

Fast and Effective Procedure Inlining*

Indiana University
Computer Science Department
Technical Report No. 484

Oscar Waddell and R. Kent Dybvig
Indiana University Computer Science Department
Lindley Hall 215
Bloomington, Indiana 47405
{owaddell,dyb}@cs.indiana.edu, 812-855-3608

May 15, 2004

Abstract

Inlining is an important optimization for programs that use procedural abstraction. Because inlining trades code size for execution speed, the effectiveness of an inlining algorithm is determined not only by its ability to recognize inlining opportunities but also by its discretion in exercising those opportunities. This paper presents a new inlining algorithm for higher-order languages that combines simple analysis techniques with demand-driven online transformation to achieve consistent and often dramatic performance gains in fast linear time. Benchmark results reported here demonstrate that this inlining algorithm is as effective as and significantly faster than offline, analysis-intensive algorithms recently described in the literature.

Keywords: procedure inlining, procedure integration, compiler optimization

1 Introduction

Inlining is an important optimization for programs that use procedural abstraction. By replacing a call with an inline expansion of the procedure body, inlining eliminates procedure linkage overhead and may uncover opportunities for further optimization. Unrestrained inlining, however, can lead to excessive code growth and excessive compile time. The effectiveness of an inlining algorithm is therefore determined not only by its ability to recognize inlining opportunities but also by its discretion in exercising those opportunities.

Recent research on procedure inlining for higher-order languages has focused on the use of offline polyvariant flow analyses to identify appropriate inlining opportunities [14, 3]. Although reported performance gains are impressive, the cost and volatility of such analyses renders the method impractical for use in production compilers. This paper presents a new inlining algorithm for higher-order languages that combines simple analysis techniques with online transformation to achieve comparable or better performance gains in fast linear time.

Several factors contribute to the speed and effectiveness of our algorithm.

- It is polyvariant. Information about the arguments at a call site is used to decide whether inlining that call is cost-effective.
- It is online. Inlining decisions are based, in part, on the effects of other optimizations performed by the algorithm: constant folding, copy propagation, and elimination of useless or unreachable code.

*This material is based on work supported in part by the National Science Foundation under grant number CDA-9312614. This report is an expanded version of a paper to be presented at the *1997 International Static Analysis Symposium*.

- It is context-sensitive. The context in which an expression occurs is used to identify legal transformations.
- It is demand-driven. The operands of binding constructs and calls are not processed until the context in which they are used has been determined.

Polyvariance permits aggressive inlining without excessive code growth. Code growth is usually limited by inlining only procedures whose size is less than a given threshold. Our polyvariant algorithm instead considers the size of the procedure when optimized for the arguments at the call site. For example, given the definition of `f` below, the algorithm may choose to inline `e2` in place of the call `(f 'key2)` even if the size of `f` as a whole exceeds the threshold.

```
(define f
  (lambda (x)
    (case x [(key1) e1] [(key2) e2] ... [(keyn) en])))
```

Online transformation provides our algorithm with more accurate information on which to base its inlining decisions. Offline algorithms, which separate inlining from subsequent optimization, must estimate the effects of these optimizations. When their estimate is too conservative, inlining opportunities are missed; when their estimate is too liberal, code size increases. Our algorithm inlines aggressively while maintaining tighter control over code growth because it uses the actual optimized code to decide when inlining a call is cost-effective. In some cases, online transformation also provides more accurate control-flow information, which helps to identify additional inlining opportunities.

Contextual information allows the algorithm to determine when an expression can be discarded or processed in a more restricted mode. For example, in the expression `(seq e1 e2)`, `e1` is processed only for its effects. When processed for effect, certain provably terminating effect-free expressions are discarded. In the expression `(if e1 e2 e3)`, `e1` is processed for its boolean value. This is useful in languages such as C and Scheme where it may be possible to determine the boolean value of an expression even when its actual value cannot be determined.

Demand-driven processing of call operands allows our algorithm to defer processing until more contextual information is available. For example, suppose `f` is `(lambda (x y) (if (< x 5) x y))`. When attempting to inline the call `(f e1 e2)`, `e1` is not processed until it is needed to evaluate the test `(< x 5)`. If `e1` reduces to `3`, the reference to `y` is unreachable and the useless expression `e2` is processed for effect only. Eliminating useless or unreachable code saves processing time and also improves the accuracy of the information collected by the algorithm. For example, if variables originally referenced only within `e2` no longer appear referenced, their bindings can be processed for effect only.

The inlining algorithm is a source-to-source transformation that operates by a mostly depth-first traversal of the source program. The algorithm collects information about the size of the residual program and the use of variables in the residual program as it is produced. This information supports inlining decisions and elimination of useless code. Termination is achieved by detecting cycles in the call graph as it is computed and by bounding the amount of effort expended on inlining at each call site in the input program. Code growth is limited by aborting an inlining attempt as soon as the size of the residual code produced in that attempt exceeds a given threshold. Polyvariance is obtained by attempting to inline all calls, while relying on the effort bound and the size threshold to abort an attempt when it becomes too costly. Applicative order and sequencing of effects are preserved by optimizing inlined calls as if they were equivalent `let` expressions. For example, suppose `f` is `(lambda (x y) e1)`. When inlining, the call `(f e2 e3)` is treated as the semantic equivalent of `(let ((x e2) (y e3)) e1)`. Values obtained by processing `e2` and `e3` are used when specializing `e1`, but effects in `e2` and `e3` are residualized at the point of call.

Inlining is copy propagation extended to `lambda` expressions. Just as copy propagation often enables constant folding, inlining enables β -reduction. The next section formally describes an inlining algorithm that exploits this observation. This algorithm is simplified in that it concentrates on recognizing inlining opportunities without attempting to select appropriately among them. Section 3 shows how the simplified algorithm is restrained to obtain the actual inlining algorithm. Section 4 presents refinements of the inlining algorithm. Section 5 provides detailed measurements of the effect of the algorithm on execution time, code size, and compile time in a production compiler. Section 6 describes related work and compares our results

$$\begin{aligned}
e ::= & (\text{const } c) \mid (\text{ref } x) \mid (\text{primref } p) \\
& \mid (\text{if } e_1 e_2 e_3) \mid (\text{seq } e_1 e_2) \mid (\text{assign } x e) \\
& \mid (\text{lambda } (x) e) \mid (\text{letrec } ((x_1 e_1) \dots (x_n e_n)) e) \\
& \mid (\text{call } e_0 e_1) \\
c \in & \text{Const}, p \in \text{Prim}, x \in \text{Var}, e \in \text{Exp}
\end{aligned}$$

Figure 1: The core language.

$$\begin{aligned}
\rho \in \text{Env} &= \text{Var} \rightarrow \text{Var} \\
\sigma \in \text{Store} &= (\text{Loc}_x \rightarrow \text{VarFlags}) \times (\text{Loc}_\gamma \rightarrow \text{ContextFlags}) \times (\text{Loc}_e \rightarrow \text{Exp} \cup \{\text{unvisited}\}) \\
\kappa \in \text{Continuation} &= \text{Exp} \rightarrow \text{Store} \rightarrow \text{Exp} \\
s, r \in \text{VarFlags} &\subseteq \{\text{ref}, \text{assign}\} \\
&\text{ContextFlags} \subseteq \{\text{inlined}\} \\
op \in \text{Operand} &::= \text{Opnd}(\text{Exp}, \text{Env}, \text{Loc}_e) \\
&\text{Var} ::= \text{Var}(\text{Identifier}, \text{Operand} \cup \{\text{null}\}, \text{VarFlags}, \text{Loc}_x) \\
\gamma \in \text{Context} &::= \text{Test} \mid \text{Effect} \mid \text{Value} \mid \text{App}(\text{Operand}, \text{Context}, \text{Loc}_\gamma) \\
l_e \in \text{Loc}_e, l_x \in \text{Loc}_x, l_\gamma \in \text{Loc}_\gamma &\quad \text{Locations}
\end{aligned}$$

Figure 2: Domain equations.

with those obtained by offline, flow-analysis-based inlining algorithms. Section 7 summarizes our work and suggests directions for future work.

2 Unrestrained Inlining Algorithm

The inlining algorithm is described as a continuation-passing style (CPS) interpreter with explicit store. The actual implementation is in direct style but otherwise follows the structure of the algorithm described here with the restraints and extensions described in Sections 3 and 4.

