

# Code-Copying Compilation in Production

An Experience Report

M. ANTON ERTL, TU Wien, Austria

BERND PAYSAN, net2o, Germany

A code-copying compiler implements a programming language by concatenating code snippets produced by a different compiler. This technique has been used in Gforth since 2003, with code snippets generated by GCC. We have solved various challenges: in particular, which code snippets can be copied and what to do about the others; and challenges posed by changes in compilers. The performance of Gforth is similar to that of SwiftForth, a commercial system with a conventional compiler; the implementation effort is comparable to 1–2 targets for SwiftForth.

## 1 INTRODUCTION

Code copying is a programming language implementation technique where the compiler of the implemented language A concatenates code snippets coming out of the compiler for language B. While there have been a number of research papers about this topic (see Section 8), we know of only one production language implementation that has used this approach for a long time: Gforth.

The present work is an experience report about the use of code copying in Gforth: How does it compare to a conventional compiler (Section 2)? Section 3 explains the concepts of code copying, while Section 4 discusses various implementation aspects. We also discuss the problems from changes in compilers (Section 5) and operating systems (Section 6) and how we overcame them.

In addition to this experience report, this paper also discusses alternative approaches (Section 7) and related work (Section 8).

The present work also appears in the EuroForth 2025 proceedings, with the same content and different formatting.

### 1.1 Is Gforth a production system?

Gforth is free software that has been developed since 1992 and first released in 1996. As it is free software, everybody can use it without contacting us, and few people do, so we do not know that much about who uses it for what purpose. However, we know that it has been used by IBM and Apple in their work on Open Firmware, and Forth, Inc. (who develop SwiftForth, but also give Forth courses) have given courses using Gforth, also in the Open Firmware context. So: Yes, Gforth is a production system.

## 2 WHY NOT JUST WRITE A CONVENTIONAL COMPILER?

One reason why people may have avoided going for a code-copying compiler is the assumption that writing a conventional compiler will produce better code, or require less effort. By “conventional” we mean that there is a large amount of hand-written architecture-specific code for each target architecture in the compiler. So before we go into details about code copying, we will address this concern.

### 2.1 Performance

Figure 1 shows the performance of the `gforth-fast` engine of Gforth<sup>1</sup> with various optimizations, of two commercial conventional Forth compilers (SwiftForth and VFX Forth), and, for of GCC-12.2

---

<sup>1</sup>Gforth also has an engine `gforth` intended for debugging. All references to Gforth performance refer to `gforth-fast`.

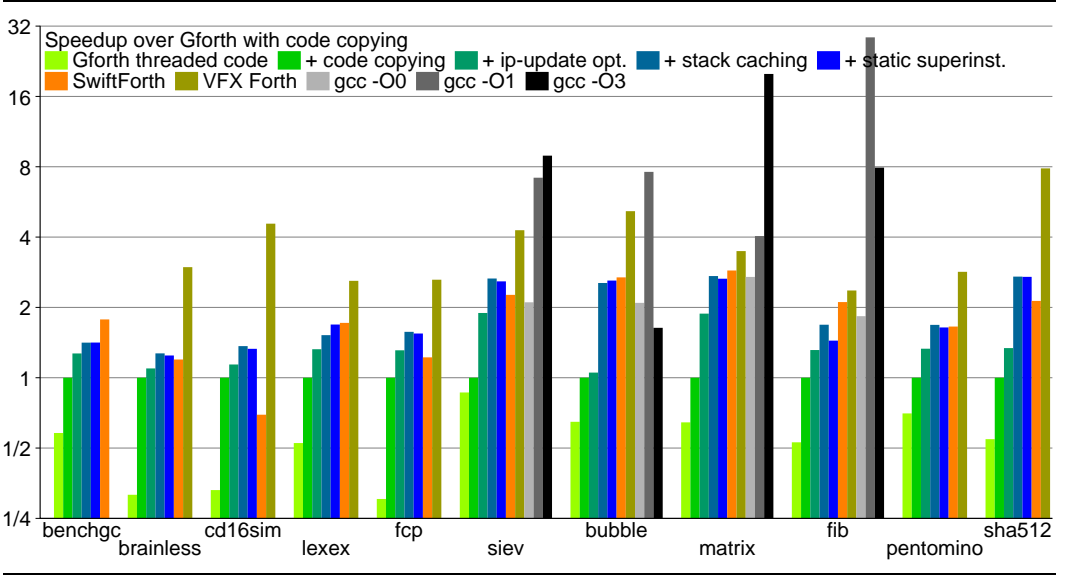


Fig. 1. Speedup factor of various systems over Gforth with code copying, on a Core i5-1135G7 (Tiger Lake)

gcc -O0, -O1, and -O3. All Forth systems use load-and-go compilers (compile time is included in the results), while GCC uses ahead-of-time compilation (only the run-time is shown in the results).

Not all benchmarks are available in C, and not all benchmarks run on all Forth systems, and the missing cases are reflected by missing bars.

The data shown is the median of 30 runs for each benchmark/system combination on a Core i5-1135G7 (Tiger Lake); each bar represents the number of cycles of Gforth with only code copying divided by the number of cycles of the system represented by the bar, i.e., the speedup of that system over Gforth with only code copying. The Gforth version used is 0.7.9\_20250817, commit 4224ab5fafea970dade64b04493ef690da8b3c32 compiled with gcc-11.4. The benchmarks are from the Forth appbench suite (benchgc-fcp), Gforth’s small (and mostly loop-dominated) benchmarks (sieve-fib), and two additional ones.

As can be seen, the performance of Gforth with all optimizations is similar to that of SwiftForth, which uses a conventional compiler, and typically around half of the performance of VFX Forth, which also uses a conventional compiler.

Before comparing Gforth with the others, let’s first take a look at the variants of Gforth, starting with the one with the best performance/effort:

**Threaded code** This is a fast interpretation technique for virtual-machine (VM) code, without any machine-code generation (see Section 3.1).

**Code copying** This method concatenates code snippets from the threaded code engine (see Section 3). It requires an estimated 500 lines of code in the Gforth source code. With this method Gforth still accesses literal data and performs control flow by accessing the VM code; it therefore also maintains a VM instruction pointer (IP), and updates it once for every VM instruction.

**IP-update optimization** This optimization reduces these IP updates. It was added by inserting 864 lines and deleting 316 lines in the Gforth source code [EP24].

**Stack caching** (actually static multi-state stack caching) eliminates many memory accesses to stack items and stack-pointer updates [EG04a, EG05]. The way this optimization is implemented in Gforth requires code copying to work.

**Static superinstructions** replace a sequence of Forth words with an optimized sequence [EGKP02]. Many of the benefits that static superinstructions have originally provided are now provided by code copying, the IP-update optimization and static stack caching; there are still cases where static superinstructions result in shorter code, but this has not led to consistent speedups in these measurements.

The code implementing stack caching and static superinstructions is quite interweaved with the rest of the code, so it is hard to give precise numbers for their size, but we estimate [Ert24] that all four optimizations combined require an estimated total of 5000 lines of code.

SwiftForth's compiler can be seen as a copy-and-patch compiler, but with the code snippets written by hand in assembly language and better resulting code than when patching using object file linkage information (see Section 7.3). SwiftForth does not have a VM interpreter substrate, and therefore does not have IP updates, so it gains the benefits of the IP update optimization without having to do anything. It deals with literal values and control flow by patching the code. SwiftForth does not perform multi-state stack caching, but it makes extensive use of static superinstructions (346 rules in 1819 lines). Overall each of the IA-32 and AMD64 targets of SwiftForth has about 7000 lines of architecture-specific code [Ert24].

Gforth with all optimizations is competitive in speed with SwiftForth, so apparently Gforth's stack caching provides enough speedup to compensate the costs that Gforth incurs for literals and control flow.

VFX Forth performs register allocation of data-stack items within a basic block, and inlines aggressively; inlining is very helpful for idiomatic Forth code, where calls and returns are the most frequent basic block boundaries. Therefore inlining also enhances the effectiveness of VFX's register allocator. The speed advantage of VFX over Gforth and SwiftForth is a result of these optimizations. In particular, for the `cd16sim` benchmark there is one call site that calls an empty definition and that is responsible for 2/3 of Gforth's run-time on this benchmark, while VFX inlines it away. We have no source code for VFX and therefore cannot report numbers about the size of its compiler. When asked about the effort to port VFX to ARM A64 (a currently ongoing project), Stephen Pelc gave the qualitative statement "far too much".

