



deGoal a tool to embed dynamic code generators into applications

Henri-Pierre Charles, Damien Couroussé, Victor Lomüller, Fernando A. Endo, Rémy Gauguey

► To cite this version:

Henri-Pierre Charles, Damien Couroussé, Victor Lomüller, Fernando A. Endo, Rémy Gauguey. deGoal a tool to embed dynamic code generators into applications. 23rd International Conference on Compiler Construction (CC), Apr 2014, Grenoble, France. springer, pp.107-112, 2014, <10.1007/978-3-642-54807-9_6>. <hal-01068248>

HAL Id: hal-01068248

<https://hal.archives-ouvertes.fr/hal-01068248>

Submitted on 25 Sep 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

deGoal a tool to embed dynamic code generators into applications

Henri-Pierre Charles Damien Couroussé Victor Lomüller
Fernando A. Endo Rémy Gauguey

September 25, 2014

Abstract

The processing applications that are now being used in mobile and embedded platforms require at the same time a fair amount of processing power and a high level of flexibility, due to the nature of the data to process. In this context we propose a lightweight code generation technique that is able to perform data dependent optimizations at run-time for processing kernels.

In this paper we present the motivations and how to use **deGoal** a tool designed to build fast and portable binary code generators called *compilettes*.

1 Introduction

Today, software development is facing two competing objectives:

- Improve programmers efficiency by using generic and expressive programming languages
- Generate efficient machine code to achieve the best execution speed and energy efficiency

These two objectives are competing because the more expressive and abstract a programming language is, the more difficult it is for a code generation tool-chain to produce efficient machine code.

We believe that code optimization, driven by run-time data characteristics, is a promising solution to tackle this issue. To achieve this, we propose **deGoal** a lightweight runtime solution for code generation.

deGoal was designed with embedded systems computing power and memory limitations in mind: our bottom-up approach allows code generation of specialized kernels, at runtime, depending on the execution context, the features of the targeted processor, and furthermore on the *data to process*: their characteristics and their values. Runtime code generation is achieved thanks to tiny *ad hoc* code generators, called *compilettes*, which will be embedded into the application

and produce machine code at runtime. *Compilettes* have only the strict necessary processing intelligence to perform the required code optimizations. As a consequence, code generation is:

1. *very fast*: 10 to 100 times faster than typical JITs or dynamic compilers which allow to use code generation inside the application and during the code execution, not in a virtual machine.
2. *lightweight*: the typical size of *compilettes* is only a few kilobytes which allows its use on constrained memory micro controllers such as the Texas Instrument MSP430 which has only 512 bytes of available memory.
Standard code generators, such as LLC of the LLVM infrastructure, have Mbytes of memory footprint, making their use impossible in this context.
3. *produce compact code*: as we are able to generate only the needed specialized code and not all variants at the same time.
4. *portable* across processor family: i.e. a *compilette* is portable on RISC platforms or on GPU platforms.
5. able to perform *cross-jit* applications, i.e. a *compilette* can run on one processor model and generate code for an other processor and download the generated code.

2 Introduction to the deGoal Infrastructure

The `deGoal` infrastructure integrates a language for kernel description and a small run-time environment for code generation. The tools used by the infrastructure are architecture agnostic, they only require a python interpreter and an ANSI C compiler. We briefly introduce the language, how applications are written with `deGoal` and how they're compiled and executed.

2.1 Kernel Description

We use a dedicated language to describe the kernel generation at runtime. This language is mixed with C code, this latter allowing to control the code generation performed in *compilettes*. This combination of C and `deGoal` code allows to efficiently design a code generator able to:

1. inject immediate values into the code being generated,
2. specialize code according to runtime data, e.g. selecting instructions,
3. perform classical compilation optimizations such as loop unrolling or dead code elimination.