The inlining algorithm, I , takes an input expression e , a context γ , an environment ρ , a continuation κ , and a store σ . The input expression is in the core language shown in Figure 1. The algorithm returns a residual expression in the same core language. Domain equations are given in Figure 2. Figure 3 presents the algorithm driver. Figure 4 shows the functions that implement the core of the algorithm. To simplify the presentation, **lambda** abstractions are restricted to a single formal parameter, and applications (**call**) are restricted to a single operand. Generalization to multiple-argument and variable-arity procedures is straightforward.

Contextual information (γ) is used to identify situations in which certain transformations are legal. Four contexts are distinguished: *test*, used when processing an expression for its boolean value, *effect*, used when processing an expression for side effects only, *value*, used when the actual value is needed, and *call*, used when processing an expression in the operator position of an application. Naturally, side effects are conserved in all contexts. Test, effect, and value contexts are represented by the nullary constructors **Test**, **Effect**, and **Value**. A call context γ is represented by the constructor **App**(op, γ_1, l_γ) where op is an **Opnd** structure (described below) for the actual parameter at the call site, γ_1 is the context of the call itself, and $\sigma(l_\gamma)$ contains a flag that is set if this call is inlined. This flag is used when residualizing a call.

The environment ρ maps source-program variables to residual-program variables and is used to rename (α -convert) bound variables to prevent inadvertent capture when code is duplicated by inlining. The initial environment ρ_0 is the identity function on variables. The store σ is a triple of functions that respectively

$$\begin{aligned}
I & : \text{Exp} \rightarrow \text{Context} \rightarrow \text{Env} \rightarrow \text{Continuation} \rightarrow \text{Store} \rightarrow \text{Exp} \\
I(\text{const } c)\gamma\rho\kappa & = \begin{cases} \kappa(\text{const void}) & \text{if } \gamma = \text{Effect} \\ \kappa(\text{const true}) & \text{if } \gamma = \text{Test and } c \neq \text{false} \\ \kappa(\text{const } c) & \text{otherwise} \end{cases} \\
I(\text{seq } e_1 \ e_2)\gamma\rho\kappa & = Ie_1 \text{Effect } \rho\lambda e'_1. Ie_2\gamma\rho\lambda e'_2. \kappa \text{seq}(e'_1, e'_2) \\
I(\text{if } e_1 \ e_2 \ e_3)\gamma\rho\kappa & = Ie_1 \text{Test } \rho\kappa_1 \\
& \text{where } \kappa_1 e'_1 = \begin{cases} Ie_2\gamma_1\rho\lambda e'_2. \kappa \text{seq}(e'_1, e'_2) & \text{if } \text{result}(e'_1) = (\text{const true}) \\ Ie_3\gamma_1\rho\lambda e'_3. \kappa \text{seq}(e'_1, e'_3) & \text{if } \text{result}(e'_1) = (\text{const false}) \\ Ie_2\gamma_1\rho\lambda e'_2. Ie_3\gamma_1\rho\kappa_2 & \text{otherwise, where} \\ \kappa_2 e'_3 = \begin{cases} \kappa \text{seq}(e'_1, e'_2) & \text{if } e'_2 = e'_3 = (\text{const } c) \\ \kappa(\text{if } e'_1 \ e'_2 \ e'_3) & \text{otherwise} \end{cases} \end{cases} \\
& \gamma_1 = \begin{cases} \text{Value} & \text{if } \gamma = \text{App}(op, \gamma_1, l_\gamma) \\ \gamma & \text{otherwise} \end{cases} \\
I(\text{assign } x \ e)\gamma\rho\kappa\sigma & = \begin{cases} Ie \text{Effect } \gamma\rho\kappa\sigma & \text{if } \rho(x) = \text{Var}(x', op, s, l_{x'}) \text{ and } \text{ref} \notin s \\ Ie \text{Value } \gamma\rho\kappa_1\sigma & \text{if } \rho(x) = \text{Var}(x', op, s, l_{x'}) \text{ and } \text{ref} \in s, \text{ where} \\ & \kappa_1 e'\sigma_1 = \kappa \text{seq}((\text{assign } x' \ e'), (\text{const } c)) \sigma_2 \\ & \sigma_2 = \sigma_1[l_{x'} \mapsto \{\text{assign}\} \cup \sigma_1(l_{x'})] \\ & c = \text{true} \text{ if } \gamma = \text{Test} \text{ else } c = \text{void} \end{cases} \\
I(\text{call } e_1 \ e_2)\gamma\rho\kappa\sigma & = Ie_1\gamma_1\rho\kappa_1\sigma_1, \text{ where} \\
& \begin{aligned} op & = \text{Opnd}(e_2, \rho, l_{e_2}) & l_{e_2} \text{ fresh} \\ \gamma_1 & = \text{App}(op, \gamma, l_{\gamma_1}) & l_{\gamma_1} \text{ fresh} \\ \sigma_1 & = \sigma[l_{e_2} \mapsto \text{unvisited}, l_{\gamma_1} \mapsto \emptyset] \\ \kappa_1 e'_1\sigma_2 & = \begin{cases} \kappa e'_1\sigma_2 & \text{if } \text{inlined} \in \sigma_2(l_{\gamma_1}) \\ \text{visit}(op, \text{Value}, \kappa_2, \sigma_2) & \text{otherwise, where} \\ & \kappa_2 = \lambda e'_2. \kappa(\text{call } e'_1 \ e'_2) \end{cases} \end{aligned} \\
I(\text{primref } p)\gamma\rho\kappa & = \begin{cases} \kappa(\text{const true}) & \text{if } \gamma = \text{Test} \\ \kappa(\text{const void}) & \text{if } \gamma = \text{Effect} \\ \kappa(\text{primref } p) & \text{if } \gamma = \text{Value} \\ \text{fold}(\text{primref } p)\gamma\rho\kappa & \text{if } \gamma = \text{App}(op, \gamma_1, l_\gamma) \end{cases} \\
I(\text{lambda } x \ e)\gamma\rho\kappa\sigma & = \begin{cases} \kappa(\text{const true})\sigma & \text{if } \gamma = \text{Test} \\ \kappa(\text{const void})\sigma & \text{if } \gamma = \text{Effect} \\ Ie \text{Value } \rho_1\kappa_1\sigma_1 & \text{if } \gamma = \text{Value, where } x \text{ abbreviates } \text{Var}(x, \text{null}, s, l_x) \\ & x' \text{ abbreviates } \text{Var}(x', \text{null}, \sigma(l_x), l_{x'}) \quad x', l_{x'} \text{ fresh} \\ & \rho_1 = \rho[x \mapsto x'] \\ & \sigma_1 = \sigma[l_{x'} \mapsto \emptyset] \\ & \kappa_1 = \lambda e'. \kappa(\text{lambda } x' \ e') \\ \text{fold}(\text{lambda } x \ e)\gamma\rho\kappa\sigma & \text{if } \gamma = \text{App}(op, \gamma_1, l_\gamma) \end{cases} \\
I(\text{ref } x)\gamma\rho\kappa\sigma & = \begin{cases} \kappa(\text{const void})\sigma & \text{if } \gamma = \text{Effect, else} \\ \kappa(\text{ref } x')\sigma_1 & \text{if } op = \text{null or assign} \in s \\ & \text{where } \text{Var}(x', op, s, l_{x'}) = \rho(x) \\ & \sigma_1 = \sigma[l_{x'} \mapsto \{\text{ref}\} \cup \sigma(l_{x'})] \\ \text{visit}(op, \text{Value}, \kappa_1, \sigma) & \text{otherwise, where } \rho(x) = \text{Var}(x', op, s, l_{x'}) \\ & \text{and } \kappa_1 = \lambda e\sigma_2. \text{copy}(\rho(x), \text{result}(e), \gamma, \kappa, \sigma_2) \end{cases}
\end{aligned}$$

Figure 3: The inlining algorithm.