VFX is faster than Gforth by typically around a factor of 2. However, it is possible to perform inlining in Gforth, too, with direct performance benefits as well as indirect benefits from better stack caching. It will be interesting to see how far Gforth (and code copying) can close the gap.

Gforth's performance with all optimizations is roughly comparable to that of `gcc -O0` on those benchmarks that are also available in C. `gcc -O1` and `gcc -O3` often produce significantly faster code; sometimes they don't, but the reasons for that are beyond the scope of this paper.

## 2.2 Portability

A major reason to avoid implementing a conventional compiler is portability/retargetability.

Gforth has supported as many architectures as we could get our hands on, as long as `gcc` and something Unix-like (e.g., Cygwin for Windows) is available on the architecture. Gforth has supported the following architectures with a code copying compiler: Alpha, ARM A32/T32, ARM A64, HPPA, IA-32, IA-64, Loongarch, SPARC, PowerPC, PowerPC64 (but we no longer can check for all architectures that they still work). Gforth supports all architectures it does not know about by falling back to threaded code, which is slower, but still works.

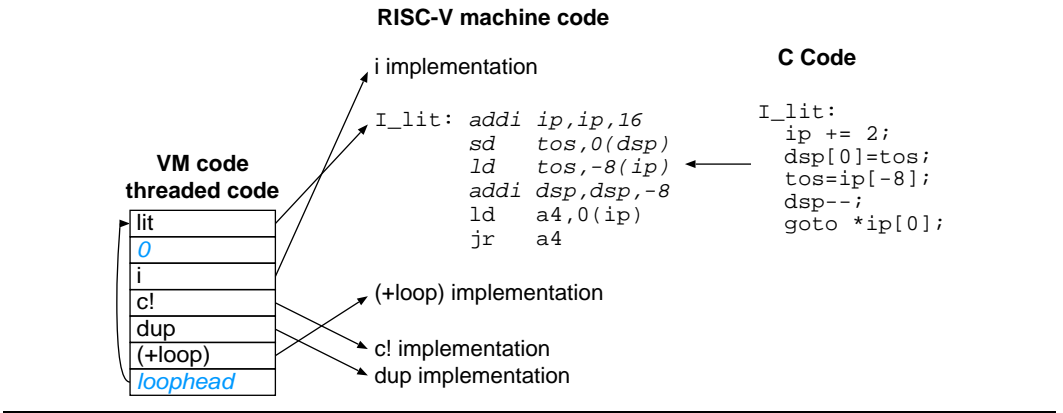


Fig. 2. Threaded-code representation of VM code. Each box is a machine word. *Slanted light blue* indicates an immediate operand of the preceding VM instruction.

In particular, when IA-64 (launched 2001) and AMD64 (launched 2003) became available to us in 2003, Gforth worked out of the box on these architectures<sup>2</sup> using the unknown-architecture support, likewise for ARM A64 in 2014 and RISC-V in 2017. A few small changes enabled code copying<sup>3</sup>, and a one-line change for configuring the number of registers for stack caching.

The benefit of code copying is that it reuses the retargeting efforts of the compiler it is based on (GCC or Clang in case of Gforth).

By contrast, SwiftForth has supported only IA-32 until the 2020s, when they started working on an AMD64 port (released on 2025-10-22). VFX has supported IA-32 initially, later ARM A32, and, also starting in the 2020s, AMD64. Both systems have interactive cross-compilers for a number of embedded targets.

The low number of desktop ports and the late support for AMD64 may be due to lack of commercial interest, but we think that the larger effort required to retarget and maintain the compiler for another architecture has something to do with it. iForth, another conventional Forth compiler, got an AMD64 port in 2009, but the IA-32 port was subsequently dropped (last release with IA-32 support in 2011).

### 2.3 Incremental development

Another benefit of code copying over writing a conventional compiler is that it can be done step-by-step: First add code copying, then add one optimization (e.g., IP-update optimization), then the next, etc., always with the fallback options of disabling the optimization or completely falling back on threaded-code interpretation.

By contrast, when coming from an interpreter, the conventional model requires a big-bang approach where a complete code generator for one target has to be developed without reusing much from an existing interpreter; and as long as you do not develop code generators for all targets, you still need to maintain the interpreter, as well as all the compiler targets. The latter will hopefully be helped by designing the compiler for retargetability, but that increases the complexity of the compiler framework.

<sup>2</sup>We added 64-bit support in 1996 while doing the Alpha port.

<sup>3</sup>For RISC-V, this was our first encounter with gcc-7 and its more aggressive code duplication (Section 5.4); we needed a little longer to find a workaround for that, but that's not specific to the architecture.

### 3 WHAT IS CODE COPYING COMPILATION?

#### 3.1 Threaded Code

The basis for Gforth’s code-copying implementation is a threaded-code interpreter [Bel73] for Gforth’s virtual machine (VM).

As a running example, we look at the VM code in Fig. 2. The first VM instruction in the example is `lit`, which has an immediate operand `0`. This VM instruction pushes its immediate operand on the data stack. It is represented by the address of the machine code that implements it; in direct-threaded code, every VM instruction is represented by the address of the machine code that implements it. In the case of `lit`, the implementation for RISC-V (RV64G) is:

```

                                # //C code
addi ip,ip,16    # ip += 2;
sd  tos,0(dsp)  # dsp[0] = tos;
ld  tos,-8(ip)  # tos = ip[-1];
addi dsp,dsp,-8 # dsp--;
ld  ca,0(ip)    # ca = ip[0];
jr  ca          # goto *ca;
```

This code uses register names that reflect their roles: `ip` is the VM instruction pointer; `tos` is the top of the data stack; `dsp` is the data stack pointer; `ca` is the code address (of the next VM instruction).

The *slanted blue* instructions are the payload which perform the actual work of the VM instruction as far as code copying is concerned. Other optimizations reduce that part further; e.g. the first instruction updates `IP`, and the `IP`-update optimization often optimizes it away.

The third instruction loads the immediate operand (`0`) from the VM code by accessing it through `IP`. This access of immediate operands and control-flow operations through `IP` is still in Gforth with all optimizations applied, and is the difference between an interpreter-based code-copying system and a copy-and-patch system (Section 7.3).

The bottom two (black) instructions perform the dispatch to the next VM instruction. The first instruction loads the machine code address of the next VM instruction, and the second instruction jumps to it.

This assembly-language code can be generated from the C code shown in the comments of the assembly language. It uses the GNU C extension “Labels as Values”,<sup>4</sup> which allows jumping to the address in `ca` with `goto *ca`<sup>5</sup>; this extension is also supported by Clang, tcc, and icc.

The other VM instruction implementations have the same pattern of *payload*, and dispatch. The last VM instruction in our example, `(+loop)` is notable: it is a VM-level conditional branch that branches back to *loophead* (given as immediate operand) or falls through to the next instruction. It is implemented with the following code

<sup>4</sup><https://gcc.gnu.org/onlinedocs/gcc/Labels-as-Values.html>

<sup>5</sup>The GCC maintainers call this a computed goto, although it is more like a Fortran assigned goto.

