Eliminating Dead Code on Recursive Data

Yanhong A. Liu^{*} and Scott D. Stoller^{*}

Computer Science Department, Indiana University, Bloomington, IN 47405 {liu,stoller}@cs.indiana.edu

Abstract. This paper describes a general and powerful method for dead code analysis and elimination in the presence of recursive data constructions. We represent partially dead recursive data using liveness patterns based on general regular tree grammars extended with the notion of live and dead, and we formulate the analysis as computing liveness patterns at all program points based on program semantics. This analysis yields a most precise liveness pattern for the data at each program point, which is significantly more precise than results from previous methods. The analysis algorithm takes cubic time in terms of the size of the program in the worst case but is very efficient in practice, as shown by our prototype implementation. The analysis results are used to identify and eliminate dead code. The general framework for representing and analyzing properties of recursive data structures using general regular tree grammars applies to other analyses as well.

1 Introduction

Dead computations produce values that never get used [1]. While programmers are not likely to write code that performs dead computations, such code appears often as the result of program optimization, modification, and reuse [40, 1]. There are also other programming activities that do not explicitly involve live or dead code but rely on similar notions. Examples are program slicing [60, 45], specialization [45], incrementalization [34, 33], and compile-time garbage collection [24, 21, 42, 57]. Analysis for identifying dead code, or code having similar properties, has been studied and used widely [8, 7, 25, 41, 1, 24, 21, 10, 26, 34, 54, 45, 33, 57]. It is essentially backward dependence analysis that aims to compute the minimum sufficient information needed for producing certain results. We call this *dead code analysis*, bearing in mind that it may be used for many other purposes.

In recent years, dead code analysis has been made more precise so as to be effective in more complicated settings [21, 10, 26, 45, 5, 33]. Since recursive data constructions are used increasingly widely in high-level languages [52, 14, 37, 3], an important problem is to identify *partially dead recursive data*—that is, recursive data whose dead parts form recursive substructures—and eliminate computations of them.¹ It is difficult because recursive data structures can be defined by the user, and dead substructures may interleave with live substructures. Several methods have been studied [24, 21, 45, 33], but all have limitations.

This paper describes a general and powerful method for analyzing and eliminating dead computations in the presence of recursive data constructions. We

^{*} The authors gratefully acknowledge the support of NSF under grant CCR-9711253 and ONR under grants N00014-99-1-0132 and N00014-99-1-0358.

¹ This is different from *partial dead code*, which is code that is dead on some but not all computation paths [26, 5].

represent partially dead recursive data using *liveness patterns* based on general regular tree grammars extended with the notion of live and dead, and we formulate the analysis as computing liveness patterns at all program points based on program semantics. This analysis yields a most precise liveness pattern for the data at each program point. The analysis algorithm takes cubic time in terms of the size of the program in the worst case but is very efficient in practice, as shown in our prototype implementation. The analysis results are used to identify and eliminate dead code. The framework for representing and analyzing properties of recursive data structures using general regular tree grammars applies to other analyses as well.

The rest of the paper is organized as follows. Section 2 describes a programming language with recursive data constructions. Section 3 defines liveness patterns that represent partially dead recursive data. Section 4 formulates the analysis as solving sets of constraints on grammars. Section 5 presents efficient algorithms for computing liveness patterns at every program point. Section 6 describes dead code elimination, our implementation, and extensions. Section 7 compares this work with related work and concludes.

2 Language

We use a simple first-order functional programming language. The expressions of the language are:

variable constructor application primitive function application conditional expression binding expression function application

Each constructor c, primitive function p, and user-defined function f has a fixed arity. If a constructor c has arity 0, then we write c instead of c(). New constructors can be declared, together with their arities. When needed, we use c^n to denote that c has arity n. For each constructor c^n , there is a primitive function c? that tests whether the argument is an application of c, and if n > 0, then for each i = 1..n, there is a primitive function $c_{\overline{i}}$ that selects the *i*th component in an application of c, e.g., $c_{\overline{2}}(c^3(x, y, z)) = y$. A program is a set of mutually recursive function definitions of the form:

$$f(v_1, \dots, v_n) \stackrel{\Delta}{=} e \tag{1}$$

together with a set of constructor declarations. Figure 1 gives some example definitions, assuming that min and max are primitive functions, and that constructors nil^0 , $cons^2$, and $triple^3$ are declared in the programs where they are used. For ease of reading, we use *null* instead of nil?, car instead of $cons_{\overline{1}}$, cdr instead of $cons_{\overline{2}}$, and 1st, 2nd, and 3rd instead of $triple_{\overline{1}}$, $triple_{\overline{2}}$, and $triple_{\overline{3}}$, respectively.