$$\begin{aligned}
fold(\text{primref } p) App(op, \gamma_1, l_\gamma) \rho \kappa \sigma &= visit(op, \text{Value}, \kappa_1, \sigma) \\
\text{where } \kappa_1 e'_1 \sigma_1 &= \begin{cases} \kappa(\text{const } c') \sigma_2 & \text{if } result(e'_1) = (\text{const } c) \text{ and } p(c) = c' \\ & \text{where } \sigma_2 = \sigma_1[l_\gamma \mapsto \{\text{inlined}\} \cup \sigma_1(l_\gamma)] \\ \kappa(\text{primref } p) \sigma_1 & \text{otherwise} \end{cases} \\
fold(\text{lambda } x e) App(op, \gamma_1, l_\gamma) \rho \kappa \sigma &= Ie \gamma_1 \rho_1 \kappa_1 \sigma_1 \\
\text{where } x &\text{ abbreviates } \text{Var}(x, \text{null}, s, l_x) \\
x' &\text{ abbreviates } \text{Var}(x', op, \sigma(l_x), l_{x'}) \quad x', l_{x'} \text{ fresh} \\
\rho_1 &= \rho[x \mapsto x'] \\
\sigma_1 &= \sigma[l_{x'} \mapsto \emptyset] \\
\kappa_1 e' \sigma_2 &= \begin{cases} visit(op, \text{Effect}, \kappa_2, \sigma_2) & \text{if } \text{ref} \notin \sigma_2(l_{x'}) \text{ and } \text{assign} \notin \sigma_2(l_{x'}) \\ visit(op, \text{Effect}, \kappa_3, \sigma_2) & \text{if } \text{ref} \notin \sigma_2(l_{x'}) \text{ and } \text{assign} \in \sigma_2(l_{x'}) \\ visit(op, \text{Value}, \kappa_3, \sigma_2) & \text{otherwise} \end{cases} \\
\kappa_2 &= \lambda e'_1 \sigma_3. \kappa seq(e'_1, e') \sigma_3[l_\gamma \mapsto \{\text{inlined}\} \cup \sigma_3(l_\gamma)] \\
\kappa_3 &= \lambda e'_1 \sigma_3. \kappa (\text{call } (\text{lambda } x' e') e'_1) \sigma_3[l_\gamma \mapsto \{\text{inlined}\} \cup \sigma_3(l_\gamma)]
\end{aligned}$$

$$copy(\text{Var}(x', op, s, l_{x'}), e, \gamma, \kappa, \sigma) = \begin{cases} I(\text{const } c) \gamma \rho_0 \kappa \sigma & \text{if } e = (\text{const } c) \\ \kappa(\text{ref } x_1) \sigma & \text{if } e = (\text{ref } x_1) \text{ and } \text{assign} \notin s_1 \\ & \text{where } x_1 \text{ abbreviates } \text{Var}(x_1, op_1, s_1, l_{x_1}) \\ fold e \gamma \rho_0 \kappa \sigma & \text{if } \gamma = App(op_1, \gamma_1, l_\gamma) \\ & \text{and } e = (\text{primref } p) \\ & \text{or } e = (\text{lambda } x_1 e_1) \\ \kappa(\text{primref } p) \sigma & \text{if } \gamma = \text{Value} \text{ and } e = (\text{primref } p) \\ \kappa(\text{const true}) \sigma & \text{if } \gamma = \text{Test} \text{ and } e = (\text{primref } p) \\ & \text{or } e = (\text{assign } x_1 e_1) \\ & \text{or } e = (\text{lambda } x_1 e_1) \\ \kappa(\text{ref } x') \sigma_1 & \text{otherwise, where } \sigma_1 = \sigma[x' \mapsto \{\text{ref}\} \cup \sigma(l_{x'})] \end{cases}$$

$$\begin{aligned}
visit(\text{Opnd}(e, \rho, l_e), \gamma, \kappa, \sigma) &= \begin{cases} Ie \gamma \rho \kappa_1 \sigma & \text{if } \sigma(l_e) = \text{unvisited} \\ & \text{where } \kappa_1 = \lambda e' \sigma_1. \kappa e' \sigma_1[l_e \mapsto e'] \\ \kappa e' \sigma & \text{otherwise, where } e' = \sigma(l_e) \end{cases} \\
seq(e_1, e_2) &= \begin{cases} e_2 & \text{if } e_1 = (\text{const void}) \text{ else} \\ (\text{seq } (\text{seq } e_1 e_3) e_4) & \text{if } e_2 = (\text{seq } e_3 e_4) \\ (\text{seq } e_1 e_2) & \text{otherwise} \end{cases} \\
result(e) &= \begin{cases} e_2 & \text{if } e = (\text{seq } e_1 e_2) \\ e & \text{otherwise} \end{cases}
\end{aligned}$$

Figure 4: The inlining algorithm (continued).

map locations l_x to subsets of *VarFlags*, locations l_γ to subsets of *ContextFlags*, and locations l_e to residual expressions in the core language or to the flag `unvisited`. As a notational convenience, we use subscripts on location labels to select the appropriate component of the store. Thus $\sigma(l_x)$ abbreviates $(\sigma \downarrow 1)(l_x)$ and $\sigma[l_\gamma \mapsto \emptyset]$ abbreviates $\langle \sigma \downarrow 1, (\sigma \downarrow 2)[l_\gamma \mapsto \emptyset], \sigma \downarrow 3 \rangle$.

Source- and residual-program variables are represented using the constructor $\text{Var}(x, op, s, l_x)$, where x is the lexical identifier corresponding to the variable, op is an operand structure, and s and $\sigma(l_x)$ are subsets of *VarFlags* that describe uses of the variable in the source and residual program, respectively. The `ref` flag supports elimination of useless assignments and useless bindings, as described in Section 2.4 and Section 2.7. The `assign` flag is used to prevent copy propagation and inlining through assigned variables. The parser introduces `Var` structures for variables and initializes the flags in the initial store, σ_0 . The final store contains flags for the residual program variables and may be used in place of σ_0 in a subsequent application of I . Where unambiguous we use x as a notational convenience for $\text{Var}(x, op, s, l_x)$.

When a variable x is bound as the formal parameter of an inlined procedure or bound by a `letrec` expression, the op field of $\text{Var}(x, op, s, l_x)$ contains an operand structure $\text{Opnd}(e, \rho, l_e)$, where e is the source expression corresponding to the actual parameter or bound value, ρ is an environment binding free variables of e , and l_e contains the residual code for e or `unvisited` if the operand has not yet been processed.

The remainder of this section explains how the algorithm handles the core forms of the language.

2.1 Constants

Handling constants is straightforward. When processed for boolean value or side effects, constant values are reduced to `true`, `false`, or `void`. This simplifies processing of conditional expressions and simplifies the construction of residual sequence expressions.

2.2 Sequences

The expression `(seq e1 e2)` discards the value of `e1` and returns the value of `e2`. Therefore, `e1` is processed in an `Effect` context, and `e2` is processed in the current context, γ . When processed for effect, constants, primitive references, `lambda` expressions, variable references, and effect-free primitive calls reduce to the constant `void`. Residual sequences (constructed by *seq*, Figure 4) are nested only in the first position. This permits constant-time access (via *result*) to the result expression of a residual sequencing form.

2.3 Conditionals

The guard of a conditional expression is processed for its boolean value using a `Test` context. When processed for boolean value, non-`false` constants, primitive references, `lambda` expressions, and certain effect-free primitive calls reduce to the constant `true`. If the *result* of the residual guard expression is a constant, the flow of control is known at compile time, so a sequence of the guard (which may contain effects) and the selected branch is residualized. Folding conditional tests eliminates unreachable code.

When the flow of control through a conditional expression cannot be determined, and the residual expressions for the two branches are equivalent in the current context, the algorithm residualizes a sequence of the test and one of the branches. For our purposes, two residual expressions are equivalent only if they are identical constants.

It may appear that this transformation is seldom useful. Indeed, a programmer is not likely to write `(if e1 false false)`. The flattening of values in `Test` and `Effect` contexts can produce such cases, however. Consider the following example.

```
(let ((x 3))
  (if (and (read) (= x 0))
      e1
      e2))
```

Copy propagation and constant folding reduce the test `(= x 0)` to `false`. With `(and (read) (= x 0))` treated as `(if (read) (= x 0) false)`, the preceding transformations reduce the entire expression to `(seq (read) e2)`.

2.4 Variable assignments

If a local variable x is unreferenced, then an assignment (**assign** x $e1$) is useless and $e1$ can be processed for effect. Because the algorithm may still be in the process of computing referenced flags for residual-program variables ($\sigma(l_{x'})$) when it encounters the assignment, it bases this decision on the conservative source-program referenced flags (s). If x is referenced in the source program, the expression $e1$ is processed in **Value** context, the **assign** flag is set in $\sigma(l_{x'})$, and an assignment to the α -converted variable x' is residualized.