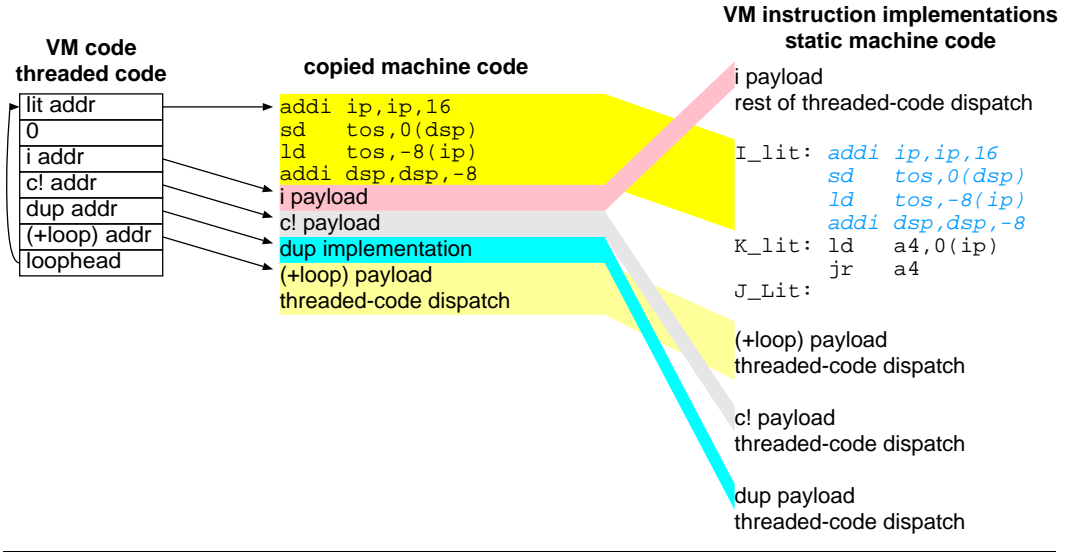


Fig. 3. Code copying.

```

addi ip,ip,16          # ip += 2;
...compute condition...
blt a5,zero,fallthrough # if (taken) {
ld ip, -8(ip)          # ip = ip[-1];
ld ca, 0(ip)           # ca = ip[0];
jr ca                  # goto *ca;
fallthrough:          # }
ld ca,0(ip)            # ca = ip[0];
jr ca                  # goto *ca;

```

If the conditional branch is taken, the new IP is loaded from the immediate operand and a dispatch is performed. It is better to have separate dispatches for the taken and the fallthrough cases for branch prediction<sup>6</sup> and because it allows to leave away the fallthrough dispatch in code-copying.

### 3.2 Code copying

Most VM instructions do not perform VM-level control flow, but just continue with the next VM instruction. Code copying copies and concatenates the machine code implementing the VM instructions, but in most cases without the dispatch code at the end. Only taken branches (i.e. VM instructions that change IP to point to some other VM instruction than the next one) need to perform a dispatch.

Figure 3 shows this for our running example. The VM code is conceptually the same as before, but for each VM instruction the machine word now points to the copied machine code instead of the original.

<sup>6</sup>Even with history-based indirect-branch prediction, branch predictors have an easier time if there are fewer targets for each indirect branch

In particular, the copied code still has the IP, which points to the threaded (VM) code, and it accesses the immediate operands `0` and `loophead` through it. The threaded code is also used on control flow: the VM-level conditional branch (`+loop`) is taken, loads the target threaded-code address `loophead` into IP, and then performs a threaded-code dispatch, which loads the code address at `loophead`, which points to the start of the concatenated code. All control flow in Gforth is performed with threaded-code dispatches in this way.

The threaded-code slots for instructions other than `lit` in this example are not accessed during execution. Gforth keeps them around to simplify the implementation.

At the end of the shown sequence the threaded-code dispatch is copied. While this is necessary for unconditional branches, it is not generally necessary for conditional branches such as (`+loop`) (as discussed above). However, the following VM instruction may make it necessary to perform a dispatch after the (`+loop`).

Code copying has also been called the `memcpy()` method [RS96], selective inlining [PR98] and (especially in Gforth) dynamic superinstructions [EG03a].

### 3.3 Benefits over threaded code

The obvious benefit of code copying is that it eliminates most threaded-code dispatches and results in straight-line execution of VM-level straight-line code, avoiding the limit of typically one taken branch per cycle. Another benefit is that the indirect branches in most of the remaining dispatches have only one target, vastly improving branch prediction accuracy in CPUs without sophisticated indirect-branch predictors, and still making life easier (and faster) for hardware with such branch predictors.

Another benefit is that code copying enables additional optimizations that require code snippets that are not represented as VM instructions (and where introducing additional VM instructions with threaded-code dispatch would make the optimization unprofitable).

E.g., the IP update optimization [EP24] leaves the IP update in front of most VM instruction implementations away and replaces it with an IP update by a larger amount for VM instructions that actually use the IP.

As another example, stack caching as implemented in Gforth inserts transitions between stack-cache states where necessary. These transitions do not have a VM instruction slot and therefore can only be inserted when code-copying is enabled. Gforth's stack-caching implementation relies on being able to insert the transitions, so stack caching is disabled when code copying is disabled [EG04a].

### 3.4 When is code copying appropriate?

The shorter the VM instruction implementations are, the larger the benefit of code copying over threaded code, because the overhead of threaded-code dispatch is relatively larger then.

Conversely, with long VM instruction implementations as in Tcl, whose VM instructions “can average hundreds of [machine] instructions” [VA04] the benefit is small, and often does not amortize the cost of copying the code or of increased I-cache misses [VA04].

Another aspect is that a compiler (to VM code) that uses more VM instructions, with each doing less, has more opportunities to optimize the VM code. This has been done for CPython recently<sup>7</sup>. With expensive VM instruction dispatch, splitting an existing VM instruction into several simpler ones increases the cost, and the optimization must be very good and must be applicable often to amortize this cost. With code-copying, the dispatch cost approaches 0, and such transformations become less of a gamble.

<sup>7</sup><https://github.com/faster-cpython/>

## 4 IMPLEMENTATION OF CODE COPYING

### 4.1 Code organization

Gforth has a big function `engine()` that contains all the code snippets (implementations of all VM instructions, and additional snippets used by optimizations), and little else.

Every code snippet has a label in front of it and behind it:

```
L_before:
  code snippet in C;
L_after:
  threaded-code dispatch;
```

You can see that more concretely in Fig. 3.

The label before it obviously points to the start of the code snippet.

Getting the right label for the end of the code snippet was initially straightforward (up to gcc-3.1), but later required extra work. If the source code falls through to the label (i.e., it does not end in an unconditional branch), like for the payload of most VM instructions in Gforth, with some extra help (see Section 5.4), the following label points right behind the code snippet, but if the code snippet cannot reach the label (e.g., because it ends in an unconditional branch, e.g. in a threaded code dispatch), gcc-3.2 and following have reordered code. We solved this problem by taking the values of all the labels, sorting them, and searching for the first label behind the label at the start of the snippet. This might include some unrelated code in cases where the code snippet does not fall through to the label, but in that case this is not a problem for correctness (but possibly for relocatability, see Section 4.5).

The function `engine()` has two code paths: the first just returns a table containing all the labels, for use in threaded-code generation and code-copying; the second starts the execution of the code by performing a threaded-code dispatch.

If code copying is disabled,<sup>8</sup> the threaded code address for each VM instruction just points to the implementation of that instruction inside `engine()`, and every threaded-code dispatch jumps around within this function.

With code copying, the first threaded code dispatch in `engine()` jumps to the copy of the VM instruction implementation and continues running there, with control-flow changes by performing a threaded-code dispatch.

### 4.2 Why does it work?

Why can we concatenate the code snippets produced in the way described above, and get code that works?

In particular, won't the register allocator have different register allocations for the different code snippets? Actually, at the start and, for fallthrough snippets, the end of the snippet, the register allocation has to be the same as at the start of every other snippet, because the compiler has to consider the possibility that every `goto *` jumps to every label whose address is taken. And the addresses of all labels before and after all code snippets are taken (to determine the code snippet address and length).

The code snippets that do not fall through end in a `goto *` in Gforth. And the register allocation at the `goto *` has to be compatible with that of all the labels whose address is taken, or it would not work even in ordinary use.

---

<sup>8</sup>Gforth option `--no-dynamic`.



More precisely, `engine()` is compiled separately from the code dealing with the threaded code, so the C compiler has to assume that every `goto *` in `engine()` can jump to any label whose address is taken.

Therefore, at a `goto *` all variables are alive (i.e., read before being overwritten) that are alive at any label whose address is taken, and each variable has to be in the same location at all those labels and all the instances of `goto *`. The code snippets that fall through to their second label are followed by a threaded-code dispatch:

```
ca = ip[0];
goto *ca;
```

so at the label between the code snippet and the dispatch, all the same variables are alive as at the `goto *`, except possibly `ca`, but that is not alive before the threaded-code dispatch, either. These variables also all have to reside at the same locations, because the `goto *` could jump to them.

### 4.3 Fallback

There are cases where certain code snippets cannot be copied (usually because they are not relocatable, see Section 4.5). How does Gforth deal with that?

Gforth falls back to plain threaded code in these cases: Append a threaded-code dispatch to the previous copied code snippet (unless the code snippet already ends with a threaded-code dispatch), and let the machine word representing the current VM instruction point to the original implementation of the VM instruction (inside `engine()`) rather than a copy). At run-time, the code performs the threaded-code dispatch, which then jumps to the original; that ends in another threaded-code dispatch, which may jump to code coming out of code-copying, or to another original implementation.

