

Automatic Parallelization for Heterogeneous Computing

Rodrigo Caetano Rocha

Prospective PhD supervisor - Dr. Zheng Wang

Abstract

1 Introduction

During the last several years, modern computers are increasingly parallel and heterogeneous [32, 30]. With parallel architectures being widely spread, parallel programming must extend beyond its traditional realm of scientific applications [24]. The ultimate goal of a programmer in a modern computing environment is to write scalable applications that can take complete advantage of all available resources, since parallelization is an important step that can have serious impact on performance and power consumption [12]. However, parallel programming is known to be difficult and error prone, imposing several challenges to the programmer, such as race conditions, deadlocks, or even non-determinism [12, 28, 27]. Thus, the difficulty of developing parallel software is a key obstacle for an average programmer to exploit the considerable computational power that different multiprocessors offer [30].

Although manual parallelization by expert programmers is known to usually results in the most efficient parallel implementation, it is a costly and time-consuming approach, besides requiring specialized knowledge about parallelism and the target architecture. As hardware parallelism increases in scale with each generation and programming costs increase, automated parallelizing technology becomes extremely attractive with the potential to greatly reduce the costs while still profiting from the parallel architecture [42, 47].

Automatic parallelization is an optimization that produces parallel code that is semantically equivalent to a sequential code received as input. This is a difficult problem since the parallel code must be correct and also perform efficiently on the target machine [49]. Despite intense research interest in the area, it has failed to deliver outside niche domains [42, 47]. Improving the performance of programs for multiprocessors presents a challenge because even when thread-level parallelism exists, it is difficult for a compiler to analyze and identify all true dependencies between potential parallel threads [18].

The traditional approach had been to concentrate on restructuring computationally intensive loops using dependence analysis to determine which loops can be parallelized. This analytical approach has only proved successful with a small and restricted-set of programs [49]. Traditional static parallelism detection techniques are not effective in finding parallelism due to limitations on the level of information that can be confidently

acquired during static analysis [42, 47]. Another limitation in traditional static parallelism relates to the combination of automatic parallelism discovery and portable mapping. Given that the number and type of processors of a parallel system is likely to change from one generation to the next, finding the right mapping for an application may have to be repeated many times throughout an application’s lifetime, hence, making automatic approaches attractive [42, 47].

Furthermore, although this area have been widely studied, there is still much to be explored regarding automatic parallelization that also optimize for energy consumption, embedded systems [13] or other emergent computing architectures [3], such as massively parallel processors including graphic processing units (GPUs), many integrated cores (MICs), and others.

Considering that most programs are executed indefinitely many more times than they are compiled, compilers can afford performing expensive optimization during compilation time. Previous work suggest different approaches for automatic parallelization based on machine learning techniques [42, 47] and also evolutionary computing [43, 44, 48, 49]. However, extensive research into automatic parallelization over the years has lead to the development of a large body of heuristic and analytical techniques which may be included into such modern approaches [49].

2 Related work

In this section we present an overview of several work related to automatic parallelization of sequential code, where we discuss several parallelization approaches for both multi-core and heterogeneous systems.

2.1 Automatic parallelism based on static analysis

Numerous early work make use of static analysis techniques in order to parallelize affine loops [22, 34, 6, 25]. Kuck et al. [22] discuss the use of data dependence graphs for several code transformations and optimizations, including loop distribution, loop fusion, loop blocking, loop interchanging, etc. Lim and Lam [25] propose an algorithm for finding the optimal affine transformations that maximize parallelism while minimizing synchronization in programs with arbitrary loop nesting and affine data accesses.

Misailovic et al. [30] present *QuickStep*. Guided by performance measurements, static analysis, loop and memory profiling information, and statistical accuracy tests on the outputs of test executions, QuickStep automatically searches this space to find a parallel program that maximizes performance subject to a specified accuracy goal.

QuickStep prioritizes the parallelization of loops that consume more execution time, where it explores the parallelization space for each loop by repeatedly applying accuracy and performance enhancing transformations. Once the final parallelization is obtained, QuickStep produces an interactive parallelization report that a developer can navigate to evaluate the acceptability of this parallelization and obtain insight into how the application responds to different parallelization strategies.