2.5 Procedure calls

The treatment of **call** illustrates the demand-driven nature of our algorithm. The argument at a call site is not processed before the body of the called procedure. Rather, when the algorithm encounters a call, it constructs an **App** context that contains the current context (the context of the call itself) and an operand structure for the argument at the call site. The expression in operator position is then processed using the newly constructed call context.

The operand structure captures the current environment so that the argument can be processed on demand. Operand structures also play a role in cycle detection described in Section 3. The residual code is cached within the operand so that an argument is processed at most once. The *visit* procedure returns the cached residual code if present. Otherwise it processes the argument in the given context, caching the result. Operands are visited in **Value** context unless they are associated with unreferenced variables, in which case they are visited for effect.

Processing of arguments is deferred until the algorithm encounters references to the corresponding formal parameters or decides to residualize the call. If a call is inlined, the residual code returned is a version of the integrated procedure's body specialized, via recursive application of the inliner, to its argument. If the call has not been inlined, the operand is processed for value using *visit*, and the call is residualized.

Processing on demand allows arguments bound to unreferenced formal parameters of an inlined procedure to be processed for effect rather than value. This often reduces the size and execution time of the residual program and may improve the accuracy of the referenced and assigned information collected for variables in the residual program. For example, if f is `(lambda (x y z) (if (> z 0) (+ z y) x))`, the inlined call `(f e1 e2 3)` reduces to `(seq e1 ((lambda (y) (+ 3 y)) e2))`. Because $e1$ is processed for effect, the residual code for the call may actually be smaller when inlined. If variable references in the source of $e1$ no longer appear in the residual code, expressions bound to those variables may be processed for effect as well.

2.6 Primitive references

The current context determines how the algorithm treats primitive references. In **Test** or **Effect** contexts it residualizes **true** or **void** because a primitive reference cannot signal an error and always has a non-false value. In a **Value** context it simply residualizes the primitive reference. In a call context, the algorithm processes the operand for its value and attempts constant folding when the result is a constant. If successful, the call is marked **inlined** and the new constant is residualized; otherwise the primitive reference is residualized.

2.7 Lambda expressions

The algorithm handles **lambda** expressions much like primitive references. In **Test** and **Effect** contexts **lambda** expressions reduce to the constants **true** and **void**, respectively. In a **Value** context the body is processed in an environment extended with a fresh variable for the formal parameter, and the **lambda** expression is reconstructed with the fresh variable and the optimized body.

In a call context **App**(op, γ, l_γ) the body is processed using the outer context of the call, γ , in an environment mapping the formal parameter x to a fresh residual-program variable x' bound to the operand op from the call site. The s field of x' is initialized with the contents of $\sigma(l_x)$ for the corresponding source-program variable x . That is, the source flags of x' inherit the residual flags of x .

If no reference to x' is residualized while processing the body, then $\text{ref} \notin \sigma_2(l_{x'})$ and the operand can be processed for effect. If x' is also unassigned, then a sequence is constructed from the residualized operand and the specialized body to ensure applicative-order sequencing of effects. Otherwise, a **let** binding (direct **lambda** call) is constructed for variables referenced or assigned in the residual code of the body. These

```

(letrec ((f (lambda (x) (if (zero? x) 1 (* x (f (- x 1))))))
  (g (lambda () (display "hello"))))
  (h (assign x (lambda () g)))
  (a (lambda (x) (b x)))
  (b (lambda (x) (a x))))
(f 5))

```

Figure 5: Sample letrec expression.

transformations effectively inline a version of the procedure specialized to the call site for which the call context was originally constructed.

2.8 Letrec expressions

Handling a **letrec** expression is similar to handling a **lambda** expression in a call context: the body is processed in an extended (circular) environment and useless bindings are pruned. A binding can be pruned only after all expressions in its scope have been processed. Therefore, after processing the body of a **letrec** expression, the algorithm must process operands of as-yet unreferenced variables, such as **g**, **h**, **a**, and **b** in Figure 5. Operands of variables not referenced in the source program are visited for effect. Processing of trivially effect-free operands, such as the **lambda** expression bound to **g**, or operands of already referenced variables such as **f**, is left for the pruning phase. Other operands are conservatively processed for value. The binding for an unassigned and unreferenced variable can then be eliminated provided its operand is effect-free, that is, it does not diverge, perform side effects, or capture a continuation. The binding for **g** is retained because a reference to **g** is residualized when the operand of **h** is conservatively processed for value and the binding for **h** is retained because its operand performs a side effect. The bindings for **a** and **b** are discarded. By visiting operands on demand, the algorithm generally obtains more complete information than would be obtained with any arbitrary processing order.

2.9 Variable references

Variable references are the focal points for copy propagation and inlining. In **Effect** context, a variable reference reduces to the constant **void**. In **Value**, **Test** and call contexts, the variable is renamed, and the operand bound to that variable, if any, is processed for value as a prelude to copy propagation and inlining. The initial processing of the operand conveniently biases the algorithm in favor of inlining at leaves of the call graph.

Copy propagation and inlining are performed by *copy*, defined in Figure 4, which takes the **Var** structure for the renamed variable, the *result* *e* of the processed operand, and the current context, continuation and store. Copies are propagated in cases one, four, and five of the definition. Inlining is attempted in the third case when *e* is a **lambda** expression processed in a call context.

3 Restraining the Algorithm

To simplify the presentation, we have not yet described the mechanisms that bound the run time of the algorithm, limit code growth due to inlining, and detect cycles. There are three mechanisms of interest.

A linear bound is enforced on the run time of the algorithm by maintaining an effort counter whenever an expression is processed more than once, i.e., for each attempt to inline an operand (see Section 2.9). The effort counter advances each time *I* is called. If the effort bound is exceeded, the inlining attempt is aborted and the call is residualized instead. Because the inlining attempt may be aborted, *fold* does not set the inlined flag until processing is complete. The effort counter is not reset for nested integrations, as this could lead to nonlinear processing time. Because the number of call sites in the original source program is fixed and because we bound the amount of effort expended processing each of these sites, the run time of the algorithm is linear in the size of the input program.

The effort counter is an intuitive, no-nonsense mechanism that permits general-purpose optimization and which allows the linear-time constant of proportionality to be scaled. Other methods for bounding the run time of the algorithm may abandon entire classes of useful optimizations. For example, if only calls to leaf routines are inlined, then programs containing η -expansions are not fully optimized.

Code growth is limited via a size counter that tracks the size of the residual code as it is produced. This counter is incremented whenever a form is residualized in the output. If the size threshold is exceeded while attempting to integrate a procedure, the inlining attempt is abandoned and the call is residualized instead. When inlining an operand bound to a variable referenced only once in the source program, no limit is imposed on the size of the residual code. Inlining procedures called only once does not directly increase the size of the residual program and, through specialization, may decrease code size.¹ While the effort counter alone prevents indefinite inline expansion, the size counter prevents excessive growth in any one part of the program.

The effort and size counters are passed as additional arguments to *I*. When constructing a call context, the current effort counter is stored in each operand structure so that the cost of processing the operand may be charged to the original call site. In particular, this allows source-program operands to be processed without an effort counter. The size counter is not stored in the operand structure. Rather, a new counter with no limit is used when the operand is first processed for value. The resulting size count is cached along with the residual code in the operand structure. The cached size is added to the current size counter if a binding to this operand is residualized.

Cycles are detected by consulting additional flags within operand structures to recognize when the value of an operand is requested during the initial processing of that operand. An “outer-pending” flag, tested and set by *copy*, is used to detect recursive references encountered while inlining a procedural operand, as in the expression `((lambda (x) (x x)) (lambda (x) (x x)))`. An “inner-pending” flag, tested and set by *visit*, is used to detect recursive references within an operand as in `(letrec ((f (lambda () (f)))) (f))`. These flags can be extended to govern loop unrolling within a recursive procedure and unfolding at the call site. When a cycle is discovered while processing the operand for a variable, a variable reference is residualized.

The cycle-detection mechanism generally interrupts recursive inlining long before the effort counter is exceeded. This is important because it prevents the algorithm from wasting effort that could otherwise be applied to inlining. To illustrate, suppose that while inlining a call to *f*, the algorithm encounters a call to *g*, where *g* is bound to a recursive procedure. If cycle-detection were disabled, then recursively visiting the operand of *g* would eventually exceed the effort limit and the entire inlining attempt for *f* would be aborted. By detecting the cycle involving *g* early, the algorithm residualizes the call to *g* and resumes its effort to inline *f*.