If other optimizations are active, the preparation for the fallback may require appending additional code. E.g., the IP needs to be up-to-date before the threaded-code dispatch, so in the presence of IP-update optimization, an IP update may be inserted before the threaded-code dispatch. Also, in Gforth the plain threaded code always expects the stack in the canonical state, so in the presence of stack caching, a transition from the current stack state to the canonical stack state may need to be inserted before the threaded-code dispatch.

Gforth may also find that it cannot copy the threaded-code dispatch. In that case it disables code copying completely and falls back to threaded code not just for individual VM instructions, but for all of them.

The option to fall back to threaded code has helped in various cases where things did not work according to our expectations (e.g., see Section 5.4). It means we always have a way to make Gforth work, albeit not as fast as we would like.

### 4.4 Instruction sets

Code copying is based on the assumption that the code snippets are independent and concatenable. At the instruction-set level this is satisfied if individual instructions are independent and concatenable. Some instruction sets have restrictions between groups of instructions. In this case a code snippet must not contain a partial group, i.e., there must not be a label within a group.

There are a few cases of such instruction-set restrictions:

**Branch delay slots** This is a misfeature of some early RISC architectures, in particular, HPPA, MIPS and SPARC: The branch instruction performs the instruction behind it before continuing at the target. This does not work with code copying if the compiler puts a label

between the branch and the instruction behind it. However, the compilers we have used (most recently gcc-14.2) do not do that.

**Load delay slots** This is a restriction of the MIPS I instruction set (eliminated in MIPS II). The instruction behind a load instruction is not allowed to read the register written by the load instruction. MIPS I also has some placement restrictions on reading and writing the hi and lo registers. Having labels right after the load or in the shadow of hi/lo reads can result in violating these restrictions in code copying. We have not tested if compilers actually place labels in a way that would lead to such violations. Instead, these concerns along with the relocatability problems (Section 4.5) and the lack of relevance of MIPS in Unix systems around 2003 were the reasons why we just configured Gforth to fall back to threaded code on MIPS (including the 64-bit MIPS port).

**Instruction groups** This is an IA-64 (aka Itanium processor family) property. Instructions within a group have restrictions on register usage that are intended to ensure that the instructions can be performed in one cycle without register renaming.<sup>9</sup> If a compiler put a label inside a group, code copying could violate these restrictions. Apparently the compilers we used (gcc-3.3, gcc-4.1.3, gcc-4.3.2) put stops (group boundaries) at labels, because in our testing IA-64 has always worked fine. If they did not, an easy fix would be to insert the stops using asm statements or at the assembly-language stage.

Based on the experiences with branch delay slots and instruction groups, it seems that gcc developers also avoid splitting groups of instructions with interdependencies by inserting a label inside these groups, but if these instruction sets still were important targets, that might change.

The problematic restrictions/features have not spread to newer architectures and all the architectures with these restrictions in general-purpose computers have been canceled in the meantime, while older or contemporary architectures without these restrictions thrive. So apparently the idea of independent, concatenable instructions has some merit, and we can expect that future instruction sets will also exhibit this property and thus support code copying.

## 4.5 Relocatability

A code snippet must be relocatable in order to be used in code copying, i.e., it must behave the same way in the original place and when copied.

*Non-relocatable code.* The main problems here are references to addresses: The code in the snippet must refer to addresses inside the snippet in a PC-relative way, and must *not* refer to addresses outside the snippet in a PC-relative way. Most architectures refer to other code addresses in a PC-relative way, so the most common reason for non-relocatability is when the VM instruction implementation performs a call to some function (e.g., for performing I/O).

Accesses to global constants or to global variables in a PC-relative way can also cause non-relocatability. Gforth avoids global variables for that reason and because of multi-threading; it stores some formerly global variables in a struct whose address is stored in a local variable inside engine(). However, computing the FP negation and the FP absolute value implicitly involve a constant that resides in memory on AMD64 (with SSE2 FP), making the implementations of these VM instructions (fnegate and fabs) non-relocatable on this architecture.

The pointer-to-struct approach could also be used for invoking functions without making the calling code non-relocatable, but for now we have not done that.

Note that asking the C compiler for position-independent code does not mean that individual code snippets are relocatable, even though the binary as a whole is, because position-independent

<sup>9</sup>Groups are often confused with bundles, which are IA-64's encoding of three instructions in 128 bits. By contrast, groups can be arbitrarily long, and can start and end somewhere in the middle of a bundle.

code may refer to code or data outside the code snippet in a PC-relative way (and usually does), while a relocatable code snippet must not do this.

*Determining relocatability.* How do we find out if a code snippet is relocatable or not? The implementations of the VM instructions actually look as follows:

```
L_skip:
    asm("SKIP4");
    asm("SKIP4");
    asm("SKIP4");
    asm("SKIP4");
L_before:
    code snippet in C
L_after:
    asm("SKIP4");
    asm("SKIP4");
    asm("SKIP4");
    asm("SKIP4");
    threaded-code dispatch
```

We compile `engine()` with these pieces to assembly language. Then we assemble the result twice: Once with `SKIP4` defined as empty string, so the `SKIP4`s assemble to nothing, and the result is as discussed earlier; and once with `SKIP4` defined as `.skip 4`, and with `engine` defined as `engine2`, so as a result the object file contains a function `engine2()` that has 16 bytes of padding before and after each code snippet.<sup>10</sup> We link both object files into the final executable. The addresses of the `L_skip` labels are taken and passed outside `engine()`, so `gcc` cannot optimize the initial skip away as dead code, and also because that usually is the next label after a threaded-code dispatch.

We now have a function `engine()` without the skips before and after the code snippets, and a function `engine2()` that has 16-byte skips before and after each code snippet. We extract the labels from each of the functions, and then compare the code snippets: If a code snippet from `engine()` contains exactly the same bytes as the corresponding code snippet from `engine2()`, then the code snippet is relocatable, otherwise it is not.

How does this work? If code from inside the code snippet references a code or data address outside the code snippet through a PC-relative address, the offset of the relative address will be different between `engine()` and `engine2()`, because the target label will be farther away in `engine2()` thanks to the skips. If there is an absolute reference (e.g., MIPS `j` instruction) to inside the code snippet, it will be different between `engine()` and `engine2()`, because the respective targets are at different addresses.

Even if the code snippet ends in an unconditional branch and the C compiler puts some other code behind that unconditional branch,<sup>11</sup> this scheme works: If the two code snippets compare equal, the code is relocatable. When used in a code-copying system, the code snippet may have some unused code behind the unconditional jump, but the generated code is still correct.

The reason for skipping 16 bytes is that this is a common code-alignment value, so the skips would not result in altered alignment (these days we ask the compiler to align to 1-byte boundaries, so skipping less might be sufficient). The reason for performing the 16-byte skip as 4 4-byte skips

<sup>10</sup>In earlier times we compiled twice rather than assembling twice, but compiling once is faster, and we do not need to worry if the two compilation runs introduce unintended differences in addition to the intended ones.

<sup>11</sup>We have not seen such an occurrence yet.

is that for some targets gcc counts the number of instructions in `asm` statements, assumes that each instruction takes at most 4 bytes, and generates code that relies on this assumption.

The absolute target addresses for the MIPS `j` and `jal` instructions have a catch: They work only for targets in the same 256MB segment of the address space. When we last looked, the functions `engine()` and `engine2()` were linked in the same 256MB segment as the functions called by some of the code snippets, and the code snippets would have been classified as relocatable. However, they were only relocatable within this 256MB segment. This is another reason why we disabled code copying for MIPS. An alternative would have been to allocate the memory for the copied code in the same 256MB segment as the original. Fortunately, among the architectures we have looked at, only MIPS has this property.

## 5 COMPILER ISSUES

In the previous section we have already mentioned a few caveats about how compilers have interfered with our initial assumptions about the generated code, and what we do about that. This section discusses additional issues.

We had quite a few problems with various gcc versions in the 2000s, and for some we found ways to deal with them, while some others were eventually fixed (after reappearing for several years). Also, the rethoric about undefined behaviour started at around that time and has spread and become more aggressive since then,<sup>12</sup> so at some point we expected to have to switch from using GNU C to assembly language as a more reliable foundation at some point [Ert14], essentially switching to a conventional compiler. But this has not happened (yet?), and actually, in the 2010s and 2020s only few new problems have appeared, and we found ways to deal with them. So GNU C seems to be a relatively stable foundation after all, once one has implemented various workarounds.

