient. The second form of the program adds a second statement and thus more dependencies.

| | |
|---|--|
| <pre> for(t = 0; t < T; t++) for(i = 1; i < N-1; i++) S1 : A[t+1][i] = 0.3333 * (A[t][i-1] + A[t][i] + A[t][i+1]); </pre> | <pre> for (t = 0; t < T; t++) for (i = 1; i < N-1; i++) S1 : B[i] = 0.3333 * (A[i-1] + A[i] + A[i+1]); for (i = 0; i < N; i++) S2 : A[i] = B[i]; </pre> |
|---|--|

(a) Perfectly nested 1-d jacobi

(b) Imperfectly nested 1-d jacobi

Figure 1: 1-d Jacobi

In the following, we will consider mainly stencil computations, or at least computations that are structured likewise : and outer “time” loop and dependencies from one time step to another, spanning along the iteration space.

2.2 Program transformations and polyhedral framework

Most highly computational algorithms are built with a set of nested loop and some statements, accessing and reading data. The program transformations we need to perform concern mostly these loops. In figure 1, we have one outer-loop iterating on time, and an inner loop on the spatial dimension.

Data dependency In the perfectly nested loop in figure 1a the dynamic instances of statement $S1$ at (i, t) reads data that has been written by $S1$ at $(t - 1, i)$, $(t - 1, i - 1)$ and $(t - 1, i + 1)$. We call this a *data dependency* : a dynamic instance of a program statement refers to data of a preceding statement. We will refer simply to a *dependency* since we're not interested in other kinds of dependencies.

2.2.1 Loop transformations

To exploit parallelism in application, we have to modify the nested-loops structure, while ensuring the program still produces the correct output. We present some correct loop transformations [2] to modify the program and produce the correct strategy.

- loop skewing can be useful to enable parallelism. Since it only changes loop bounds and alters accordingly the use of index variables, it is always legal to perform skewing on loops.
- loop tiling is useful to improve cache performance. When performing loop tiling on an outer loop (containing other loops), we divide the iteration space into chunks and iterate over these parts.
- loop distribution (or loop splitting, loop fission) can be used to distribute a single loop into different loops. This can create perfect nested loops, or reduce the dependences in the nested loops. The inverse transformation is loop fusion.

Loop tiling and distribution legality is constrained by the dependences. We will speak of tile dependencies when we perform tiling, and see that loop transformations alter these dependencies.

| | |
|--|---|
| <pre> for (t = 0; t < T; t++) for(i = t + 1; i < t+N-1; i++) S1 : B[i - t] = 0.3333 * (A[i-1 - t] + A[i - t] + A[i+1 - t]); </pre> | <pre> for(Tt = 0; tT < T; Tt+=64) for(Ti = 0; Ti < N; Ti+=64) for (t = Tt; t < Tt + 64; t++) for(i = Ti; i < Ti + 64; i++) S1 : B[i] = 0.3333 * (A[i-1] + A[i] + A[i+1]); </pre> |
|--|---|

(a) Skewed 1-d jacobi

(b) Tiled 1-d jacobi

Figure 2: 1-d jacobi after loop transformations

Yet these transformations can be expressed in a more high level model, which will also allows to find suitable transformations for parallelization. Programming manually transformed loops is a very tedious work and error prone. During my internship, I programmed various tilings using loop transformations, but the goal today in high performance computing is to build an automatic system using the model described in the following paragraph. Though I didn't use it myself except for some experiments with Pluto [CITE THE ARTICLE REFERRING TO PLUTO], I couldn't dissociate the optimization of iteration spaces and tilings from this field of research, also very important in the team where I was working.

2.2.2 Polyhedral model

Simple transformations are easy to achieve without excessively complexifying the original program. But when loop nests become more complicated, a mathematical framework can prove useful to express loop skewing, tiling, distributing etc. We will express the properties of our stencil applications in a linear-algebraic representation of programs, the polyhedral model, and provide a brief description of it. Given a representation of a program, it will be easier to express further transformations and show their validity.

This model has been proved very useful in automatic parallelization at compile time [3], because it provides many abstractions to reason about program transformations.

Polyhedral representation of programs A program is a set of dynamic instances of statements. Each instance S is defined by its iteration vector \vec{i} containing the values

of the indexes of the loops containing the statement. When inner loop bounds depend on outer indexes values, the set of iterations vectors representing different instances of a statement define a polytope D_s (see figure 3a), the domain of the statement, with m_s its dimensionality. In this model, we consider only convex set as iteration spaces.

Polyhedral dependencies The difficulty of parallelizing stencil applications arises when considering the dependencies between dynamic instances of the statements. Dependencies are affine and exact, and only carried by the time loops in our stencil applications. The data dependence graph is a directed multi-graph, with each node S_i representing a statement and each edge $e \in E$ representing a dependence between statements S_i and S_j . e can be represented here a constant distance vector by analyzing the source and target iterators in the domain D_s .

In 1-d jacobi we have only 3 dependencies $\vec{d}_1 = (-1, 1)$, $\vec{d}_2 = (0, 1)$ and $\vec{d}_3 = (1, 1)$ (3b)

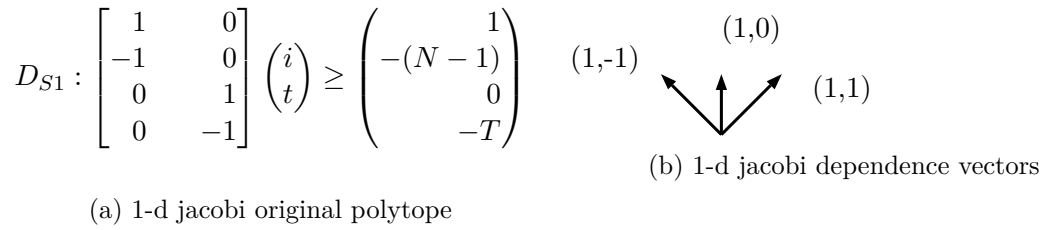


Figure 3: Elements of the polyhedral representation of 1-d jacobi

2.3 Tiling

Tiling in stencil applications is a key transformation to improve locality, memory reuse and parallelism. There has been a considerable amount of research in this domain, and we present here a summary of previous work.