Polyvariance is achieved by attempting to integrate a procedure at each of its call sites and relying on the size counter to abort any attempts for which the residual code produced would exceed the size limit. Thus, the algorithm bases the inlining decision at each call site on the size of the procedure when specialized to its arguments. In contrast, a monovariant approach bases the decision solely on the size of the procedure without taking into account the arguments. Polyvariance thus permits more aggressive inlining without an increase in the size threshold, resulting in better performance with tighter control over code growth.

4 Extensions

4.1 Primitive handlers

The algorithm in Section 2 folds primitive calls only when all of the operands at the call site residualize with constant result expressions. In fact, by exploiting algebraic properties of primitives, constant folding may be possible even when some operands are not constants. For example, `(member e1 '())` is processed as if it were the equivalent expression `(seq e1 false)`. This provides *member* with the same polyvariant inlining benefits that would automatically result had it been defined as a user procedure.

Our implementation currently performs only a few such transformations. We consider adding a new primitive handler only if: (1) it eliminates the need to process an operand for value, or (2) the transformation

¹Code size increases only if procedural arguments passed to a called-once procedure are inlined.

exposes opportunities for further folding.² The `member` example above satisfies both conditions: `e1` can be processed for effect, and exposing the constant `false` may allow folding of a conditional test.

4.2 Pruning letrec bindings

After processing the body of a `letrec` expression, the algorithm must process operands of as-yet unreferenced variables, such as `g`, `h`, `a`, and `b` in Figure 5. In the simplified algorithm, these operands are conservatively visited for value unless trivially effect-free. Using an effort counter, the algorithm can instead visit operands of as-yet unreferenced variables for effect. To support this, operand structures are extended so that cached residual code produced in effect context can be distinguished from that produced in value context. An operand is processed for effect, and the resulting code cached, only if the operand has not already been processed. If the effort counter is exceeded, the operand is visited conservatively for value. When an operand is visited for value, cached residual code is returned only if it was obtained in value context. Caching ensures that operands of as-yet unreferenced variables are processed at most twice: once for effect using an effort counter, and if necessary, once for value.

4.3 Improving accuracy

In Figure 3, the third case of the rule for variable references visits the operand in value context as a prelude to copy propagation or inlining. When a reference is processed in test context, immediately processing its operand for value may reduce accuracy. Consider the difference between `(if (cons e1 e2) 1 2)` and `(let ((x (cons e1 e2))) (if x 1 2))`. In the first expression, the call to `cons` is processed in test context. Because `cons` cannot return `false`, `e1` and `e2` can be processed for effect, and the conditional test folds. In the second expression, processing the operand of `x` for value causes the call to `cons` to be residualized with `e1` and `e2` processed for value. This prevents the conditional expression from folding and may cause otherwise useless bindings referenced within `e1` and `e2` to be retained.

To solve this problem, the rule for variable references is modified so that when processing a reference in test context, its operand is visited in test context using an effort counter. An effort counter must be used because it may be necessary to revisit the operand in value context if the variable is referenced elsewhere in the program. If the effort limit is exceeded, the operand is processed conservatively in value context. Operand structures are extended to distinguish residual code produced in test context from that produced in value context. When an operand is visited for value, cached residual code is returned only if it was produced in a value context. With these changes in place, the second example above folds to `1` just as the first did.

In the definition of *fold* in Figure 4, the continuation κ_3 constructs the equivalent of a `let` binding when the formal parameter of a called `lambda` expression appears to be referenced or assigned. Where the referenced or assigned information is conservative, this can be improved upon by two simple transformations. First, when the residual code for `e2` is `(const c)`, `(let ((x e1)) e2)` can be treated instead as the equivalent expression `(seq e1 (const c))` because `x` is clearly not used even if it appeared to be while processing `e2`. Similarly, when the residual code for `e3` is `(const c)`, `(let ((x e1)) (seq e2 e3))` can be treated instead as `(seq (let ((x e1)) e2) (const c))`. With these two transformations in place, the following example reduces to `12` instead of a `let` binding with `12` as its body. This allows the constant to be propagated into other contexts where it may lead to constant folding.

```
(let ((f (lambda (g) (g g))))
  (let ((z (lambda (n) (if (< n 10) (lambda (x) 5) (lambda (x) 15)))))
    (let ((y (lambda (m) (m 4))))
      (let ((x (y z)))
        (+ 7 (f x))))))
```

4.4 Inlining recursive procedures

A procedure is inlined by setting a size counter and an effort counter (if not already set) then processing the procedure body in an environment mapping the formal parameters to fresh variables containing operand

²For example, strength reductions that simply replace one call with another are left for a later pass of the compiler.

structures for the actual parameters at the call site. Calls to non-recursive procedures are inlined when the size of the resulting residual code does not exceed the size limit and the effort limit is not exceeded. Calls to recursive procedures are inlined only if the above conditions hold and (1) no recursive calls appear in the residual code, or (2) the procedure body can be specialized to the actual parameters at the call site. Many calls satisfying condition (1) are already inlined as a natural consequence of the algorithm already presented. In this section we describe extensions which allow the algorithm to discover additional cases in which calls satisfy condition (1) and which allow the algorithm to inline calls satisfying condition (2).

To illustrate the first condition, consider the factorial function defined below.

```
(define f
  (lambda (x)
    (if (zero? x)
        1
        (* x (f (- x 1))))))
```

Using the cycle detection mechanisms described in Section 3, the algorithm can recognize that a variable is bound to a recursive operand when that operand is initially processed for value. When attempting to inline a call to the recursive procedure bound to **f**, we first extend the environment so that **f** maps to a fresh variable, **f'**, with no operand, then process the procedure body in the call context. Because **f'** has no operand, recursive calls within the body are forced to residualize (see the last line in definition of *copy*, Figure 4). If, after processing the body, **f'** is unreferenced, then no recursive calls were residualized. That is, the call folded completely, as happens for (**f** 0), satisfying condition (1) above.

If, after processing the body, **f'** is referenced but the fresh formal parameters are not referenced, then copy propagation or inlining successfully replaced references to the formals within the body. This often means the operands of recursive calls to **f** are known. For example, processing the call (**f** 5) in this manner yields (*** 5 (f' 4)**). If some of the operands at the original call site have constant residual expressions, a size counter is set along with modest effort and unfolding counters and we attempt to completely unfold the call. If these limits are not exceeded and no recursive calls remain in the resulting residual code, condition (1) is satisfied. For example, the call (**f** 5) reduces to 120.

If condition (1) cannot be met, a call to a recursive procedure may still be inlined if the procedure body can be specialized to the arguments at the call site. For example, consider the definition of **fold** below.

```
(define fold
  (lambda (f x b null? car cdr)
    (if (null? x)
        b
        (f (car x) (fold f (cdr x) b null? car cdr)))))
```

When the operand bound to **fold** is initially processed, we record the residual code for the operands at each recursive call. When attempting to specialize **fold** we collect the set of formal parameters that are invariant with respect to copy propagation and inlining. For our purposes, a formal parameter is invariant if it is unassigned and if for all recursive calls the corresponding actual parameter is merely a reference to that formal parameter. In the case of **fold**, all formal parameters except **x** are invariant in this sense. The set of invariant formals can be computed in linear time because the residual code for each of the arguments at the recursive call sites is tested at most once. Once computed, this set is cached within the operand structure.

Specialization is attempted if some formals are invariant and the residual expressions for their operands at the original call site are not variable references. This prevents trivial specialization. To specialize the procedure, we extend the environment with formal parameters mapped to fresh variables, set effort and size counters and process the body. The fresh formal parameters contain operands only if the corresponding formal parameters are marked invariant. This ensures that the resulting optimized code is valid for all iterations. If the size and effort limits are not exceeded, the original call is replaced by a call to a local **letrec** binding for the specialized procedure body.