### 5.1 Code reordering

When we started, gcc arranged the basic blocks in source order. This changed with gcc-3.2. This has an effect on how we find the next label (Section 4.1). But we also saw cases where the compiler moved basic blocks from between `L_before` and `L_after` to outside these labels, which caused problems.

To avoid such problems, we tried to have only straight-line code in the VM instruction implementations. We extracted loops and most `if`-statements into functions that are compiled separately, and the VM instruction implementation only contains a call to this function. This costs a little performance (from the function call as well as turning the VM instruction implementation into non-relocatable code on most architectures), but fortunately the VM instruction affected by this are executed relatively rarely.

However, conditional VM branches are executed frequently, and in the ideal case they contain a conditional branch, in the following form (also seen for `(+loop)` in Section 3):

---

<sup>12</sup><http://blog.lvm.org/2011/05/what-every-c-programmer-should-know.html>

```

    ... skips ...
L_before:
    ... stack handling etc. ...
    if (VM_branch_taken)
        ip = ip[-1]; /*VM-branch target*/
        threaded-code dispatch;

L_after:
    ... skips ...
    threaded-code dispatch;

```

Ideally such VM-instruction implementations are compiled such that the basic blocks in the machine code are in the same order as in the source code, so that the code controlled by the `if` is between `L_before` and `L_after`, and the second threaded-code dispatch can be left away by code-copying in the usual case. For now, `gcc` does it that way for our code. But if `gcc` ever started changing this, a possible way to steer it back on the right path may be to use `__builtin_expect(VM_branch_taken, ...)` instead of just `VM_branch_taken`.

## 5.2 Code alignment

Compilers insert padding to align branch targets to instruction-fetch boundaries or cache-line boundaries. In particular, they do this for branch targets behind unconditional branches and loop heads.

When code copying, the padding inserted for the original code is often inappropriate for the target code. Therefore, we suppress this padding by compiling `engine()` with the options `-falign-labels=1 -falign-jumps=1`.

Instead, our code-copying implementation performs its own alignment (but on 2007-era processors where we measured the effects, the effects were in the noise).

## 5.3 Code deduplication

Starting with `gcc-3.0`, `gcc` started to compile all the `goto *` instances to an unconditional jump to one instance of an indirect branch. The reason for this probably was to reduce the control-flow edges in the data-flow analysis, for  $m$  `goto *` and  $n$  labels from  $nm$  to  $n + m$ .

In a number of `gcc` versions (up to the early `gcc-4.x` releases), `gcc` then did not eliminate the unconditional jump afterwards, with some versions eliminating them and some versions regressing, but eventually the `gcc` maintainers managed to make the unconditional-branch elimination stick, for our code.

So if that is a solved problem, why do we mention it here? We occasionally see this problem reappear in some form, so it's not completely gone.

E.g., when we managed to extend stack-caching support on AMD64 to three registers, we found that on AMD64 `gcc` compiled the `goto *` to an unconditional branch to common code that contains a lot of register shuffling (with no overall effect) and finally the indirect branch. Apparently the register shuffling made the common code so long that the branch-elimination heuristic decided not to eliminate the branch.

Fortunately, we found out that the register shuffling (and, consequently, the unconditional branch) go away with the compilation option `-fno-tree-vectorize`. Apparently without this option `gcc` tries to vectorize loads and stores of adjacent values, and is less precise in the data flow analysis for that than for individual values, leading to the register shuffling.

For the problems in the gcc-3.x and 4.x era, Gforth contains a workaround that has just one threaded-code dispatch and jumps there from all the VM instruction implementations. Gforth has labels before and after this dispatch, and because there is only one, gcc does not deduplicate it; this allows Gforth to use it as a code snippet that is appended whenever a threaded-code dispatch is needed.

In order to work with this workaround and still be relocatable, we implemented conditional VM branches to just set the IP on a taken branch, and then continue through `L_after` to the dispatch code. This results in worse code than we would have liked, but it was the best that was possible on these compiler versions. This approach remains an option when building Gforth,

## 5.4 Code duplication

On our first encounter with gcc-7, we found that the generated code looked as a straightforward compiler would generate for:

```
L_skip:
    ... skipping ...
    code snippet in C;
    threaded-code dispatch;
L_before:
    code snippet in C;
    threaded-code dispatch;
L_after:
    threaded-code dispatch;
```

I.e., gcc-7 duplicated code reached by jumping to a label and the same code being reached in a straight-line way. This may be a useful optimization, but it means that our code snippets now contain the dispatch code, which is contrary to our intentions.

We found the following workaround: In order to convince gcc that this code duplication does not pay off, after each label we insert 8 asm statements, each containing a comment with a text unique to that label (so gcc hopefully will not try to deduplicate the code). Currently this is enough to convince gcc to avoid the code duplication

## 5.5 Register allocation

Virtual machines have a number of “registers”, which are implemented in C code as C (local) variables. At least for the frequently-used variables, it would help performance if they were allocated to real-machine registers.

Up to and including gcc-9, we explicitly assigned registers to several of these variables on many platforms with GNU C’s feature “Explicit Register Variables”. In gcc-10 and later, disabling the explicit register variables produced better results than enabling them.

With either approach, we have the following problem: In the Gforth engine, gcc only used callee-saved registers for these variables. With explicit register variables, because gcc does not accept caller-saved registers for those. But if left to itself, gcc does not use caller-saved variables, either, because `engine()` contains about 100 VM instruction implementations that perform calls, and these calls apparently cause the compiler to avoid using caller-saved registers for these variables, especially for those that are used in < 100 VM instructions, such as the return-stack pointer of Gforth. A problem here is that gcc does not know that VM instructions that access the return stack are used frequently, while VM instructions that perform calls tend to be used rarely. This is a problem even for architectures like Alpha that have a lot of registers in principle, but a calling convention with relatively few callee-saved registers.

For being able to use additional registers for stack caching without spilling other VM registers, we use the following observation: All VM instruction implementations that contain a call only use the canonical state with one stack item in a register, due to non-relocatability. So additional stack cache registers are dead at the end of these VM instruction implementations, and there is no reason to preserve these registers across the calls. But how do we tell gcc about that?

```
L_skip:
    ... skipping ...
L_before:
    code snippet containing a call;
    asm("" : "=X"(spb));
    asm("" : "=X"(spc));
L_after:
    threaded-code dispatch;
```

The empty asm statements right before L\_after claim to overwrite spb and sbc (the variables holding the additional stack-cache items in some stack-cache states). Therefore, these variables are dead at the call and do not need to be preserved. This means that this VM instruction implementation is no hindrance to allocating spb and spc in a caller-saved register. And indeed, one of these variables is allocated by gcc in a caller-saved register.

Another way to influence the register allocator that we have not used is the GNU C extension “Label Attributes” (available since gcc-5). We can declare the VM instruction implementations with calls as being cold, and/or declare frequently-used VM instruction implementations to be hot by following the label with an attribute:

```
L_skip:
    ... skipping ...
L_before: __attribute__((cold));
    code snippet containing a call;
L_after:
    threaded-code dispatch;
```

With that, the register allocator is hopefully more willing to use caller-saved registers for local variables of the VM.

## 5.6 Cache consistency

Many architectures do not guarantee cache consistency between data and instruction caches, and require a special piece of code between generating code and executing code; this incantation typically consists of a few lines of architecture-specific (or, on some architectures worse, implementation-specific or OS-specific) code, and for a long time has been the only non-portable part of Gforth’s code copying implementation. Gcc-4.3 introduced `__builtin__clear_cache()`, which would eliminate this last piece of non-portability. We use `__builtin__clear_cache()` on RISC-V.

Unfortunately, `__builtin__clear_cache()` is not implemented correctly on at least PowerPC64.<sup>13</sup> We have switched Gforth back to using architecture-specific implementations of this functionality (except on RISC-V). When implementing your own code-copying compiler, check if `__builtin__clear_cache()` is compiled to non-empty code on each architecture that requires special code to make the caches consistent. If it compiles to non-empty code, that code will hopefully be correct.

<sup>13</sup>[https://gcc.gnu.org/bugzilla/show\\_bug.cgi?id=93811](https://gcc.gnu.org/bugzilla/show_bug.cgi?id=93811)

Another problem with such architectures is multi-threading: The code-generating thread must ensure that the D-cache lines are written to a common memory, and then the code-executing threads must invalidate these regions in the I-cache (to get rid of stale I-cache lines); due to prefetching and branch prediction, this may even be necessary if code in the address range has never been executed.