Automatic parallelization based on static analysis is challenging and limited in the type of code that they can parallelize [8, 19, 11]. In order to identify parallel regions of code they must use complex interprocedural analyses to prove that the code has no data dependencies for all possible inputs. Typical code targeted by these compilers consists

of nested loops with affine array accesses written in a language with limited aliasing. Large systems written in modern languages usually contain multiple modules and memory aliasing, which makes them not amenable to automatic parallelization. Furthermore, code whose memory access patterns are indeterminable at compile time due to dependence on program inputs can be impossible for these compilers to parallelize [8].

2.2 Thread-level speculative parallelism

In contrast to automatic parallelization that aims at generating a parallel code, thread-level speculation (TLS) approaches performs sequential code parallelization during runtime in a speculative manner. It consists of selecting regions of code to execute in parallel, relying on the hardware to detect dependence violations and re-execute the conflicting sections such that sequential execution semantics is preserved [18, 11, 50].

Although TLS overcomes limitations intrinsic with conservative compile-time automatic parallelizing tools by extracting parallel threads optimistically and only ensuring absence of data dependence violations at runtime. TLS usually requires a specialized hardware for monitoring memory accesses made by the parallel threads and for restarting any thread that attempt to violate true dependencies [21, 18, 50]. The overheads associated with maintaining speculative state is a significant barrier for the use of TLS [51].

While TSL is usually performed mainly by hardware, there also exists software-based speculative parallelism, such as Softspec, which is the combination of a compiler and runtime system that can target program binaries [8].

2.3 Machine learning based approach

Wang and O’Boyle [46] propose a machine learning based approach for mapping fine-grain parallelism of OpenMP programs to multi-cores, namely, Intel Xeon and the Cell processors. They use an artificial neural network (ANN) model to predict the scalability of a program and select the optimal thread number for it. Afterwards, a multi-class Support Vector Machine (SVM) model is used to select a scheduling policy, such as cyclic, dynamic, guided or block scheduling.

The approach proposed by Tournavitis et al. [42] and Wang et al. [47], integrates profile-driven parallelism detection and machine learning based mapping in a single framework for semi-automated parallelism, where the user is expected to approve those loops where parallelization is likely to be beneficial, but correctness cannot be proven conclusively. They use profiling data to extract actual control and data dependencies and enhance the corresponding static analyses with dynamic information. Subsequently, they apply a previously trained prediction mechanism to each parallel loop candidate and decide if and how the parallel mapping should be performed. Finally, they generate parallel code using standard OpenMP annotations.

2.4 Automatic parallelism based on evolutionary computing

Walsh and Ryan propose *Paragen* [43, 44], a precursor for the use of evolutionary computing in automatic parallelization. Paragen makes no attempt to analyse the original code, instead it attempts to directly performs both loops and instructions parallelization. The correctness of the resulting parallel code are determined by the fitness function based

on the outputs produced by the original sequential code and the one produced by the parallel code.

Williams and Williams [48, 49] also propose an evolutionary approach for automatic parallelization. This work suggests two genetic representation. The *gene-transformation* where each gene represents a loop transformation (such as loop-skewing, loop-interchange, loop-fusion, loop-splitting, loop-reversal, etc.) and the chromosome has variable-length and is composed by an ordered sequence of transformations to be applied. The *gene-statement* where each gene represents a single instruction in the source code. For the gene-statement representation, no crossover operation is performed and the loop transformations are encoded as mutation operators with predefined probabilities, where the length of the chromosome can change depending on the mutations applied. The fitness function perform an analysis of the resulting code of each chromosome generating an estimated execution-time. They conclude that the gene-statement representation is better than the gene-transformations representation.

Evolutionary based approaches have an optimistic parallelization strategy, experimentally validating the best individuals. Although this approach is not limited by dependence analysis techniques, the evolutionary approach produces a probabilistic result, only with statistical guarantees. In Section 2.7 we present several other work that discuss probabilistic and approximate solutions for a variety of problems.

