

POLITECNICO DI MILANO  
Scuola di Ingegneria Industriale e dell'Informazione  
Dipartimento di Elettronica, Informazione e Bioingegneria  
Master Degree In Computer Science and Engineering



**POLITECNICO**  
**MILANO 1863**

Thesis  
Title

Advisor: Prof. Giovanni AGOSTA

Thesis by:  
Pietro Ghiglio Matr. 920491

Academic Year 2019–2020



*To someone very special...*



# Acknowledgments

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.



# Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.





# Sommario

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.



# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Background</b>	<b>3</b>
1.1 LLVM . . . . .	3
1.1.1 LLVM-IR . . . . .	3
1.1.2 SSA and Phi nodes . . . . .	3
1.1.3 Class hierarchy . . . . .	4
1.1.4 LLVM Metadata . . . . .	4
1.1.5 LLVM Passes . . . . .	5
1.2 How LLVM handles debug information . . . . .	5
1.2.1 Metadata classes . . . . .	5
1.2.2 Transformation passes guidelines . . . . .	5
1.3 Program Instrumentation . . . . .	7
1.4 Debugging . . . . .	8
1.4.1 Debug information . . . . .	8
1.4.2 DWARF format . . . . .	9
<b>2 State of the art</b>	<b>11</b>
2.1 Overview . . . . .	11
2.2 Simulation . . . . .	12
2.3 Direct measuring . . . . .	14
2.4 Performance counters . . . . .	14
2.5 Instruction Level Energy modeling . . . . .	15
2.5.1 Characterization of an ISA Energy model . . . . .	16
2.5.2 Why employing an ISA energy model . . . . .	16
2.5.3 Producing an ISA energy model . . . . .	16
2.5.4 Extensions . . . . .	17
2.6 Source Code-level visualization . . . . .	21
2.7 Final Remarks . . . . .	21

<b>3</b>	<b>—</b>	<b>23</b>
3.1	Instrumentation . . . . .	23
3.2	Source level visualization . . . . .	23
3.3	Compiler optimizations . . . . .	23
<b>Conclusions</b>		<b>25</b>
<b>Bibliography</b>		<b>25</b>
<b>A</b>	<b>First appendix</b>	<b>31</b>
<b>B</b>	<b>Second appendix</b>	<b>33</b>
<b>C</b>	<b>Third appendix</b>	<b>35</b>
<b>Index</b>		<b>35</b>

# List of Figures

1.1	Example of optimization with merged debug location . . . . .	7
-----	--	---



# List of Tables





# List of Algorithms



# Introduction

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc

elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

# Chapter 1

## Background

### 1.1 LLVM

The LLVM Project [20] is a collection of modular and reusable compiler toolchain technologies. It is built around an intermediated representation called LLVM-IR, and provides a set of APIs to interact with it. LLVM provides an optimizer that works on the intermediate representation, and also several code generation helpers that allow to target all the main hardware architectures.

#### 1.1.1 LLVM-IR

The LLVM-IR is a language that resembles a generic assembly language, while also providing some high level features such as unlimited registers, explicit stack memory allocation and pointer deferentation. This allows LLVM-IR to be both the ideal target for high-level language developers, that do not have to worry about architecture specific details, and also the ideal source language for compiler back-end developers, that have to implement only a translator from LLVM-IR to their target architecture's assembly language, without worrying about high-level language features.

The LLVM-IR is accesible in three formats: in-memory represantation, that allows manipulation through the LLVM APIs, binary format, used by many LLVM tools, and the human-readable textual format, that can also very conveniently be parsed by means of the APIs.

#### 1.1.2 SSA and Phi nodes

The LLVM-IR is by definition in SSA (Static Single Assignment) form. The SSA form requires a variable to be assigned only once, and requires every variable to be defined before its uses. It is called static because it does not take into account dynamic (related to the program's runtime) considerations. For instance, an

assignment in a loop counts always as one assignment, even if at runtime it will be performed several times.

—esempio

It is always clear which definition to use, unless a basic block has multiple predecessors. In that case it is necessary to add phi nodes that carry the information to disambiguate the uses at runtime.

—esempio

### 1.1.3 Class hierarchy

The class hierarchy defined in the LLVM APIs consists of hundreds of classes, a complete and exhaustive view is given by the LLVM Doxygen Documentation. The main components of the hierarchy are:

- Module: the entire program/compile unit. Contains the global values of the program: mainly the global variables and the functions.
- Function: a function in the compile unit, contains mainly a set of arguments and its control flow graph in the form of a set of basic blocks.
- Basic Block: a set of instructions with no branches between them.
- Instruction: An instruction of the IR.

Another key class in the LLVM class hierarchy is the Value class. It represents anything that has a type and can be used as an operand to an instruction: function arguments, constants, instructions, basic blocks and functions are all Values. A Value also carries information of what other Values it uses, and what other Values use it.

### 1.1.4 LLVM Metadata

The LLVM-IR allows metadata to be attached to Instructions, Functions, Global Variables or Modules. Metadata can convey extra information about the code to the optimizers and code generator. The main use of metadata is debug information, but they may also carry information about loop boundaries or other assumption that are useful during the various stages of the compilation process.

