

UNIVERSITÉ de VERSAILLES ST-QUENTIN EN YVELINES

**Habilitation à diriger des recherches**

Discipline :

**INFORMATIQUE**

Présentée par :

**Sid-Ahmed-Ali TOUATI**

Titre :

**Méthodes d'optimisations de programmes bas niveau**

***On Backend Code Optimisation Methods***

Soutenue le mercredi 30 juin 2010 devant le jury suivant :

Pr.	Jens KNOOP	Université technique de Vienne	Rapporteur
Pr.	Jagannathan RAMANUJAM	Université de Louisiane	Rapporteur
Pr.	Denis TRYSTRAM	Institut polytechnique de Grenoble	Rapporteur
Dr.	Christian BERTIN	Directeur de centre d'expertise chez STMicroelectronics	Examinateur
Dr.	Alain DARTE	École normale supérieure de Lyon	Examinateur
Pr.	William JALBY	Université de Versailles St-Quentin en Yvelines	Examinateur
Pr.	Pascal SAINRAT	Université Paul Sabatier de Toulouse	Président

Thèse d'habilitation préparée au sein des laboratoires PRISM et de l'INRIA-Saclay



# Préface

*“... Was this is an example concerning foundational research based on rigorous mathematical and logical modelling and reasoning. I would like to single out as another example of the uniqueness and quality of the research of the applicant an example from another pole of the methodological spectrum applied by the applicant in his research, concerning practical experiments conducted in the search of hard problems, i.e. an example from the engineering pole of research. ... Moreover, I would like to add that the research work of the applicant is methodological sound, that it presents new and important scientific insights in theory and practice, and clearly demonstrate the ability of the applicant to contribute to the development and advancement of the field. ...”*

Professor Jens KNOOP  
Head of the institute of computer science  
Technical university of Vienna (Austria).

*“... The habilitation thesis and Dr. Touati’s work over the last ten years demonstrate that he is highly motivated and that he has the important ability to quickly grasp the central issues, and ask and answer the most important questions concerning research problems. His strong mathematical and reasoning abilities have added to his broad background for continued important research work. ...”*

Professor Jagannathan RAMANUJAM  
John E. and Beatrice L. Ritter distinguished professor  
Louisiana State University (USA)

*“... J’ai apprécié l’effort de Sid TOUATI pour fournir un document complet et cohérent. Il est clair dans l’ensemble. Les résultats techniques sont introduits avec précision. La lecture est parfois aride, mais le texte est ponctué de nombreuses figures très claires. Les nombreux résultats sont obtenus par des techniques variées qui montrent un large spectre de connaissances et une bonne culture informatique. Sid TOUATI est aujourd’hui un spécialiste reconnu du domaine de la compilation de bas niveau sur les processeurs modernes. Le travail produit est de qualité, ce qui est attesté par de nombreuses publications (8 journaux internationaux parmi les meilleurs du domaine et une dizaine de conférences de bon niveau). Le travail est bien équilibré entre des analyses conceptuelles parfois sophistiquées (complexité, analyses d’algorithmes, preuves d’optimalité, heuristiques) et des réalisations expérimentales (production de codes en logiciels libres, disponibles sur le site web de Sid TOUATI). ...”*

Professor Denis TRYSTRAM  
Distinguished professor at Grenoble Institute of Technology (France)

I am sincerely grateful to this honourable committee for providing me the title of “habilité à diriger des recherches”. Many thanks to professors TRYSTRAM, RAMANUJAM and KNOOP for their valuable reviews.

My first thinking after getting this title goes to my parents that were unable to attend today.

Then, I would like to thank professor William JALBY for his full support to me at the university. I must never forget the support of doctor Albert COHEN at INRIA.

I sincerely thank doctor Christine EISENBEIS, my first PhD advisor, because she helped me to go into this difficult but exciting field of research.

Thank you doctor Alain DARTE for your detailed remarks every-time and everywhere to improve the quality of my research.

At the end, I am grateful to all my students that worked and are still working hard to finish their PhD.

After one decade of continuous effort in publishing, teaching, implementing, experimenting, advising, preparing projects, travelling, discussing, studying, thinking and reading, we arrive to the most popular time of a defence, which is the time of celebrating our success.

