

GRAMMAR-BASED COMPRESSION OF INTERPRETED CODE

Automatically designing and implementing compact interpretable bytecodes.

Programs that must run at or near top speed must use native machine code, but some programs have more modest performance requirements. For example, a cellular telephone handset might include software that reacts to key-strokes and spends most of its time waiting for the next input. An interpreted code can run quickly enough for such applications and can be smaller than machine code. Witness the success of the language Forth in embedded systems.

Good encodings are, however, difficult to design. They must anticipate the sequences of operations that programmers will use most often and assign the shortest codes to the most common sequences.

Typical programs can be analyzed for statistics to guide the design. Indeed, the designer of a compact representation may target a single program and design a compact language exclusively for that program. Of course, designing a language for every program is too labor-intensive to be done by

hand. It requires both automation and a different interpreter for each compacted program, which can also be expensive. A better solution may be to design an interpreter for a set of programs and use one interpreted language for all.

Our focus here is on the automatic design and implementation of compact interpretable bytecodes. The objective is a form that is compact for a set of sample programs and for other programs with similar characteristics. The

key to designing such compact bytecodes is to identify frequently occurring patterns of program constructs, and to replace them with a single interpreted construct. This process is unlike the Huffman fixed-to-variable length coding, which encodes single symbols using a variable number of bits, and more like Tunstall variable-to-fixed length coding, which encodes multiple symbols as a single, fixed-size codeword [12].

Representation

We could start our search for frequently occurring patterns with programs represented in a high-level language (such as

**PROGRAMMING
languages normally obey
a grammar, and this
restriction can help
compression. It is not
necessary to represent
invalid programs,
which confers an
immediate advantage
over general-purpose
compression schemes
that must compress
everything.**

C++), an intermediate code (such as bytecode or expression trees), or the instruction set of a target machine. Since our goals include direct interpretation, starting with a high-level language can be problematic. Few high-level languages are directly interpreted, so our compacting interpreter would itself need to produce a lower-level representation for

interpretation. At the other end of the spectrum, we may start with machine code and have the interpreter translate its compact representation into instructions that can be directly executed. It can, however, be tricky for the interpreter to maintain control of the execution of the program in this case. Thus, most systems rely on a compiler front-end to produce an intermediate code that can be interpreted.

Some systems create specialized instructions for common sequences in postfix bytecode [3, 7]. Others operate on the corresponding expression trees [6, 9]. For example, a simple expression, such as $1+(2 \times 3)$, translates into the tree `AddInt(1, MulInt(2, 3))`. Proebsting's greedy heuristic looks for the most frequent parent/child pair in all the expression trees and creates a new instruction or "superoperator" for that pair of operations. For example, if multiplication by two is the most common pair, then a new, unary instruction `MulInt(2,*)` replaces all multiplications by two. After replacement, our example expression would use only two operands and two operations, rather than the original three operands and two operations. This process may be repeated to obtain a desired interpreter language size. Ernst et al. describe a similar method that enhances a RISC-like virtual machine instruction set [4].

A variation on the superoperator theme represents the program as a dictionary and a skeleton [8]. The dictionary contains the enhanced instructions while the skeleton consists of a sequence of original machine code instructions and references (calls) to the enhanced instructions in the dictionary. The skeleton, in this case, acts as an interpreter.

The preceding methods and the one described here [5] all use addressable codes, typically bytecodes. A notable alternative interprets Huffman-coded instructions. Such encodings traditionally require unaligned, bit-level decoding, but methods have been devised to decode them more quickly [2, 7].

Grammar-based Techniques

Programming languages normally obey a grammar, and this restriction can help compression. It is not necessary to represent invalid programs, which confers an immediate advantage over general-purpose compression schemes that must compress everything. A grammar also categorizes program elements, and compressors can exploit this labeling by specializing an encoding for each element type.