Metadata can either be a simple string attached to an instruction, or they can be a Metadata Node (MDNode). MDNodes can reference each other and are specified by other classes in the LLVM APIs. See section 1.2 or the LLVM Language Reference [12] for more details.

### 1.1.5 LLVM Passes

LLVM passes are where most of the interesting parts of the compiler exist. Passes perform the transformations and optimizations that make up the compiler, they build the analysis results that are used by these transformations, and they are, above all, a structuring technique for compiler code.

Passes are categorized in two ways: by the granularity at which they operate, and by the fact that they perform changes on the module or not.

By the first categorization, passes are identified as:

- Module Passes: operate on an entire Module.
- Function Passes: operate on a single Function.
- Loop Passes: operate only on loops.
- Region Passes: operate on subsets of Basic Blocks of a Function, with a single entry point and a single exit point.

By the second categorization, passes are identified as:

- Analysis Passes: passes that only perform an analysis of the given entity, without modifying it.
- Transformation Passes: passes that may modify the given entity. They exploit the results of the Analysis Passes, often (but not only) in order to perform optimizations: they may add, remove, move or replace instructions and basic blocks, with the ultimate goal of improving performances or reduce the size of the binary.

Passes may depend on other passes, for instance a pass that performs an optimization may require the results of a pass that performs a specific analysis. They are therefore handled by a Pass Manager that schedules the passes, ensuring that all the dependencies for a pass are met before executing it.

## 1.2 How LLVM handles debug information

### 1.2.1 Metadata classes

### 1.2.2 Transformation passes guidelines

As we've seen in section 1.1.5, during the compilation a module may undergo some changes: instructions may be removed, moved, merged together, and replaced with new instructions, all in order to improve the performances of the resulting program.

These transformations have the side effect of obfuscating the correspondence between source code and binary code: before the optimization occurs, debug information provides a very clear, one-to-many relation between source location and LLVM-IR instructions. But as the module progresses into the optimization pipeline, it becomes more and more difficult to maintain this relation.

In general it is not possible to map unambiguously source locations to optimized code, but the LLVM project provides a set of guidelines that specify how to correctly update debug info when implementing transformation passes [7].

Here we provide a short summary of such guidelines<sup>1</sup>, highlighting some behaviors that, even when following them, lead to a loss of information regarding source-binary mapping. These behaviors are not bugs or mistakes of the people who provided the guidelines, but are instead related to the fact that they want to provide a debugging experience as close as possible to the one that a user would have while debugging the unoptimized code.

The guiding principles for a developer that wants to update debug info are the following:

1. Do not provide misleading information: a developer should not speculate, and providing no information is better than providing wrong information that may lead a developer to wrong considerations about the behavior of his program.
2. Provide as much information as possible: when it's not misleading, information should be preserved.

In order to achieve this, when choosing what to do with the debug information of a given instruction, a developer has three alternatives:

- Preserve the original location.
- Merge two locations: two debug locations can be merged together. Location merge is performed by computing the intersection of the two locations: the resulting location will contain only the information that the original two had in common.
- Delete the location.

Locations can be safely preserved when the modified instruction either remains in the same basic block, or its basic block is folded into a predecessor that branches unconditionally. For instance, an optimization that replaces the instruction `add x, x` with a binary shift to the left (`shl x, 1`) can safely keep the location of the original `add`.

Location should be merged when two instructions are replaced with a new instruction. An example of that is figure 1.1, in which the two stores can be merged

---

<sup>1</sup>Provided at a speech and the 2020 LLVM Conference by Adrian Pranti and Vedant Kumar



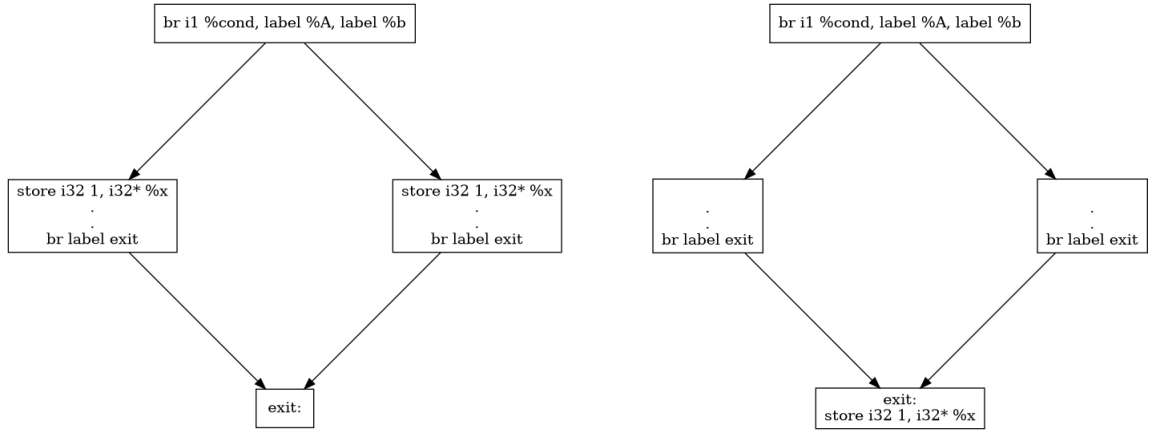


Figure 1.1: Example of optimization with merged debug location

into a new one, inserted in the exit basic block: the new instruction effectively replaces the original two, and therefore its location will be the merged location of the old ones.