2.5 Automatic parallelization for heterogeneous systems

Most studies focus on automatic parallelization for CPU only. However, there have been some recent work [23, 3, 1, 17] toward parallelization for heterogeneous architectures, usually GPUs.

Leung et al. [23] extend the Java virtual machine to detect affine loops, in plain Java bytecodes, which can be parallelized and executed more quickly on the GPU rather than on the CPU. They present a cost model that balances the higher raw performance of the available GPU against the cost of transferring input and output data between main memory and GPU memory. There are several obstacles that need to be overcome when parallelizing Java bytecode language, such as unstructured control flow, the lack of multi-dimensional arrays, the precise exception semantics, and the proliferation of indirect references.

Baskaran et al. [2, 3] propose a framework for automatic parallelization and optimization of affine loops on GPUs. A polyhedral compiler model of data dependence is used to perform program transformation for efficient data access from GPU global memory. Their framework performs source-to-source transformations, generating CUDA from C code, optimized for efficient data access, including tiling strategies. They use the PluTo [7] polyhedral parallel tiling infrastructure and the polyhedral code generator called CLooG [4], which transforms a polyhedral representation of a program and affine scheduling constraints into final loop code.

Parallelizing techniques has several limitations, one such limitation is regarding indirect array accesses. Govindarajan and Anantpur [17] present an algorithm for executing on GPU, loops with cross iteration dependencies due to indirect memory accesses. The algorithm described is able to compute memory dependencies at runtime, where dependence computation and actual computation are performed in pipeline.

2.6 Evolutionary computing applied to other compiler optimizations

Stephenson et al. [41] make use of genetic programming to optimize the priority functions associated with two compiler heuristics, namely, predicated hyperblock formation and register allocation, achieving significant speedups over standard baselines.

Schulte et al. [38] propose a post-compilation optimization approach that leverages ideas from evolutionary computation (population-based stochastic search), mutational robustness (random mutations yield independent implementations of the same specification), profile-guided optimization (performance measured on workloads) and relaxed notions of program semantics (customizing software to particular runtime goals and environments provided by a software engineer). The proposed approach is a steady-state evolutionary algorithm, which maintains a population of candidate optimizations (assembly programs), using randomized operators to generate variations, and selecting those that improve an objective function (the power model) while retaining all required functionality expressed in a test suite. Their results show that the evolutionary approach is effective on multiple architectures, finding hardware-specific optimizations, with an average reduction of 20% in energy consumption.

2.7 Performance versus accuracy trade-off

Several work discuss probabilistic and approximate architectures in order to design hardwares with improved performance regarding power consumption, execution time, physical space and architectural simplicity [36, 35, 26, 20]. In contrast with traditional approaches, approximate architectures does not attempt to correct the errors introduced by components which are susceptible to perturbations. Palem et al. suggest their work has the potential to address challenges and impediments to Moore’s law arising from material properties and manufacturing difficulties, indicating that their shifts from the current-day deterministic design paradigm to statistical and probabilistic designs of the future [36, 35]. Palem et al. demonstrate a XOR gate with a $3\times$ reduction in energy consumption for a 2.3% decrease in the probability of correctness. Chakrapani et al. [10] present a mathematically rigorous foundation characterizing the trade-off between the energy consumed and the quality of solutions.

Some recent work also propose numerous code transformations aiming at the trade-off between performance and accuracy. For certain types of applications, accuracy-aware program transformations may generate error-tolerant code that may perform better than the original code [52]. The main transformations are: (i) *substitution transformations* replace parts of a program with code that computes approximations of the output computed by the original parts but with less computational overhead; (ii) *sampling transformations* produces a transformed code that performs the same computation as the original code but only on a subset of the elements obtained by some sampling policy; (iii) *loop perforation* is another type of program transformation which transforms a loop into a new loop where only a subset of the original loop iterations is performed.