Tiling legality In figure 1b we have an example of imperfectly-nested loops. It raises more problems when searching for valid tiling hyperplanes. Using the polyhedral framework, we have conditions provided in prior research, expressed in Theorem 1 of [4]. A legal tiling is a tiling that ensures that we can construct a total order between tiles for execution (there is no tiles that depend on each other) and provided inter-tile dependencies are satisfied, a tile can execute atomically.

Statement-wise hyperplane 1. A hyperplane for a statement S_i is of the following form:

$$\phi_{S_i}(\vec{x}) = \vec{h} \cdot \vec{x} + h_0$$

where $\vec{x} \in D_s$ is an instance of the statement S_i , h_0 is the translation and \vec{h} is the hyperplane normal vector.

Legal tiling hyperplane 1. Let ϕ_{S_i} be a one-dimensional affine transform for statement S_i . For $\{\phi_{S_0}, \dots, \phi_{S_k}\}$ to be a legal tiling hyperplane, the following should hold for each edge e from S_i to S_j :

$$\phi_{S_j}(t) - \phi_{S_i}(s) \geq 0, P_e$$

In our example with only one statement, a legal tiling hyperplane with normal vector $\vec{h} = (h_i, h_t)$ must satisfy the following conditions :

$$\vec{h} \cdot \vec{d}_1 = h_t - h_i \geq 0, \vec{h} \cdot \vec{d}_2 = h_t \geq 0, \vec{h} \cdot \vec{d}_3 = h_i + h_t \geq 0$$

Here rectangular tilings are not legal . Parallel computations should minimize inter-processors communication and rectangular tilings are not satisfactory when dependence vectors have negative components along the space dimension.

Tiling for locality The spatial dimensions of most problems involving stencil computations are usually quite large, often larger than the cache sizes or other on-chip memory size. Parallelism is present along the spatial dimension, but tiling is necessary to improve data locality during the computation and the different memory speeds. The objective here is to fit the tile spatial domain into on-chip memory, potentially at different levels. The tiling concerns here the inner loops.

Tiling for reuse The goal is to work around the problem of limited memory bandwidth by using present memory as much as possible before releasing it to memory levels with longer access delays. Many computations are throttled by this issue, an efficient schedule along the time dimension can dramatically improve the performances. Tiling only the inner loops to provide locality is not sufficient. But transforming the nested loops including the outer one can enhance memory reuse between different tiles.

Tiling for communication-free parallelism In current multiprocessor architectures, communication between processing units is costly and has to be avoided. Tiling for parallelism implies finding calculations that can be executed at the same time and without data from other concurrent calculations.

A tiling enables a tile-wise *concurrent start* if all that tiles along one dimension of the iteration space can be started concurrently. If tiles carry dependences along one dimension, they can't be started concurrently along this dimension.

When dependences span the entire iteration space, there is often a trade-off between locality, memory reuse and and communication-free parallelism.

2.3.1 Tiling techniques in previous research work

Skewed tiling Skewing the iteration space is a standard tiling technique for time-iterated stencils along with rectangular space tiling. Visually speaking, figure 4 presents the shape of the tiles, parallelograms along the time dimension, and regular and legal tiling long the spatial dimensions. This improves memory reuse but creates inter-tiles dependencies along the space dimension, restricting parallelism to tile wavefronts along diagonals, and forbidding a concurrent start.

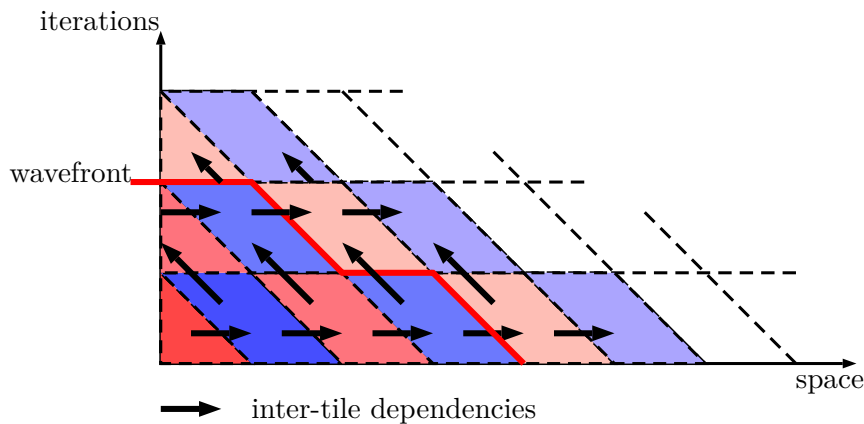


Figure 4: Skewed tiling : inter-tiles dependences and wavefront

Overlapped tiling Using overlapped tiling enables a concurrent start among the spatial dimension and removes synchronization between tiles. But it also causes a significant amount of redundant computation, especially if the time dimension is important. Moreover, redundant computations often need to be stored in shared memory, causing another overhead. Figure 5 shows schematically the shape of the tile along on spatial dimension and the time dimension. The shape of the tile is defined by the dependencies projected on the spatial dimension. Thus, overlapped tiling can be very efficient on low-order stencils and few iterations.

Hierarchical overlapped tiling [5] has proven efficient for stencil computations by balancing communication overhead and redundant computation cost.