In all the cases in which the previous rules do not apply, locations should be dropped. In particular, they should be dropped whenever an instruction is moved from a basic block with multiple predecessors, to one of the predecessors. This is done to avoid situations in which, while debugging, the program seems to have taken a branch in a conditional, while the actual conditions are not the one that would have resulted in the branch being taken.

Dropping locations and merging locations is a very reasonable course of action when dealing with debugging: they lead to a debugging experience that is as close as possible to the not optimized one. But they also lead to a loss of information in the source-binary mapping: when two locations are merged, we will most likely lose the information on the original source code lines, as they probably will not be equal, and when a location is dropped we will of course lose the information it carried.

We have therefore developed a methodology that allows to propagate debug locations through the optimization pipeline, while also bringing to the developers a view of the optimizations performed on their program, so that they can understand how it has been optimized by the compiler.

### 1.3 Program Instrumentation

Instrumenting a program means to insert additional code that was not originally in the program's source, typically in order to produce additional information (regarding some functional or non functional properties) during the program's runtime. It can be performed directly on the source code, on the executable binary, or during the compilation. An example of instrumentation are the many sanitizers that are part

of the LLVM project, they make runtime checks about memory and thread safety.

LLVM provides some helper classes to perform instrumentation on an IR Module, and in general a user may define his own transformation pass that inserts new code into the program being compiled.

Instrumentation often introduces a performance overhead, due of course to the fact that more instructions are executed while the program is running, so it is usually performed only during the development stage of an application.

## 1.4 Debugging

A debugger is a computer program used to test and debug other programs. It allows a programmer to run the target program in controlled conditions, pause the program's execution, check the state of variables and more.

### 1.4.1 Debug information

The main functionality of a debugger, over which more advanced features can be built, are setting break points and accessing the content of a variable defined in the source code.

This is achieved by means of debug information: information stored by the compiler in the program's executable, with the purpose of providing a correspondence between source level entities (variable, source code locations, data types) and low level entities (assembly instructions and memory locations).

The format used to store them may vary with the compiler/operating system used, but the stored information are mainly:

- Definition of the data types employed in the program and their layout in memory, both language-defined (eg. `int`, `float`, `unsigned` in C) or user defined (eg. C structs or C++ classes).
- Mapping between variables defined in the source code and memory locations in which they are stored. This allows a debugger to output the value of a variable given its name.
- Mapping between source code locations and assembly instruction. This allows the debugger to pause the program's execution when a given source code location is reached.

These information are useful not only for debugging purposes, they may also be employed by any other tool that requires a mapping between source code and binary executable, such as a profiler or a test coverage tool, that may be able to annotate the source code with the information that they have gathered.

### 1.4.2 DWARF format

The DWARF format [19] is a debugging file format used by many compilers and debuggers to support source-level debugging. It is designed to be extensible with respect of the source language, and to be architecture and operating system independent.

The main data structure used to store debug information is the DIE (Debug Information Entry). DIEs are used to describe both data types and variables, and can reference each other creating a tree structure.

Another data structure that is very useful for our purposes is the Line Number Table: it contains the mapping between memory addresses of the executable code, and the source line corresponding to those addresses.

Each row of the table contains the following fields:

- Address: the program counter value of a machine instruction.
- Line: the source line number.
- Column: the column number within the line.
- File: an integer that identifies the source file.
- Statement: boolean indicating if the current instruction is the beginning of a statement.
- Block: boolean indicating if the current instruction is the beginning of a basic block.

And other fields that are described in the DWARF documentation.



## Chapter 2

# State of the art

### 2.1 Overview

This chapter will provide an overview of the state of the art methods to measure, estimate and visualize the energy consumption of software.

In general, there is no unique solution to this problem: the proposed techniques vary by both the applicative domain, the properties of the result and the procedure with which the result is obtained.

For the applicative domains, we have identified three main cases:

- Embedded systems: embedded systems are often employed as sensors in contexts where the only source of electricity is their own battery. Therefore, energy consumption has always been a concern of both hardware and software developers.
- Smart phones: similarly to embedded systems, smart phones have to rely on their own battery to operate. The current rate of battery improvement is around 5% a year, but the workloads that smart phones have to withstand increases by an order of magnitude every 5 years [14]. This means that energy consumption has to be tackled also from the software perspective.
- Multicore CPU: energy consumption is not a big concern from the point of view of PC users. But it is a primary concern in large datacenters where heat dissipation requires good engineering solutions, and whose impact on global CO2 emission and energy consumption is non-negligible.

In terms of the properties of the result of the measurement/estimate, solutions differ by their granularity: some provide a single quantity (the total amount of energy consumed by a program), other have a finer grain, allowing to attribute energy measures to either source code entities or hardware components.