Misailovic et al. [31] provide a probabilistic approach based that guarantees for the accuracy produced by loop perforation, while reducing execution time and resource consumptions. The transformations are analyzed based on the perforation noise, defined by the difference between the result that the perforated computation produces and the result that the original computation produces. The perforation noise is then used to produce

the probabilistic guarantee of the transformation’s accuracy.

Sidiroglou et al. [39] show that perforating appropriately selected loops can produce significant performance gains (up to a factor of seven reduction in overall execution time) in return for small (less than 10%) accuracy losses.

3 Research directions

It is widely known that detecting parallelism with static compilers in large, general programs is challenging [19, 11]. In order to overcome the limitations of automatic parallelization based on static analysis for general multiprocessors, we intend to explore optimistic and evolutionary approaches aided by profiling and machine learning techniques. We aim at approaches where it is possible to parallelize aggressively rather than conservatively [11, 48, 49]. In this section we suggest some possible research directions and their respective challenges.

3.1 An evolutionary approach

Automatic parallelization using evolutionary computing has an optimistic approach, where the code is parallelized more aggressively. Each individual of a generated population can then be scored based on profiling techniques. While this approach is not limited by dependence analysis mechanisms, several challenges need to be addressed, including: (i) define a genetic representation for the source code, the parallel transformations and optimizations; (ii) define a fitness function that score each individual considering the trade-off between performance and accuracy; (iii) measure, with statistical confidence, the accuracy of the parallelized source code; (iv) efficiently profile each individual.

Efficiently profiling each individual is important for reducing the latency involved in the overall parallelization process using the evolutionary approach. Some of the strategies for efficiently profiling the program corresponding to each individual include: simulating or actually executing the program; executing a transformed version of the individual’s program after performing sampling or loop perforations [39, 31, 52]; estimating the performance without actually executing the program [14, 15, 16].

Besides the traditional *Darwinian evolution model* for randomly creating the initial population, mutations and crossovers, it is also possible to use the *learnable evolution model* [29, 9], where the individuals are *genetically engineered*. In the *learnable evolution model*, individuals can be created guided by static analysis, profiling, or machine learning.

3.2 A machine learning based approach

In addition to being used for mapping parallelism to multi-core processors [46, 42, 47], machine learning can also be used for parallelizing a sequential code. Similar to the evolutionary approach, once we have optimistically identified all once the optimization and parallelization search space has been identified in an optimistic manner, by analyzing the sequential code, machine learning techniques, such as ANN [33, 45], can be used for automatically selecting the best parallelization.

Several work propose solutions using neural networks for a variety optimization problems in general [33, 45, 40], which can be applied either using supervised or unsupervised learning [5].

3.3 Solving dependence violations during runtime

By modeling the automatic parallelization problem as an optimization problem, defining a parallelization search space in an optimistic manner, both the evolutionary and the machine learning based approaches will produce probabilistic solutions.

If a probabilistic solution is unacceptable, we can use speculative techniques for detecting dependence violations during runtime for re-execution of conflicting sections providing precise solutions. This dependence monitoring mechanism can be implemented entirely in software, similar to Softspec [8], allowing this technique to be used in general parallel hardware. Although this mechanism guarantee correctness, it adds some overhead to the overall execution, which needs to be properly balanced for performance.

It is also possible to integrate techniques from dynamic parallelization, such as those used by the *hybrid analysis* [37]. *Hybrid analysis* combines static and runtime analysis of memory references into a single framework. In the context of heterogeneous systems, Govindarajan and Anantpur [17] describe an algorithm which is able to compute memory dependencies at runtime, where dependence computation and actual computation are pipelined on a GPU.

3.4 Parallelization for heterogeneous systems and energy consumption

Although automatic parallelization have been widely studied, there is still much to be explored regarding energy consumption, embedded [13] or emergent computing architectures [23, 3, 1, 17], such as massively parallel processors including graphic processing units (GPUs) or many integrated cores (MICs).

Several new challenges are added when parallelizing for heterogeneous systems. In addition to identifying parallelizable loops, it must also decide, for each loop, if the loop parallelization will be profitable, and if so, where it should be executed, i.e., whether to execute it on the CPU, the GPU, or any other available parallel processor [42, 47, 23]. In a heterogeneous system, besides the raw performance of each processor, synchronization, data communication or even energy consumption must also be considered when evaluating how the parallel mapping should be performed.