Until now we have ignored this problem, and relied on our luck. Typically Gforth programs only start subthreads after finishing compiling the source code (and thus code generation), which may explain why we have not seen any problems from that. A system with on-demand code generation (the narrow meaning of JIT) may be more likely to encounter such problems, however.

## 5.7 Spectre

GCC offers mitigations against Spectre v2 [KHF<sup>+</sup>19]. While all of these mitigations are expensive, because they disable indirect-branch prediction, the option `-mindirect-branch=thunk-inline` is less expensive than `-mindirect-branch=thunk`, because the latter makes the code snippets non-relocatable, so every VM instruction performs an indirect branch, while with the former option the relocatability of the code snippets is not affected, resulting in fewer indirect branches and therefore less slowdown.

On a Ryzen 3900X, we see slowdowns by a factor of 2.1–7.6 from using `-mindirect-branch=thunk-inline` and slowdown factors of 7.5–18.1 from using `-mindirect-branch=thunk`.

However, if you want to implement your programming language with Spectre mitigations, you will prefer approaches such as copy-and-patch compilation that avoid performing so many indirect branches. You will also want to use mitigations against other Spectre vulnerabilities (e.g., speculative load hardening [ZBC<sup>+</sup>23] against Spectre v1), which will introduce additional slowdowns for any approach, but unfortunately, these other mitigations require more work than just setting a C compiler flag.

## 5.8 Control-flow protection

There are exploit techniques such as return-oriented and jump-oriented programming that work by returning or jumping to arbitrary code. To make it more difficult to use these techniques, architectures and compilers offer ways to check that branches and returns only jump to targets that the compiler had in mind. E.g., gcc with the option `-fcf-protection=full` inserts an `endbr64` instruction at every indirect-branch target (i.e., every label in `engine()`), and the CPU can be told to report an error on an indirect branch to some other code. `Endbr64` is an AMD64 instruction, some other architectures have similar features.

This works with code copying: It copies the `endbr64` instruction to those places that the dispatch code will later indirect-branch to (and to additional places).

We use `-fcf-protection=none` in Gforth, however, because Gforth offers enough gadgets<sup>14</sup> already at the intended targets of indirect branches: All the VM instructions; moreover, Gforth and its VM is a low-level language that allows arbitrary memory access within the process. So a Gforth program that is exposed to untrusted input has to successfully defend against an attacker at the front line (source-level bounds checks etc.) and cannot make life harder for the attacker who has breached the front-line defense.

However, if your language is better suited to defense-in-depth, you can enable `-fcf-protection=full`, and they will work with code copying. This feature may cost a little

<sup>14</sup>In the context of return-oriented and jump-oriented programming, a gadget is a machine-code sequence that an attacker may want to return/jump to.



performance, though: All the `endbr64` instructions need to be decoded and executed. In a small experiment with Gforth on a Ryzen 8700G (Zen4), we saw an increase in instruction count by a factor 1.45 and an increase in cycle count by a factor 1.04 from `-fcf-protection=full`. Narrower processors may see a bigger slowdown (the instructions per cycle on Zen4 increased from 3.83 to 5.34). VM implementations with more machine instructions per VM instruction will see a smaller effect.

## 5.9 Clang

Clang supports “Labels as Values”, and Gforth is built with clang on platforms where GCC is not available. However, using Clang poses a number of problems:

- Clang wants to understand the assembly language in `asm` statements, and stops compiling when it sees `asm("SKIP4")`. One can work around that, and that is done in the ports that need clang, but we have not done that for the experiments on Debian Linux in the following.
- Clang takes much longer than gcc to compile Gforth’s `engine()` and also needs more memory. As an example, for `gforth-ipc` (an indirect-threaded-code Gforth without code copying nor other optimizations, and therefore without `SKIP4`), on a Ryzen 5800X `gcc-12.2` takes 3s and 346MB to compile `engine()`, while `clang-14.0.6` takes 699s and 5603MB. For `engine()` for `gforth-fast` (with all optimizations enabled), clang takes 3399s and 18264MB before it stops compiling because of `SKIP4` (gcc takes 26s and 1804MB).
- Clang generates a lot of register and memory shuffling code, similar to what we have seen with `gcc-3.0`. As a result, running the small benchmarks on Clang-compiled `gforth-ipc` executes 6.4 times more AMD64 instructions than on GCC-compiled `gforth-ipc` and consumes 4.2 times more Ryzen 5800X cycles.

As a result, Gforth selects GCC whenever it can. We expect that the clang compilation speed will be a problem for other code-copying compilers. The bad code generation may be less pronounced in language implementations that rely less on copy propagation than Gforth. Clang may be more viable when using tail calls instead of using one function and “Labels as Values” (see Section 7.1).

## 6 OS ISSUES

Over the years operating systems have restricted executing dynamically-generated code more and more. In the beginning, all memory was allocated with read, write, and execute (RWX) permissions; later, `malloc()` only allocated RW memory, and one has to use `mmap()` to get RWX memory.

Recently, some operating systems (in particular MacOS on Apple silicon) do not serve `mmap()` calls that ask for RWX memory (this restriction is also known as `W^X`). This is a problem for all systems with run-time code generation, not just code-copying compilers, but, e.g., Java JITs as well. For a single-threaded language implementation, one can `mprotect()` the memory to W when generating the code, and to X when executing it, but that does not work for multi-threaded code, unless you want to start a new page whenever you generate a new piece of code.

MacOS provides a MacOS-specific API for JIT compilers that supports switching the memory into W in the code-generating thread and keeping it X in the other threads, and Bernd Paysan has actually invested the time to use this API.

Several of the BSDs also has `W^X` by default, but allows to mark binaries such that RWX works. The command for marking the binary is short, but specific to the BSD variant.<sup>15</sup>

An approach that may work without special APIs is to have the code generation in one process and the execution in a different process, both mapping the same memory, but with different permissions. Another option may be to map the same memory within one process twice, at one

<sup>15</sup>[https://www.reddit.com/r/BSD/comments/10isrl3/notes\\_about\\_mmap\\_mprotect\\_and\\_wx\\_on\\_different\\_bsd/](https://www.reddit.com/r/BSD/comments/10isrl3/notes_about_mmap_mprotect_and_wx_on_different_bsd/)

address range with W permission, and at the other address range with X permission. We have not tried either approach.

If all else fails or you don't want to jump through the hoops that these operating systems put up, code-copying based on threaded code always allows you to fall back to plain threaded code, which works fine on operating systems with the W^X restriction. E.g., Gforth-0.7 (which was not specifically designed for this circumstance) automatically falls back to plain threaded code on MacOS on Apple silicon: the `mmap()` call for allocating the code memory fails, so Gforth-0.7 falls back to using `malloc()`, and because that does not produce executable memory on modern OSs, Gforth-0.7 turns off dynamic code generation.

## 7 ALTERNATIVE APPROACHES

In this section we describe approaches that are interesting but that are not implemented in production Gforth.

### 7.1 Tail calls

Instead of putting all VM-instruction implementations in one function and using `goto *` for threaded-code dispatch, one can also put each VM instruction implementation in a separate function and use optimized tail-calls for threaded-code dispatch, as follows:

```
typedef void (*vm_inst)(void **ip, long *dsp, long tos);
```

```
void lit(void **ip, long *dsp, long tos)
{
    ... payload including ip update ...;
    (*((vm_inst *)ip)[0])(ip, dsp, tos);
}
```

The last line of the function performs the threaded-code dispatch. The tail-call must be optimized into a jump, otherwise the C stack grows and eventually overflows. When we first considered this approach [Ert95], GCC did not tail-call optimize such code, but in the meantime it does, as does Clang [XK21]; Clang even provides a way to require that a call is tail-call-optimized, and will report an error if it cannot meet this requirement.

The VM registers are passed as parameters, at least as long as the calling convention supports passing them in machine registers. With gcc, additional VM registers could be stored in global explicit register variables; on AMD64 this results in 12 general-purpose and 8 floating-point registers available for VM registers. Clang does not support explicit register variables, but it supports using a calling convention for these functions and calls that uses as many registers as possible for parameter-passing.

So for dealing with VM registers efficiently, one has to pass VM-registers in parameters or keep them in global register variables with compiler-dependent and ABI-dependent code, but that is a relatively small effort.