The procedure adopted to obtain the result are widely different, [15] provides an overview of some techniques, and groups them in simulation-based and measurement-based. Simulation based technique require a model of the target architecture, and, provided a segment of binary code, perform a cycle-accurate simulation of the events that occur during the program's run.

Measurement-based techniques, instead, can be further sub-grouped in roughly three categories:

- Direct measurement: the measure can be taken by plugging the device to an instrument that allows to measure the current or power absorbed by the device while running the program.
- Performance-counter based: some hardware architectures provide special register that store information about the energy consumption, and a set of APIs that allow to read their contents.
- Modeling based: some solutions propose to model mathematically the energy behavior of the target architecture, perform some experiments in order to estimate the parameters of the model, and use the model in order to obtain an estimate of the consumed energy.

The dimensions that we have indicated are not completely orthogonal. In particular, modeling based approaches usually target embedded devices, since they have simpler underlying architectures that are inherently easier to analyze.

Performance counter based method, instead, are bound to specific architectures that provide such counters, such as Intel's Running Average Power Limit (RAPL), for their Sandy Bridge architectures, the Intel System Management Controller (SMC) for the Xeon Phi, or NVidia's NVidia Management Library (NVML), that allows to obtain energy consumption of their GPUs.

## 2.2 Simulation

Simulation based methods are methods that provide an estimate of the energy consumption by running the assembly code of the program in an architectural simulator. Said simulator must also have been provided with some energy/power model of the target.

The first proposal for such a technique has been published by Tiwari et al. in 2000 [4]. The simulator that they developed, Wattch, is the first simulator to operate at the architectural level: it does not require the full RTL design (the Verilog of the target architecture), but relays instead on a more high-level description of the CPU.

Given such a description, that includes functional units, caches, register files, memories, TLB and other components, Wattch employs a parameterizable power

model that, through a cycle-accurate simulation, outputs an estimate of the consumed energy. The simulation is run by interfacing with the SimpleScalar [1] architectural simulator.

Being more high-level than RTL-based simulations, Wattch is faster and doesn't rely on the Verilog description of the target (usually not disclosed by companies), to the detriment of the accuracy of the result, since it does not model in full detail the entire logic of the target.

Despite being faster than RTL-based simulators, running a binary file with Wattch is several order of magnitudes (around 10000 slower, as reported in [2]) slower than running the actual program, even if the latter has been instrumented.

Since it's original publication, the original work on Wattch has been expanded in several ways. The main contribution in that sense has been given by Li et al., who developed a completely new power simulator, McPat [9], offering more modern and advanced features than Wattch.

It provides the support to compute power-area integrated metrics (energy-delay-area product), models static, dynamic, and short-circuit power dissipation (whereas Wattch only modeled dynamic power dissipation), allows to model multicore architectures, that have become increasingly widespread, and also provides an XML interface to the simulator, that allows McPat to be ported to different performance simulators.

Validating the correctness of the output of these tools is not an easy task: they provide a very fine grained output (the power/energy estimates for each of the CPU's sub-components), but hardware manufacturers often do not disclose design data with such a level of detail. In [22], the authors, that work for the IBM corporation, have access to such data, and therefore they can provide a more insightful validation of the estimates emitted by McPat. They conclude that, while the procedure employed by McPat to obtain such results is sound, the power models that it exposes are often incomplete, too high-level, or represent an implementation of the structure that differs from the core at hand. The authors provide also some guidelines to improve power modeling accuracy, but ultimately state that academic researches would greatly benefit from the availability of validated power models for contemporary commercial chips, emitted by the hardware producers themselves.

To conclude this section: simulation based power models are very interesting as they can characterize a hardware architecture with great detail, but their (slow) simulation speed, and some concerns regarding their accuracy, make them not practical for software developers: they are more suited to hardware/compiler developers that want to characterize the power/energy behavior of a target architecture, not to programmers that want to characterize the energy behavior of software.

## 2.3 Direct measuring

Directly measuring the current drawn by a device during the program's execution is the method that provides the greatest accuracy, but it's also the one that has the greatest "overhead" for a developer that wants to assess the energy consumption of his software.

Given its accuracy, it is often used as "ground truth" when evaluating the performances of other methods (simulations, performance counter or modeling).

Experimental setups may differ, depending on the target architecture and the tools at disposal, but they usually consist in a measuring point placed between the device and the power supply. For example in [17], they state that their experimental setup consist in a precision current-sense amplifier that amplifies the voltage drop across a shunt resistor, the output signal is then sampled by an Analog to Digital Converter, and sent to a PC.

Given the measured current,  $I$ , and the supply voltage  $V_{cc}$  the power drawn by the running target is given by  $P = I \times V_{cc}$ . The total energy consumed is given by  $E = P \times T$ , where  $T$  is the running time, which can be further decomposed in  $T = N \times \tau$ , where  $N$  is the number of clock cycles taken by the program, and  $\tau$  is the clock period [21].