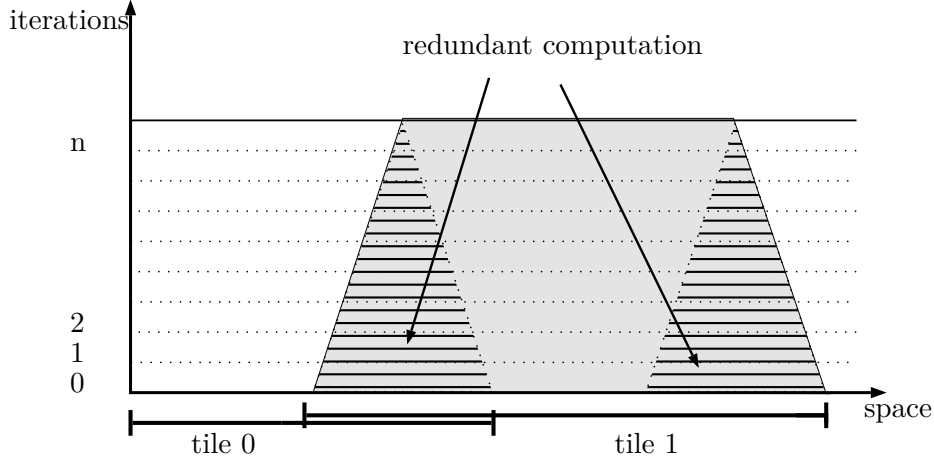


Figure 5: Overlapped tiling

Split tiling Split tiling [6] does not make redundant computation while allowing a concurrent start along the space dimension. In figure 6, the lower darker tiles are executed concurrently in a first step, then the lighter tiles, mapped onto threads so to reuse a maximum amount of memory. The vertical arrows represent the step by step calculation : between each level and between light and dark tiles, we have a barrier. However, we could schedule the computation differently, using tasks where we wait for dependences to be satisfied. The new tiling, diamond tiling, shows a more natural space partition to use a schedule without space-wide barriers.

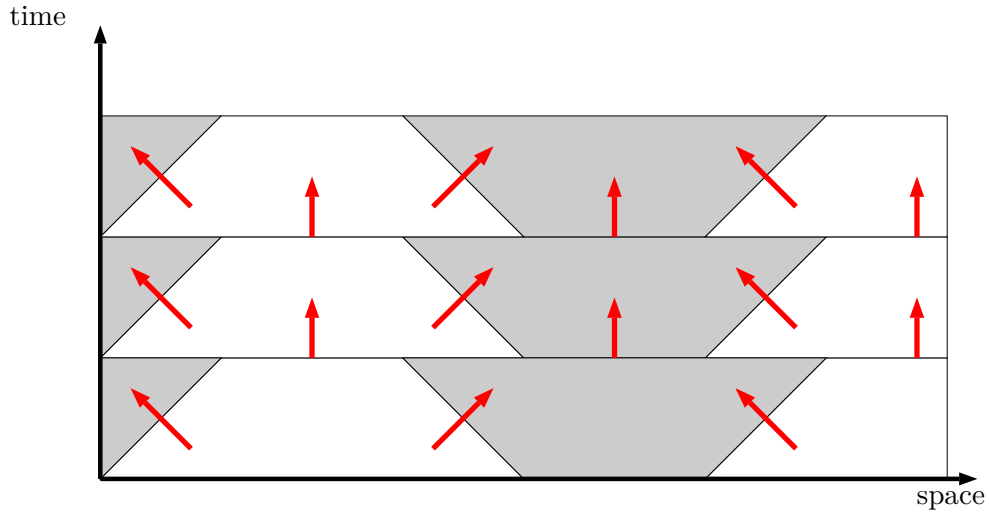


Figure 6: Split tiling and inter-tile dependencies

Diamond tiling Diamond tiling [1] allows a concurrent start along space dimensions.

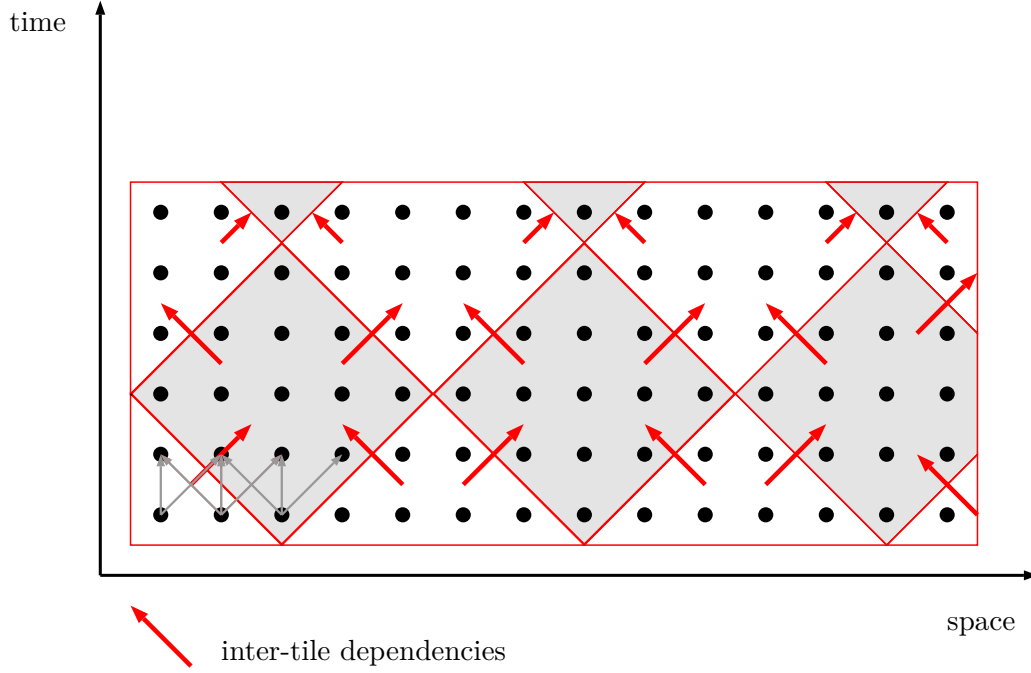


Figure 7: Diamond tiling with inter-tile dependencies

2.4 Scheduling

The second part of this work was to show that the task approach to stencil computations can be as efficient as the distributed loops. An efficient task scheduler will also avoid the programmer the additional effort of thinking about how to help a static loop schedule to help improve reuse.

2.4.1 Loop scheduling in OpenMP

In the exploitation of loop level parallelism, scheduling the iterations can be a difficult task for the programmer. In the shared-memory programming model OpenMP, there is typically two scheduling strategies : static or dynamic. For loops with good balances among iterations,

2.4.2 Tasks and libkpn

In [7] G. Kahn defines a language for parallel programming, viewing parallel programs as a networks of computation units, the tasks. These wrap parts of the program meant for concurrent execution.