`deGoal` uses a pseudo-assembly language whose instructions are similar to a neutral RISC-like instruction set. The goal is to achieve:

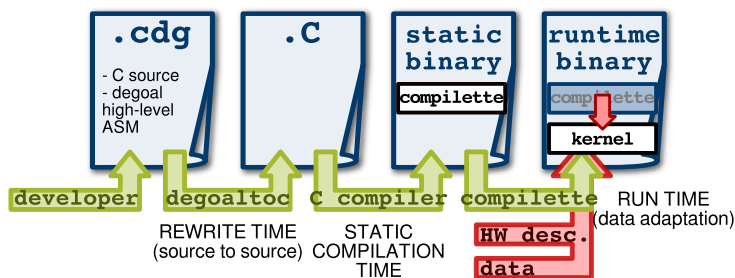


Figure 1: deGoal work flow: from the writing of applications source code to the execution of a kernel generated at run-time

Architecture	Port status	SIMD support	Instruction bundling
ARM Thumb-2 (+NEON/VFP)	✓	✓	N/A
STxP70 (STHORM)	✓	N/A	✓
PTX (NVIDIA)	✓	✓	N/A
ARM32	✓	✗	N/A
MSP430	✓	N/A	N/A
K1	✓	✓	✓
MIPS	Ⓒ	✗	N/A

Table 1: deGoal support status

- A rich instruction set focused on vectorial and multimedia instructions.
- The capability to use the run-time information in the specialized code.
- Cross-platform code generation: the architecture targeted by the *complette* may also be different from the architecture on which the code generator runs.
- Fast code generation, thanks to the “multiple time” compilation scheme. The intermediate representation (IR) is processed at static compile time. At run time the application has only to generate binary code mixed with data.

2.2 Compilation Chain

To achieve fast code generation, the compilation chain is split into several steps that run at different “times” (Figure 1).

Static compilation The *compilette* is rewritten into a standard C file by the `degoaltoc` tool. The C version of the *compilette* is then statically compiled and linked to the targeted architecture `deGoal` back-end using the C compiler of the target platform.

Runtime The application first invokes the *compilette* to generate the machine code of the kernel, once the optimizing data from the execution context are available. The kernel can then be executed as a standard procedure.

Given that the back-end is composed of portable C functions, our compilation chain is able to generate cross-platform code.

2.3 Run-time

At runtime, the *compilette* generates code according to run-time data and environment (rightmost block on Figure 1). At this time, registers are allocated, instructions scheduled and bundled (for VLIW architectures).

3 Current Status

Table 1 details the current support status of `deGoal` (MIPS is a work in progress). The column “SIMD support” shows the ability to take advantage of hardware vectors efficiently. The last column indicates if `deGoal` is able to generate code for VLIW processors.

The basic infrastructure is licensed under a BSD style license but all specific developments are restricted to their respective owners.

4 Related works

There is an extensive amount of literature about approaches related to our work with `deGoal`. Other works are related :

Java JIT mix interpretation and dynamic compilation for hotspots. Such techniques usually require large memory to embed JIT framework, and performance overhead. Some research works have tried to tackle these limitations: memory footprint can be reduced to a few hundreds of KB [6], but the binary code produced often presents a lower performance because of the smaller amount of optimizing intelligence embedded in the JIT compiler [12].

Java JITs are unable to directly take data value as parameters. They use indirect hotspot detection by tracing the application activity at runtime.

In `deGoal`, the objective is to reduce the cost incurred by runtime code generation. Our approach allows to generate code at least 10 times faster than traditional JITs: JITs hardly go below 1000 cycles per instruction

generated while we obtain 25 to 80 cycles per instruction generated on the STxP70 processor.

LLVM [9] (Low Level Virtual Machine) is a compilation framework that can target many architectures, including x86, ARM or PTX. One of its advantages is the unified internal representation (LLVM IR) that encode a virtual low-level instruction with some high-level information embedded on it. Various tools were built on top of it.

In **deGoal**, we don't use IR at run-time, we keep only calls (with parameters) to binary code generators.

Partial evaluation Our approach is similar to partial evaluation techniques [4, 8], which consists in pre-computing during the static compilation passes the maximum of the generated code to reduce the run-time overhead. At run-time, the finalization of the machine code consists in: selecting code templates, filling pre-compiled binary code with data values and jump addresses.

Using **deGoal** we compile statically an *ad hoc* code generator (the *complette*) for each kernel to specialize. The originality of our approach relies in the possibility to perform run-time instruction selection depending on the *data* to process [2].

DyC [7] is a tool that creates code generators from an annotated C code. Like ``C`, it adds some tokens such as `@` to evaluate C expressions and inject the results as an immediate value into the machine code.

deGoal is different from DyC because the parameters given to the binary run-time generators can drive specialized optimization such as loop-unrolling or vectorizers.