Discourse of Sid-Ahmed-Ali TOUATI

June 30th, 2010

Associate professor at the University of Versailles Saint-Quentin en Yvelines

## Résumé

Ce manuscrit synthétise plus d'une décennie de notre recherche académique sur le sujet d'optimisation de codes bas niveau, dont le but est une intégration dans un compilateur optimisant ou dans un outil d'optimisation semi-automatique. Dans les programmes bas niveau, les caractéristiques du processeur sont connues et peuvent être utilisées pour générer des codes plus *en harmonie* avec le matériel.

Nous commençons notre document par une vue générale sur le problème d'ordonnancement des phases de compilation. Actuellement, des centaines d'étapes de compilation et d'optimisation de codes existent; un problème fondamental et ouvert reste de savoir comment les combiner et les ordonner efficacement. Pour pallier rapidement cette difficulté, une stratégie du moindre effort consiste à appliquer une compilation itérative en exécutant successivement le programme avant de décider de la technique d'optimisation de code à employer et avec quels paramètres. Nous prouvons que l'approche de compilation itérative ne simplifie pas fondamentalement le problème, et l'utilisation de modèles statiques de performances reste un choix raisonnable.

Un problème classique de conflit entre deux étapes de compilation est celui qui lie l'allocation de registres et l'ordonnancement d'instructions. Nous montrons comment gérer efficacement cet antagonisme en séparant les contraintes de registres des contraintes d'ordonnancement d'instructions. Cela est possible grâce à la notion de saturation en registres (RS), qui est le besoin maximal en registres pour tous les ordonnancements possibles d'un graphe. Nous apportons une contribution formelle et une heuristique efficace, qui permettent la détection de contraintes de registres toujours vérifiées; ils peuvent par conséquent être négligées.

Nous introduisons la plate-forme SIRA, qui permet de garantir l'absence de code de vidage avant l'ordonnancement d'instructions. SIRA est un modèle basé sur la théorie des graphes permettant de borner le besoin maximal en registres pour tout pipeline logiciel, sans altérer, si possible, le parallélisme d'instructions. SIRA modélise les contraintes cycliques des registres dans différentes architectures possibles : avec plusieurs types de registres, avec tampons ou files d'attente, et avec des bancs de registres rotatifs. Nous apportons une heuristique efficace qui montre des résultats satisfaisants, que ce soit comme outil indépendant, ou comme passe intégrée dans un vrai compilateur.

Dans le contexte des processeurs exhibant des retards d'accès aux registres (VLIW, EPIC, DSP), nous attirons l'attention sur le problème qui peut survenir lorsque les contraintes de registres sont traitées avant l'ordonnancement d'instructions. Ce problème est la création de circuits négatifs ou nuls dans le graphe de dépendances de données. Nous montrons comment éliminer ces circuits indésirables dans le contexte de SIRA.

SIRA définit une relation formelle entre le nombre de registres alloués, le parallélisme d'instructions et le facteur de déroulage d'une boucle. Nous nous basons sur cette relation pour écrire un algorithme optimal qui minimise le facteur de déroulage tout en sauvegardant le parallélisme d'instructions et en garantissant l'absence de code de vidage. D'après nos connaissances, ceci est le premier résultat qui démontre que le compactage de la taille de code n'est pas un objectif antagoniste à l'optimisation des performances de code.

L'interaction entre la hiérarchie mémoire et le parallélisme d'instructions est un point central si l'on souhaite réduire le coût des latences d'opérations de chargement. Premièrement, notre étude pratique avec des micro-benchmarks montre que les processeurs superscalaires ayant une exécution dans le désordre ont un *bug* de performances dans leur mécanisme de désambiguïcation mémoire. Nous montrons ensuite qu'une vectorisation des opérations mémoire résoud ce problème pour des codes réguliers. Deuxièmement, nous étudions l'optimisation de préchargement de données pour des codes VLIW embarqués irréguliers.

Finalement, avec l'arrivée des processeurs multicœurs, nous observons que les temps d'exécution des programmes deviennent très variables. Afin d'améliorer la reproductibilité des résultats expérimentaux, nous avons conçu le *Speedup-Test*, un protocole statistique rigoureux. Nous nous basons sur des tests statistiques connus (tests de Shapiro-Wilk, F de Fisher, de Student, de Kolmogorov-Smirnov, de Wilcoxon-Mann-Whitney) afin d'évaluer si une accélération observée du temps d'exécution médian ou moyen est significative.

