

FACULDADE DE ENGENHARIA DA UNIVERSIDADE DO PORTO

# Autotuning Parallel Application in Heterogeneous Systems

João Alberto Trigo de Bordalo Morais



**FEUP** FACULDADE DE ENGENHARIA  
UNIVERSIDADE DO PORTO

Mestrado Integrado em Engenharia Informática e Computação

Supervisor: Jorge Manuel Gomes Barbosa

February 10, 2017



# **Autotuning Parallel Application in Heterogeneous Systems**

**João Alberto Trigo de Bordalo Morais**

Mestrado Integrado em Engenharia Informática e Computação

February 10, 2017



# Abstract

Nowadays computational platforms have been evolving to the high computational power direction, however it requires a lot of energy to achieve such high performance with single but powerful processing unit. To manage this energy cost and keep with high performance, computers are built under the assumption of heterogeneous systems, in other words, computers that have different kind of processing units with different functions, such as CPU, GPU, Xeon Phi and FPGA. So, developers should take advantage of parallel activity and scheduling tasks by using the various parts of the heterogeneous systems.

Now the problem is how to efficiently achieve the highest performance possible when running software applications by taking the most advantage of such heterogeneous systems and keeping the energy cost at the minimum level without jeopardizing the application performance and its results. Overall, the problem consists in the coexistence work of multicore specs, its parallelism and its shared cache problems; CPU and GPU parallelism and scheduling tasks; performance; and energy costs.

For this problem's solution is expected to find/create an autotuner, or at least a concept proof, that can achieve the best performance in a software application by enhancing the application's code automatically in a level that takes the best benefit of the available hardware without elevated energy costs. To do so, after creating its code, the developer runs the autotuner and it will enhance, automatically, the code to get the best performance.

This kind of solution requires some validation process and metrics to make sure that it is doing its work and with proper results. To do so, the idea of the process' validation is going to be about comparing the behaviour of three different codes: a version of a serialized code; a version of the same code but with an expert manually paralleling it; and a version of the serialized code but automatically parallelized. The metrics that will be used to compare these three code versions are the following: processing power; execution time; number of memory accesses; and energy consuming.

With this solution, applications will achieve its highest performance possible in an automatic way and developers will have less burdened about creating parallel code, consequently, saving them time.



# Resumo

Atualmente as plataformas computacionais têm vindo a evoluir na direção do elevado poder computacional, no entanto, estas requerem uma quantidade enorme de energia para atingir elevado desempenho individualmente. De modo a gerir este custo energético e manter a elevada performance, os computadores são construídos sobre a assunção de sistemas heterogéneos, isto é, computadores compostos por diferentes tipos de unidades de processamentos com diferentes funcionalidades, como por exemplo, CPU, GPU, Xeon Phi e FPGA. É neste sentido que os programadores devem tirar proveito de atividade paralela e escalonamento de tarefas recorrendo às várias partes que compõem o sistema heterogéneo.

O problema incide sobre como atingir de forma eficiente o maior desempenho possível quando se corre uma aplicação de software, tirando o maior proveito dos sistemas heterogéneos e mantendo o nível de custo energético o mais baixo possível sem prejudicar o resultado e o desempenho da aplicação.

Para solucionar este problema é esperado encontrar/criar um autotuner, ou pelo menos uma prova de conceito, que consegue atingir o melhor desempenho numa aplicação de software, aprimorando automaticamente o código da aplicação a um nível que take o melhor proveito do hardware disponível sem custos elevados de energia. Para tal, após o código criado, o programador correrá o autotuner e este irá aprimorar, automaticamente, o código para atingir o melhor desempenho.

Este tipo de solução requer um processo de validação e métricas para assegurar que se está a fazer o trabalho corretamente e com resultados aceitáveis. Para tal, a ideia da validação do processo consiste em comparar o comportamento de três diferentes códigos: uma versão sequencial de um código; a versão deste mesmo código mas paralelizada por um perito; e a versão do código sequencial mas paralelizado automaticamente. As métricas que serão utilizadas para comprar estas três versões de código são as seguintes: poder de processamento; tempo de execução; número de acessos a memória; e custo energético.

Com esta solução, as aplicações conseguiram atingir o seu melhor desempenho possível de forma automática e sobrecarregando menos os programadores a criarem código paralelo o que, consequentemente, poupar-lhes-á tempo.





# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Motivation and Goal . . . . .	2
1.3	Structure of the Report . . . . .	2
<b>2</b>	<b>Achieving the Highest Processing Power</b>	<b>3</b>
2.1	Introduction . . . . .	3
2.2	Using Computers' Heterogeneous Components . . . . .	3
2.2.1	OpenCL . . . . .	4
2.2.2	StarPU . . . . .	4
2.2.3	Twin Peaks . . . . .	5
2.3	Using Code Parallelization . . . . .	5
2.3.1	OpenMP . . . . .	6
2.3.2	Kremlin . . . . .	6
2.3.3	Kismet . . . . .	6
2.4	Overview . . . . .	7
<b>3</b>	<b>Matrix Multiplication</b>	<b>9</b>
3.1	Introduction . . . . .	9
3.2	Using Computers' Heterogeneous Components . . . . .	9
3.2.1	OpenCL . . . . .	9
3.3	Overview . . . . .	9
<b>4</b>	<b>Methodology</b>	<b>11</b>
4.1	Introduction . . . . .	11
4.2	Research Method . . . . .	12
4.2.1	Deep learn on Kremlin's usage . . . . .	13
4.2.2	Kremlin's application in specific code samples . . . . .	14
4.2.3	Code parallelization with Kremlin's data . . . . .	14
4.2.4	Manually Code parallelization . . . . .	14
4.2.5	Results analysis . . . . .	15
4.3	Data collection from executed experiences . . . . .	15
4.4	Data analysis method . . . . .	15
4.5	Data validation . . . . .	15
4.6	Experimental environment setup . . . . .	16