Thompson and Booth describe how to use a probabilistic grammar for a context-free language to encode strings from the language [11]. One of their techniques, termed derivation encoding [10], repre-

sents a program by the sequence of grammar rules used to derive it (in a leftmost or rightmost derivation) from the starting symbol of the grammar. Thompson and Booth suggest using a Huffman code, based on the probabilities of the grammar rules, to encode the rule choices.

Another grammar-based encoding method—parsing encoding—represents a program as the sequence of choices made by a parser as it processes the program. A top-down parser makes essentially the same choices as the derivation encoding, but a bottom-up or shift-reduce parser is different. The parser is typically a push-down automaton, and the choices it makes are which action (shift or reduce) to perform and which state to transition to.

Cameron [1] demonstrated the power of derivation encoding by using a probabilistic grammar to obtain a derivation along with its probability. He then encoded the derivation using an arithmetic encoder and was able to compress programs to almost 10% of their original size.

These methods do not produce interpretable results. The compressed form of the program must be decompressed and compiled before execution.

Grammar Rewriting

How can we exploit the compression potential of grammar-based methods in a language an interpreter can decode without decompressing it first? One solution [5] starts with some representative sample programs and a grammar for the original (uncompressed) instruction set. Each program when parsed using the grammar yields a parse tree that describes a derivation of the program (see Figure 1). The list of rules used in the derivation forms a byte-code encoding of the program. The compressor transforms the grammar so that it parses the same language but uses fewer derivation steps, at least for the sample programs. The revised grammar defines a bytecode that will be smaller for the sample and for

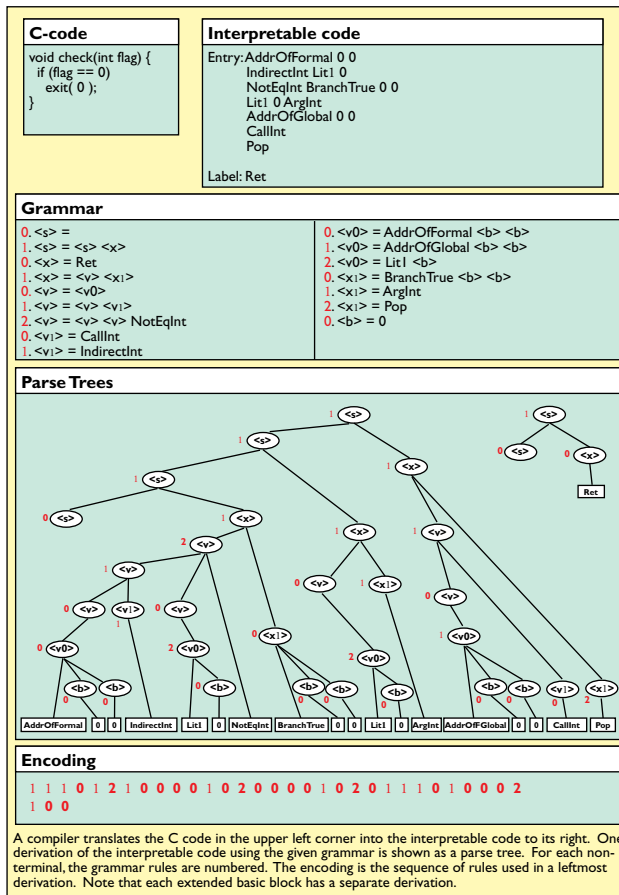


Figure 1. Using a grammar to encode interpretable code.

Each rule number, in conjunction with the current non-terminal, acts as an instruction for our interpreter. We can see for the grammar in Figure 1 that there are three instructions: 0, 1, and 2. Unless we encode each rule number as a byte, this is not, in general, a very practical code for interpretation. In order to create a practical and concise encoding of the program, we modify the grammar so that each non-terminal has close to 256 rules. The modification process takes two rules, $A \rightarrow \alpha B \beta$ and $B \rightarrow \gamma$, and adds to the grammar a third rule, $A \rightarrow \alpha \gamma \beta$, where A and B are non-terminals and α , β , and γ are strings of terminals and non-terminals. We call this process inlining a B rule into an A rule. Inlining doesn't change the language accepted by the grammar. However, it shortens the sequence of rules (the derivation) needed to express some programs, and it increases the number of rules for some non-terminals.