A parallel program is an oriented graph, where nodes are the tasks (or processes) and edges communication pipes. There is additional edges that represent input and output of the program. Each process computes a sequential program, receiving data only from its incoming edges and issuing data on his outgoing lines. In the following paragraph, we distinguish the *dataflow dependencies* introduced by Kahn and the *creation dependencies*. The first imply a partial execution order. A task cannot execute before all its incoming dependencies haven't been satisfied. Also a node cannot wait or produce data on more than one channel simultaneously.

Creation dependencies represent the parent-child relations between task. There is also an order between sibling task since the spawning or creation of child tasks is sequentially ordered in the parent task.

We distinguish here the task graph, containing task and pipes, from the creation tree.

In previous implementations of task parallelism languages such as Cilk and OpenMP 4, dataflow dependencies are restricted to tasks that are directly related in the creation tree or sibling tasks, depending on the model. On the contrary, Libkpn allows arbitrary dependency creations between tasks .

Two models can be distinguished in previous languages and frameworks. We illustrate these two models in figure 8. Dashed arrows represent the creation tree, whereas full arrows are dataflow dependencies. In OpenMP and SMPs (figure 8a, the dependencies are built on sequentially ordered memory descriptors. A master task (top dot) creates all its children sequentially. The execution order of its children is determined by the data dependencies.

In the Cilk-like models of multithreaded computations, the creation tree is not limited to one task sequentially spawning the task graph, since each process can spawn children. The limitation here is in the fact that a parent can only wait on its child. In this model of multithreaded computation in figure 8b, threads (diamonds) are non-blocking function and cannot wait but instead spawn successor threads (dashed arrows). Then the successors wait for the data dependency to be satisfied. The procedures enclosing the threads (dashed rectangles) are sequentially executed parts of the program, similar to tasks. In Libkpn,

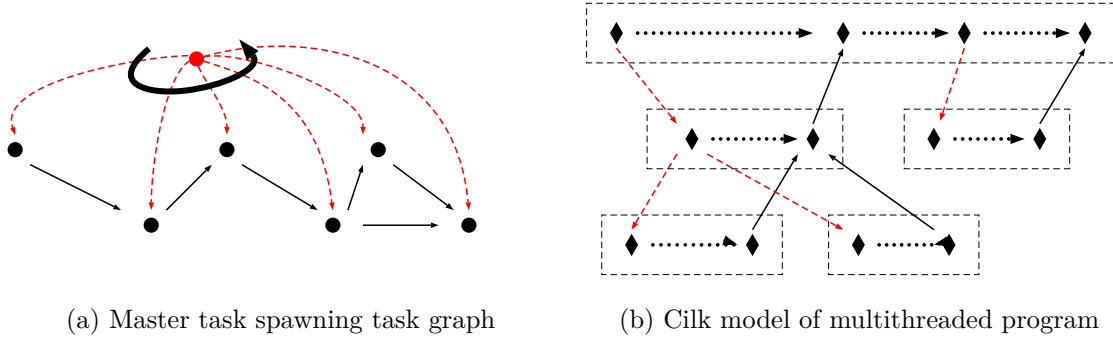


Figure 8: Different models for task parallelism languages

the model doesn't restrict the creation of arbitrary dataflow dependencies. This can result in cyclic task graphs or error-prone implementations, but any process network can be implemented. Data consumer-producer relationships are represented as first-class objects that can be passed between tasks. A channel between two tasks can be created without knowing the producer or the consumer. In figure 9, when the bottom task is executed, the consumer task may not have been created. The dependencies here do not follow the two-branch creation tree either

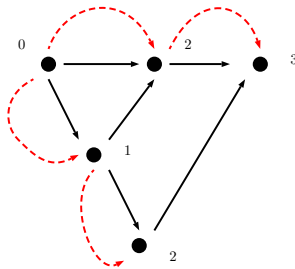


Figure 9

2.4.3 Dynamic task scheduling

While some systems offer compile-time scheduling [8], a natural approach to task-based parallelism is scheduling at runtime to adapt to dynamic data dependencies and unknown hardware platforms. Tasks can be in three states : either running, ready when all the dependencies have been satisfied, or stalled, when the tasks waits for data. Libkpn's relies on a work-stealing scheduler with additional mechanisms. The work-stealing algorithm [9] manages the transfer between the set of ready and running tasks. The additional mechanism is needed to deal with the stalled tasks.

During its execution, a task can :

- spawn another task.
- die, when its execution terminates.
- stall, if it waits for a dependency to be satisfied.
- enable another task, by satisfying the dependency on which it stalled.

Work-stealing scheduler The work-stealing scheduler outlined here schedules the execution of a set of tasks with a dependency graph on a parallel computer. Each processor has a ready deque containing tasks ready to be executed, with a top and a bottom. Tasks can be inserted at the bottom of the deque, and removed at either end. A processor treats its ready deque as a call stack : it removes the task at the end to get work, but other processors can steal tasks from the deque by removing the top task.

A processor P begins to work by pulling the bottom task T_a of its ready deque, and executes it until T_a spawns, stalls, dies or enables another task :

- if T_a **spawns** a task T_b , then T_a is placed at the bottom of the ready deque of P and P starts executing T_b .
- when T_a **stalls**, P checks its ready deque and begins work on the bottom task. If the queue is empty, it tries stealing the top task of a randomly chosen processor.
- when T_a **dies**, it follows the same rule as when T_a stalls.
- if T_a **enables** another task T_b , T_b is placed at the bottom of the ready deque of P .

Here when a processor attempts to steal work from another processor, it chooses his victim randomly. The attempt to steal can either fail, if the deque of the victim is empty, or succeed. If it fails, the processor chooses another victim at random, and attempts stealing again, until it finds work.

This algorithm ensures some execution properties with bounds on memory and time. If we call t_1 the minimum serial execution time of the multithreaded computation on a single processor and t_∞ the minimum execution time on an infinite number of processors, the execution time of the multithreaded program on p processors is bounded by $\frac{t_1}{p} + O(\text{Infty})$. The space requirement is also bounded by $p * S_1$ with S_1 the space required by the computation on a single processor. This means that there is no space overhead coming from the multi-processor execution, and the time overhead is bounded.