## CONTENTS

<b>5</b>	<b>Resultls and Discussion</b>	<b>17</b>
5.1	Kremlins reports . . . . .	17
5.2	Original Vs Manual vs Kremlin . . . . .	17
<b>6</b>	<b>Conclusion</b>	<b>19</b>
	<b>References</b>	<b>21</b>
<b>7</b>	<b>Appendices</b>	<b>23</b>
7.1	Developed code . . . . .	23
7.1.1	Original Matrix Multiplication (Mult) . . . . .	23
7.1.2	Original Matrix Multiplication By line (MutLine) . . . . .	24
7.1.3	Manual Matrix Multiplication (Mult) . . . . .	25
7.1.4	Manual Matrix Multiplication By line (MultLine) . . . . .	26
7.1.5	Kremlin Matrix Multiplication (Mult) . . . . .	27
7.1.6	Kremlin Matrix Multiplication By line (MultLine) . . . . .	28
7.2	Kremlin's Reports . . . . .	29
7.2.1	Kremlin report for Matrix Multiplication, Mult version . . . . .	29
7.2.2	Kremlin report for Matrix Multiplication, MultLine version . . . . .	30

# List of Figures

2.1	The OpenCL platform model and the OpenCL memory model . . . . .	4
3.1	The OpenCL platform model and the OpenCL memory model . . . . .	10
4.1	Followed up methodology . . . . .	12



# Chapter 1

## 2 Introduction

### 4 1.1 Context

Previously, computer systems were built to maximize their processing power in compactness and  
6 individually because programs were developed with a sequential approach. With the advance  
in microchips' technology, computers increased their processing capacity per volume, however  
8 some issues arose, such as high energy cost, high temperature and low equipment durability. To  
solve these issues some measures needed to take place in order to make computers systems more  
10 reliable, durable, efficient, and powerful.

Recently, the computing industry has moved away from exponential scaling of clock frequency  
12 toward chip multiprocessors in order to better manage trade-offs among performance, energy effi-  
ciency, and reliability [Dat08]

14 Combining different computer processing components, such as CPU, GPU Xeon Phi and  
FPGA, in a single computer system removed some heavy burden in the main processing core,  
16 making the computer system with better performance and reliable. However some concerns arose:  
how to properly use these components without jeopardizing the computer system and application  
18 performance. Some processing components can handle specif jobs better then others and com-  
bined the computer can achieve a whole new performance level; for instance, the use of a GPU  
20 together with a CPU to accelerate deep learning algorithm, analytics, and engineering applica-  
tions [Nvi], however this kind of utility is not yet well optimized and its utility is only recently  
22 emerging.

- mencionar multiplicação de matrizes, -dois algoritmos ligeiramente diferentes para a multi-  
24 plicação

## 1.2 Motivation and Goal

My motivation for this thesis is to advance a little further on the field of the automatic code parallelization and replace the manual parallelization labor because it requires a lot of time and effort to achieve significant performance. 2 4

## 1.3 Structure of the Report

This report is divided in three more chapters. The next one is called *Achieving the Highest Processing Power*, and it is related to the state of the art of my thesis' scope. In this chapter there are three sections. The first section is related to the context of the state of the art in the field. The other two sections are two different but complementary approaches which help and describe the state of the art. 6 8 10

The third chapter addresses the problem involved in my thesis and how I propose to solve it, including the approach, the methodology and solution's validation. The last chapter includes final consideration related to the work developed so far in *Preparação da Dissertação* course, expected results with my proposed solution and work plan to develop the solution. In the end of this report there are the references used to develop this report. 12 14

## Chapter 2

# Achieving the Highest Processing Power

The introduction describes a brief overview about each content of each chapter that this report is made up with. This chapter will focus on the state of the art in how to achieve the highest processing power. Related work and already known technologies are the main point in this chapter.

### 2.1 Introduction

Following the context introduced in the previous chapter, the idea of having different processing components in a computer system doesn't improve the applications performance on its own. This is where the developers' work is crucial to take advantages of such different systems. The developers' work is to schedule the application's tasks to the different components so that these components can work simultaneously, avoiding overheads caused by their parallel activity, accessing memory at the wrong moment, memory conflicts, task dependency, wrong application's results compared with the sequential application. [LR]

As mentioned previously, trying to create parallelized code can arise many problems and must be handled so the applications don't lose their functionalities. In order to do so, it requires a lot of time and effort to make it correctly parallelized. So trying to make code parallelization automatic is the next step in the direction of taking the most advantage of heterogeneous systems which, consequently, improves applications performance.

This chapter is divided in two parts: one part will focus in the system's heterogeneity, how they can be used in favor of enhancing performance; and the main point of the other part is taking advantage of parallel activity by transforming sequential code into parallelized code.

### 2.2 Using Computers' Heterogeneous Components

Technologies and frameworks in this field have been developed in order to manipulate and control efficiently the different processing components. The main goal of this technologies is to optimize

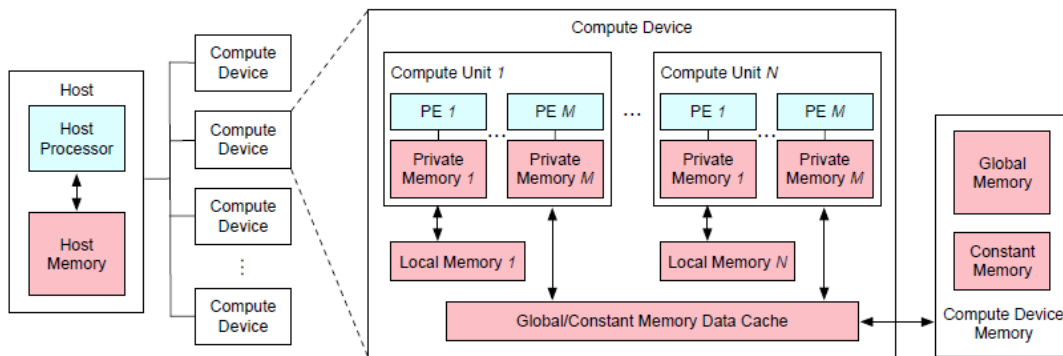


Figure 2.1: The OpenCL platform model and the OpenCL memory model

application parallelization; application memory management; application workload; application scheduling queue and application kernel dimension.

An interesting fact is that the following software/frameworks that will be present are built, as its bases, under OpenCL programing language due to the fact that this language's use is targeted to heterogeneous parallel programing with CPUs and GPUs.

## 2.2.1 OpenCL

OpenCL is a programming language for heterogeneous parallel programing targeted to CPUs, GPUs and other processors [She]. In a small brief, this language is designed to take advantage of different types of processors and facilitates heterogeneous computing integration in applications' code. The user programs in a virtual platform and the source code that has been developed there is compatible for any system that supports OpenCL. Additionally, OpenCL allows users to control the applications' tunning parallelism through its hardware abstraction. In figure 3.1 there is an idea of the OpenCL plataform model and memory model for a better understanding of this hardware abstractions that was previously mentioned.