References

- [1] Mehdi Amini, Béatrice Creusillet, Stéphanie Even, Ronan Keryell, Onig Goubier, Serge Guelton, Janice O. McMahon, François-Xavier Pasquier, Grégoire Péan, Pierre Villalon, and Others. Par4all: From convex array regions to heterogeneous computing. In *IMPACT & HiPEAC*, 2012.
- [2] Muthu Manikandan Baskaran, Uday Bondhugula, Sriram Krishnamoorthy, J. Ramanujam, Atanas Rountev, and P. Sadayappan. A compiler framework for optimization of affine loop nests for gpgpus. In *ICS*, pages 225–234. ACM, 2008.
- [3] Muthu Manikandan Baskaran, J. Ramanujam, and P. Sadayappan. Automatic c-to-cuda code generation for affine programs. In *CC/ETAPS*, pages 244–263. Springer-Verlag, 2010.

- [4] Cédric Bastoul. Code generation in the polyhedral model is easier than you think. In *PACT*, pages 7–16. IEEE, 2004.
- [5] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, Université De Montréal, and Montréal Québec. Greedy layer-wise training of deep networks. In *NIPS*. MIT Press, 2007.
- [6] Bill Blume, Rudolf Eigenmann, Keith Faigin, John Grout, Jay Hoeflinger, David Padua, Paul Petersen, Bill Pottenger, Lawrence Rauchwerger, Peng Tu, and Stephen Weatherford. Polaris: The next generation in parallelizing compilers. In *LCPC*, pages 10–1. Springer-Verlag, 1994.
- [7] Uday Bondhugula, Albert Hartono, J. Ramanujam, and P. Sadayappan. A practical automatic polyhedral parallelizer and locality optimizer. In *PLDI*, pages 101–113. ACM, 2008.
- [8] Derek Bruening, Srikrishna Devabhaktuni, and Saman Amarasinghe. Softspec: Software-based speculative parallelism. In *FDDO*. ACM, 1998.
- [9] Guido Cervone, Ryszard S. Michalski, Kenneth K. Kaufman, and Liviu A. Panait. Combining machine learning with evolutionary computation: Recent results on LEM. In *IWML*, pages 41–58, 2000.
- [10] Lakshmi N.B. Chakrapani, Kirthi Krishna Muntimadugu, Avinash Lingamneni, Jason George, and Krishna V. Palem. Highly energy and performance efficient embedded computing through approximately correct arithmetic: A mathematical foundation and preliminary experimental validation. In *CASES*, pages 187–196. ACM, 2008.
- [11] M.K. Chen and K. Olukotun. The jrpm system for dynamically parallelizing java programs. In *ISCA*, pages 434–445. IEEE Comp. Soc., 2003.
- [12] Johan Cockx, Bart Vanhoof, Richard Stahl, and Patrick David. System and method for automatic parallelization of sequential code, 2010. US Patent 7,797,691.
- [13] D. Cordes, P. Marwedel, and A. Mallik. Automatic parallelization of embedded software using hierarchical task graphs and integer linear programming. In *CODES+ISSS*, pages 267–276. IEEE/ACM, 2010.
- [14] T. Fahringer. Estimating and optimizing performance for parallel programs. *Computer*, 28(11):47–56, 1995.
- [15] T. Fahringer, A. Pozgaj, J. Luitz, and H. Moritsch. Evaluation of P3T+: a performance estimator for distributed and parallel applications. In *IPDPS*, pages 229–234. IEEE, 2000.
- [16] Thomas Fahringer. *Automatic Performance Prediction of Parallel Programs*. Springer Publishing Company, Incorporated, 1st edition, 2011.
- [17] R. Govindarajan and Jayvant Anantpur. Runtime dependence computation and execution of loops on heterogeneous systems. In *CGO*, pages 1–10. IEEE/ACM, 2013.