Eliminating unused invariant bindings is analogous to pruning unused **letrec** bindings. Processing of recursive call operands is deferred until the rest of the body has been processed. The operands of variables marked referenced or assigned at this point are visited for value. Operands for as-yet unreferenced variables are visited as described in Section 4.2. Bindings for unused formals are then eliminated from the specialized

Benchmark	Lines	Description
boyer	528	Logic programming benchmark originally by Bob Boyer
dynamic	1,503	Henglein’s dynamic type inferencer [13]
graphs	558	Counts directed graphs with distinguished root and k vertices, each having out-degree at most 2
lattice	214	Enumerates the lattice of maps between two lattices
matrix	585	Tests whether matrix is maximal among all matrices obtained by reordering rows and columns
maze	787	Hexagonal maze generator by Olin Shivers
nbody	1,534	Implementation [17] of Greengard multipole algorithm [12]
splay	971	Implementation of splay trees
conform	395	Type checker by Jim Miller
earley	455	Earley’s parser by Marc Feeley
em-fun	416	Jeffrey Siskind’s functional implementation of an EM clustering algorithm
em-imp	404	Siskind’s imperative implementation of an EM clustering algorithm
interpret	842	Marc Feeley’s Scheme interpreter evaluating <code>tak1</code>
peval	496	Feeley’s simple Scheme partial evaluator
simplex	182	Simplex algorithm

Table 1: Description of benchmarks. Size excludes blank lines and comments.

procedure body, and each recursive call is rewritten as a sequence that evaluates operands of unused formals for effect, then makes the recursive call passing operands only for used formals.

Specializing `(fold * x 1 zero? (lambda (x) x) (lambda (x) (- x 1)))` in this manner yields the same code as the earlier definition of factorial, modulo renaming. Although we do not explicitly perform specialization for mutually recursive functions, inlining often converts mutual recursion into direct recursion [15]. The resulting direct recursion can then be specialized as described above.

5 Performance

We have implemented the inlining algorithm as an additional pass within *Chez Scheme*, a commercial optimizing compiler for Scheme. Although described in Section 2 as an interpreter written in continuation-passing style (CPS), the inlining algorithm is implemented in direct style, relying on the applicative-order semantics of Scheme to sequence effects. When a size or effort counter exceeds its bound, the inlining attempt is abandoned via non-local exit to the point at which the counter was set, and a call is residualized instead. Because the continuations which provide these non-local exits are never reinstated, simple one-shot continuations suffice [8]. Operations on sets of flags are implemented efficiently as operations on bit-vectors. When creating a fresh variable, for example, residual-variable flags are copied to source-variable flags by a single logical shift operation. We do not CPS-convert input programs; doing so would not materially simplify the algorithm and could exaggerate the benefits obtained by inlining [1].

Table 1 describes the benchmarks used to evaluate the inlining algorithm. The benchmarks in the first part of the table were provided by Jagannathan and Wright [14] and have already been processed by their local simplification pass, which performs many of the same optimizations as our inliner. To facilitate comparison with their results, performance data was collected on a 150 MHz SGI MIPS R4400 workstation.

Table 2 shows the performance increase obtained when the benchmarks are optimized using different size and effort bounds. With the exception of `nbody`, all programs show consistent and significant speedups. Profiling reveals that `nbody` spends most of its time indexing three levels deep into a structure of nested arrays. Thus, the variability we observe is likely due to cache effects. When run on the Intel Pentium Pro, which better tolerates cache effects, `nbody` shows consistent speedups ranging from 1.06 to 1.11. As is to be expected, increasing the effort limit generally improves performance somewhat. Both `graphs` and `peval` improve substantially at higher effort limits.

Inlining more than doubles the performance of `lattice`, `graphs`, and `conform`. A large percentage of runtime

Size limit:	Effort limit 100						Effort limit 1000					
	10	20	30	40	50	60	10	20	30	40	50	60
boyer	1.03	1.04	1.04	1.04	1.04	1.04	1.03	1.04	1.06	1.06	1.07	1.06
dynamic	1.05	1.07	1.07	1.08	1.06	1.05	1.08	1.10	1.05	1.05	1.10	1.07
graphs	1.79	1.78	1.78	1.78	1.78	1.77	1.76	2.25	2.24	2.25	2.25	1.59
lattice	4.30	4.37	4.39	4.40	4.39	4.39	4.30	4.57	4.48	4.57	4.56	4.57
matrix	1.29	1.27	1.27	1.29	1.30	1.29	1.31	1.28	1.30	1.31	1.32	1.30
maze	1.45	1.38	1.40	1.35	1.40	1.42	1.46	1.39	1.40	1.42	1.41	1.39
nbody	.95	.97	.98	1.04	1.04	1.05	.92	1.00	.98	.93	1.04	1.03
splay	1.26	1.29	1.30	1.29	1.28	1.28	1.32	1.20	1.33	1.31	1.38	1.35
conform	2.41	2.50	2.62	2.70	2.69	2.66	2.44	2.52	2.47	2.68	2.67	2.69
earley	1.10	1.10	1.08	1.09	1.09	1.10	1.09	1.11	1.12	1.12	1.12	1.12
em-fun	1.15	1.18	1.05	1.02	1.09	1.02	1.20	1.15	1.07	1.15	1.15	1.14
em-imp	1.16	1.10	1.06	1.11	1.15	1.06	1.16	1.09	1.06	1.13	1.05	1.06
interpret	1.07	1.03	1.09	1.09	1.10	1.10	1.08	.99	1.10	1.07	1.07	1.07
peval	1.22	1.07	1.05	1.07	1.07	1.06	1.07	1.21	1.23	1.22	1.23	1.24
simplex	1.32	1.32	1.30	1.32	1.32	1.32	1.31	1.32	1.31	1.32	1.30	1.30

Table 2: Performance increase from inlining with different size and effort bounds.

Size limit:	Effort limit 100						Effort limit 1000					
	10	20	30	40	50	60	10	20	30	40	50	60
boyer	.97	.97	.99	1.00	1.00	1.00	.97	.99	.99	1.01	1.01	1.01
dynamic	.87	1.18	1.20	1.20	1.31	1.31	.86	1.26	1.30	1.37	1.55	1.57
graphs	.72	.71	.71	.71	.71	.71	.80	.85	.85	.85	.85	.90
lattice	1.13	.91	.91	.91	.91	.91	1.13	.89	.89	.89	.89	.89
matrix	.96	1.01	1.03	1.03	1.03	1.03	.94	.99	1.01	1.00	1.00	1.00
maze	.74	.74	.78	.78	.78	.78	.74	.74	.78	.78	.78	.78
nbody	1.20	1.22	1.25	1.25	1.25	1.25	1.24	1.27	1.33	1.34	1.40	1.46
splay	1.00	1.00	1.02	1.06	1.06	1.06	.99	1.01	1.12	1.12	1.10	1.12
conform	.73	.80	.82	.86	.86	.86	.72	.75	.76	.80	.80	.80
earley	.91	.92	.96	.96	.96	.96	.82	.85	.89	.89	.91	.94
em-fun	.88	.92	1.07	1.09	1.09	1.09	.88	.91	1.15	1.17	1.13	1.12
em-imp	.87	.89	.89	.92	.92	.92	.85	.89	.88	.90	.90	.90
interpret	1.19	1.26	1.31	1.31	1.31	1.31	1.17	1.23	1.28	1.29	1.29	1.35
peval	1.01	1.04	1.04	1.04	1.04	1.04	1.00	1.03	1.03	1.03	1.03	1.03
simplex	.86	.86	.86	.86	.86	.86	.86	.86	.86	.86	.90	.90

Table 3: Ratio of code size to original code size with different size and effort bounds. Numbers less than 1 indicate a decrease in code size.