The OpenCL's programs has two parts: the compute kernels that are executed depending the number of processing devices; and the program that will be run. The program creates a set of commands and puts them in a queue for each device, additionally, to manage the execution of each kernels, additional commands are queued in the different kernels. When the computation is finished, the result data, from the previous kernels activity, return back to the original program.

## 2.2.2 StarPU

StartPU is a software tool with the purpose for programmers to use the computing power available in CPUs and GPUs, wihtout needing to care about if their programs are adapted to a specific machine and its processing components. [Sta] In fact, StartPU is a run-timpe support library that provides scheduling applications-provided tasks on heterogeneous environments, such as CPUs



and GPUs. Additionally, it comes with programming language support, for the programming C language extensions and for OpenCL.

Programs submit computational tasks, with CPU and/or GPU implementations, and StarPU schedules these tasks and associated data transfers on available CPUs and GPUs. The data that a task manipulates are automatically transferred among accelerators and the main memory, so that programmers are freed from the scheduling issues and technical details associated with these transfers.

StarPU takes particular care of scheduling tasks efficiently, using well-known algorithms from the literature (Task Scheduling Policy). In addition, it allows scheduling experts, such as compiler or computational library developers, to implement custom scheduling policies in a portable fashion (Defining A New Scheduling Policy).

### 2.2.3 Twin Peaks

"Software platform that enables applications originally targeted for GPUs to be executed efficiently on multicore CPUs", mentioned by Jayanth Gummaraju, Laurent Morichetti, Michael Houston, Ben Sander, Benedict R. Gaster, Bixia Zheng, in the paper *Twin peaks: a software platform for heterogeneous computing on general-purpose and graphics processors* [Gum10]. This is a small definition of the Twin Peaks' job. The aim of this software is, firstly, to program applications using an API written in OpenCL; secondly, to compile the applications code to, for instance, add syntactic and semantic checks to make sure that the kernels meet the OpenCL requirements; and execute applications in the heterogeneous environment using CPUs and GPUs.

## 2.3 Using Code Parallelization

Great advances have been made in the code parallelization. However, currently this kind of practice (the code parallelization) mostly is done by programmers and it requires a lot of effort, time and knowledge. It requires knowledge in the best practices related to what should and can't be parallelized, good knowledge on the code: its functionalities and its correct outputs because without these knowledges the chances to parallelize code correctly would be low since it is important to know if, firstly, is possible to parallelize and if, secondly, the parallelization doesn't jeopardize the programs results, outcomes and performance; to sum up, it requires time and effort to get a deep understanding of the code and to try if the code is correctly parallelized. [Jeo]

Since this practice is very costly, although grants great results at performance levels, this field has been developing ways to have results less costly, mostly in effort and time-consuming. These developments created tools to help programmers develop parallelized code, using OpenMP directives, or software tools which recommend possible parallelized regions and its theoretical speed up gain, with Kremlin, or even a way to estimate how much can a program be parallelized, with Kismet software. [GJ]

The following software tools that will be presented have, as its base support, OpenMP directives to help in parallelizing code, or at least, to measure performance.

### 2.3.1 OpenMP

OpenMP was designed to be a flexible standard, easily implemented across different platforms. the main objectives are: control structure, the data environment, synchronization, and the runtime library.

In terms of how it really does its job, OpenMP was designed to exploit certain characteristics of shared-memory architectures. The ability to directly access memory throughout the system, combined with fast shared memory locks, makes shared-memory architectures best suited for supporting OpenMP. In practice, OpenMP is a set of compiler directives and callable runtime library routines that extend Fortran (and separately, C and C++) to express shared-memory parallelism. [Nc98]. To be more precise, OpenMP provides standard environment variables to accompany the runtime library functions where it makes sense and to simplify the start-up scripts for portable applications. This helps application developers who, in addition to creating portable applications, need a portable runtime environment. OpenMP has been designed to be extensible and evolve with user requirements. The OpenMP Architecture Review Board was created to provide long-term support and enhancements of the OpenMP specifications.

### 2.3.2 Kremlin

The true purpose of Kremlin lies in asking the following question: "What parts of this program should I spent time parallelizing?" [Par]. So, in overall, Kremlin profiles a serial program and tells the programmer not only what regions should be parallelized, but also the order in which they should be parallelized to maximize the return on their effort. Giving a non parallelized code, Kremlin guides the programmer how to achieve better performance in its program though parallelization by presenting a list of code regions that could be parallelized. this list contains a plan that will minimize the number of regions that must be parallelized to maximize the programs performance, though parallelization.

At the core of the Kremlin system is a heavyweight analysis of a sequential program's execution that is used to create predictions about the structure of a hypothetical, optimized parallel implementation of the program. These predictions incorporate both optimism and pessimism to create results that are surprisingly accurate. [GJLT11]

Overall, Kremlin is an automatic tool that, given a serial version of a program, will make recommendations to the user as to what regions (e.g. loops or functions) of the program to attack first. [GJL<sup>+</sup>12]

### 2.3.3 Kismet

Opposed to Kremlin, Kismet helps mitigate the risk of parallel software engineering by answering the question, "What is the best performance I can expect if I parallelize this program?" [Par]. Kismet profiles serial programs and reports the upper bound on parallel speedup based on the program's inherent parallelism and the system it will be running on.

Kismet performs dynamic program analysis on an unmodified serial version of a program to determine the amount of parallelism available in each region(e.g. loop and function) of the program. Kismet then incorporates system constraints to calculate an approximate upper bound on the program's attainable parallel speedup. [Tay]

In order to estimate the parallel performance of a serial program, Kismet uses a parallel execution time model. Kismet's parallel execution time model is based on the major components that affect parallel performance, including the amount of parallelism available, the serial execution time of the program, parallelization platform overheads, synchronization and memory system effects which contribute in some cases to super-linear speedups.

## 2.4 Overview

As mentioned before, the previously presented software tools, for both cases (using computers' heterogeneous components and using code parallelization) have their base support even being a programming language, for OpenCL, or a set of compile directives, for OpenMP. Those software tools have improved applications performance somehow, which is already good. However, looking as a software that can do everything on its own, with the minimum programmer's input, in other words, that can do things almost automatically, none of them can make it. The only software tool that is close to that automation is Kremlin because it gives what a developer should do in their code in order to increase its efficiency and performance.