5 Application domain examples

As examples, here are some references of work in different application domains where *complettes* have been used:

Specialized memory allocator : memory allocators are specialized depending on the data size they can handle. *Lhuilier et al* [10] built an example with a very low memory footprint, able to adapt itself to the size of data set.

Hardware support thanks to the fast generation code scheme we are able to generate specialized code which run faster with a low overhead. We have used this support in

Mono-core specialization in an MPSoC context where each node is able to generate an optimized version of a matrix multiplication function [5].

GPU code specialization on an NVIDIA GPU we have developed a “cross-JIT” approach where a CPU generate a specialized GPU code depending on data sets [3].

Microcontroller with no floating point support, we are able to generate specialized floating point routine on the fly[1]

Video compression need specialized code depending on data sets as shown in [11]. `deGoal` can be used in this domain.

Thanks to the low memory footprint of both code generator and generated code, `deGoal` is perfect for embedded systems.

This article is only a functional tutorial, results with discussions about acceleration and produced code size can be found in the following bibliography.

References

- [1] Aracil, C., Couroussé, D.: Software acceleration of floating-point multiplication using runtime code generation. In: Proceedings of the 4th International Conference on Energy Aware Computing. Istanbul, Turkey (Dec 2013)
- [2] Charles, H.P.: Basic infrastructure for dynamic code generation. In: Charles, H.P., Clauss, P., Ptrot, F. (eds.) Workshop ”Dynamic Compilation Everywhere”, in conjunction with the 7th HiPEAC conference. Paris, France (january 2012), <https://sites.google.com/site/dynamiccompilationeverywhere/home/accepted-papers>
- [3] Charles, H.P., Lomüller, V.: GPU design pattern, chap. Data Size and Data Type Dynamic GPU Code Generation. Saxe-Coburg publications (2012)
- [4] Consel, C.: A general approach for run-time specialization and its application to c. pp. 145–156. ACM Press (1996)
- [5] Couroussé, D., Lomüller, V., Charles, H.P.: Introduction to dynamic code generation – an experiment with matrix multiplication for the sthorm platform. In: Smart Multicore Embedded Systems, chap. 6, pp. 103–124. Springer Verlag (2013)
- [6] Gal, A., Eich, B., Shaver, M., Anderson, D., Mandelin, D., Haghghat, M.R., Kaplan, B., Hoare, G., Zbarsky, B., Orendorff, J., Ruderman, J., Smith, E.W., Reitmaier, R., Bebenita, M., Chang, M., Franz, M.: Trace-based just-in-time type specialization for dynamic languages. In: Proceedings of the 2009 ACM SIGPLAN conference on Programming language design and implementation. pp. 465–478. PLDI ’09, ACM, New York, NY, USA (2009), <http://doi.acm.org/10.1145/1542476.1542528>

- [7] Grant, B., Mock, M., Philipose, M., Chambers, C., Eggers, S.J.: DyC: an expressive annotation-directed dynamic compiler for C. *Theor. Comput. Sci.* 248(1-2), 147–199 (Oct 2000), [http://dx.doi.org/10.1016/S0304-3975\(00\)00051-7](http://dx.doi.org/10.1016/S0304-3975(00)00051-7)
- [8] Jones, N.D.: An introduction to partial evaluation. *ACM Comput. Surv.* 28, 480–503 (September 1996), <http://doi.acm.org/10.1145/243439.243447>
- [9] Lattner, C.: LLVM: An Infrastructure for Multi-Stage Optimization. Master’s thesis, Computer Science Dept., University of Illinois at Urbana-Champaign, Urbana, IL (2002)
- [10] LHUILLIER, Y., COUROUSSÉ, D.: Embedded system memory allocator optimization using dynamic code generation. In: *Inproceedings of the 1st workshop on "Dynamic compilation everywhere" DCE2012* (January 2012)
- [11] Sajjad, K., Tran, S.M., Barthou, D., Charles, H.P., Preda, M.: A global approach for mpeg-4 avc encoder optimization. *14th Workshop on Compilers for Parallel Computing* (2009)
- [12] Shaylor, N.: A just-in-time compiler for memory-constrained low-power devices. In: *Java VM’02*. pp. 119–126. USENIX Association, Berkeley, CA, USA (2002)