In [5], instead, they employed a power meter located between the target's power socket and the A/C outlet, in order to establish the energy consumption of servers running Intel multicore CPUs.

As we said, the main drawback of directly measuring the energy consumption is the fact that a whole experimental setup is required, with appropriate tools that a software developer may not even have at his disposal. Another drawback of this approach is that it provides only a raw quantity (program X during this run consumed Y Joule), but a software developer may also desire some clues about which source-code entities lead to that energy consumption.

## 2.4 Performance counters

Performance counters are a feature of some hardware architectures. These architectures expose some registers that store information about the energy consumption of a program (among other metrics). In the following section we will give an overview of Intel's RAPL as an example of such a feature.

RAPL provides a set of counters providing energy and power information. It is not an analog power meter, but rather uses a software power model that estimates power consumption by means of performance counters and hardware power models [18]. The RAPL counters can be accessed by a user by reading appropriate files on the target machine.



A typical usage of RAPL is to read the energy measures, perform a task, and then read again the energy measures, taking the difference between the two readings as the estimate of the energy consumed by the task.

RAPL provides a fine-grained view of the energy consumption, with respect of the hardware components, offering separate estimates for each of the following domain:

- Package: the whole CPU.
- Core: the central components of the CPU, such as ALU, FPU and L1 and L2 cache.
- Uncore: components that are shared between cores, such as L3 cache and the memory controller.
- DRAM: the main memory.

A typical use case of RAPL is given in [11], in which the authors first provide a java library that allows to easily access the RAPL energy estimates, and then use said library to benchmark several data access and data organization patterns, providing some guidelines for application-level energy optimizations. They follow the typical usage pattern of measurement  $\rightarrow$  task  $\rightarrow$  measurement.

Regarding RAPL's accuracy, there are some discording opinions: in [16], where RAPL is introduced to the general public, they state that the prediction provided by their power model matches actual measurements, showing high correlation between the two. This is partially disproved in [5], where the authors compare results obtained with RAPL to results obtained via direct measurement. They show that the average error between RAPL's estimate and the actual energy consumption ranges from 8% to 73%. The discrepancy between the estimate and the ground truth is also non-constant: it varies greatly depending on both the performed task and the configuration of the machine. They even show that optimizing an application using data collected from RAPL as benchmark leads to an effective increase in the total (directly measured) consumed energy.

To summarize this section, RAPL offers a very interesting set of features: provides a fine granularity in terms of hardware components, can be accessed from source code, allowing to profile arbitrary regions of code, and has very little overhead for the programmer (he just has to add the calls to RAPL where he is interested), but it suffers of accuracy problems. This means that, at least, it is not a good candidate to be used as ground truth to validate other estimation methods.

## 2.5 Instruction Level Energy modeling

Given a target's Instruction Set Architecture (ISA), an energy model is a model of the energy consumed by each instruction. They have been introduced in 1996 by

Tiwari et al. [21].

### 2.5.1 Characterization of an ISA Energy model

The main components of an energy model are:

- Instruction base cost ( $B_i$ , for each instruction  $i$ ): the cost associated with the basic processing needed to execute an instruction.
- Effect of circuit state ( $O_{i,j}$ , for each pair of instruction  $i, j$ ): the cost of the switching activity resulting from executing two consecutive instructions differing one from another.
- Other inter-instruction effects ( $E_k$ , for each additional effect  $k$ ): any other effect that can occur in real program, such as stalls or cache misses.

Given these components and a program  $P$ , the total energy consumed by it,  $E_p$ , is given by:

$$E_p = \sum_i (B_i \times N_i) + \sum_{i,j} (O_{i,j} \times N_{i,j}) + \sum_k E_k$$

Where  $N_i$  is the number of occurrences of instruction  $i$ , and  $N_{i,j}$  is the number of times there has been a switch from instruction  $i$  to instruction  $j$ .

### 2.5.2 Why employing an ISA energy model

The most common way to describe a processor's power consumption is through the average power consumption.

This single number may not provide enough information to characterize the energy consumed by a program running on the target processor: different programs may employ the functional units of the CPU in different ways, leading to different measurements at equal running time.

ISA Energy Models offer a more detailed view of the energy profile of the target architecture. They therefore allow to identify variations of consumed energy from one program to another, and may also guide decision of both humans (hardware/software design) and software (compilers or operating systems).

### 2.5.3 Producing an ISA energy model

Energy models can be produced through an experimental procedure.

In order to obtain instruction base costs, a program consisting of a large loop of a repeated instruction is written. Then one can measure the average current drawn by the processor while executing the program,  $\hat{i}$ , and multiply it by the supply voltage  $V_{cc}$ , obtaining the base energy consumption.

Instruction may also be grouped together, since instruction with similar functionality will have similar base cost.

In order to obtain the circuit state effects, loop of pairs of instruction are required. The difference between the instruction's base costs and the average current measured provides the circuit state overhead.

A similar approach can be employed to obtain the costs of other inter-instruction effects: writing large loops in which the examined effect occurs several times, measuring the average current and subtracting the costs that are already known (base costs and circuit state).