Both approaches, using computers' heterogeneous components and using code parallelization, have the role to answer the state of the art premise: "achieving the highest processing power".

## Achieving the Highest Processing Power

## Chapter 3

# 2 Matrix Multiplication

4 The introduction describes a brief overview about each content of each chapter that this report  
is made up with. This chapter will focus on the state of the art in how to achieve the highest  
6 processing power. Related work and already known technologies are the main point in this chapter.

### 3.1 Introduction

### 8 3.2 Using Computers' Heterogeneous Components

#### 3.2.1 OpenCL

### 10 3.3 Overview

As mentioned before, the previously presented software tools, for both cases (using computers'  
12 heterogeneous components and using code parallelization) have their base support even being a  
programming language, for OpenCL, or a set of compile directives, for OpenMP. Those software  
14 tools have improved applications performance somehow, which is already good. However, looking  
as a software that can do everything on its own, with the minimum programmer's input, in other  
16 words, that can do things almost automatically, none of them can make it. The only software tool  
that is close to that automation is Kremlin because it gives what a developer should do in their  
18 code in order to increase its efficiency and performance.

Both approaches, using computers' heterogeneous components and using code parallelization,  
20 have the role to answer the state of the art premise: "achieving the highest processing power".

## Matrix Multiplication

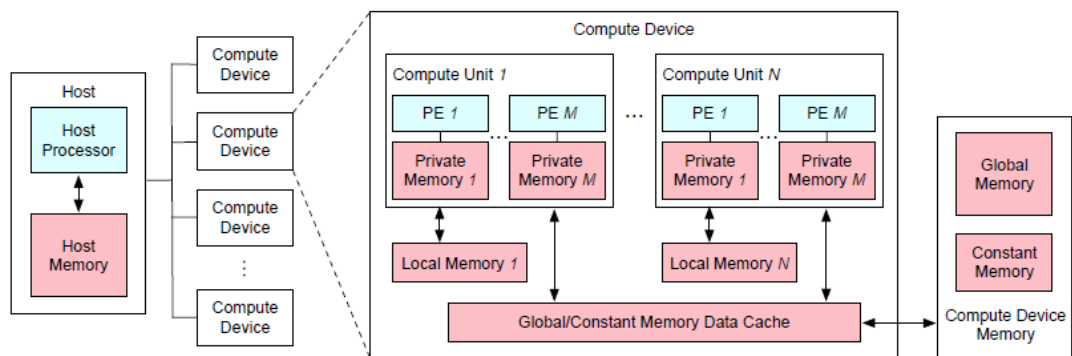


Figure 3.1: The OpenCL platform model and the OpenCL memory model

# Chapter 4

## 2 Methodology

### 4 4.1 Introduction

According to the state of the art presented in chapter two, there are many means to, in some kind of automatic way, improve an applications performance. During my research, my focus was to find ways to automatically enhance the execution time in applications and programs. For this propose, Kremlin had a crucial impact in other to understand the viability of automatically parallelise code.

To study the utility and impact of automatic tools, the matrix multiplication algorithm will be used as a reference to make the performance comparison between original algorithm, an expert manually parallelising the original algorithm and using the Kremlin's indication to parallelise the original algorithm.

To increase the credibility of this experiment, two similar algorithms for the matrix multiplication were used. As mentioned and explained in the chapter two, there is the traditional way of multiplying square matrices, naming as a quick reference *Mult* algorithm, see in the appendix's list 7.1 this algorithm implementation, written in C++ programming language; and the optimized algorithm that multiplies each element from the first matrix with the correspondent line of this matrix element but for the second matrix, naming this algorithm as *MultLine*, see in the appendix's list 7.2 this algorithm implementation, written in C++ programming language. These algorithms differs from one another in the variables preparation and the order of the loops, which differs in the memory access. The *MultLine* algorithm is an optimized version for matrix multiplication because it takes advantages of what is preloaded in cache and starts pre-calculating the intermediate values that will lead to the final and correct result of the multiplication, which means that the computer won't need to load unnecessary values to cache memory and/or will need afterwards.

Several experiments were conducted to understand the influence of Kremlin's indications versus code being manually parallelized by an expert. The data's length, in this case, the matrix size; the number of threads used and if the code was parallelized were the used metrics to evaluate the results, based on a comparison of the execution time, changing these variables.

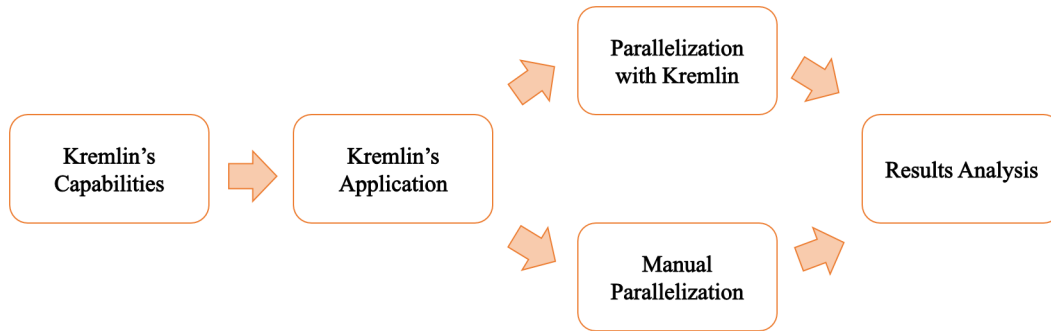


Figure 4.1: Followed up methodology

In this chapter it is explained the methodology and the steps followed that guided to report in the Results and Discussion chapter the results and conclusions obtained from the obtained outcomes coming from the experiences. This chapter also includes detailed information of the acquired data from the conducted experiences, as in, how it is obtained and its meaning; also includes the methods that were used to analyse the obtained data and the reason behind those methods; how the data was validated in order to verify its correctness, accuracy and reliability; and, in the end, it includes the setup and material used to conduct such experiments.

## 4.2 Research Method

In the figure 4.1 is outlined the steps that were followed to study the impact of the code being automatic parallelised. This methodology has five states. Firstly and using simple applications, an evaluation was made to Kremlin in order to understand how to use this software tool and evaluate the results that Kremlin can achieve, for instance, if it has similar results comparing with an expert parallelizing manually the same code. After this, Kremlin will be applied to a set of codes with specific characteristics.