**Mots clés :** parallélisme d'instructions, ordonnancement d'instructions, allocation de registres, saturation en registres, pipeline logiciel, programmation linéaire, programmation linéaire en nombres entiers, hiérarchie mémoire, évaluation des performances des programmes, compilation, optimisation de code.

## Abstract

This manuscript is a synthesis of our research effort since one full decade on the topic of low level code optimisation, devoted to an integration in a compiler backend or in a semi-automatic optimisation tool. At the backend level, processor characteristics are known and can be used to generate codes using the underlying hardware more efficiently.

We start our document by a global view on the phase ordering problem in optimising compilation. Nowadays, hundreds of compilation passes and code optimisation methods exist, but nobody knows exactly how to combine and order them efficiently. Consequently, a best effort strategy consists in doing an iterative compilation by successively executing the program to decide about the passes and optimisation parameters to apply. We prove that iterative compilation does not fundamentally simplify the problem, and using static performance models remains a reasonable choice.

A well known phase ordering dilemma between register allocation and instruction scheduling has been debated for long time in the literature. We show how to efficiently decouple register constraints from instruction scheduling by introducing the notion of register saturation (RS). RS is the maximal register need of all the possible schedules of a data dependence graph. We provide formal methods for its efficient computation, that allows to detect obsolete register constraints. Consequently, they can be neglected from the instruction scheduling process.

In order to guarantee the absence of spilling before instruction scheduling, we introduce the SIRA framework. It is a graph theoretical approach that bound the maximal register need for any subsequent software pipelining, while saving instruction level parallelism. SIRA model periodic register constraints in the context of multiple register types, buffers and rotating register files. We provide an efficient heuristic that show satisfactory results as a standalone tool, as well as an integrated compilation pass inside a real compiler.

In the context of processors with architecturally visible delays to access registers (VLIW, EPIC, DSP), we highlight an open problem that arises when register constraints are handled before instruction scheduling. This problem is the creation of non-positive cycles inside data dependence graphs. We show how to remove these undesirable cycles in the context of SIRA.

SIRA defines a formal relationship between the number of allocated registers, the instruction level parallelism and the loop unrolling factor. We use this relationship to write an optimal algorithm that minimises the unrolling factor while saving instruction level parallelism and guaranteeing the absence of spilling. As far as we know, this is the first result in the literature proving that code size compaction and code performance are not antagonistic optimisation objectives.

The interaction between memory hierarchy and instruction level parallelism is of crucial issue if we want to hide or to tolerate load latencies. Firstly, we practically demonstrate that superscalar out-of-order processors have a performance bug in their memory disambiguation mechanism. We show that a load/store vectorisation solves this problem for regular codes. For irregular codes, we study the combination of low level data pre-loading and prefetching, designed for embedded VLIW processors.

Finally, with the introduction of multicore processors, we observe that program execution times may be very variable in practice. In order to improve the reproducibility of the experimental results, we design the *Speedup-Test*, which is a rigorous statistical protocol. We rely on well known statistical tests (Shapiro-wilk's test, Fisher's F-test, Student's t-test, Kolmogorov-Smirnov's test, Wilcoxon-Mann-Whitney's test) to evaluate if an observed speedup of the average or the median execution time is significant.

**Keywords :** Instruction-Level Parallelism, Instruction Scheduling, Register Allocation, Register Saturation, Software Pipelining, Linear Programming, Integer Linear Programming, Memory Hierarchy, Program Performance Evaluation, Compilation, Code Optimisation