With the tail-calling approach, there is a fixed allocation of VM registers to machine registers, either coming from the position in the parameter list, or from the explicit register allocation.

We expect that the VM instruction implementations can be compiled faster and with less memory with the tail-calling approach, because the compiler will hopefully not try to perform data-flow analysis between the functions, while it tries to do it when the implementations are all contained in one function. We can then squander the compilation speed gain on introducing more code snippets, for various optimization purposes (Xu and Kjolstad report using 98831 code snippets [XK21]).

Another benefit is that we should see no or little of the register-and-memory shuffling that we see with Clang, or with gcc without `-fno-tree-vectorize`.

So far you have only seen how tail calls can be used to implement threaded code. How can it be used for code-copying compilation?

In order to do that, we need a way to get rid of the dispatch part of the implementation. Unfortunately, compilers tend to mix the instructions from the payload part with those from the dispatch part; just inserting a label between them will not work, because there is nothing that jumps to this label. Maybe an asm statement can be made to act as a barrier, but preliminary experiments failed to produce satisfying results.

One way that may be more promising is to have, in addition to functions that end in a threaded-code dispatch (to have a fallback option), variants intended only for code copying that end in a direct [XK21] or indirect tail-call without threaded-code dispatch. On many architectures this is just one instruction, that must be last in the function. However, there are exceptions: Some architectures have delayed branches (HPPA, MIPS, SPARC); some architectures require two instructions for indirect branches (PowerPC, IA-64). In some programming models, a direct jump to a function is expressed as an indirect jump to a target loaded from the global offset table (GOT), and as a result the direct jump also is expressed with more than one instruction.

Once we have solved the problem of keeping the payload separate from the tail call, how do we know where the tail call starts so that we can use the code between the start of the function and this instruction as code snippet? Xu and Kjolstad extract the function size (and the code) from the object file (see Section 7.2), and apparently use their own architecture-specific knowledge about the size of the last instruction to determine where it starts. A way to determine the size of this last instruction may be to have a function that performs only this tail-call, and look at its size.

## 7.2 Snippets from object files

Gforth extracts code snippets from the executable at run-time and has some startup overhead while it examines all the code snippets for relocatability and performs its table setup.

An alternative is to extract code snippets from object files [NHCL98, XK21] at system build time using the Binary File Descriptor library (GNU BFD). One advantage of this approach is that the object file contains additional information, such as the function size, or linkage information for symbols external to the object file.

## 7.3 Copy-and-patch compilation

Gforth accesses immediate operands and control-flow information through IP. This requires a register for IP, results in less efficient accesses to immediate operands and less efficient control flow than with ordinary compilers, and requires keeping the VM code around.

An alternative is to have code snippets that contain dummy immediate arguments and perform control flow directly to dummy targets, and then patch the constants or target addresses in these code snippets with the actual values, resulting in copy-and-patch compilation.

One approach for copy-and-patch compilation has been based on using the linkage information in object files [NHCL98, TCL<sup>+</sup>00, XK21]. References to external symbols are used for patchable immediate operands and patchable control-flow targets. The linkage information describes where to patch and how to patch (e.g., absolute or relative address). This requires some architecture/ABI-specific work, but ABIs have a finite number of relocation types (e.g., 52 in the AMD64 ABI [LMG<sup>+</sup>]) and only a few are actually used in the code snippets.

However, by referring to an external symbol the copy-and-patch compiler usually cannot patch the immediate operand of instructions like RISC-V's `addi`. The external symbol is a 64-bit (or

32-bit) value, while the immediate operand of `addi` is 12 bits long, so the addition of a constant (whatever its size) is compiled to several instructions.

Another approach is to start with code snippets delimited by labels in one C function, like Gforth’s code copying uses, but perform patching in addition [VA04, EG04b].

We implemented copy-and-patch compilation for Gforth in a prototype for IA-32 and PowerPC using the latter approach [EG04b]. This work was based on Gforth’s approach of extracting code snippets from the executable at system startup time. The `engine()` function was compiled thrice, twice with the same immediate arguments, and once with different immediate arguments. The first two versions were compared to determine relocatability, the third version was compared to find out the placeholders of the immediate arguments.

This approach can make use of the RISC-V `addi` instruction, but needs to fall back to code that uses several instructions when the immediate operand becomes too large. It needs quite a bit of knowledge about the instruction encodings, in particular, the sizes of the immediate-operand fields. We considered determining the encoding and size by varying the immediate operands a lot more, but did not implement that idea; dealing with each architecture manually is probably less work.

We originally intended to turn this copy-and-patch compiler into a production engine for Gforth, but in those years several GCC releases resulted in falling back to threaded code, so the copy-and-patch approach looked too brittle, and we let it bit-rot. Later, the rethoric by the advocates of C code without undefined behaviour kept the distrust in GCC high. If we had continued to maintain this engine, maybe we could now report on its success and the hurdles we had to overcome. Or maybe it would have been a bridge to far.

## 8 RELATED WORK

GCC-2.0 (released February 1992) introduced “Labels as Values”, which not only proved useful for implementing threaded code (we started the Gforth project[Ert93] in July 1992), but also for compiling by copying compiler-generated code snippets between two labels, with all the code snippets being within a function. This method was first outlined by Rossi and Sivalingam [RS96, Section 2.5], who refer to an unpublished discussion between Xavier Leroy and Kenneth Oksanen. Piumarta and Riccardi provided a more elaborate treatment [PR98], with deduplication of code sequences.

Ertl and Gregg implemented code-copying in Gforth, and in the beginning the main benefit was in indirect branch prediction accuracy [EG03a, EG03b, CEG07]; it turned out that leaving away deduplication (or conversely, introducing replication, as we framed it) helped the branch predictors at the time. Indirect branch predictors have improved a lot in general-purpose processors [RSS15], but code copying still provides a good speedup.<sup>16</sup>

Once you have code copying, you can eliminate instruction-pointer (IP) updates, either by leaving away the unneeded VM instruction slots [PR98], or by replacing several IP updates with a combined one [EP24]. While IP updates play a minor role for performance on CPUs from the 2000s, they can be the decisive bottleneck on loop-dominated benchmarks in the 2020s.

Another optimization that was facilitated by code copying is multi-state stack caching [Ert95, EG04a, EG05].

Tempo is a partial evaluator that uses code copying and patching by extracting information from object files [NHCL98]; Tempo was later used to specialize an interpreter into a compiler [TCL\*00].

Iliasov [Ili03] describes a copy-and-patch compiler with a minimal patching component: Only literals need to be patched; control flow is performed by performing indirect jumps to addresses provided as literals.

<sup>16</sup>See Section 2.1 and <http://www.complang.tuwien.ac.at/anton/interpreter-branch-pred.txt>.

QEMU is a full-system emulator. It is a production system with a long history, and has many more users than Gforth. QEMU can emulate machines with a different instruction set than the host machine. It uses dynamic translation techniques for that, originally implemented in its Dyngen component [Bel05] using code-copying and patching, similar to what we described in Section 7.2 and 7.3. But Dyngen uses ordinary functions, not tail-calling functions, and has to get rid of the function prologue and epilogue. Dyngen is gcc-3.x-specific, and it apparently was too difficult to adapt it to newer gcc versions or other compilers, so it was replaced with TCG in QEMU-0.10.0 released in 2009. TCG is based on QOP by Paul Brook, who described it as “Hand written code generator”<sup>17</sup>, so TCG probably is not based on copying and pasting compiler-generated code.

In Gforth we have dealt with changes in GCC by finding workarounds, or, for versions where we were not successful, by falling back to threaded code. Another approach is to actually define the properties that a compiler’s code generation should have to support code copying; then modify a compiler to provide those properties (when asked for it), and report an error if it fails to provide the properties. This approach has been explored by Prokopski and Verbrugge [PV07, PV08], but their patches have not been integrated into GCC.

Several code-copying JavaVM implementations have been implemented, among them SableVM [GH03] and the Cacao interpreter [ETK06]. A particular challenge solved by these implementations was quickening of VM instructions, where VM instructions rewrite themselves into faster code on first execution. SableVM stopped being maintained after the research project ended (last release 2007). The Cacao interpreter bit-rotted while the main thrust of Cacao continued to use conventional code generation technology.