- [18] Lance Hammond, Mark Willey, and Kunle Olukotun. Data speculation support for a chip multiprocessor. *SIGOPS Oper. Syst. Rev.*, 32(5):58–69, 1998.
- [19] Ken Kennedy and John R. Allen. *Optimizing Compilers for Modern Architectures: A Dependence-based Approach*. Morgan Kaufmann Publishers Inc., 2002.
- [20] Christoph M. Kirsch and Hannes Payer. Incorrect systems: It’s not the problem, it’s the solution. In *DAC*, pages 913–917. ACM, 2012.
- [21] Venkata Krishnan and Josep Torrellas. Hardware and software support for speculative execution of sequential binaries on a chip-multiprocessor. In *ICS*, pages 85–92. ACM, 1998.
- [22] D. J. Kuck, R. H. Kuhn, D. A. Padua, B. Leasure, and M. Wolfe. Dependence graphs and compiler optimizations. In *POPL*, pages 207–218. ACM, 1981.
- [23] Alan Leung, Ondřej Lhoták, and Ghulam Lashari. Automatic parallelization for graphics processing units. In *PPPJ*, pages 91–100. ACM, 2009.
- [24] Jiangtian Li, Xiaosong Ma, K. Singh, M. Schulz, B.R. de Supinski, and S.A. McKee. Machine learning based online performance prediction for runtime parallelization and task scheduling. In *ISPASS*, pages 89–100. IEEE, 2009.
- [25] Amy W. Lim and Monica S. Lam. Maximizing parallelism and minimizing synchronization with affine transforms. In *POPL*, pages 201–214. ACM, 1997.
- [26] Avinash Lingamneni, Kirthi Krishna Muntimadugu, Christian Enz, Richard M. Karp, Krishna V. Palem, and Christian Piguet. Algorithmic methodologies for ultra-efficient inexact architectures for sustaining technology scaling. In *CF*, pages 3–12. ACM, 2012.
- [27] Michael McCool, James Reinders, and Arch Robison. *Structured parallel programming: patterns for efficient computation*. Elsevier, 2012.
- [28] Michael D. McCool. Structured parallel programming with deterministic patterns. In *Proceedings of the 2Nd USENIX Conference on Hot Topics in Parallelism, Hot-Par’10*, pages 5–5, Berkeley, CA, USA, 2010. USENIX Association.
- [29] Ryszard S. Michalski. Learnable evolution model: Evolutionary processes guided by machine learning. *Machine Learning*, 38(1-2):9–40, 2000.
- [30] Sasa Misailovic, Deokhwan Kim, and Martin Rinard. Parallelizing sequential programs with statistical accuracy tests. *ACM Trans. Embed. Comput. Syst.*, 12(2):88:1–88:26, 2013.
- [31] Sasa Misailovic, DanielM. Roy, and MartinC. Rinard. Probabilistically accurate program transformations. In *Static Analysis*, volume 6887 of *Lecture Notes in Computer Science*, pages 316–333. Springer, 2011.
- [32] Siddharth Mohanty and Murray Cole. Autotuning wavefront abstractions for heterogeneous architectures. In *WAMCA*, pages 42–47. SBC, 2012.