Which rules should we inline? The goal of the inlining is to produce a grammar that provides short derivations for programs. Starting with a derivation of a program using the original grammar, the best single inline that we could perform is the most frequently occurring pair of rules; one used to expand a non-terminal on the right-hand side of the other. If this pair

similar programs.

A derivation transforms the grammar's start non-terminal into the final program by expanding the leftmost non-terminal at each step using some grammar rule. Our encoding is the list of the rules used in the derivation. Each rule is represented as an index: the i th rule for a non-terminal is represented as the index i .

Note that a separate derivation is generated for each basic block. Keeping the derivations separate allows direct interpretation of the encoding. When the interpreter encounters a control transfer, it knows the encoding at the target starts with a rule to expand the start non-terminal.

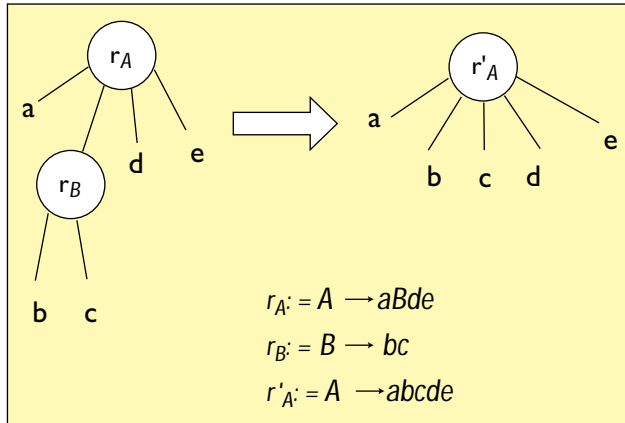


Figure 2. Edge contraction.

were used m times in the derivation, inlining would decrease the derivation length by m rules.

We can view this process as operating on the forest of parse trees obtained from parsing the original, uncompressed sample program using the original grammar. The parse produces a forest since we restart the parser from the start non-terminal at every potential branch target (that is, `Label`). For our purposes, a parse tree is a rooted tree in which each internal node is labeled with a rule and each leaf with a terminal symbol. The root is labeled with a rule for the start non-terminal. In general, an internal node that is labeled with a rule $A \rightarrow a_1 a_2 \dots a_k$ (where a_i is a terminal or non-terminal symbol) has k children. If a_i is a non-terminal then the i th child (from the left) is labeled with a rule for non-terminal a_i . If a_i is a terminal then the i th child is a leaf labeled with a_i . The program appears as the sequence of labels at the leaves of the parse trees in the forest, reading from left to right. A leftmost derivation is simply the sequence of rules encountered in a preorder traversal of each parse tree in the forest.

The inlining of one rule r_B into another rule r_A creates a new rule r'_A whose addition to the grammar permits a different (and shorter) parse of the program. One such new parse can be obtained by contracting every edge from a node labeled r_A to a node labeled r_B in the original forest—meaning the children of r_B become the children of r_A —and relabeling the node labeled r_A with the new rule r'_A (see Figure 2). If the number of edge contractions is m , the resulting forest has m fewer internal nodes and thus represents a derivation that is shorter by m steps.

To construct an expanded grammar, we parse a sample program (or a set of sample programs) using the original grammar and obtain a forest of parse trees. We then inline the pair of rules at the endpoints of the most frequent edge in the forest, contract all occurrences of this edge, add the new inlined rule to the grammar, and repeat. We stop creating rules for a non-terminal once it has 256 rules and thus create one bytecoded instruction set for each non-terminal. A grammar with several non-terminals will thus result in programs that interleave bytecodes for several different non-terminals, but the interpreter always knows how to decode the next byte because the context defines the non-terminal used to parse the next byte.