This language has call-by-value semantics. Well-defined expressions evaluate to constructed data, such as cons(3, nil). We use \perp to denote the value of undefined (non-terminating) expressions; an expression must evaluate to \perp if any of its subexpressions evaluates to \perp . Since a program can use data constructions

| $\begin{array}{ll} minmax(x) & : \text{ compute min and max for all suffixes of } x \\ minmax(x) & \stackrel{\Delta}{=} \text{ if } null(x) \text{ then } nil \\ \text{ else if } null(cdr(x)) \text{ then } \\ & cons(triple(car(x), car(x), car(x)), nil) \\ \text{ else let } v = minmax(cdr(x)) \text{ in } \\ & cons(triple(car(x), car(x), car(x)), nil) \\ \end{array}$ | $\begin{array}{rl} len(x) & : \text{ compute length of } x \\ len(x) & \stackrel{\Delta}{=} \mathbf{if} \; null(x) \; \mathbf{then} \; 0 \\ & \mathbf{else} \; 1 + len(cdr(x)) \end{array}$ |
|--|---|
| $\min(car(x), 2nd(car(v)))), \ \max(car(x), 3rd(car(v)))), \ v)$ | odd(x) : get elements of x at odd positions even(x) : get elements of x at even positions |
| $\begin{array}{ll} list second(x) &: \text{ list the second element in each triple in } x \\ list second(x) &\stackrel{\Delta}{=} \textbf{if } null(x) \textbf{ then } nil \\ \textbf{else } cons(2nd(car(x)), list second(cdr(x))) \end{array}$ | $odd(x) \stackrel{\Delta}{=} \begin{array}{l} \text{if } null(x) \text{ then } nil \\ \text{else } cons(car(x), \\ even(cdr(x))) \end{array}$ |
| $\begin{array}{ll} getmin(x) & : \text{ compute the min elements for all suffixes of } x \\ getmin(x) & \stackrel{\Delta}{=} listsecond(minmax(x)) \end{array}$ | $even(x) \stackrel{\Delta}{=} \begin{array}{l} \mbox{if } null(x) \mbox{ then } nil \\ \mbox{else } odd(cdr(x)) \end{array}$ |

Fig. 1. Example function definitions.

 $c(e_1, ..., e_n)$ in recursive function definitions, it can build data structures of unbounded sizes, i.e., sizes not bounded in any way by the size of the program but determined by particular inputs to the program.

There can be values, which can be subparts of constructed data, computed by a program that are not needed in obtaining the output of the program. To improve program efficiency, we can eliminate such dead computations and use a special symbol _ as a placeholder for their values. A constructor application does not evaluate to _ even if some arguments evaluate to _. A primitive function application (or a conditional expression) must evaluate to _, if not \perp , if any of its subexpressions (or the condition, respectively) evaluates to _. Whether a function application (or a binding expression) evaluates to _ depends on the values of the arguments (or the bound variable, respectively) and how they are used in the function definition (or the body, respectively).

Dead code may exist in a program especially when only certain parts of its result or intermediate results are needed. Such parts can be specified by a user or determined by how these results are used in computing other values, e.g., by how the value of *minmax* is used in computing *getmin*, in which case all the max operations are dead.

3 Liveness patterns

We represent partially dead recursive data using liveness patterns. A liveness pattern indicates which parts of data must be dead and which parts may be live. D indicates that a part must be completely dead, and L indicates that a part may be completely live. Partial liveness is represented using constructors. For example, cons(D, L) indicates a cons structure with a definitely dead head and a possibly live tail. Also, nil() indicates the liveness pattern corresponding to a nil structure, so there is no confusion between a liveness pattern and a data value. A liveness pattern is a function; when applied to a data value, it returns the data with the live parts unchanged and the dead parts replaced by $_$. For example,

 $cons(D, cons(L, D)) \left(cons(0, cons(1, cons(2, nil))) \right) = cons(_, cons(1, _)).$

Formally, liveness patterns are domain projections [48, 17], which provide a clean tool for describing substructures of constructed data by projecting out the parts that are of interest [56, 27, 39, 45, 33]. Let X be the domain of all possible values computed by our programs, including \perp and values containing $_$. We define an ordering \sqsubseteq on X, where we read $x_1 \sqsubseteq x_2$ as " x_1 is more dead than x_2 ": for all x in X, $\perp \sqsubseteq x$, and for two values x_1 and x_2 other than \perp ,

$$x_1 \sqsubseteq x_2 \quad \text{iff} \ x_1 = _, \text{ or } x_1 = x_2, \text{ or} \\ x_1 = c(x_{11}, \dots, x_{1n}), \ x_2 = c(x_{21}, \dots, x_{2n}), \text{ and } x_{1i} \sqsubseteq x_{2i} \text{ for } i = 1..n.$$

$$(2)$$

A liveness pattern over X is a function $\pi : X \to X$ such that $\pi(x) \sqsubseteq x$ and $\pi(\pi(x)) = \pi(x)$ for all $x \in X$. L is the identity function: L(x) = x. D is the absence function: $D(x) = _$ for all $x \neq \bot$, and $D(\bot) = \bot$. $c^n(\pi_1, ..., \pi_n)$ is the function:

$$c^{n}(\pi_{1},...,\pi_{n})(x) = \begin{cases} c^{n}(\pi_{1}(x_{1}),...,\pi_{n}(x_{n})) \text{ if } x = c^{n}(x_{1},...,x_{n}) \\ \bot & \text{otherwise} \end{cases}$$
(3)

Grammar-based liveness patterns. We represent liveness patterns as grammars. For example, the grammar $S \rightarrow nil() | cons(D, S)$ projects out a list whose elements are dead but whose spine is live. It generates the set of sentences $\{nil(), cons(D, nil()), cons(D, cons(D, nil())), ...\}$. Applying each element to a given value, say, cons(2, cons(4, nil)), yields $\bot, \bot, cons(_, cons(_, nil)), \bot, ...,$ in which $cons(_, cons(_, nil))$ is the least upper bound.