# Contents

<b>1 Prologue</b>	<b>11</b>
1.1 Introduction . . . . .	11
1.2 Inside this Manuscript . . . . .	12
<b>2 Phase Ordering in Optimising Compilation</b>	<b>15</b>
2.1 Introduction to the Phase Ordering Problem . . . . .	15
2.2 Background on Phase Ordering . . . . .	16
2.2.1 Performance Modelling and Prediction . . . . .	16
2.2.2 Some Attempts in Phase Ordering . . . . .	17
2.3 Towards a Theoretical Model for Phase Ordering Problem . . . . .	18
2.3.1 Decidability Results . . . . .	19
2.3.2 Another Formulation of the Phase Ordering Problem . . . . .	20
2.4 Examples of Decidable Simplified Cases . . . . .	21
2.4.1 Models with Compilation Costs . . . . .	21
2.4.2 One-pass Generative Compilers . . . . .	22
2.5 Compiler Optimisation Parameters Space Exploration . . . . .	24
2.5.1 Towards a Theoretical model . . . . .	24
2.5.2 Examples of Simplified Decidable Cases . . . . .	25
2.6 Conclusion on Phase Ordering . . . . .	27
<b>3 The Register Need</b>	<b>29</b>
3.1 Data Dependence Graph and Processor Models . . . . .	29
3.2 The Acyclic Register Need . . . . .	30
3.3 The Periodic Register Need . . . . .	32
3.3.1 Software Pipelining, Periodic Scheduling, Cyclic Scheduling . . . . .	32
3.3.2 The Circular Lifetime Intervals . . . . .	33
3.4 Computing the Periodic Register Need . . . . .	34
3.5 Some Results on the Periodic Register Need . . . . .	37
3.5.1 Minimal Periodic Register Need vs. Initiation Interval . . . . .	37
3.5.2 Computing the Periodic Register Sufficiency . . . . .	37
3.5.3 Stage Scheduling under Register Constraints . . . . .	38
3.6 Conclusion on the Register Requirement . . . . .	41
<b>4 The Register Saturation</b>	<b>43</b>
4.1 Motivations on the Register Saturation Concept . . . . .	43
4.2 Computing the Acyclic Register Saturation . . . . .	45
4.2.1 Characterising the Register Saturation . . . . .	46
4.2.2 Efficient Algorithmic Heuristic for RS Computation . . . . .	48
4.2.3 Experimental Efficiency of GREEDY-K . . . . .	50
4.3 Computing the Periodic Register Saturation . . . . .	51
4.4 Conclusion on the Register Saturation . . . . .	54

<b>5 Spill Code Reduction</b>	<b>55</b>
5.1 Introduction on Register Constraints in Software Pipelining . . . . .	55
5.2 Related Work in Periodic Register Allocation . . . . .	56
5.3 SIRA: Schedule Independant Register Allocation . . . . .	57
5.3.1 Reuse Graphs . . . . .	57
5.3.2 DDG Associated to Reuse Graph . . . . .	58
5.3.3 Exact SIRA with Integer Linear Programming . . . . .	60
5.3.4 SIRA with Fixed Reuse Edges . . . . .	61
5.4 SIRALINA: An Efficient Polynomial Heuristic for SIRA . . . . .	62
5.5 Experimental Results with SIRA . . . . .	65
5.6 Conclusion on Spill Code Reduction . . . . .	66
<b>6 Exploiting the Register Access Delays</b>	<b>67</b>
6.1 Problem Description of DDG Cycles with Non-positive Distances . . . . .	67
6.2 Eliminating Non-Positive Cycles . . . . .	68
6.3 Experimental Results on Eliminating Non-Positive Cycles . . . . .	71
6.4 Conclusion on Non-Positive Cycles Elimination . . . . .	72
<b>7 Loop Unrolling Degree Minimisation</b>	<b>75</b>
7.1 Introduction . . . . .	75
7.2 Unroll Degree Minimisation of Unscheduled Loops . . . . .	77
7.2.1 Problem Description of Unroll Factor Minimisation for Unscheduled Loops . . . . .	77
7.2.2 Algorithmic Solution for Unroll Factor Minimisation: Single Register Type . . . . .	78
7.2.3 Solution for LCM Problem . . . . .	79
7.2.4 Unroll Factor Minimisation in Presence of Multiple Register Types . . . . .	81
7.2.5 Solution for Minimal Loop Unrolling . . . . .	85
7.3 Unroll Degree Minimisation of Scheduled Loops . . . . .	86
7.4 Experimental Results . . . . .	87
7.5 Conclusion on Loop Unroll Degree Minimisation . . . . .	88
<b>8 Memory Hierarchy Effects and ILP</b>	<b>91</b>
8.1 Problem of Memory Disambiguation at Runtime . . . . .	91
8.1.1 Introduction . . . . .	91
8.1.2 Related Work . . . . .	92
8.1.3 Experimental Environment . . . . .	93
8.1.4 Experimentation Methodology . . . . .	93
8.1.5 Experimental Study of Cache Behavior . . . . .	94
8.1.6 The Effectiveness of Load/Store Vectorisation . . . . .	97
8.1.7 Conclusion on Memory Disambiguation Mechanisms . . . . .	99
8.2 Data Preloading and Prefetching . . . . .	100
8.2.1 Introduction . . . . .	100
8.2.2 Related Work . . . . .	100
8.2.3 Problems of Optimising Cache Effects at the Instruction Level . . . . .	102
8.2.4 Target Processor Description . . . . .	103
8.2.5 Our Methodology of Instruction-Level Code Optimisation . . . . .	104
8.2.6 Experimental Results . . . . .	108
8.2.7 Conclusion on Pre-fetching and Pre-Loading . . . . .	108
<b>9 Statistical Performance Analysis</b>	<b>111</b>
9.1 Code Performance Variation . . . . .	111
9.2 The Speedup-Test Protocole . . . . .	112
9.2.1 The Observed Speedups . . . . .	112
9.2.2 The Speedup of the Observed Average Execution Time . . . . .	114
9.2.3 The Speedup of the Observed Median Execution Time, as well as Individual Runs	115
9.3 Discussion and Conclusion on the Speedup-Test . . . . .	117