Occasionally, a rule for a non-terminal may be subsumed by a new rule. That is, after the addition of the new rule, the first rule is no longer used in the derivation. If the unused rule is one that was added via inlining, we are free to remove it from the grammar.

(We cannot, however, remove any of the original grammar rules, or we risk changing the grammar’s language.) In our current implementation, we remove unused inlined rules in order to decrease the size of the

| input | original | compressed | | | |
|-------|-----------|----------------|-------|----------------|-------|
| | | trained on gcc | | trained on lcc | |
| | | bytes | ratio | bytes | ratio |
| gcc | 1,423,370 | 471,111 | 33% | 577,814 | 41% |
| lcc | 199,497 | 75,077 | 38% | 57,722 | 29% |
| gzip | 47,066 | 19,466 | 41% | 19,706 | 42% |
| 8q | 436 | 138 | 32% | 152 | 35% |

Benchmark program sizes before and after compaction.

removal may cause some non-terminals to have fewer than 256 rules. The implementation could be made to respond with more inlining, but the number of reductions is typically small, and the incremental value of the next inlining step drops with time, so the added complexity might not pay off.

This construction procedure is greedy; it always inlines the most frequent pair of rules. This is a heuristic solution to the problem of finding a set of rules to add to the grammar that permits the shortest derivation of the sample program. We rely on this heuristic since finding an exact solution is, unfortunately, NP-hard.

The resulting expanded grammar is ambiguous, since we leave the original rules in the grammar. We can use any valid derivation, but the size of the representation is the number of rules in the derivation, so compression demands a minimum length derivation. We use Earley’s parsing algorithm, slightly modified, to obtain a shortest derivation for a given sequence. The derivation is then the compressed bytecode representation of the program and is suitable for interpretation.

The Interpreter

This system has two interpreters. The initial interpreter accepts the initial, uncompressed bytecode. The initial interpreter and the expanded grammar form the raw material from which the system builds the second interpreter, which accepts compressed bytecode.

At the core of the initial interpreter is a routine comprised of a single C switch:

```
void interpret1(
    unsigned char op, istate *istate
) {
    switch (op) { ... }
}
```

The routine accepts a single, uncompressed operator and pointer to a structure that records the state of an interpreter. The interpreter state could be maintained as variables local to a single interpretation routine, but it was helpful to be able to change the state from multiple routines.

The preceding switch has one case for each instruction in the initial instruction set, and the cases manipulate a small execution stack. Stack elements use a union of the basic machine types. For example, the case for `AddInt` pops two elements, adds them as integers, and pushes the result:

```
case AddInt:
    stack = istate->stack;
    a = stack[istate->top--].i;
    b = stack[istate->top--].i;
    stack[++istate->top].i = a + b;
    return;
```

The base interpreter, which is called `interp`, simply calls `interpret1` repeatedly. The second interpreter, which interprets compressed bytecodes, introduces another level of interpretation between `interp` and `interpret1`:

```
void interp(istate *istate) {
    while (1)
        interpNT(istate, NT_start);
}
```

`interpNT` adds an argument that identifies a non-terminal and thus which of several specialized byte-coded instruction sets to use. `interpNT` fetches the next bytecode, which, with the given non-terminal, identifies the rule for the next derivation step. A table encodes for each rule the sequence of terminals and non-terminals on the rule's right-hand side. `interpNT` advances left-to-right across this right-hand side. When it encounters a terminal symbol, it calls `interpret1` to execute that symbol. When it encounters a non-terminal, it calls itself recursively,

HOW can we exploit the compression potential of grammar-based methods in a language an interpreter can decode without decompressing it first? One solution starts with some representative sample programs and a grammar for the original (uncompressed) instruction set.

with the given non-terminal to define the new context and new specialized bytecode.