Formally, the grammars we use for describing liveness patterns are regular tree grammars [16], which allow bounded, and often precise, representations of unbounded data [23, 38, 39, 2, 51, 9, 45]. A regular-tree-grammar-based liveness pattern G is a quadruple $\langle \mathcal{T}, \mathcal{N}, \mathcal{P}, S \rangle$, where \mathcal{T} is a set of terminal symbols including L, D, and all possible constructors c, \mathcal{N} is a set of nonterminal symbols N, \mathcal{P} is a set of productions of the form:

$$N \to D, \quad N \to L, \quad \text{or} \quad N \to c^n(N_1, \dots, N_n),$$
(4)

and nonterminal S is the start symbol. So, our liveness patterns use general regular tree grammars [23, 16, 9] extended with the special constants D and L. The language \mathcal{L}_G generated by G is the set $\{\pi \in \mathcal{T}^* \mid S \xrightarrow{*}_G \pi\}$ of sentences. The projection function that G represents is:

$$G(x) = \sqcup \{ \pi(x) \mid \pi \in \mathcal{L}_G \}$$
(5)

where \sqcup is the least upper bound operation for \sqsubseteq . It is easy to see that G(x) is well-defined for all $x \in X$. We overload L to denote the grammar that generates sentence L, and overload D to denote the grammar that generates only sentence D. For ease of presentation, when no confusion arises, we write grammars in compact forms. For example, $\{S \to nil() | cons(N, S), N \to triple(D, L, D)\}$, where | denotes alternation, projects out a list whose elements are triples whose first and third components are dead.

We extend regular tree grammars to allow productions of the form:

$$N \to N', \quad N \to c_i(N'), \quad \text{or} \quad N \to [N'] R'$$
 (6)

for R' of the form $L, c^n(N_1, ..., N_n)$, or N'', and we define:

$$c_{\overline{i}}(\pi) = \begin{cases} L & \text{if } \pi = L \\ \pi_i & \text{if } \pi = c^n(\pi_1, \dots, \pi_n) \\ D & \text{otherwise} \end{cases} \quad \text{and} \quad [\pi]\pi' = \begin{cases} D & \text{if } \pi = D \\ \pi' & \text{otherwise} \end{cases}$$
(7)

These extended forms are for convenience later; the selector form in the middle of (6) is the same as that first used by Jones and Muchnick [23], and the conditional form on the right of (6) is for similar purposes as those used in several other analyses, e.g., the operator \triangleright used by Wadler and Hughes for strictness analysis [56]. Given an extended regular tree grammar G that contains productions of the forms in (4) and (6), we can construct a regular tree grammar G'that contains only productions of the form (4) such that $\mathcal{L}_G = \mathcal{L}_{G'}$, i.e., G' and G represent the same projection function; an algorithm is given in Section 5.

When using a grammar G, what matters is the projection function that G represents. In fact, different grammars can represent the same projection function. A basic idea of this work is to capture the information of interest—liveness patterns—using grammars that are constructed based on program semantics and then simplify the grammars to equivalent grammars in simpler forms where, in particular, the only grammar that represents D is $\{S \rightarrow D\}$.

We define an ordering \leq on regular-tree-grammar-based liveness patterns. For two grammars G_1 and G_2 , we define:

$$G_1 \leq G_2 \quad \text{iff} \quad \forall \pi_1 \in \mathcal{L}_{G_1}, \ \exists \pi_2 \in \mathcal{L}_{G_2}, \ \pi_1 \leq \pi_2, \tag{8}$$

where for two sentences π_1 and π_2 , we define (overloading \leq):

$$\pi_1 \leq \pi_2$$
 iff $\pi_1 = D$, or $\pi_2 = L$, or
 $\pi_1 = c(\pi_{11}, ..., \pi_{1n}), \ \pi_2 = c(\pi_{21}, ..., \pi_{2n}), \ \text{and} \ \pi_{1i} \leq \pi_{2i} \ \text{for} \ i = 1..n.$

For convenience, we define $G_1 \ge G_2$ if and only if $G_2 \le G_1$. It is easy to see that

if
$$G_1 \leq G_2$$
, then $\forall x, G_1(x) \sqsubseteq G_2(x)$. (10)

This means that if $G_1 \leq G_2$, then G_1 projects out values that are more dead than those that G_2 projects out. This is a basis of our correctness proof. The converse is not true, e.g., $G_1 = \{S \rightarrow cons(L, L)\}$ and $G_2 = \{S \rightarrow cons(L, D) \mid cons(D, L)\}$ form a counterexample for the converse, but this converse is not used. Note that this ordering on grammars does not form a complete partial order.