<b>10 Epilogue</b>	<b>121</b>
10.1 Problem of Instruction Selection . . . . .	121
10.2 Perspectives on Code Optimisation for Multi-Core Processors . . . . .	122
10.3 General Conclusion . . . . .	122
<b>A Benchmarks Presentation</b>	<b>125</b>
A.1 Qualitative Benchmarks Presentation . . . . .	125
A.2 Quantitative Benchmarks Presentation . . . . .	126
A.3 Changing the Architectural Configuration of the Processor . . . . .	130
<b>B Experiments on Register Saturation</b>	<b>131</b>
B.1 The Acyclic Register Saturation . . . . .	131
B.1.1 On the Oprimal RS Computation . . . . .	131
B.1.2 On the Accuracy of GREEDY-K Heuristic vs. Optimal RS . . . . .	131
B.1.3 GREEDY-K Execution Times . . . . .	133
B.2 The Periodic Register Saturation . . . . .	133
B.2.1 Optimal PRS Computation . . . . .	135
B.2.2 Approximate PRS Computation with Heuristic . . . . .	136
<b>C Experiments on SIRA</b>	<b>139</b>
C.1 Efficiency of SIRALINA on Standalone DDG . . . . .	139
C.1.1 Naming conventions for register optimisation orders . . . . .	139
C.1.2 Experimental efficiency of SIRALINA . . . . .	139
C.1.3 Measuring the Increase of the MII . . . . .	140
C.1.4 Efficiency of SIRALINA Execution Times . . . . .	140
C.2 Efficiency of SIRALINA plugged Inside Industrial Compiler . . . . .	140
C.2.1 Static Performance Results . . . . .	145
C.2.2 Execution Time Performance Results . . . . .	147
<b>D Experiments on Non-Positive Cycles Elimination</b>	<b>151</b>
D.1 Experimental Setup . . . . .	151
D.1.1 Heuristics Nomenclature . . . . .	151
D.1.2 Empirical Efficiency Measures . . . . .	151
D.2 Comparison of the Heuristics Execution Times . . . . .	152
D.2.1 Time to Minimise Register Pressure for a fixed $II$ . . . . .	152
D.3 Convergence of the Proactive Heuristic (Iterative SIRALINA) . . . . .	154
D.4 Qualitative Analysis of the Heuristics . . . . .	154
D.4.1 Number of Saved Registers . . . . .	154
D.4.2 Proportion of Sucess . . . . .	157
D.4.3 Increase of the MII when Sucess . . . . .	157
D.5 Conclusion on Non-Positive Cycles Elimination Strategy . . . . .	157
<b>E Experiments on Unroll Degree Minimisation</b>	<b>159</b>
E.1 Standalone Experiments with Single Register types . . . . .	159
E.1.1 Experiments with Unscheduled Loops . . . . .	159
E.1.2 Results on Randomly Generated DDG . . . . .	159
E.1.3 Experiments on Real DDG . . . . .	160
E.1.4 Experiments with Scheduled Loops . . . . .	162
E.2 Experiments with Multiple Register Types . . . . .	165
<b>F Experiments on Preloading and Prefetching</b>	<b>169</b>
<b>G Synthèse des travaux de recherche en français</b>	<b>173</b>