Before applying Kremlin, I manually parallelised the same sample of code in order to evaluate the results and, afterwards, compare with the Kremlin output. Since these states (the experiences with Kremlin and the manual code parallelization) required several attempts there were transitions between manual and Kremlin states.

Finally, in the last state, after several attempts and tuning exercises applied to both code cases (manual and kremlin), data was collected from this experience to evaluate and validate its correctness in order to conclude how helpfull can automatic parallelization can be.

To sum up, this methodology as three main stages: learn and evaluate Kremlin's uses and results; finding the tuning parameter through several attempts using Kremlin's outputs and manually parallelize the application's code; and, in the end, compare and analyze the results in every



attempt to take conclusions;

## 2 4.2.1 Deep learn on Kremlin's usage

Firstly, and according to all the presented tools/frameworks mentioned in the second chapter, *Achieving the Highest Processing Power*, in the *Using Code Parallelization* section 2.3, Kremlin was chosen because it presented the best results, easy usage and accessibility comparing to the other presented ones regarding the way the tool/framework could automatically parallelise code 2.3.2.

Kremlin is a tool that indicates, for a serial program, which block can be parallelised and teorical calculated values, such as, overall speedup; self parallelism for each block; the ideal time reduced for each block, in percentage; the actual time reduced for each block, in percentage; and the block coverage considering the whole program, in percentage. The way this tool was used is as it follows: first, an object file, \*.o extension, is required from the compilation of a serial code. Afterwards, it is time to use the Kremlin's compiler with the generated object file so that it can profile the application. In order to do so, Kremlin's compiler runs the program as it is supposed to work. Now that the profiling is done, Kremlin generates the indications that should be followed to parallelise de provided serial code. It also includes the blocks that can be parallelised and the impact of this theoretical parallelization with the calculations done during the profiling. Since this parallelization report is done, the program has interpret it, confront with the code an apply it.

The Kremlin's usage seems easy, linear and fast forward, however it has some limitations that I experienced during the learning of Kremlin's capabilities: Kremlin's requires a specific environment mentioned in the Kremlin's repository [Kre]. It requires several software, libraries, compilers installations and a modern Unix operative system as its base, such as MAC OS, RHEL 7 or other Linux distribution compatible with the software specification required. Additionally, when installing the Kremlin's tool, some minor fixed are required in order to successfully install.

From the experiences that I have been through, Kremlin has another limitation: it can't compile and profile all kind of programs: it can only profile programs that use C/C++ as its programming language; programs that take advantage of data structures from the *Standard Library*, such as, stack, list, priority queue, queue, list, hash table, map, multimap, heap, etc., since it doesn't recognize these structures; another Kremlin's limitations is its capability of compiling programs that have a deep function call level greater that seven. By deep function call level I mean the depth a function has starting from the *main* function until it is called, like a tree function call tree. For instance: in a program there is the *main()* function, a first level, that calls a *foo1()* function, and this function calls a *foo2()*, that this calls a *foo3()* function, and so on. In this case, the depth of *foo3()* function is four. Another small issue that kremlin's tool has is the definition of the iterator variable used in the *for*'s loops must be defined outside of the loop, as it is in C programming language.

### 4.2.2 Kremlin's application in specific code samples

After all the experiences made in the previous state and as mentioned in the introduction of this chapter, the matrix multiplication algorithm was used to see the potentialities of Kremlin's compiler to profile and indicate the regions that can be parallelised. So, Kremlin was used in two similar, relatively in the code structure, matrix multiplication codes. The reason behind the choice was because these two versions of the algorithm are really close to one another, which means that the testing environment is similar to one another, consequently, the results should be similar.

### 4.2.3 Code parallelization with Kremlin's data

Kremlin's tool just points the regions/blocks where the program can be parallelised. In both code samples there are various numbers of inner *for* loops, for the *Mult* code there are three inner *for* loops, which one of them has a degree of three and the rest a degree of two; and for the *MultLine* code there are four inner *for* loops, which one of them has a degree of three and the rest a degree of two as well. At this time, after reading the report provided by Kremlin's tool, the developer must locate the loops, apply which loop should be parallelised, if it should be, and in case of inner *for* loops, what loop should be parallelised using the OpenMP *pragma* directives.

In my case, I followed all the instructions provided by Kremlin, located all the *for* loops blocks indicated by kremlin's tool and applied the OpenMP *pragma* directives.

Following the two reports, 7.7 7.8, and looking at the code's structure for both codes, it can be divided in 2 bigger parts: the *for* loops used for matrices initialization and the *for* loop for the matrix multiplication. With this information, code understanding and using an expert knowledge, the code parallelization was done.

### 4.2.4 Manually Code parallelization

In order to not be manipulated by the Kremlin's indications, the both codes were previously manually parallelised, this way it was guaranteed that the expert parallelization wasn't bias nor influenced.

For this parallelization, as mentioned before, it requires knowledge in, firstly, matrix multiplication algorithm; code understanding; best practice in what can and can't be parallelised, taking into account the overhead that could occur; and understand the thread behaviour in order to make it do the proper job without jeopardizing the programs outputs and/or possible caused overhead.

Analysing the code, only the *for* loop for the matrix multiplication was parallelised and applied the OpenMP *pragma* directives applied to the innermost *for* loop. In this case, each code as a slight difference because for the *Mult* code each value of the result matrix but be calculated individually, so each thread is responsible for it and must treat that value as a private variable that isn't shared by the other threads. In the opposite, and since the *MultLine* code calculates the values by adding the multiplication to the respective matrix's cell, each thread don't need to have their own private variable.

The bigger part of the code responsible for the matrices initialization wasn't parallelised, unlike in kremlin's case, because the gain would be noticeable on a large matrix size or it even could cause thread trampling, which could lead an overhead increase.

### 4.2.5 Results analysis

To obtain the final execution time of each implementation (Original Matrix Multiplication 7.1, Original Matrix Multiplication by line 7.2, Manual Matrix Multiplication by an expert 7.3, Manual Matrix Multiplication by line by an expert 7.4, Kremlin Matrix Multiplication 7.7 and Kremlin Matrix Multiplication by line 7.6), these six implementations suffered many modifications and tweaks since this process is a try-error until it is found the believed best parallelization. It is hardly possible to parallelise a whole program at the first try.

After compiling all these implementations and registering all the execution time for different matrix sizes and number of threads, in this case not applied to the original codes, this data was organized so it could be used to compare results and conclude about the performed experiences.