Size limit:	Effort limit 100						Effort limit 1000					
	10	20	30	40	50	60	10	20	30	40	50	60
boyer	.03	.04	.03	.04	.04	.04	.02	.08	.07	.08	.09	.09
dynamic	.33	.48	.52	.50	.54	.53	.34	.56	.57	.68	.81	1.06
graphs	.08	.08	.09	.09	.08	.07	.11	.13	.14	.13	.15	.14
lattice	.08	.08	.09	.08	.07	.08	.09	.11	.12	.10	.11	.12
matrix	.11	.11	.11	.11	.13	.11	.14	.18	.20	.21	.21	.21
maze	.11	.11	.13	.12	.13	.12	.11	.12	.13	.13	.14	.13
nbody	.20	.20	.20	.22	.20	.22	.29	.47	.48	.51	.57	.57
splay	.10	.10	.10	.09	.10	.11	.13	.15	.13	.13	.14	.16
conform	.08	.12	.13	.14	.14	.15	.24	.27	.28	.32	.34	.31
earley	.13	.15	.15	.15	.15	.17	.26	.31	.40	.41	.44	.45
em-fun	.17	.19	.20	.21	.21	.22	.40	.45	.66	.70	.68	.64
em-imp	.16	.18	.18	.18	.19	.18	.34	.45	.46	.50	.54	.51
interpret	.28	.29	.33	.34	.34	.34	.29	.31	.47	.51	.56	.55
peval	.08	.08	.09	.12	.13	.13	.17	.19	.19	.19	.22	.19
simplex	.05	.06	.06	.06	.05	.05	.09	.10	.10	.14	.14	.13

Table 4: Run time of the algorithm (in seconds) with different size and effort bounds.

calls made by `lattice` (73%) and `graphs` (43%) are to a user-defined `memv` procedure that searches a list for a given key. With a size limit of 4 our algorithm completely unfolds calls to `memv` in both programs because the target list happens to be a short constant at each call site for `memv`. The speedup of 4.57 for `lattice` at size limit 20 represents an additional speedup of 1.15 beyond the speedup of 3.97 due to unfolding calls to `memv` at size limit 4. The speedup of 2.25 for `graphs` at size limit 20 represents an additional speedup of 1.43 beyond the speedup of 1.57 due to unfolding calls to `memv` at size limit 4. The algorithm cannot unfold any of the recursive calls within `conform`. The large speedup is due simply to the fact that `conform` contains many small procedures that can be inlined.

Table 3 shows the effect of different size and effort bounds on the size of the resulting object code. The table reports the ratio of code size after optimization to the original code size. For many programs, inlining actually decreases code size. For most programs, code size increases gradually as the size limit increases. This confirms the intuition that the size counter is an effective mechanism for regulating code growth.

Table 4 shows the run time of the algorithm for different size and effort bounds. The algorithm completes both analysis and transformation of each program (except `dynamic`) in less than one second. At the highest settings for size and effort limits, `dynamic` is processed in just over one second. Increasing the effort limit by an order of magnitude typically only doubles the run time of the algorithm. When the size limit is varied, the run time of the algorithm increases roughly in proportion to the increase in code size.

6 Related work

Jagannathan and Wright [14] use a polyvariant flow analysis to identify inlining opportunities and to estimate the size of a procedure when specialized for a particular call site. Their approach is offline: all inlining decisions are made prior to any transformation. Our algorithm combines a less accurate analysis (effectively sub-OCFA [4]) with online transformation and polyvariant specialization. The time required by their flow analysis varies widely and can be excessive—110 seconds to analyze `dynamic`—rendering their method impractical for use in a production compiler.

Ashley [3] evaluated the effectiveness of four different flow analyses for inlining. The analyses range from a fast analysis less accurate than OCFA to a polyvariant analysis similar to 1CFA. His inlining algorithm is based on that of Jagannathan and Wright but exposes more opportunities for inlining.

Table 5 provides a rough comparison of the best speedups obtained by the flow-directed inlining algorithm of Jagannathan and Wright with the best speedups obtained by our algorithm (labeled POLI, for *polyvariant online linear inliner*). The table also compares processing time and the ranges of object code size for the

Performance Increase				Processing Time (seconds)		Code Size Ratios		
Baseline	block	non-block		POLI	PCFA	block	non-block	
Algorithm	POLI	POLI	PCFA			POLI	POLI	PCFA
boyer	1.07	1.34	1.72	.09	1.4	.97 – 1.01	.82 – .86	.95 – 1.10
dynamic	1.10	1.20	1.20	.56	110.3	.86 – 1.57	.78 – 1.43	1.52 – 2.48
graphs	2.25	2.45	1.30	.13	.3	.71 – .90	.65 – .82	.79 – .87
lattice	4.57	6.40	1.27	.11	.2	.89 – 1.13	.72 – .92	1.02 – 3.09
matrix	1.32	1.35	1.11	.21	.4	.94 – 1.03	.86 – .94	.93 – 1.21
maze	1.46	1.88	1.67	.11	.5	.74 – .78	.48 – .51	.89 – .98
nbody	1.04	1.38	1.11	.22	2.6	1.20 – 1.46	1.08 – 1.31	1.25 – 2.62
splay	1.38	1.35	1.20	.14	4.0	.99 – 1.12	.94 – 1.06	.94 – .95

Table 5: Performance increase, processing time and range of code size ratios for our algorithm (POLI) and the polyvariant flow-directed inlining of Jagannathan and Wright (PCFA). The columns labelled “non-block” indicate results relative to a baseline that is not block compiled. All programs were compiled and run using *Chez Scheme*. All data was collected on the 150 MHz SGI MIPS R4400.

two methods. The comparison is complicated by the fact that their optimizer effectively block compiles the benchmarks, yet their baseline for comparison is not block-compiled. Enabling block compilation reduces baseline object-code size and improves baseline performance considerably. To compensate for the difference in benchmarking methodology we present both our absolute results and the results we obtain relative to a baseline that is not block compiled. It appears that we achieve equivalent or better speedups for all benchmarks except *boyer* and better code size ratios for all benchmarks except *splay*, yet our algorithm completes both analysis and transformation in far less time than their analysis alone.

Table 6 compares the performance increase, processing time, and code size ratios for Ashley’s inlining algorithm justified by sub-OCFA and 1CFA analyses with the results obtained by our algorithm with effort limits of 100 and 1000. All data was generated using the same compiler and benchmarking methodology. In particular, the baseline and optimized programs were block-compiled in all cases. The data was collected on virtually identical hardware (our machine has half the physical memory and second-level cache). All results in the table reflect a size threshold of 30. Our algorithm achieves better or equivalent speedups for virtually all of the benchmarks and typically completes both analysis and transformation in less time than is required for flow analysis alone. The volatility of the 1CFA analysis is highlighted by the fact that it exhausted virtual memory (128Mb core, 100Mb swap) while attempting to analyze *interpret*. Our code size ratios are consistently better than those obtained using either flow analysis. This supports our claim that online transformation provides more accurate size information on which to base inlining decisions.

The intuitive appeal of an offline flow-directed approach is that the results of flow analysis might be used to justify other optimizations in addition to inlining. In fact, Ashley found that inlining violates the soundness of flow information about procedures which could otherwise be used by procedure call optimizations in later passes of the compiler. This means that flow analysis must be repeated after inlining.

Our one-pass algorithm differs significantly from the iterated multi-pass algorithm used by Appel in the SML/NJ compiler [1]. Appel uses an offline method in which a data gathering pass is run alternately with a transformation pass that performs several optimizations including constant folding, β -contraction, uncurrying, and inlining of functions called only once. This transformation phase is iterated until the number of contractions found is below some threshold. Next, inline expansion is performed using several heuristics to estimate the cost and benefit of inlining. To improve these estimates, this “round” of optimization is iterated beginning with the contraction phase. Although our algorithm can be iterated, we have found that few programs actually benefit from additional passes.

Appel and Jim recently described a linear-time algorithm that performs constant folding, dead-variable elimination, and inlining of functions called only once [2]. They have implemented an $O(n^2)$ variant of the algorithm to replace the contraction phase of the SML/NJ compiler described above. The new algorithm is online in that usage counts for variables are updated as the optimizations are performed. The algorithm

	Performance Increase				Processing Time (seconds)				Code Size Ratio			
	<0cfa	1cfa	100	1000	<0cfa	1cfa	100	1000	<0cfa	1cfa	100	1000
boyer	1.01	1.03	1.01	1.02	.01	.07	.02	.03	1.47	2.33	.98	1.01
dynamic	1.11	1.05	1.10	1.10	.32	2.09	.20	.23	3.56	4.72	1.46	1.61
graphs	1.20	1.27	1.70	2.18	.05	.20	.03	.05	.73	.82	.68	.72
lattice	1.03	1.33	3.92	4.06	.01	.19	.03	.04	1.20	1.51	.94	.92
matrix	1.05	1.09	1.22	1.21	.05	3.02	.06	.09	1.11	1.23	1.04	1.02
maze	1.27	1.27	1.28	1.29	.05	.19	.03	.04	.88	.89	.75	.75
nbody	1.06	1.10	1.09	1.06	.13	1.98	.08	.19	1.29	1.37	1.19	1.25
splay	1.09	1.11	1.23	1.26	.06	.95	.02	.05	1.42	1.46	1.07	1.18
conform	1.72	1.75	1.94	1.95	.04	.31	.05	.12	1.06	1.19	.86	.81
earley	1.10	1.11	1.07	1.10	.07	.45	.06	.15	1.06	1.13	1.00	.94
interpret	1.12	—	1.14	1.15	.39	—	.13	.19	2.06	—	1.47	1.45
peval	1.15	1.15	1.14	1.15	.07	.41	.04	.07	1.20	1.31	1.04	1.03
simplex	1.39	1.39	1.43	1.43	.03	.19	.03	.04	.89	.92	.81	.81

Table 6: Comparison of performance increase, processing time, and ratio of code size to original code size for two flow-directed algorithms and our algorithm run with effort limits of 100 and 1000. All data was gathered on the 200 MHz Intel Pentium Pro. All programs were compiled and run using *Chez Scheme*. Processing time for <0cfa and 1cfa cases includes analysis time only, i.e., it does not reflect the cost of inlining or subsequent simplification. A size threshold of 30 was used in all cases.

is demand-driven in that bodies of binding forms are processed before processing the bindings. The new algorithm typically performs as many contractions in one pass as their earlier offline algorithm performs in three. They do not consider the more difficult problem of inlining functions called more than once.