Notation. We use CON to denote the grammar that projects any constructor but none of its arguments: letting \mathcal{T}_c be the set of all possible constructors,

$$CON = \langle \mathcal{T}_c \cup \{D\}, \{S\}, \{S \to c^n(\overbrace{D, \dots, D}^n) \mid c^n \in \mathcal{T}_c\}, S \rangle^2$$
(11)

Given a grammar $G = \langle \mathcal{T}, \mathcal{N}, \mathcal{P}, S \rangle$, we use $con_{c,i}(G)$ to denote using G as the *i*th component of a c^n structure $(i \leq n)$ whose other components are dead: assuming S' is a nonterminal not used in G,

$$\underline{con_{c,i}(G)} = \langle \mathcal{T}, \ \mathcal{N} \cup \{S'\}, \ \mathcal{P} \cup \{S' \to c^n(\overbrace{D,...,D}^{i-1}, S, \overbrace{D,...,D}^{n-i})\}, \ S'\rangle,$$
(12)

² For convenience, we fold D into the right sides of the productions.

and we use $sel_{c,i}(G)$ to denote the part of G corresponding to the *i*th component of a c^n structure $(i \leq n)$: assuming S' is a nonterminal not in used in G,

$$sel_{c,i}(G) = \langle \mathcal{T}, \ \mathcal{N} \cup \{S'\}, \ \mathcal{P} \cup \{S' \to c_i(S)\}, \ S' \rangle.$$
 (13)

For example, if *nil* and *cons* are all possible constructors where a liveness pattern *CON* is used, as we assume for functions *len*, *odd*, and *even*, then $CON = \{S \rightarrow nil() | cons(D, D)\}$. If $G = \{S \rightarrow L\}$, then $con_{cons,1}(G) = \{S' \rightarrow cons(S, D), S \rightarrow L\}$ and $sel_{cons,1}(G) = \{S' \rightarrow car(S), S \rightarrow L\}$. Finally, we define a conditional:

$$cond(G_1, G_2) = \begin{cases} D & \text{if } G_1 = D \\ G_2 & \text{otherwise.} \end{cases}$$
(14)

4 Analysis of liveness patterns using constraints

Dead code analysis computes liveness patterns associated with values at *program points*, such as function definitions, parameters, and subexpressions. We develop such a backward dependence analysis. Given liveness patterns associated with certain program points, it computes liveness patterns at all program points, so that the liveness specified by the given liveness patterns is guaranteed. The basic idea is that a liveness pattern associated with a program point is constrained by liveness patterns associated with other points based on the semantics of the program segments involved.

Sufficiency conditions. The resulting liveness patterns must satisfy two kinds of sufficiency conditions. First, the resulting grammar at a program point must project out values that are more live than required by the given grammar (if any) associated with that point. Precisely, at each subexpression e where a liveness pattern is given, if the given grammar is G_0 , and the resulting grammar is G, then $G_0(e) \sqsubseteq G(e)$ for all values of the free variables in e. Second, the resulting grammars must satisfy the constraints determined by the program semantics. Precisely, assume a resulting grammar is associated with each parameter and each subexpression of all function definitions. Let Ge denote that grammar G is associated with e. Then (1) the liveness patterns at function parameters must be sufficient to guarantee the liveness pattern at the function return, i.e., for each definition of the form $f(G^{I}v_1, \ldots, G^{In}v_n) \triangleq Ge$, the following sufficiency condition must be satisfied for all values v_1, \ldots, v_n :

$$G(f(v_1, ..., v_n)) \sqsubseteq f(G_1(v_1), ..., G_n(v_n))$$
(15)

and (2) the liveness pattern at each subexpression must be sufficient to guarantee the liveness pattern at the enclosing expression, i.e., for each subexpression that is of a form in the left column below, the corresponding sufficiency condition in the right column must be satisfied for all values of the free variables in the subexpression:

 $\begin{array}{lll} G^{i}c(^{G_{1}i}e_{1},...,^{G_{n}i}e_{n}) & G(c(e_{1},...,e_{n})) \sqsubseteq c(G_{1}(e_{1}),...,G_{n}(e_{n})) \\ G^{i}c_{\overline{i}}(^{G'}e') & G(c_{\overline{i}}(e')) \sqsubseteq c_{\overline{i}}(G'(e')) \\ G^{i}c_{\overline{i}}(^{G'}e') & G(c_{\overline{i}}(e)) \sqsubseteq c_{\overline{i}}(G'(e')) \\ G^{i}q(^{G_{1}i}e_{1},...,^{G_{n}i}e_{n}) & \text{if } q \text{ is } p \text{ other than } c_{\overline{i}} \text{ or } c_{1}^{2} \\ & G(q(e_{1},...,e_{n})) \sqsubseteq q(G_{1}(e_{1}),...,G_{n}(e_{n})) \\ G^{i}\text{if } G^{1}ie_{1} \text{ then } G^{2}ie_{2} \text{ else } G^{3}ie_{3} & G(\text{if } e_{1} \text{ then } e_{2} \text{ else } e_{3}) \sqsubseteq \text{ if } G_{1}(e_{1}) \text{ then } G_{2}(e_{2}) \text{ else } G_{3}(e_{3}) \\ G^{i}\text{let } u = ^{G_{1}i}e_{1} \text{ in } ^{G_{2}i}e_{2} & G(\text{let } u = e_{1} \text{ in } e_{2}) \sqsubseteq \text{ let } u = G_{1}(e_{1}) \text{ in } G_{2}(e_{2}) \\ G^{i}f(^{G_{1}i}e_{1},...,^{G_{n}i}e_{n}) & G(f(e_{1},...,e_{n})) \sqsubseteq f(G_{1}(e_{1}),...,G_{n}(e_{n})) \end{array}$

Note that no approximation is made in these conditions. For example, the condition of a conditional expression does not have to be evaluated, so it does not have to be associated with L. In particular, the liveness patterns associated with a function application are not related to liveness patterns associated with the definition of the function, and thus, different applications of the same function may require different parts of the function definition to be live. For example, consider functions f and g below:

$$f(x,y) \stackrel{\Delta}{=}$$
 if $x > 0$ then x else y $g(z) \stackrel{\Delta}{=} {}^{G_0:} f({}^{L:}1, {}^{D:}z) + f({}^{L:}0, {}^{L:}z)$

Given $G_0 = L$ at the definition of g, the liveness patterns associated with all program points, where the liveness patterns not explicitly written are all L, satisfy the sufficiency conditions. Note that the two calls to f need different parts of f to be live.

Grammar constraints. Given liveness patterns associated with certain subexpressions, we construct a set of constraints on the resulting liveness patterns that guarantee the sufficiency conditions. First, at each subexpression e where a liveness pattern is given, if the given grammar is G_0 , and the resulting grammar is G, then we construct $G \ge G_0$. Second, for a function definition of the form $f({}^{G_1}v_1,...,{}^{G_n}:v_n) \triangleq e$, we construct, for i = 1..n and for each occurrence ${}^{G'_i}v_i$ in e, the constraint

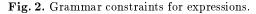
$$G_i \ge G'_i \tag{16}$$

and, for each subexpression of e that is of a form in the left column of Figure 2, we construct the corresponding constraints in the right column. These constraints make approximations while guaranteeing the sufficiency conditions, as explained below.

Formula (16) for function definitions requires that the liveness pattern at formal parameter v_i be greater than or equal to the liveness patterns at all uses of v_i . Rule (7) for function calls requires that, for all non-dead calls of the same function, the liveness patterns at the arguments be greater than or equal to the liveness patterns at the corresponding formal parameters, and that the liveness pattern for the return value of the function be greater than or equal to the liveness patterns at all calls. Thus, if a function call is dead, then all its arguments are also dead, even though the formal parameters of the function may not be dead due to other calls to the same function.

Other constraints are based on the semantics of each construct locally. Rules (1)-(3) handle data constructions. Rule (1) says that liveness pattern at a component of a construction must be no less than the corresponding component in

```
\begin{array}{ll} (1) \ \ ^{G_{i}}c(^{G_{1}}i_{e_{1}},...,^{G_{n}}i_{e_{n}}) & G_{i} \geq sel_{c,i}(G) \ \ \text{for} \ \ i=1..n \\ (2) \ \ ^{G_{i}}c_{i}^{-}(^{G'}i'_{e'}) & G' \geq cond(G,con_{c,i}(G)) \\ (3) \ \ ^{G_{i}}c_{i}^{-}(^{G'}i'_{e'}) & G' \geq cond(G,CON) \\ (4) \ \ ^{G_{i}}q(^{G_{1}}i_{e_{1}},...,^{G_{n}}i_{e_{n}}) & \text{if} \ q \ \text{is} \ p \ \text{other than} \ \ c_{i}^{-} \ \text{or} \ c_{i}^{2} \\ & G_{i} \geq cond(G,L) \ \text{for} \ \ i=1..n \\ (5) \ \ ^{G_{i}}i_{e_{1}} \ \text{then} \ \ ^{G_{2}}i_{e_{2}} \ \text{else} \ \ ^{G_{3}}i_{e_{3}} & G_{1} \geq cond(G,L), \ G_{2} \geq G, \ \ G_{3} \geq G \\ (6) \ \ ^{G_{i}}\text{let} \ \ u = \ ^{G_{1}}i_{e_{1}} \ \text{in} \ \ ^{G_{2}}i_{e_{2}} & G_{1} \geq cond(G,G') \ \text{for each occurrence of} \ \ ^{G'_{i}}u \ \text{in} \ e_{2}, \ \ G_{2} \geq G \\ (7) \ \ ^{G_{i}}f(^{G_{1}}i_{e_{1}},...,^{G_{n}i}e_{n}) \ \ \text{where} \ f(\frac{^{G'_{1}}i_{v_{1}},...,\frac{^{G'_{n}}i_{v_{n}}}{G_{i}} \geq cond(G,G'_{i}) \ \text{for} \ i = 1..n, \ \ G'_{2} \geq G \\ (7) \ \ ^{G_{i}}f(^{G_{1}}i_{e_{1}},...,^{G_{n}i}e_{n}) \ \ \text{where} \ f(\frac{^{G'_{1}}i_{v_{1}},...,\frac{^{G'_{n}}i_{v_{n}}}{G_{i}} \geq cond(G,G'_{i}) \ \text{for} \ i = 1..n, \ \ G'_{2} \geq G \\ \end{array}
```