## 4.3 Data collection from executed experiences

-from kremlins' usage -from kremlins output when using matrix multiplication -times from original code -Times from kremlins matrix multiplication -times from manually code parallelization -variables: Parallelized vs not parallelized, matrix size, n° of threads)

In order to achieve such performance in the conditions mentioned in the problem's section, I purpose an autotuner or a concept proof that will enhance the application. For that, this autotuner will receive the program's original source code and through parallel optimization a new code will emerge. This code's modification will increase applications performance without jeopardizing the application's outcome.

## 4.4 Data analysis method

-3 versions of the same code (original, manual code parallelization, kremlins parallelization) - Comparing execution time for each case -Table and graphic analysis -Impact for each concrete case (increase or decrease in execution time) To create this autotuner, the first step is to identify what parameters exist to tune and what impact they have in the application. For this matter, with the help of Kremlin and manual expert parallelization applications will increase its performance and parameters will be found. As Kremlin detect possible parallelized code and, additionally, measures the applications speedup with such modifications, and with manual expert parallelization, there will be a confrontation with these results and find what is better between these two scenarios. With this confrontation, and with different use cases, the tuning parameters will be found and the autotuner will be made.

## 4.5 Data validation

Comparing with teoric and expected results vs experimental results To validate the whole methodology process, not only the autotuner itself but to compare the results between Kremlin outcomes and the code manually being parallelized, for evaluation metrics will be used: system energy consumptions when running the application; applications execution time; number of memory accesses and cache misses on the application; and the processing power measured by the number of instructions per secs. These evaluation metrics will be compared in three different cases: the original sequential code; manually paralleled code; and "automatic" paralleled code. In this last case it can be with Kremlin or with autotuner, depending in which of the methodology's state I am currently in.

To measure the above mentioned evaluation metrics, some libraries will be used: to measure energy consumptions RAPL will be used [MMH]; to measure application execution time OpemMP will be used [Nc98]; to measure number of memory accesses, cache misses and processing power PAPI [MMH] will be used.

Finally, two use cases will be used to validate this solution: a biopharmaceutical HPC application for accelerating drug discovery; and a self-adaptive navigation system to be used in smart cities. These use cases will be the applications that, the autotuner will try to achieve the highest processing power without jeopardizing the applications' outcome and having a low energy consumption.

## 4.6 Experimental environment setup

-Kremlins setup -matrix mul setup

## **Chapter 5**

# **<sup>2</sup> Results and Discussion**

### **5.1 Kremlins reports**

### **<sup>2</sup> 5.2 Original Vs Manual vs Kremlin**

## Results and Discussion

## **Chapter 6**

## **4 Conclusion**

## Conclusion



## References

- [Dat08] Stencil computation optimization and auto-tuning on state-of-the-art multicore architectures. *2008 SC - International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2008*, 2008. Cited on page 1.
- [GJ] Saturnino Garcia Jr. *A Practical Oracle for Sequential Code Parallelization*. Cited on page 5.
- [GJL<sup>+</sup>12] Saturnino Garcia, Donghwan Jeon, Christopher Louie, Michael Bedford, and San Diego. the Kremlin Oracle for the Kremlin Open - Source Tool Helps Programmers By Automatically Identi-. pages 42–53, 2012. Cited on page 6.
- [GJLT11] Saturnino Garcia, Donghwan Jeon, Chris Louie, and Michael Bedford Taylor. Kremlin : Rethinking and Rebooting gprof for the Multicore Age. 2011. Cited on page 6.
- [Gum10] Twin peaks: a software platform for heterogeneous computing on general-purpose and graphics processors. *Proceedings of the 19th international conference on Parallel architectures and compilation techniques - PACT '10*, page 205, 2010. Cited on page 5.
- [Jeo] Donghwan Jeon. *Parallel Speedup Estimates for Serial Programs*. Cited on page 5.
- [Kre] kremlin: like gprof but for parallelization. <https://bitbucket.org/elsaturnino/kremlin>. Accessed: 2017-04-17. Cited on page 13.
- [LR] Changmin Lee and Won W Ro. for Parallel Processing on CPU / GPU Hybrids. Cited on page 3.
- [MMH] H Marcus, V Marcus, and H Hermann. Measuring Energy Consumption for Short Code Paths Using RAPL. Cited on page 16.
- [Nc98] S Ilicon G Raphics I Nc. OpenMP : An Industry-. pages 46–55, 1998. Cited on pages 6 and 16.
- [Nvi] What is gpu-accelerated computing? <http://www.nvidia.com/object/what-is-gpu-computing.html>. Accessed: 2017-02-09. Cited on page 1.
- [Par] Saturnino (sat) garcia. <http://home.sandiego.edu/~sat/>. Accessed: 2017-02-12. Cited on page 6.
- [She] Jie Shen. *Efficient High Performance Computing on Heterogeneous Platforms*. Cited on page 4.
- [Sta] *StarPU Handbook*. Cited on page 4.
- [Tay] Michael Bedford Taylor. Kismet : Parallel Speedup Estimates for Serial Programs. pages 519–536. Cited on page 7.

## REFERENCES

## Chapter 7

# Appendices

## 7.1 Developed code

### 7.1.1 Original Matrix Multiplication (Mult)

Listing 7.1: Matrix Multiplication original algorithm, written in C++

```
10
1  double OnMult(int m_ar, int m_br)
12 {
13     double Time1, Time2;
14     double temp;
15     int i, j, k;
16     double *pha, *phb, *phc;
17
18     //Matrixes Memory allocation
19     pha = (double *)malloc((m_ar * m_ar) * sizeof(double));
20     phb = (double *)malloc((m_ar * m_ar) * sizeof(double));
21     phc = (double *)malloc((m_ar * m_ar) * sizeof(double));
22
23     //Starting counting time
24     Time1 = omp_get_wtime();
25
26     //Loading matrix values
27     for(i=0; i<m_ar; i++)
28         for(j=0; j<m_ar; j++)
29             pha[i*m_ar + j] = (double)1.0;
30
31     for(i=0; i<m_br; i++)
32         for(j=0; j<m_br; j++)
33             phb[i*m_br + j] = (double)(i+1);
34
35     //Matrix Multiplication
36     for(i=0; i<m_ar; i++)
37     {
38         for(j=0; j<m_br; j++)
39         {
40             temp = 0;
41             for(k=0; k<m_ar; k++)
42             {
43                 temp += pha[i*m_ar+k] * phb[k*m_br+j];
44             }
45         }
46     }
47 }
```