The main disadvantage of this approach is that several different programs must be written: for an ISA with  $n$  instructions,  $\mathcal{O}(n)$  programs are required to produce base costs and  $\mathcal{O}(n^2)$  for circuit state effects.

Estimation of other inter-instruction effects also gets more difficult as the complexity of the architecture increases.

On the other hand, this approach has the big advantage of not requiring a model of the circuit of the target processor, information that is often not disclosed by the manufacturing companies.

In [13], the same authors of [21] employ their technique to model the instruction level energy consumption of a Digital Signal Processing (DSP) embedded system. They describe their experimental setup, consisting in a standard, off-the-shelf, dual-slope integrating digital ammeter connected between the power supply and the pins of the DSP chip. They exploit this power model in order to design a scheduling algorithm that minimizes the total energy consumed.

In this work, they also highlight some practical issues regarding the methodology employed to construct the energy model: impact of operand values and tables size. For the impact of operand values, they propose to make the measurement using a wide a range of operand values, and averaging the consumption values.

By table size, instead, they mean that the number of experiments that need to be carried out to model all the instruction can be overwhelming, and so they propose to group instructions by similar functionality, assigning an average cost to each group, thus reducing the number of needed experiments. The variation of the energy cost of instructions grouped in this way is around 5%.

#### 2.5.4 Extensions

The original work of Tiwari et al. has been extended in several ways through the years, both in terms of the complexity of the model, and in terms of how said model has been exploited.

In [8], Eder et al. characterize the energy behavior of a multithreaded architecture. They target processor is the XMOS XS1-L. In this processors threads are

executing in a round robin fashion, this makes program execution time-deterministic and allows to easily model the multithreaded behavior.

The key difference between their model and Tiwari's model is that they take into account the energy cost of context switching. Therefore, according to their model, the energy consumption of a program,  $E_p$ , is given by:

$$E_p = P_{base}N_{idle}T_{clk} + \sum_{i=1}^{N_t} \sum_{i \in ISA} ((M_t P_i O + P_{base}) N_{i,t} T_{clk})$$

Where  $T_{clk}$  is the clock period,  $P_{base}$  is dissipated power when the processor is idle,  $N_{idle}$  is the number of clock cycles in which the processor was idle,  $N_t$  is the maximum number of running threads,  $M_t$  is a multiplier that accounts for the round-robin execution in the pipeline, and, for each instruction  $i$ ,  $P_i$  is the dissipated power and  $N_{i,t}$  is the amount of times instruction the instruction has been executed in thread  $t$ , finally,  $O$  is the inter-instruction overhead, that they assume to be constant.

In order to perform their experiments, they have designed a software suite that allows to automatically generate benchmarks used to characterize the energy model, loading them on the target and monitor their execution. Tests are generated only for instructions with no effects on control flow and no non-deterministic timing. They obtain execution statistics by hardware simulation (this can be replaced by profiling). They observed that the number of operands has significant impact on power consumption, while data width has an also an impact on power consumption, but way lower. They also choose to generalize inter-instruction overhead, observing that it exhibits little variance between different couples of instructions.

For instructions that cannot be directly tested, they propose two solutions: either group instructions by number of operands, and put the untestable instr in the appropriate group or assign default cost to untestable instr.

They conclude by stating that the model in which instructions are grouped by number of operands performs worse than the model in which instruction are considered individually: the first one exhibits an average error of 16%, the latter 7%). Both the models provide a consistent underestimation.

lee: energy model of risc architectures, targeting embedded systems. combining empirical method and statistical analysis. developed a model whose unknown are estimated tanks to data from empirical observations, through linear regression. they say that modeling in tiwari's way is too simplistic, since it relays only on average

current. First they model the energy consumption: pipelined processor,  $e(X,Y)$  is the energy when  $X$  is executed after  $Y$  (this allows to consider switching activity). it depends on several variables ( $v$ ), they consider  $f(v_x,v_y)$ : the variation on the variable  $v$  between  $x$  and  $y$ , which depends on the hamming distance between  $v_x$  and  $v_y$ . they test several sets of programs: in the first set, they test same instructions, same operands different location and determine the impact of the instruction fetch stage. in the second set, they derive instruction base cost in each of the pipeline stages.

nunez-yanez: system level energy/power modeling (on chip system: memory, cpu, gpu; no camera/display). they require RTL design information (often not available), and then through linear regression they produce a power model of the whole SoC. They provide high detail, being able to characterize power consumption of different components (cpu, ram, cache), but the requirement of the RTL design is too strict.

[2]: the first effort that we have found in both estimating energy consumption and mapping it to source code level entities. The analysis is performed at level of the parse-tree. The parse tree of a C program is decorated by associating a cost-contribution (atom) to each node. The authors have also introduced *kernel instructions* as a form of target independent assembly instructions, to which energy costs can be assigned. Energy cost are estimated through least square fitting, obtained by comparing to the output of the ARMulator instruction set simulator. Following the grammar's rule of the C language, rules to combine instruction costs are defined. The parse tree is then instrumented in order to produce a trace during the program's execution.