the liveness pattern at the result of the construction. As a special case of (1), for any constructor of arity 0, no constraint is added. Rule (2) requires that, if the result of a selection by c_i is not dead, then the argument be as live as a construction using c whose *i*th component is as live as the result of the selection. Rule (3) says that, if the result of an application of a tester is not dead, then the liveness pattern at the argument needs to project out the outermost constructor but none of the components. Rule (4) says that, if the result of a primitive operation is not dead, then each argument must be live. If we assume that primitive functions are defined only on primitive types such as Boolean and integer, then we could use CON in place of L in the constraints. Rule (5) requires that the condition be live if the result of the entire conditional expression is not dead, and that both branches be as live as the result. Again, we could use CON in place of L as a sufficient context for e_1 ; furthermore, if e_2 equals e_3 , in fact as long as $G(e_2)$ equals $G(e_3)$, then we could use D in place of L as a sufficient context for e_1 and thus no constraint for G_1 would be needed. Rule (6) is similar to a function call, since it equals an application of $\lambda u.e_2$ to e_1 . It requires that the defining expression e_1 be as live as all its uses in the body, and that the body be as live as the result.

We can show by a standard inductive argument that the constraints for each construct guarantee sufficient information, and thus an inductive proof shows that the sufficiency conditions are satisfied.

5 Construction and simplification of liveness patterns

We describe a straightforward method for building minimum grammars that satisfy the above constraints; these grammars may contain productions of the extended forms in (6). Then, we simplify the grammars by eliminating extended forms; this makes explicit whether the grammar associated with a program point equals dead.

Constructing the grammars. Let \mathcal{T}_c be the set of all possible constructors in the program. Let \mathcal{N}_0 be the set of nonterminals used in the given liveness patterns associated with selected subexpressions. We associate a unique nonterminal, not in \mathcal{N}_0 , with each parameter and each subexpression of all function definitions.

Then we add productions using these terminals and nonterminals. Finally, the resulting grammar at a program point is formed by using these terminals, non-terminals, and productions, and by using the nonterminal associated with that point as the start symbol.

We add two kinds of productions. For each subexpression e where a grammar G_0 is given, let N_0 be the start symbol of G_0 , and let N be the nonterminal associated with e. We add $N \to N_0$ as well as all productions in G_0 . Second, for each function definition $f({}^{N_1}v_1, ..., {}^{N_n} : v_n) \triangleq e$, we add, for each i = 1..n and for each occurrence ${}^{N'_i}v_i$ in e, the production

$$N_i \to N'_i$$
 (17)

and, for each subexpression of e that is of a form in the left column of Figure 3, the corresponding productions in the right column.

| $(1) \ ^{N} c (\ ^{N_{1}} e_{1},, \ ^{N_{n}} e_{n})$ | $N_i \rightarrow c_i^-(N)$ for $i = 1n$ |
|--|---|
| (2) ${}^{N}c_{\overline{i}}({}^{N'}e)$ | $N' \rightarrow [N] c^n (D,, D, N, D,, D)$ with $i - 1$ D's before N and $n - i$ D's after N |
| $(3)^{-N} c?(^{N'}e)$ | $N' \to [N] c^n(D,, D)$ with $n \ D$'s for each possible constructor c^n |
| (4) ${}^{N}q({}^{N_1}e_1,,{}^{N_n}e_n)$ if q is p other t | |
| | $N_i \rightarrow [N] L$ for $i = 1n$ |
| (5) N if $^{N_{1}:}e_{1}$ then $^{N_{2}:}e_{2}$ else $^{N_{3}:}e_{3}$ | $N_1 \rightarrow [N] L, N_2 \rightarrow N, N_3 \rightarrow N$ |
| (6) ^N let $u = {}^{N_1}e_1$ in ${}^{N_2}e_2$ | $N_1 \rightarrow [N] N'$ for each occurrence of $N'^{:}u$ in $e_2, N_2 \rightarrow N$ |
| (7) ${}^{N}f({}^{N_1}e_1,,{}^{N_n}e_n)$ where $f({}^{N_1'}v_1,,{}^{N_n}e_n)$ | |
| | $N_i \rightarrow [N] N_i'$ for $i = 1n, N' \rightarrow N$ |
| | |

Fig. 3. Productions added for expressions.