## Appendices

```
33     }
34     phc[i*m_ar+j]=temp;
35 }
36 }
37
38 //Stopping time
39 Time2 = omp_get_wtime();
40
41 //Freeing memory used for matrixes
42 free(pha);
43 free(phb);
44 free(phc);
45
46 return Time2 - Time1;
47 }
```

### 7.1.2 Original Matrix Multiplication By line (MutLine)

Listing 7.2: Matrix Multiplication by line original algorithm, written in C++

```
1 double OnMultLine(int m_ar, int m_br)
2 {
3     double Time1, Time2;
4     double temp;
5     int i, j, k;
6     double *pha, *phb, *phc;
7
8     //Matrixes Memory allocation
9     pha = (double *)malloc((m_ar * m_ar) * sizeof(double));
10    phb = (double *)malloc((m_ar * m_ar) * sizeof(double));
11    phc = (double *)malloc((m_ar * m_ar) * sizeof(double));
12
13    //Starting counting time
14    Time1 = omp_get_wtime();
15
16    //Loading matrix values
17    for(i=0; i<m_ar; i++)
18        for(j=0; j<m_ar; j++)
19            pha[i*m_ar + j] = (double)1.0;
20
21    for(i=0; i<m_br; i++)
22        for(j=0; j<m_br; j++)
23            phb[i*m_br + j] = (double)(i+1);
24
25    for(i=0; i<m_ar; i++)
26        for(j=0; j<m_ar; j++)
27            phc[i*m_ar + j] = (double)0.0;
28
29
30    //Matrix Multiplication
31    for(i=0; i<m_ar; i++)
32    {
33        for( k=0; k<m_ar; k++)
34        {
35            for( j=0; j<m_br; j++)
```

```

35         {
36             phc[i*m_ar+j] += pha[i*m_ar+k] * phb[k*m_br+j];
37         }
38     }
39 }
40
41
42 //Stopping time
43 Time2 = omp_get_wtime();
44
45 //Freeing memory used for matrixes
46 free(pha);
47 free(phb);
48 free(phc);
49
50 return Time2 - Time1;;
51 }

```

## 7.1.3 Manual Matrix Multiplication (Mult)

Listing 7.3: Matrix Multiplication manually parallelised using OpenMP library, written in C++

```

26 double OnMultThreading(int m_ar, int m_br, int x)
27 {
28     double Time1, Time2;
29     double temp;
30     int i, j, k;
31     double *pha, *phb, *phc;
32
33     //Matrixes Memory allocation
34     pha = (double *)malloc((m_ar * m_ar) * sizeof(double));
35     phb = (double *)malloc((m_ar * m_ar) * sizeof(double));
36     phc = (double *)malloc((m_ar * m_ar) * sizeof(double));
37
38     //Starting counting time
39     Time1 = omp_get_wtime();
40
41     //Loading matrix values
42     for(i=0; i<m_ar; i++)
43         for(j=0; j<m_ar; j++)
44             pha[i*m_ar + j] = (double)1,0;
45
46     for(i=0; i<m_br; i++)
47         for(j=0; j<m_br; j++)
48             phb[i*m_br + j] = (double)(i+1);
49
50     //Matrix Multiplication
51     for(i=0; i<m_ar; i++)
52     {
53         for( j=0; j<m_br; j++)
54         {
55             temp = 0;
56             #pragma omp parallel for reduction(+:temp) num_threads (x)
57             for( k=0; k<m_ar; k++)
58             {

```

2

4

## Appendices

```

33         temp += pha[i*m_ar+k] * phb[k*m_br+j];
34     }
35     phc[i*m_ar+j]=temp;
36 }
37 }
38
39 //Stoping time
40 Time2 = omp_get_wtime();
41
42 //Freeing memory used for matrixes
43 free(pha);
44 free(phb);
45 free(phc);
46
47 return Time2 - Time1;
48 }

```

### 7.1.4 Manual Matrix Multiplication By line (MultLine)

Listing 7.4: Matrix Multiplication by line manually parallelised using OpenMP library, written in C++

```

1 double OnMultLineThreading(int m_ar, int m_br,int x)
2 {
3     double Time1, Time2;
4     int i, j, k;
5     double *pha, *phb, *phc;
6
7     //Matrixes Memory allocation
8     pha = (double *)malloc((m_ar * m_ar) * sizeof(double));
9     phb = (double *)malloc((m_ar * m_ar) * sizeof(double));
10    phc = (double *)malloc((m_ar * m_ar) * sizeof(double));
11
12    //Starting counting time
13    Time1 = omp_get_wtime();
14
15    //Loading matrix values
16    for(i=0; i<m_ar; i++)
17        for(j=0; j<m_ar; j++)
18            pha[i*m_ar + j] = (double)1,0;
19
20    for(i=0; i<m_br; i++)
21        for(j=0; j<m_br; j++)
22            phb[i*m_br + j] = (double)(i+1);
23
24    for(i=0; i<m_ar; i++)
25        for(j=0; j<m_ar; j++)
26            phc[i*m_ar + j] = (double)0,0;
27
28
29    //Matrix Multiplication
30    for(i=0; i<m_ar; i++)
31        {
32            for( k=0; k<m_ar; k++)

```

## Appendices

```
32     {
33         #pragma omp parallel for num_threads (x)
34         for( j=0; j<m_br; j++)
35         {
36             phc[i*m_ar+j] += pha[i*m_ar+k] * phb[k*m_br+j];
37         }
38     }
39 }
40
41
42 //Stopping time
43 Time2 = omp_get_wtime();
44
45 //Freeing memory used for matrixes
46 free(pha);
47 free(phb);
48 free(phc);
49
50 return Time2 - Time1;
51 }
```