The final cost is obtained by combining the costs of the of the atoms and the data from the output of the instrumentation, providing a view that maps to each node in the parse tree (which corresponds to a source level entity such as an operation, a function call or an assignment) its contribution to the overall energy cost.

This approach relays on analyzing the parse tree. this allows source code visualization but binds to the source language: in order to change source language, it would be required to perform a complete analysis of the grammar rule of the new

language. Also, source level entities related to the grammar of a language may not have a trivial correspondence to assembly instructions, which makes changing source language even harder.

[3]: they propose a methodology, based on LLVM-IR analysis, to estimate a program's energy consumption, and to propagate their analysis at source code level. They provide a technique to understand how LLVM-IR instruction are related to the target's assembly instruction, then, by means of an instrumentation that outputs a trace of the basic blocks executed during a run of a program, they gather data regarding the dynamic behavior of the program. Finally, given a vector of energy costs for each assembly instruction, they are able to characterize the energy cost of a program by summing the energy costs of the executed basic blocks, while the energy cost of each basic block is obtained by summing the energy costs of all the assembly instructions corresponding to LLVM-IR instructions contained in the basic block.

Their technique to understand the LLVM-IR to assembly mapping is based on statistical analysis. Given a dataset of several LLVM-IR programs, they compile them to obtain assembly programs for the target architecture. Then they establish correlation between LLVM-IR instructions and assembly instructions, formulation a *non-negative least square* problem. This analysis has to be performed every time one wants to target a different architecture.

They also require the energy cost of each assembly instruction to be known. They acknowledge the fact the this information is often not disclosed by the manufacturers, and state that it may be approximated by a linear function of the clock cycles, since the current absorbed by each clock cycle exhibits very little variance. This a very interesting claim since data about the clock cycles taken by each instruction is often available, and it may allow a very easy way to provide an estimate of the energy cost of each assembly instruction.

[6]: They provide a LLVM-IR to assembly mapping technique that differs from [3], based on debug information and disassembly of the binary. Given this mapping, they provide a methodology to statically estimate the worst case energy consumption

(WCEC) of a program, both at LLVM-IR level or the assembly level, and they also employ a profiling technique similar to [3] in order to obtain an energy estimation given a run of the program.

Their LLVM-IR to assembly mapping technique is based on an LLVM pass that replaces the line number information, contained in the LLVM Debug Info classes, with an unique id of the instruction being considered. Then, after that the modified LLVM module has been compiled, by disassembling the binary and parsing the line table, for each entry in the line table, there will a pair  $\langle addr, id \rangle$ , such that the assembly instruction at the address  $addr$ , can be mapped to the LLVM instruction with identifier  $id$ .

This procedure by itself does not suffice in providing a complete mapping: for efficiency reasons, some assembly instruction do not have a corresponding entry in the line table, but they can be safely mapped to the same LLVM instruction of the last previous assembly instruction with an entry in the line table.

Their statical, worst-case energy consumption estimation is based on the Implicit Path Enumeration Technique (IPET) [10]. It requires the program's CFG, annotated with information regarding properties of the dynamic behavior of the program, such as loop bounds or mutually exclusive conditions. Given these annotations, that must be specified by the user, the problem can be formulated as an Integer Linear Programming problem, whose cost function is  $\sum_{i=0}^N c_i \times x_i$ , where, for each basic block  $i$ ,  $c_i$  indicates the cost of executing the basic block, and  $x_i$  indicates the number of executions of the basic block. Their profiling technique, instead, consist in an instrumentation that outputs the identifier of the basic block every time the basic block is run.

rieger-survey: alcuni tool commerciali, altra letteratura. molto android/java (ex. powertutor) roth: piattaforma hardware per creare energy models + xml dell'energy model per integrarlo nel loro compilatore, il quale fornisce una stima worst case. pereira: source level view, statistical method to provide ranking, no energy estimation, just visualization.