It is easy to show that the resulting grammars satisfy the grammar constraints in Figure 2 and thus give sufficient information at every program point. To show this, simply notice that the productions in Figure 3 can be obtained from the constraints in Figure 2 by replacing G with N and \geq with \rightarrow , and by replacing grammar operations with the corresponding productions based on definitions: $sel_{c,i}(G)$ with $c_{\overline{i}}(N)$, $con_{c,i}(G)$ with c(D, ..., D, N, D, ..., D), CON with c(D, ..., D) for all c, and $cond(G_1, G_2)$ with $[N_1]N_2$. Thus, each production constructed here guarantees exactly a corresponding grammar constraint in Figure 2 simply by definitions. Furthermore, the resulting grammars are minimal among all solutions that use the same set of nonterminals, and they give minimum sufficient information. To see this, notice that a smaller grammar at any point would make the nonterminal at that point correspond to a smaller grammar than the grammar generated by the right hand side(s) of the nonterminal, violating the corresponding grammar constraints.

Let *n* denote the size of the program. Assume that the maximum arity of constructors, primitive functions, and user-defined functions is bounded by a constant. Since a constant number of productions are added at each program point, the above construction takes O(n) time.

Example 1. For functions len, odd, and even in Figure 1, the nonterminals labeling the program points and the added productions are shown in Figure 4. For example, we have $N_{12} \rightarrow [N_{13}] cons(N_{13}, N_0)$, where $N_0 \rightarrow D$ is the last production on the last line. It means that N_{12} is conditioned on N_{13} : if N_{13} is D, so is N_{12} , otherwise N_{12} is $cons(N_{13}, N_0)$, i.e., it projects out a cons structure, the first component of which is projected out by N_{13} . Suppose we need the result of

 $\begin{array}{rcl} len(^{N_{29}}x) & \stackrel{\Delta}{=} & ^{N_{28}}\text{if} & ^{N_{27}}null(^{N_{26}}x) & \textbf{then} & ^{N_{25}}0 \\ & & \textbf{else} & ^{N_{24}N_{23}}1 + ^{N_{22}}len(^{N_{21}}cdr(^{N_{20}}x)) \end{array}$ $\begin{array}{rcl} odd(^{N_{19}}x) & \stackrel{\Delta}{=} & ^{N_{18}} \, \mathbf{if} \, ^{N_{17}} null(^{N_{16}}x) \, \mathbf{then} \, ^{N_{15}} nil \\ & \mathbf{else} \, ^{N_{14}} cons(^{N_{13}} car(^{N_{12}}x),^{N_{11}} even(^{N_{10}} cdr(^{N_{9}}x))) \end{array}$ $\begin{array}{rcl} even(^{N_8}x) & \stackrel{\Delta}{=} & ^{N_7} \operatorname{if} {}^{N_6}null(^{N_5}x) \operatorname{then} {}^{N_4}nil \\ & & \\ & & \\ & & \\ \operatorname{else} {}^{N_3}odd(^{N_2}cdr(^{N_1}x)) \end{array}$ $\begin{array}{l} N_{29} \rightarrow N_{26}, \; N_{29} \rightarrow N_{20}, \; N_{26} \rightarrow [N_{27}] \; cons(N_0,N_0), \\ N_{26} \rightarrow [N_{27}] ril(), \; N_{27} \rightarrow [N_{28}]L, \; N_{25} \rightarrow N_{28}, \\ N_{23} \rightarrow [N_{24}]L, \; N_{20} \rightarrow [N_{21}] \; cons(N_0,N_{21}), \; N_{21} \rightarrow [N_{22}]N_{29}, \; N_{28} \rightarrow N_{22}, \; N_{22} \rightarrow [N_{24}]L, \; N_{24} \rightarrow N_{28}, \\ \end{array}$ $\begin{array}{l} N_{19} \rightarrow N_{16}, \ N_{19} \rightarrow N_{12}, \ N_{19} \rightarrow N_{9}, \ N_{16} \rightarrow [N_{17}] \ cons(N_0, N_0), \ N_{16} \rightarrow [N_{17}] \ nil(), \ N_{17} \rightarrow [N_{18}] L, \\ N_{15} \rightarrow N_{18}, \ N_{12} \rightarrow [N_{13}] \ cons(N_{13}, N_0), \ N_{13} \rightarrow car(N_{14}), \ N_{9} \rightarrow [N_{10}] \ cons(N_0, N_{10}), \ N_{10} \rightarrow [N_{11}] N_8, \\ N_7 \rightarrow N_{11}, \ N_{11} \rightarrow cdr(N_{14}), \ N_{14} \rightarrow N_{18}, \end{array}$ $\begin{array}{l} N_8 \rightarrow N_5, \ N_8 \rightarrow N_1, \ N_5 \rightarrow [N_6] \ cons(N_0, N_0), \ N_5 \rightarrow [N_6] \ nil(), \ N_6 \rightarrow [N_7] \ L, \ N_4 \rightarrow N_7, \\ N_1 \rightarrow [N_2] \ cons(N_0, N_2), \ N_2 \rightarrow [N_3] \ N_{19}, N_{18} \rightarrow N_3, \ N_3 \rightarrow N_7, \ N_0 \rightarrow D \end{array}$

Fig. 4. Productions constructed for the example functions.

len; we add $N_{28} \rightarrow L$, since N_{28} corresponds to the return value of len. Suppose we need the result of odd; we add $N_{18} \rightarrow L$. Suppose we need to know whether the result of odd is nil or cons; we add $N_{18} \rightarrow nil(), N_{18} \rightarrow cons(D, D)$.