Libkpn scheduler Libkpn's scheduler is a work-stealing algorithm with some optimizations. The inner mechanisms of the scheduler haven't been described above, but in multiprocessor environment, task-scheduler synchronization is needed when any user-space execution hands back control to the scheduler. In this case Libkpn relies on relaxed C11

primitives. The internal algorithm, managing the tasks cycles is correct but yields better performance since it gets rid of most of the heavy acquire/release lock synchronizations that would be used in a naive implementation. The scheduling algorithm itself is based on work-stealing with some other heuristics concerning task-stealing to improve reuse. A processor can for example steal either from the top or from the bottom of a ready deque of another processor. Since my internship is not focused on the scheduling but in implementing applications that use this scheduler, I will not dwell on this subject detailed in **** PAPER? ****. But the understanding of the tool I was going to use seemed quite necessary to catch the differences between Libkpn and other tools.

3 Harris Corner Detection

In the first part of my internship, I had to explore the possibility of using tasks in OpenMP for Harris Corner Detection pipeline [10]. Corner detection is widely used in computer graphics, and as a comparison point we took a state-of-the-art compiler and domain specific language for image processing pipelines [11].

3.1 Harris Corner Detection's pipeline and parallelization opportunities

The algorithm of Harris Corner Detection is an image processing pipeline divided in several elementary steps involving simple computations or stencil computations :

- Compute I_x and I_y derivatives of the image along x and y directions.
- Compute the product of derivatives at every pixel I_{xx} , I_{xy} and I_{yy} .
- Smooth the image, calculating at each point the mean of the neighboring points : S_{xx} , S_{xy} , S_{yy} .
- The output is defined as follows :

$$H = \begin{bmatrix} S_{xx} & S_{xy} \\ S_{xy} & S_{yy} \end{bmatrix}, \text{ Harris} = \text{Det}(H) - k(\text{Trace}(H)), k \in \mathbb{R}$$

The pipeline graph in figure 10 shows how the different steps are linked together, and the parallelization opportunities in the pipeline. For example, steps I_{xx} , I_{xy} and I_{yy} can be executed in parallel and only depend on I_x and I_y . In computer vision, the size of

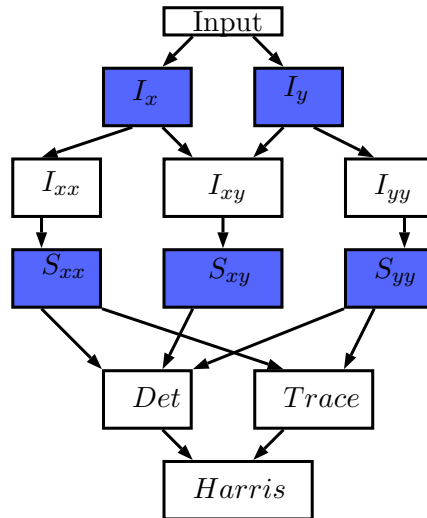


Figure 10: Harris-corner detection image processing pipeline

the images is often far greater than the sizes of fast memory. To make an efficient use of on-chip memory and to expose more parallelism, tiling the spatial dimension is also useful. Because of the stencils in the pipeline, there is need for communication between tile, thus a careful scheduling.

In the implementation from Polymage, the pipeline parallelism is not exploited. The tool uses a heuristic for grouping different stages of the pipeline and use overlapping tiles in the computation.

In the task approach, the goal is to get rid of redundant computation and to expose more parallelism in the pipeline.

3.2 Algorithm with tasks

In order provide more parallelism in the pipeline, we implemented Harris Corner Detection using both the tiling extracted from the output of Polymage and task for the stages of the pipeline. But some steps are executed sequentially in the pipeline (I_{xx} and S_{xx} , I_{yy} and S_{yy} , etc.), therefore we grouped them into one task. We also grouped the three last stages, given their weak computational weight to balance the different tasks. In each task (I_{xx}/S_{xx}) we fuse the loops to produce only two nested loops with one statement. The task graph in figure 11 shows the modified pipeline.

For each spatial tile in the image, we create the set of tasks corresponding to the pipeline

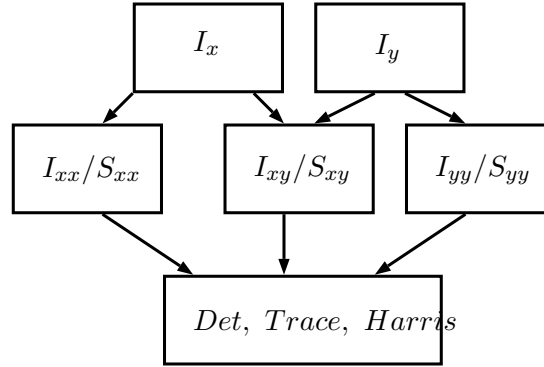


Figure 11: Harris Corner Detection with fused stages

with their dependencies, both intra-tile and inter-tile. The need of synchronization between tasks of different tiles comes from the absence of overlapping. It is the scheduler's job to use memory efficiently in consideration of these dependences.

3.3 Results

The benchmarks were run on an Intel Xeon E5-2630 with 12 cores (24 hyperthreads) and 15Mb cache, and a dual-core Intel Core i7-4510U with 4Mb cache, and we used *gcc 4.9.1* and *icc 15.1* to compare versions from different compilers with 3 levels of optimization.

The algorithm using tasks doesn't perform as well as the optimized version from Polymage. The gain from the absence of redundant computation, and a potential good reuse of memory in task scheduling doesn't compensate the loss due to more synchronization and fused-loops distribution. But the excellent performance of Polymage, nearly scaling linearly with the number of cores (see [11] for more results) also relies on good compiler optimizations of the heavily computational step resulting from loop-fusion.

We tried to measure performance without advanced compiler optimizations, using *-O2* for *gcc* and *icc*. This disables automatic loop vectorization. Since most of the time is spent in the computational fused loops, performance is badly reduced and parallel implementations of Harris Corner Detection are slower than the sequential with automatic vectorization.