- [33] Mattais Ohlsson, Carsten Peterson, and Bo Söderberg. Neural networks for optimization problems with inequality constraints: The knapsack problem. *Neural Comput.*, 5(2):331–339, 1993.
- [34] David A. Padua, Rudolf Eigenmann, Jay Hoeflinger, Paul Petersen, Peng Tu, Stephen Weatherford, and Keith Faigin. Polaris: A new-generation parallelizing compiler for MPPs. Technical report, In CSRD Rept. No. 1306. Univ. of Illinois at Urbana-Champaign, 1993.
- [35] Krishna Palem and Avinash Lingamneni. What to do about the end of moore’s law, probably! In *DAC*, pages 924–929. ACM, 2012.
- [36] Krishna V. Palem, Lakshmi N.B. Chakrapani, Zvi M. Kedem, Avinash Lingamneni, and Kirthi Krishna Muntimadugu. Sustaining moore’s law in embedded computing through probabilistic and approximate design: Retrospects and prospects. In *CASES*, pages 1–10. ACM, 2009.
- [37] Silvius Rus, Lawrence Rauchwerger, and Jay Hoeflinger. Hybrid analysis: Static & dynamic memory reference analysis. *International Journal of Parallel Programming*, 31(4):251–283, 2003.
- [38] Eric Schulte, Jonathan Dorn, Stephen Harding, Stephanie Forrest, and Westley Weimer. Post-compiler software optimization for reducing energy. In *ASPLOS*, pages 639–652. ACM, 2014.
- [39] Stelios Sidiroglou-Douskos, Sasa Misailovic, Henry Hoffmann, and Martin Rinard. Managing performance vs. accuracy trade-offs with loop perforation. In *ES-ECP/FSE*, pages 124–134. ACM SIGSOFT, 2011.
- [40] Kate A. Smith. Neural networks for combinatorial optimization: A review of more than a decade of research. *INFORMS Journal on Computing*, 11(1):15–34, 1999.
- [41] Mark Stephenson, Una-May O’Reilly, MartinC. Martin, and Saman Amarasinghe. Genetic programming applied to compiler heuristic optimization. In *Genetic Programming*, volume 2610 of *Lecture Notes in Computer Science*, pages 238–253. Springer, 2003.
- [42] Georgios Tournavitis, Zheng Wang, Björn Franke, and Michael F. P. O’Boyle. Towards a holistic approach to auto-parallelization: integrating profile-driven parallelism detection and machine-learning based mapping. In *PLDI*, pages 177–187. ACM SIGPLAN, 2009.
- [43] Paul Walsh and Conor Ryan. Automatic conversion of programs from serial to parallel using genetic programming - the paragen system. In *PARCO*, pages 415–422, 1995.
- [44] Paul Walsh and Conor Ryan. Paragen: A novel technique for the autoparallelisation of sequential programs using gp. In *Proceedings of the 1st Annual Conference on Genetic Programming*, pages 406–409, Cambridge, MA, USA, 1996. MIT Press.

- [45] Baoyun Wang, Heng Dong, and Zhenya He. A chaotic annealing neural network with gain sharpening and its application to the 0/1 knapsack problem. *Neural Processing Letters*, 9(3):243–247, 1999.
- [46] Zheng Wang and Michael F.P. O’Boyle. Mapping parallelism to multi-cores: A machine learning based approach. In *Proceedings of the 14th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP ’09, pages 75–84, New York, NY, USA, 2009. ACM.
- [47] Zheng Wang, Georgios Tournavitis, Björn Franke, and Michael F. P. O’boyle. Integrating profile-driven parallelism detection and machine-learning-based mapping. *ACM Trans. Archit. Code Optim.*, 11(1):2:1–2:26, 2014.
- [48] Kenneth P Williams and Shirley A Williams. Genetic compilers: A new technique for automatic parallelisation. *ESPPE*, pages 27–30, 1996.
- [49] Kenneth P. Williams and Shirley A. Williams. Two evolutionary representations for automatic parallelization. In *GECCO*, volume 2, pages 1429–1436, Orlando, Florida, USA, 1999. Morgan Kaufmann.
- [50] Peng Wu, Arun Kejariwal, and Călin Caşcaval. Compiler-driven dependence profiling to guide program parallelization. In *LCPC*, volume 5335 of *Lecture Notes in Computer Science*, pages 232–248. Springer Berlin Heidelberg, 2008.
- [51] Paraskevas Yiapanis, Demian Rosas-Ham, Gavin Brown, and Mikel Luján. Optimizing software runtime systems for speculative parallelization. *ACM Trans. Archit. Code Optim.*, 9(4):39:1–39:27, 2013.
- [52] Zeyuan Allen Zhu, Sasa Misailovic, Jonathan A. Kelner, and Martin Rinard. Randomized accuracy-aware program transformations for efficient approximate computations. In *POPL*, pages 441–454. ACM, 2012.