Performance

The table here reports the size of several bytecode sequences as compressed by our method. Each input was compressed twice, with grammars generated

from two different training sets, namely the compilers `lcc` and `gcc`. Predictably, `lcc` and `gcc` each compress somewhat better with their own grammar, but the other inputs compress about as well with either grammar.

The interpreters are small: 4,029 bytes for the initial, uncompressed bytecode and 13,185 for the bytecode generated from the `lcc` training set. Thus adding 9,156 bytes to the interpreter saves roughly 900KB in the bytecode for `gcc`. The grammar occupies 8,751 bytes and thus accounts for most of the difference in interpreter size.

The initial, uncompressed bytecode takes roughly 70 times longer to execute than native machine code, and the compressed bytecode adds another factor of two, but trade-offs favored size at every turn. Interpreter state is stored on the stack to simplify implementation, but it could be moved to registers. Also, double interpretation could be eliminated by hard-coding a switch for the compressed bytecode, which would suit systems that burn the interpreter into a cheap ROM but that download bytecode into scarcer RAM.

For calibration and as a very rough bound on what might be achievable with good, general-purpose data compression, `gzip` compresses the original,

uncompressed inputs described previously to 31–44% of their original size, with the larger inputs naturally getting the better ratios. Thus the compressed bytecode is competitive with `gzip` despite operating with an additional constraint, namely support for direct interpretation's random access. For example, `gzip` is free to exploit redundant patterns that span basic blocks, where our bytecode compressor must stop and discard all contextual information at every branch target. ■

REFERENCES

1. Cameron, R.D. Source encoding using syntactic information models. *IEEE Transactions on Information Theory* 34, 4 (1988), 843–850.
2. Choueka, Y., Klein, S.T., and Perl, Y. Efficient variants of Huffman codes in high level languages. In *Proceedings of the 8th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval* (1985), 122–130.
3. Clausen, L., Schultz, U., Consel, C., and Muller, G. Java bytecode compression for low-end embedded systems. *ACM TOPLAS* 22, 3 (May 2000), 471–489.
4. Ernst, J., Evans, W., Fraser, C.W., Lucco, S., and Proebsting, T.A. Code compression. In *Proceedings of the ACM SIGPLAN'97 Conference on Programming Language Design and Implementation*, 358–365.
5. Evans, W. and Fraser, C.W. Bytecode compression via profiled grammar rewriting. In *Proceedings of the ACM SIGPLAN'01 Conference on Programming Language Design and Implementation*, 148–155.
6. Hoogerbrugge, J., Augusteijn, L., Trum, J., and van de Weil, R. A code compression system based on pipelined interpreters. *Software-Practice and Experience* 29, 11 (1999), 1005–1023.
7. Latendresse, M. Automatic generation of compact programs and virtual machines for Scheme. In *Proceedings of the Workshop on Scheme and Functional Programming 2000*, 45–52.
8. Liao, S., Devadas, S., and Keutzer, K. A text-compression-based method for code size minimization in embedded systems. *ACM Transactions on Design Automation of Electronic Systems* 4, 1 (Jan. 1999), 12–38.
9. Proebsting, T.A. Optimizing an ANSI C interpreter with superoperators. In *Proceedings of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Jan. 1995), 322–332.
10. Stone, R.G. On the choice of grammar and parser for the compact analytical encoding of programs. *Computer Journal* 29, 4 (1986), 307–314.
11. Thompson, R.A. and Booth, T.L. Encoding of probabilistic context-free languages. In Z. Kohavi and A. Paz, Eds, *Theory of Machines and Computations*. Academic Press, 1971.
12. Tunstall, B.P. *Synthesis of Noiseless Compression Codes*. Ph.D. dissertation, Georgia Institute of Technology, 1967.

WILLIAM S. EVANS (will@cs.ubc.ca) is an assistant professor of computer science at the University of British Columbia.

CHRISTOPHER W. FRASER (cwfraser@microsoft.com) is a senior researcher in the programming language systems group at Microsoft Research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