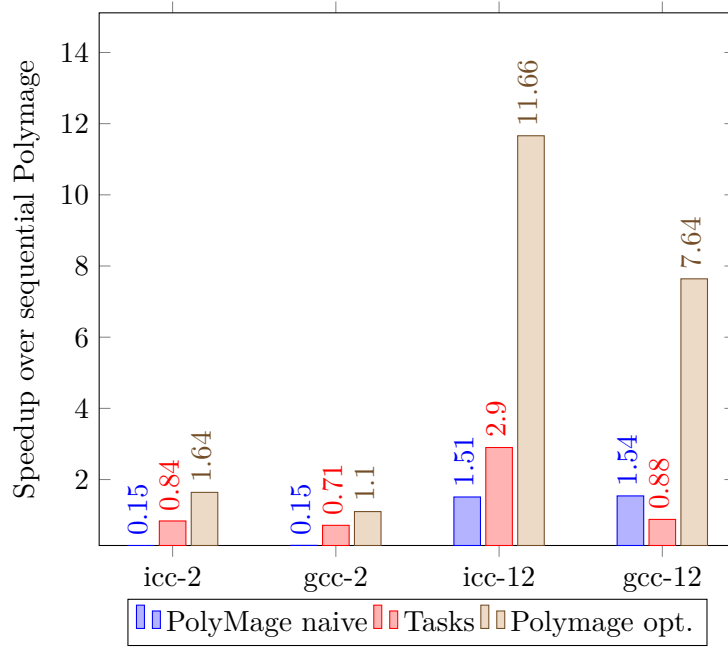


Figure 12: Speedup relative to sequential Polymage, on machines with 2 and 12 cores, with gcc4.9.1 and icc 15.1

But the version using OpenMP tasks performs better without compiler vectorization than Polymage version using heavily fused loops.

The different implementations are available on github at :
<https://github.com/victornicolet/harris-corner-implementations> ,
 with both the sources from Polymage and the algorithm using tasks.

4 1-d Jacobi

Time iterated stencils can be good candidates for an efficient use of tasks. Here we are interested in small iteration numbers, and we implement a specific tiling strategies in this case.

The OpenMP implementations are available at
<http://www.github.com/victornicolet/tiling-strategies>,
 along with some results.

4.1 Half-diamonds tiling

Diamond tiling has proven an efficient tiling strategy since it provides a concurrent start, locality and possible reuse. In our case, we consider only cases when the number of iterations is low. The iteration space is tiled using hyperplanes $(1, 1)$ and $(-1, 1)$ as shown in figure 13. We easily show the legality of this tiling since with the provided hyperplanes. For each statement S with the tiling hyperplanes $h_1 = (1, 1)$ and $h_2 = (-1, 1)$ we have $h_1.S - h_2.S = 2 * i > 0$ and for each dependency $d \in \{(-1, 1), (0, 1), (1, 1)\}$ we have $h_1.d \geq 0$ and $h_2.d \geq 0$

We have implemented a first naive version, using a global barrier between lower and upper tiles, and static scheduling.

Grouping tiles take advantage of the memory architecture and fast L1 caches, it can be more interesting to adjust the size of the chunks in the loops, and the execution schedule.

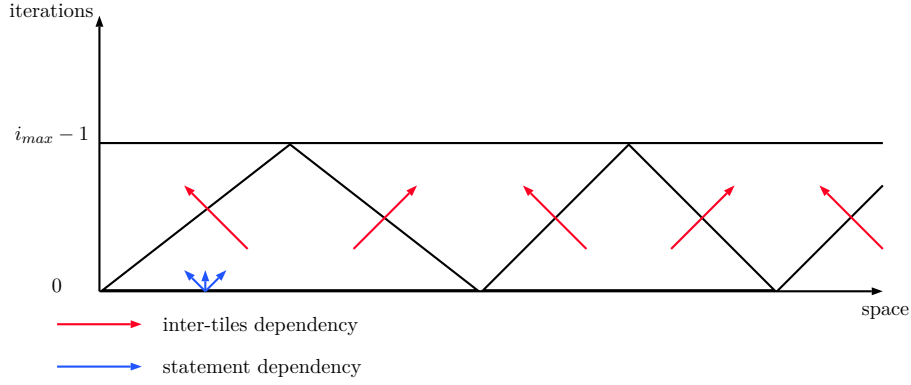


Figure 13: Half-diamonds tiling for low-iterations stencil computations

In the previous version, if the data size is significantly greater than the lower caches, the memory of the first lower tiles will be flushed out the caches before being reused by the upper tiles. To solve this problem, we group loop parts together, adjusting their size proportionally to the number of physical cores of the platform. The goal is to take advantage of static scheduling, which distributes evenly chunks of the loop across processing units to specify exactly how much computation we want to distribute. The group size is calculated depending on the number of processors n_{cores} , the size of the L1 cache and the size of the tile in memory (the size of the array containing the temporary elements between lower and upper tiles) :

$$Group\ size = n_{cores} * (L1\ cache\ size) / (tile\ size\ in\ memory)$$

Results The tiling for low iteration numbers without grouping tiles has proven efficient, scaling well with the number of processors. But the expected gain with grouping tiles hasn't been observed. The version with grouped tiles performs similarly to the simpler version, with a constant overhead (see figure 14). The schedule on the implementation with a global barrier between lower and upper tiles is already taking advantage of reuse through a good mapping of tile on the threads.

We still need to determine why grouping the tile into a size fitting the fastest cache level doesn't improve the performance, since it should improve reuse. Varying the group size doesn't affect performance either, even if the tile scheduling on the different processor is as expected : in each group upper tiles are executed on the same core as the lower tiles at their right.

4.2 A task-based approach

Using a task based approach, we wanted to show that the scheduler can improve data reuse in caches as much as the static scheduling with grouping. We constitute tasks containing one lower tile and the upper tile immediately to its left. The choice of taking the left tile will reduce task stalling. Indeed, the data dependency will probably be satisfied by the left task when the current task will begin to compute its upper tile. In figure 15 we summarize the task graph construction: each task spawns its successor, then computes the lower tile and waits at the barrier that the left task has finished to compute the upper tile. The upper half-tiles at the borders of the domain are treated separately.