Bondorf uses an abstract interpreter written in continuation-passing style (CPS) to improve the accuracy of binding-time analysis in an off-line partial evaluator [7]. The CPS structure of the interpreter enables him to move static contexts into binding forms, thereby reducing the need to CPS-convert source programs as part of binding-time improvement. We achieve a similar effect in direct style via our contexts, which abstract the relevant information about the continuation in which expression values are used.

Biswas [6] describes a demand-driven set-based analysis and uses this to obtain a polynomial-time ($O(n^3)$) algorithm for dead-code elimination in a purely functional higher-order language. No empirical evidence is given to support the effectiveness of the optimization.

Dean and Chambers [11] describe an online approach in which the cost and benefit of inlining at a call site are estimated by examining the code produced when the routine is tentatively inlined and optimized for the call site. To amortize the high cost of these inlining trials, they cache the results in a persistent database indexed by a description of the static information at the call site that justified optimizations that made inlining acceptable. By using inlining trials instead of offline scoring estimates, they were able to reduce overall compile time and code size with minimal loss of run-time performance. Our algorithm can be viewed as incorporating a lightweight form of inlining trials.

Wegman and Zadeck [16] present a fast global constant propagation algorithm and show how it can be extended to perform many of the specializations needed when procedures are inlined. Their algorithm propagates constants through the static single-assignment (SSA) graph [10] using a demand-driven traversal of the control-flow graph that facilitates elimination of unreachable code which, in turn, improves the accuracy of the information collected. The inlining algorithm they describe does not address key issues such as identifying inlining candidates, limiting code growth, and achieving termination in the presence of recursion. Their algorithm is not online, polyvariant, or context-sensitive in the sense we have described.

Ball describes an analysis that determines which parameters contribute to the value of the expressions in a procedure body [5]. When constant parameters are available at a call site he uses the parameter dependency information to guide inlining decisions with an estimate of code savings and performance gain that is based on predictions about the impact of subsequent optimizations. We suspect the speed of our algorithm could be improved by using the type-group analysis of Dean and Chambers, or the more direct parameter dependency

analysis of Ball, to exploit polyvariance only where it is likely to be useful.

7 Conclusion

We have presented a procedure inlining algorithm that is both fast and effective. The algorithm improves performance dramatically for several benchmark programs without sacrificing performance on any. Moreover our algorithm is fast enough to process programs of several thousand lines in under one second on available hardware. The run time of our algorithm is linear in the size of the input program, which is important for scalable optimizations. Because it is both fast and effective, the algorithm is well-suited for use in a production compiler. The algorithm also provides a competent basis for comparison that can be used to isolate the additional benefit, if any, obtained through more involved analysis techniques. Finally, we have presented evidence that online techniques are faster than offline, analysis-intensive techniques, yet are equally effective for inlining.

There are several directions we intend to pursue as future work. Inlining decisions might be improved by investing more effort in frequently executed parts of the program. Burger [9] has developed a framework in which profiling data can be collected and used to dynamically recompile programs as they run. When recompiling, his implementation currently performs only low-level optimizations such as block reordering. We plan to extend the scope of recompilation to include front-end optimizations such as inlining.

Because inlining can simplify control flow and eliminate polymorphism, we intend to see whether the run time of a polyvariant flow analysis can be improved by running our algorithm as a pre-pass. We also want to determine what additional benefit can be obtained by extending our algorithm to exploit the results of flow analysis. We plan to study these effects using the polyvariant flow analysis of Ashley [4].

Our implementation inlines procedures with free variables only when those variables can be eliminated during optimization or when the scope of the variables includes the call site. In the second case, inlining may add new free variables to closures. Although inlining often decreases and never increases the total amount of allocation for the benchmarks we tried, we believe there are cases where increasing closure size can impair performance. We plan to see whether performance can be improved by investigating the interaction between inlining, closure conversion, and lambda-lifting.

Acknowledgements: Mike Ashley and Bob Burger contributed to the development of the inlining algorithm. Mike Ashley, Carl Bruggeman, Bob Burger, and Mark Leone provided helpful comments on drafts of this paper. We thank Suresh Jagannathan and Andrew Wright for providing us with details on their inlining algorithm and the benchmark code they used to test it, and Andy Hanson for allowing us to run benchmarks on his SGI workstation.

References

- [1] Andrew W. Appel. *Compiling with Continuations*. Cambridge University Press, 1992.
- [2] Andrew W. Appel and Trevor Jim. Shrinking lambda expressions in linear time. To appear in *Journal of Functional Programming*.
- [3] J. Michael Ashley. The effectiveness of flow analysis for inlining. To appear in *Proceedings of the 1997 ACM SIGPLAN International Conference on Functional Programming*.
- [4] J. Michael Ashley. A practical and flexible flow analysis for higher-order languages. In *Proceedings of the ACM Symposium on Principles of Programming Languages*, pages 184–194, 1996.
- [5] J. Eugene Ball. Predicting the effects of optimization on a procedure body. *SIGPLAN Notices*, 14(8):214–220, August 1979. *Proceedings of the ACM SIGPLAN '79 Symposium on Compiler Construction*.
- [6] Sandip K. Biswas. A demand-driven set-based analysis. In *Proceedings of the ACM Symposium on Principles of Programming Languages*, pages 372–385, 1997.

- [7] Anders Bondorf. Improving binding times without explicit CPS-conversion. In *Proceedings of the 1992 ACM Conference on LISP and Functional Programming*, pages 1–10, 1992.
- [8] Carl Bruggeman, Oscar Waddell, and R. Kent Dybvig. Representing control in the presence of one-shot continuations. In *Proceedings of the ACM SIGPLAN '96 Conference on Programming Language Design and Implementation*, pages 99–107, 1996.
- [9] Robert G. Burger. *Efficient Compilation and Profile-Driven Recompilation in Scheme*. PhD thesis, Indiana University, 1997.
- [10] R. Cytron, J. Ferrante, B. K. Rosen, M. N. Wegman, and F. K Zadeck. Efficiently computing static single assignment form and the control dependence graph. Technical Report RC 14756, IBM, 1991.
- [11] Jeffrey Dean and Craig Chambers. Towards better inlining decisions using inlining trials. In *Proceedings of the 1994 ACM Conference on LISP and Functional Programming*, pages 273–282, 1994.
- [12] L. Greengard. *The Rapid Evaluation of Potential Fields in Particle Systems*. ACM Press, 1987.
- [13] Fritz Henglein. Global tagging optimization by type inference. In *Proceedings of the 1992 ACM Conference on LISP and Functional Programming*, pages 205–215, 1992.
- [14] Suresh Jagannathan and Andrew Wright. Flow-directed inlining. In *Proceedings of the ACM SIGPLAN '96 Conference on Programming Language Design and Implementation*, pages 193–205, 1996.
- [15] Owen Kaser, C. R. Ramakrishnan, and Shaunak Pawagi. On the conversion of indirect to direct recursion. *ACM Letters on Programming Languages and Systems*, 2(4):151–164, mar 1993.
- [16] M. N. Wegman and F. K. Zadeck. Constant propagation with conditional branches. *ACM Transactions on Programming Languages and Systems*, 3(2):181–210, 1991.
- [17] F. Zhao. An $O(n)$ algorithm for three-dimensional n-body simulations. Master’s thesis, Massachusetts Institute of Technology, 1987.