### 7.1.5 Kremlin Matrix Multiplication (Mult)

Listing 7.5: Matrix Multiplication with Kremlin's indications for parallelization, written in C++

```
28
1  double OnMultKremlin(int m_ar, int m_br, int x)
2  {
3      double Time1, Time2;
4      double temp;
5      int i, j, k;
6      double *pha, *phb, *phc;
7
8      //Matrixes Memory allocation
9      pha = (double *)malloc((m_ar * m_ar) * sizeof(double));
10     phb = (double *)malloc((m_ar * m_ar) * sizeof(double));
11     phc = (double *)malloc((m_ar * m_ar) * sizeof(double));
12
13     //Starting counting time
14     Time1 = omp_get_wtime();
15
16     //Loading matrix values
17     for(i=0; i<m_ar; i++)
18         #pragma omp parallel for num_threads (x)
19         for(j=0; j<m_ar; j++)
20             pha[i*m_ar + j] = (double)1,0;
21
22     for(i=0; i<m_br; i++)
23         #pragma omp parallel for num_threads (x)
24         for(j=0; j<m_br; j++)
25             phb[i*m_br + j] = (double)(i+1);
26
27
28     //Matrix Multiplication
29     for(i=0; i<m_ar; i++)
```

2

4

## Appendices

```

30 {   for( j=0; j<m_br; j++)
31     {   temp = 0;
32         #pragma omp parallel for reduction(+:temp) num_threads (x)
33         for( k=0; k<m_ar; k++)
34             {
35                 temp += pha[i*m_ar+k] * phb[k*m_br+j];
36             }
37         phc[i*m_ar+j]=temp;
38     }
39 }
40
41 //Stopping time
42 Time2 = omp_get_wtime();
43
44 //Freeing memory used for matrixes
45 free(pha);
46 free(phb);
47 free(phc);
48
49 return Time2 - Time1;
50 }

```

### 7.1.6 Kremlin Matrix Multiplication By line (MultLine)

Listing 7.6: Matrix Multiplication by line with Kremlin's indications for parallelization, written in C++

```

1 double OnMultLineKremlin(int m_ar, int m_br, int x)
2 {
3     double Time1, Time2;
4     int i, j, k;
5     double *pha, *phb, *phc;
6
7     //Matrixes Memory allocation
8     pha = (double *)malloc((m_ar * m_ar) * sizeof(double));
9     phb = (double *)malloc((m_ar * m_ar) * sizeof(double));
10    phc = (double *)malloc((m_ar * m_ar) * sizeof(double));
11
12    //Starting counting time
13    Time1 = omp_get_wtime();
14
15    //Loading matrix values
16    for(i=0; i<m_ar; i++)
17        #pragma omp parallel for num_threads (x)
18        for(j=0; j<m_ar; j++)
19            pha[i*m_ar + j] = (double)1.0;
20
21    for(i=0; i<m_br; i++)
22        #pragma omp parallel for num_threads (x)
23        for(j=0; j<m_br; j++)
24            phb[i*m_br + j] = (double)(i+1);
25
26    for(i=0; i<m_ar; i++)

```



## Appendices

```
26      #pragma omp parallel for num_threads (x)
28      for(j=0; j<m_ar; j++)
29          phc[i*m_ar + j] = (double)0.0;
30
31      //Matrix Multiplication
32      for(i=0; i<m_ar; i++)
33      {      for( k=0; k<m_ar; k++)
34          {
35              #pragma omp parallel for num_threads (x)
36              for( j=0; j<m_br; j++)
37              {
38                  phc[i*m_ar+j] += pha[i*m_ar+k] * phb[k*m_br+j];
39              }
40          }
41      }
42  }
43
44      //Stoping time
45      Time2 = omp_get_wtime();
46
47      //Freeing memory used for matrixes
48      free(pha);
49      free(phb);
50      free(phc);
51
52      return Time2 - Time1;
53  }
```

## 7.2 Kremlin's Reports

### 7.2.1 Kremlin report for Matrix Multiplication, Mult version

Listing 7.7: Kremlin's indication of the blocks that should be parallelised and theoretical variables that where calculated for Mult algorithm version

```
36 1 Speedup: 3.65
37 Serial : 3289
38 Parallel: 901
39
40
41 5 [ 0] TimeRed(4)=66.38%, TimeRed(Ideal)=70.96%, Cov=88.51%, SelfP=5.05, DOALL
42 LOOP matrixmul.cpp [ 148 - 181]: OnMult
43 FUNC matrixmul.cpp [ 142 - 142]: OnMult called at file matrixmul.cpp, line 397
44
45 9 [ 1] TimeRed(4)=3.10%, TimeRed(Ideal)=3.37%, Cov=4.13%, SelfP=5.44, DOALL
46 LOOP matrixmul.cpp [ 149 - 167]: OnMult
47 FUNC matrixmul.cpp [ 142 - 142]: OnMult called at file matrixmul.cpp, line 397
48
49 13 [ 2] TimeRed(4)=3.10%, TimeRed(Ideal)=3.37%, Cov=4.13%, SelfP=5.44, DOALL
50 LOOP matrixmul.cpp [ 149 - 161]: OnMult
51 FUNC matrixmul.cpp [ 142 - 142]: OnMult called at file matrixmul.cpp, line 397
```

## 7.2.2 Kremlin report for Matrix Multiplication, MultLine version

6

Listing 7.8: Kremlin's indication of the blocks that should be parallelised and theoretical variables that where calculated for MultLine algorithm version

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22	<pre> ..... Speedup: 3.52 Serial  : 2334 Parallel: 663  [ 0] TimeRed(4)=63.01%, TimeRed(Ideal)=63.20%, Cov=84.02%, SelfP=4.03, DOALL     LOOP matrixmul.cpp [ 213 - 247]: OnMultLine     FUNC matrixmul.cpp [ 207 - 207]: OnMultLine called at file matrixmul.cpp, line 400  [ 1] TimeRed(4)=2.86%, TimeRed(Ideal)=2.96%, Cov=3.81%, SelfP=4.45, DOALL     LOOP matrixmul.cpp [ 213 - 235]: OnMultLine     FUNC matrixmul.cpp [ 207 - 207]: OnMultLine called at file matrixmul.cpp, line 400  [ 2] TimeRed(4)=2.86%, TimeRed(Ideal)=2.96%, Cov=3.81%, SelfP=4.45, DOALL     LOOP matrixmul.cpp [ 213 - 231]: OnMultLine     FUNC matrixmul.cpp [ 207 - 207]: OnMultLine called at file matrixmul.cpp, line 400  [ 3] TimeRed(4)=2.86%, TimeRed(Ideal)=2.96%, Cov=3.81%, SelfP=4.45, DOALL     LOOP matrixmul.cpp [ 213 - 225]: OnMultLine     FUNC matrixmul.cpp [ 207 - 207]: OnMultLine called at file matrixmul.cpp, line 400 </pre>	8 10 12 14 16 18 20 22
---	--	---