Data reuse can be exploited at the border between upper and lower tiles, by executing neighboring tasks on the same processor. But parallelism shouldn't suffer and the computation load has to be balanced between processing units.

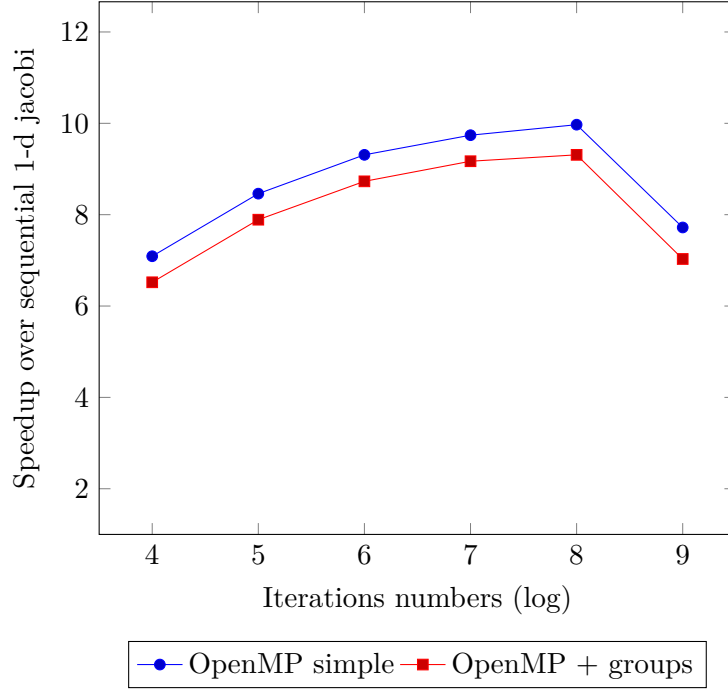


Figure 14: Speedup relative to sequential version, on Intel Xeon E5 with 12 cores

4.3 First results

This first version uses fixed-shape tiles : half diamonds with fixed hyperplanes, only the width varies depending on the number of iterations. It doesn't perform well compared to the statically scheduled OpenMP reference we use with the same tiling. On small iteration numbers the task version is slower than the sequential version. It can be explained by the granularity of the execution, the tasks are very small given the small iterations numbers. The time spent in the computation is mostly scheduling: the runtime manages many tasks, which terminate very rapidly.

The performance of the scheduler measured with this implementation with such very small tasks was still better than OpenMP4's using tasks pragmas.

4.4 Tile shape determination at runtime

The previous results showed that the runtime was underperforming because of the high granularity of the task graph. Since our tiling was designed for small iteration numbers, the tile sizes were relatively small in terms of memory occupation and computational weight.

With tiles of variable size shaped at runtime, depending on the number of iterations, we can increase the weight of this tasks. Similarly to the computation of group size in the previous section, the tile size depends on the size of the cache to ensure a good locality, and that sufficient time is spent in the task. We vary the orientation of the oblique hyperplanes with a slope between 1 (as in the previous version) and a value so that there is enough tiles to be distributed among the working threads. Tiling hyperplanes are now $h_1 = (-a, 1)$ and $h_2 = (a, 1)$ with $1/a = n_{cores} * (L1\ cache\ size) / (tile\ size\ in\ memory)$. This tiling is still legal and valid with $a \leq 1$, the scalar product between dependencies and hyperplanes normal vectors is positive or zero.

The only change using this tiling is the amount of computation in each tile. We can fear that using this skewed tiling space with large strides in the loop could lead to worse compiler automatic vectorization.

This technique combined with the task approach to this problem has given better results, in

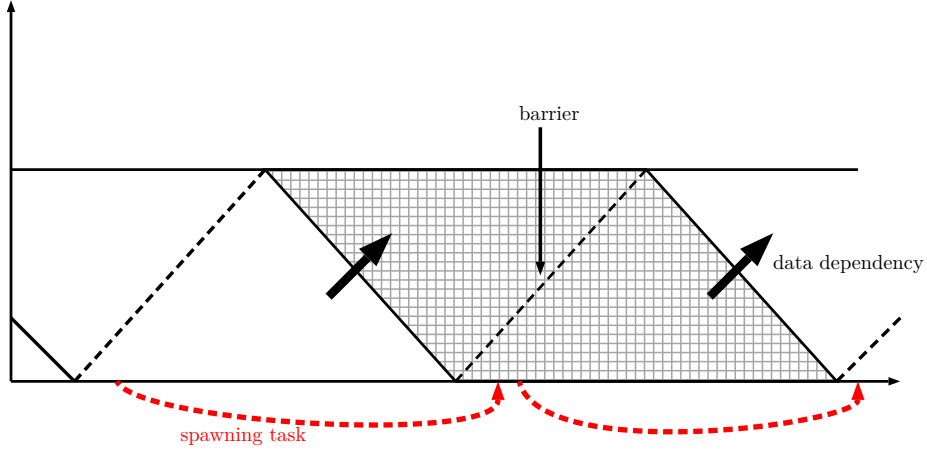


Figure 15: Task dependencies, synchronization points and creation

par with OpenMP on small problems but performing better with increasing sizes, showing speedups better than linearly-increasing with the number of cores. On the 12-core machine with a problem size of 32 MB, we reached speedups up to 20 times better than sequential, 3 times better than OpenMP. The results are shown in figure 16.