Maxine is a Java VM implementation with two-level compilation (baseline and optimizing compiler), where the baseline compiler is a copy-and-patch compiler that uses templates written in Java and where the code is generated by the optimizing compiler (which uses conventional compiler techniques) or by HotSpot [WHV<sup>+</sup>13].

Xu and Kjolstad implement two copy-and-patch compilers: One that directly compiles from the abstract syntax tree (AST) without going through a VM and one for WebAssembly. Their technique works by having each code snippet (called stencil in the paper) in a tail-calling function with references to external symbols as placeholders for patching, and extracting the code snippets from object files. They use 1666 code snippets for the WebAssembly compiler, and 98831 code snippets for the AST compiler; the latter is notable, because it is beyond practical for the technique where all code snippets are in one function.

## 9 CONCLUSION

Code-copying compilers make retargeting of the compiler much easier by using code snippets coming from a different compiler. Gforth demonstrates that code-copying without patching can produce code with similar performance as a compiler with a hand-written architecture-specific code generator. Gforth has used code copying since 2003, on many architectures, and has dealt with many GCC versions in those years. If all else fails, Gforth can fall back to threaded code, but it usually does not have to.

Copy-and-patch compilation promise an improvement in performance over copying without patching (as in Gforth) at a moderate increase in architecture-specific code. However, while there have been a number of publications about this technology, no production system is known to us that currently uses it.

<sup>17</sup><https://qemu-devel.nongnu.narkive.com/bCtjCaPs/hand-written-code-generator-2>

## REFERENCES

- [Bel73] James R. Bell. Threaded code. *Communications of the ACM*, 16(6):370–372, 1973.
- [Bel05] Fabrice Bellard. QEMU, a fast and portable dynamic translator. In *Freenix Track of Usenix Annual Technical Conference*, pages 41–46, 2005.
- [CEG07] Kevin Casey, M. Anton Ertl, and David Gregg. Optimizing indirect branch prediction accuracy in virtual machine interpreters. *ACM Transactions on Programming Languages and Systems*, 29(6):37:1–37:36, October 2007.
- [EG03a] M. Anton Ertl and David Gregg. Optimizing indirect branch prediction accuracy in virtual machine interpreters. In *SIGPLAN Conference on Programming Language Design and Implementation (PLDI’03)*, 2003.
- [EG03b] M. Anton Ertl and David Gregg. The structure and performance of *Efficient* interpreters. *The Journal of Instruction-Level Parallelism*, 5, November 2003. <http://www.jilp.org/vol5/>.
- [EG04a] M. Anton Ertl and David Gregg. Combining stack caching with dynamic superinstructions. In *Interpreters, Virtual Machines and Emulators (IVME ’04)*, pages 7–14, 2004.
- [EG04b] M. Anton Ertl and David Gregg. Retargeting JIT compilers by using C-compiler generated executable code. In *Parallel Architecture and Compilation Techniques (PACT’04)*, pages 41–50, 2004.
- [EG05] M. Anton Ertl and David Gregg. Stack caching in Forth. In M. Anton Ertl, editor, *21st EuroForth Conference*, pages 6–15, 2005.
- [EGKP02] M. Anton Ertl, David Gregg, Andreas Krall, and Bernd Paysan. vmgen — a generator of efficient virtual machine interpreters. *Software—Practice and Experience*, 32(3):265–294, 2002.
- [EP24] M. Anton Ertl and Bernd Paysan. The Performance Effects of Virtual-Machine Instruction Pointer Updates. In Jonathan Aldrich and Guido Salvaneschi, editors, *38th European Conference on Object-Oriented Programming (ECOOP 2024)*, volume 313 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 14:1–14:26, Dagstuhl, Germany, 2024. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.
- [Ert93] M. Anton Ertl. A portable Forth engine. In *EuroFORTH ’93 conference proceedings*, Mariánské Lázně (Marienbad), 1993.
- [Ert95] M. Anton Ertl. Stack caching for interpreters. In *SIGPLAN Conference on Programming Language Design and Implementation (PLDI’95)*, pages 315–327, 1995.
- [Ert14] M. Anton Ertl. How to get rid of C. In *30th EuroForth Conference*, pages 63–65, 2014.
- [Ert24] M. Anton Ertl. Interpreter vs. compiler performance at run-time. In *Tagungsband des Jahrestreffens 2024 der GI-Fachgruppe “Programmiersprachen und Rechenkonzepte”*, INSIGHTS — Schriftenreihe der Fakultät Technik, pages 7–12, 2024.
- [ETK06] M. Anton Ertl, Christian Thalinger, and Andreas Krall. Superinstructions and replication in the Cacao JVM interpreter. *Journal of .NET Technologies*, 4:25–32, 2006. Journal papers from .NET Technologies 2006 conference.
- [GH03] Etienne Gagnon and Laurie Hendren. Effective inline-threaded interpretation of Java bytecode using preparation sequences. In *Compiler Construction (CC ’03)*, volume 2622 of *LNCS*, pages 170–184. Springer, 2003.
- [Ili03] Alex Iliasov. Templates-based portable just-in-time compiler. *SIGPLAN Notices*, 38(8):37–43, August 2003.
- [KHF<sup>+</sup>19] Paul Kocher, Jann Horn, Anders Fogh, Daniel Genkin, Daniel Gruss, Werner Haas, Mike Hamburg, Moritz Lipp, Stefan Mangard, Thomas Prescher, Michael Schwarz, and Yuval Yarom. Spectre attacks: Exploiting speculative execution. In *40th IEEE Symposium on Security and Privacy (S&P’19)*, 2019.
- [LMG<sup>+</sup>] H.J. Lu, Michael Matz, Milind Girkar, Jan Hubička, Andreas Jaeger, and Mark Mitchell, editors. *System V Application Binary Interface — AMD64 Architecture Processor Supplement (With LP64 and ILP32 Programming Models)*.
- [NHCL98] François Noël, Luke Hornof, Charles Consel, and Julia L. Lawall. Automatic, template-based run-time specialization: Implementation and experimental study. In *IEEE International Conference on Computer Languages (ICCL ’98)*, pages 123–142, 1998.
- [PR98] Ian Piumarta and Fabio Ricciardi. Optimizing direct threaded code by selective inlining. In *SIGPLAN ’98 Conference on Programming Language Design and Implementation*, pages 291–300, 1998.
- [PV07] Gregory B. Prokopski and Clark Verbrugge. Towards GCC as a compiler for multiple VMs. In *Proceedings of the GCC Developers’ Summit*, pages 117–129, 2007.
- [PV08] Gregory B. Prokopski and Clark Verbrugge. Compiler-guaranteed safety in code-copying virtual machines. In *Compiler Construction (CC’08)*, pages 163–177. Springer LNCS 4959, 2008.
- [RS96] Markku Rossi and Kengatharan Sivalingam. A survey of instruction dispatch techniques for byte-code interpreters. Technical Report TKO-C79, Faculty of Information Technology, Helsinki University of Technology, May 1996.
- [RSS15] Erven Rohou, Bharath Narasimha Swamy, and André Seznec. Branch prediction and the performance of interpreters — don’t trust folklore. In *Code Generation and Optimization (CGO)*, 2015.

- [TCL<sup>+</sup>00] Scott Thibault, Charles Consel, Julia L. Lawall, Renaud Marlet, and Gilles Muller. Static and dynamic program compilation by interpreter specialization. *Higher-Order and Symbolic Computation*, 13(3):161–178, September 2000.
- [VA04] Benjamin Vitale and Tarek S. Abdelrahman. Catenation and specialization for Tcl virtual machine performance. In *IVME '04 Proceedings*, pages 42–50, 2004.
- [WHV<sup>+</sup>13] Christian Wimmer, Michael Haupt, Michael L. Van De Vanter, Mick Jordan, Laurent Daynès, and Douglas Simon. Maxine: An approachable virtual machine for, and in, Java. *ACM Transactions on Architecture and Code Optimization*, 9(4):30:1–30:24, January 2013.
- [XK21] Haoran Xu and Fredrik Kjolstad. Copy-and-patch compilation. *Proc. ACM Program. Lang.*, 5(OOPSLA):136:1–136:30, October 2021.
- [ZBC<sup>+</sup>23] Zhiyuan Zhang, Gilles Barthe, Chitchanok Chuengsatiansup, Peter Schwabe, and Yuval Yarom. Ultimate SLH: Taking speculative load hardening to the next level. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 7125–7142, Anaheim, CA, August 2023. USENIX Association.