## 2.6 Source Code-level visualization

## 2.7 Final Remarks





# Chapter 3

---

## 3.1 Instrumentation

## 3.2 Source level visualization

## 3.3 Compiler optimizations



# Conclusions



# Bibliography

- [1] T. Austin, E. Larson, and D. Ernst. “SimpleScalar: an infrastructure for computer system modeling”. In: *Computer* 35.2 (2002), pp. 59–67. DOI: 10.1109/2.982917.
- [2] C. Brandolese. “Source-Level Estimation of Energy Consumption and Execution Time of Embedded Software”. In: *2008 11th EUROMICRO Conference on Digital System Design Architectures, Methods and Tools*. 2008, pp. 115–123. DOI: 10.1109/DSD.2008.43.
- [3] C. Brandolese, S. Corbetta, and W. Fornaciari. “Software energy estimation based on statistical characterization of intermediate compilation code”. In: *IEEE/ACM International Symposium on Low Power Electronics and Design*. 2011, pp. 333–338. DOI: 10.1109/ISLPED.2011.5993659.
- [4] David Brooks, Vivek Tiwari, and Margaret Martonosi. “Wattch: A Framework for Architectural-Level Power Analysis and Optimizations”. In: *Proceedings of the 27th Annual International Symposium on Computer Architecture*. ISCA '00. Vancouver, British Columbia, Canada: Association for Computing Machinery, 2000, pp. 83–94. ISBN: 1581132328. DOI: 10.1145/339647.339657. URL: <https://doi.org/10.1145/339647.339657>.
- [5] Muhammad Fahad et al. “A Comparative Study of Methods for Measurement of Energy of Computing”. In: *Energies* 12 (June 2019). DOI: 10.3390/en12112204.
- [6] Kyriakos Georgiou et al. “Energy Transparency for Deeply Embedded Programs”. In: *ACM Trans. Archit. Code Optim.* 14.1 (Mar. 2017). ISSN: 1544-3566. DOI: 10.1145/3046679. URL: <https://doi.org/10.1145/3046679>.
- [7] *How to Update Debug Info: A Guide for LLVM Pass Authors*. URL: <http://www.llvm.org/docs/HowToUpdateDebugInfo.html>.
- [8] Steve Kerrison and Kerstin Eder. “Energy Modeling of Software for a Hardware Multithreaded Embedded Microprocessor”. In: *ACM Trans. Embed. Comput. Syst.* 14.3 (Apr. 2015). ISSN: 1539-9087. DOI: 10.1145/2700104. URL: <https://doi.org/10.1145/2700104>.

- [9] Sheng Li et al. “McPAT: An Integrated Power, Area, and Timing Modeling Framework for Multicore and Manycore Architectures”. In: *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*. MICRO 42. New York, New York: Association for Computing Machinery, 2009, pp. 469–480. ISBN: 9781605587981. DOI: 10.1145/1669112.1669172. URL: <https://doi.org/10.1145/1669112.1669172>.
- [10] Yau-Tsun Steven Li and Sharad Malik. “Performance Analysis of Embedded Software Using Implicit Path Enumeration”. In: *Proceedings of the ACM SIGPLAN 1995 Workshop on Languages, Compilers and Tools for Real-Time Systems*. New York, NY, USA: Association for Computing Machinery, 1995. ISBN: 9781450373081. DOI: 10.1145/216636.216666. URL: <https://doi.org/10.1145/216636.216666>.
- [11] Kenan Liu, Gustavo Pinto, and Yu Liu. “Data-Oriented Characterization of Application-Level Energy Optimization”. In: Apr. 2015. DOI: 10.1007/978-3-662-46675-9\_21.
- [12] *LLVM Language Reference Manual*. URL: <https://llvm.org/docs/LangRef.html>.
- [13] Mike Tien-Chien Lee et al. “Power analysis and minimization techniques for embedded DSP software”. In: *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 5.1 (1997), pp. 123–135. DOI: 10.1109/92.555992.
- [14] Jose Nunez-Yanez and Geza Lore. “Enabling accurate modeling of power and energy consumption in an ARM-based System-on-Chip”. In: *Microprocessors and Microsystems* 37.3 (2013), pp. 319–332. ISSN: 0141-9331. DOI: <https://doi.org/10.1016/j.micpro.2012.12.004>. URL: <http://www.sciencedirect.com/science/article/pii/S0141933113000021>.
- [15] Felix Rieger and Christoph Bockisch. “Survey of Approaches for Assessing Software Energy Consumption”. In: CoCoS 2017. Vancouver, BC, Canada: Association for Computing Machinery, 2017. ISBN: 9781450355216. DOI: 10.1145/3141842.3141846. URL: <https://doi.org/10.1145/3141842.3141846>.
- [16] E. Rotem et al. “Power-Management Architecture of the Intel Microarchitecture Code-Named Sandy Bridge”. In: *IEEE Micro* 32.2 (2012), pp. 20–27. DOI: 10.1109/MM.2012.12.
- [17] Mikko Roth, Arno Luppold, and Heiko Falk. “Measuring and Modeling Energy Consumption of Embedded Systems for Optimizing Compilers”. In: *Proceedings of the 21st International Workshop on Software and Compilers for Embedded Systems*. SCOPES ’18. Sankt Goar, Germany: Association for Computing Machinery, 2018, pp. 86–89. ISBN: 9781450357807. DOI: 10.1145/3207719.3207729. URL: <https://doi.org/10.1145/3207719.3207729>.

- [18] *Running Average Power Limit*. URL: <https://01.org/blogs/2014/running-average-power-limit-%E2%80%93-rapl>.
- [19] *The DWARF Debugging Standard*. URL: <http://dwarfstd.org/>.
- [20] *The LLVM Compiler Infrastructure*. URL: <https://llvm.org/>.
- [21] V Tiwari et al. “Instruction level power analysis and optimization of software”. In: vol. 13. Feb. 1996, pp. 326–328. ISBN: 0-8186-7228-5. DOI: 10.1109/ICVD.1996.489624.
- [22] S. L. Xi et al. “Quantifying sources of error in McPAT and potential impacts on architectural studies”. In: *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*. 2015, pp. 577–589. DOI: 10.1109/HPCA.2015.7056064.





## Appendix A

### First appendix



## Appendix B

### Second appendix



## Appendix C

### Third appendix