We have bad performance with small iteration sizes and small space sizes. This may be coming from the lack of parallelism with a too-small number of tasks because the slope of the sides of the tiles was not bounded : for small iterations, the tile is very wide and it could span a large part of the space. The final program at the moment I am

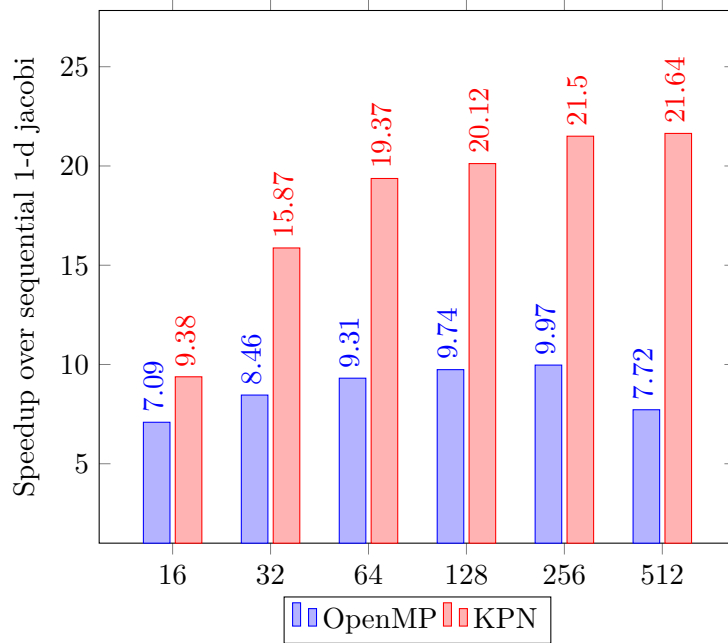


Figure 16: Speedup relative to sequential 1d-jacobi on 12 cores, with different iteration sizes on a 32 Mb problem

writing this report has still a few errors concerning the correctness of the output, but they do not affect performance in any way. They come from hand-tuning the loops and the tiling, a process that we seek to automate through polyhedral code generation tools and compilation. The main goals here were to show a well performing tiling for time-iterated stencils, and validate the approach of implementing the algorithm using lightweight task and dynamic scheduling instead of the classical approach with static scheduling through

loop-chunk distribution.

5 Conclusion and future work

During my internship, I have had the opportunity to apply and augment my knowledge in computer science. I have learned a lot about the subjects mentioned in this report, from the tilings and the polyhedral compilation problems, to the Kahn's parallel tasks network. During the first month of my internship, I mostly studied the topics I didn't master through a bibliographic research on both recent and older publications. The first application didn't provide satisfactory results but prepared me for the second subjects, especially on mastering profiling and debugging tools adapted with parallel applications. Then after a seminar from Uday Bondhugula, we discussed the possibility of taking advantage from more parallelism in one of the applications presented, Harris Corner Detection, without convincing results. Then my focus moved on to a more simple example, 1-d jacobi, using a partly new tiling strategy for small iterations and implementing an application for Likbkpn, the library developed at Parkas. At the time of writing this report, my internship is not yet finished. There is plenty applications to implement to show the strength of the library but since I have few time left, I will concentrate on improving the current example.

Optimization is a very experimental field, and I learned to use various tools to validate or refute the hypothesis made during the process.

6 References

References

- [1] V. Bandishti, I. Pananilath, and U. Bondhugula, “Tiling Stencil Computations to Maximize Parallelism,” in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC ’12. Los Alamitos, CA, USA: IEEE Computer Society Press, 2012, pp. 40:1–40:11. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2388996.2389051>
- [2] D. F. Bacon, S. L. Graham, and O. J. Sharp, “Compiler Transformations for High-performance Computing,” *ACM Comput. Surv.*, vol. 26, no. 4, pp. 345–420, Dec. 1994. [Online]. Available: <http://doi.acm.org/10.1145/197405.197406>
- [3] C. Bastoul, “Code Generation in the Polyhedral Model Is Easier Than You Think,” in *Proceedings of the 13th International Conference on Parallel Architectures and Compilation Techniques*, ser. PACT ’04. Washington, DC, USA: IEEE Computer Society, 2004, pp. 7–16. [Online]. Available: <http://dx.doi.org/10.1109/PACT.2004.11>
- [4] U. Bondhugula, M. Baskaran, S. Krishnamoorthy, J. Ramanujam, A. Rountev, and P. Sadayappan, “Automatic Transformations for Communication-Minimized Parallelization and Locality Optimization in the Polyhedral Model,” in *Compiler Construction*, ser. Lecture Notes in Computer Science, L. Hendren, Ed. Springer Berlin Heidelberg, 2008, no. 4959, pp. 132–146.
- [5] X. Zhou, J.-P. Giacalone, M. J. Garzaran, R. H. Kuhn, Y. Ni, and D. Padua, “Hierarchical Overlapped Tiling,” in *CGO ’12 Proceedings of the Tenth International Symposium on Code Generation and Optimization*, 2012, pp. 207–218.
- [6] T. Grosser, A. Cohen, P. H. J. Kelly, J. Ramanujam, P. Sadayappan, and S. Verdoolaege, “Split Tiling for GPUs: Automatic Parallelization Using Trapezoidal Tiles,” in *Proceedings of the 6th Workshop on General Purpose Processor Using Graphics Processing Units*, ser. GPGPU-6. New York, NY, USA: ACM, 2013, pp. 24–31. [Online]. Available: <http://doi.acm.org/10.1145/2458523.2458526>
- [7] K. Gilles, “The semantics of a simple language for parallel programming,” in *In Information Processing’74: Proceedings of the IFIP Congress*, 1974, pp. 471–475.
- [8] W. Thies, M. Karczmarek, and S. Amarasinghe, “Streamit: A language for streaming applications,” in *Compiler Construction*. Springer, 2002, pp. 179–196.
- [9] R. D. Blumofe and C. E. Leiserson, “Scheduling Multithreaded Computations by Work Stealing,” *J. ACM*, vol. 46, no. 5, pp. 720–748, Sep. 1999. [Online]. Available: <http://doi.acm.org/10.1145/324133.324234>
- [10] C. HARRIS, “A combined corner and edge detector,” in *Proc. 4th Alvey Vision Conf.*, 1988, pp. 147–151.
- [11] R. T. Mullapudi, V. Vasista, and U. Bondhugula, “PolyMage: Automatic Optimization for Image Processing Pipelines,” in *Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS ’15. New York, NY, USA: ACM, 2015, pp. 429–443. [Online]. Available: <http://doi.acm.org/10.1145/2694344.2694364>
- [12] A. Duran, X. Teruel, R. Ferrer, X. Martorell, and E. Ayguade, “Barcelona OpenMP Tasks Suite: A Set of Benchmarks Targeting the Exploitation of Task Parallelism in OpenMP,” in *International Conference on Parallel Processing, 2009. ICPP ’09*, Sep. 2009, pp. 124–131.