Simplifying the grammars. The grammars obtained above may contain productions of the extended forms in (6) and thus be difficult to understand and use. We simplify the grammars by removing extended forms using an iterative algorithm given in Figure 5. After the simplification, nonterminals that do not appear on

input: a grammar $\langle \mathcal{T}, \mathcal{N}, \mathcal{P}, S \rangle$

```
/* assume R is of the form L or c(N_1, ..., N_n), and R' is of the form L, c(N_1, ..., N_n), or N'' */
repeat
```

if \mathcal{P} contains $N \to N'$ and $N' \to R$, then add $N \to R$ to \mathcal{P} ;

if \mathcal{P} contains $N \to c_{\overline{i}}(N')$ and $N' \to L$, then add $N \to L$ to \mathcal{P} ; if \mathcal{P} contains $N \to c_{\overline{i}}(N')$ and $N' \to c^n(N_1, ..., N_n)$, then add $N \to N_i$ to \mathcal{P} ; if \mathcal{P} contains $N \to [N']R'$ and $N' \to R$, then add $N \to R'$ to \mathcal{P} ;

until no more productions can be added;

remove all productions of the extended forms from \mathcal{P} ;

return simplified grammar $\langle \mathcal{T}, \mathcal{N}, \mathcal{P}, S \rangle$

Fig. 5. Algorithm for simplifying the grammars.

the left side of a production with L or $c(N_1, ..., N_n)$ on the right side are implied

to derive only D. We can read off the grammar at any function parameter or subexpression by starting at the associated nonterminal and collecting all productions whose left sides are reachable from this start symbol. The correctness of this algorithm is based on the definitions of the extended forms and can be proved in a similar way to when only the selector form is used [23].

Nonterminals are associated with program points, so there are O(n) of them. Each step adds a production of the form $N \to L$, $N \to c(N_1, ..., N_k)$, or $N \to N'$. Since each right side of the form $c(N_1, ..., N_k)$ is among the right sides of the originally constructed grammar, there are at most O(n) of them. Thus, for each nonterminal, at most O(n) productions are added. So totally at most $O(n^2)$ productions are added. Adding a production has O(n) cost to check what other productions to add. Thus, the overall simplification takes $O(n^3)$ time. Although this appears expensive, the analysis is very fast in practice, as shown by our prototype implementation.

Example 2. Suppose we need the result of *len* and therefore added $N_{28} \rightarrow L$; we obtain the productions

 $\begin{array}{l} N_{29} \rightarrow nil(), \ N_{29} \rightarrow cons(N_0,N_0), \ N_{29} \rightarrow cons(N_0,N_{21}), \ N_{28} \rightarrow L, \ N_{27} \rightarrow L, \ N_{26} \rightarrow nil(), \\ N_{26} \rightarrow cons(N_0,N_0), \ N_{25} \rightarrow L, \ N_{24} \rightarrow L, \ N_{23} \rightarrow L, \ N_{22} \rightarrow L, N_{21} \rightarrow nil(), \ N_{21} \rightarrow cons(N_0,N_0), \\ N_{21} \rightarrow cons(N_0,N_{21}), \ N_{20} \rightarrow cons(N_0,N_{21}) \end{array}$

Suppose we need the result of *odd* and therefore added $N_{18} \rightarrow L$; we obtain the productions

 $\begin{array}{l} N_{19} \rightarrow cons(N_0,N_0), \ N_{19} \rightarrow nil(), \ N_{19} \rightarrow cons(N_{13},N_0), \ N_{19} \rightarrow cons(N_0,N_{10}), \ N_{18} \rightarrow L, \\ N_{17} \rightarrow L, \ N_{16} \rightarrow nil(), \ N_{16} \rightarrow cons(N_0,N_0), \ N_{15} \rightarrow L, \ N_{14} \rightarrow L, \ N_{13} \rightarrow L, \ N_{12} \rightarrow cons(N_{13},N_0), \\ N_{11} \rightarrow L, \ N_{10} \rightarrow nil(), \ N_{10} \rightarrow cons(N_0,N_0), \ N_{10} \rightarrow cons(N_0,N_2), \ N_{9} \rightarrow cons(N_0,N_{10}), \\ N_8 \rightarrow cons(N_0,N_0), \ N_8 \rightarrow nil(), \ N_8 \rightarrow cons(N_0,N_2), \ N_7 \rightarrow L, \ N_6 \rightarrow L, \ N_5 \rightarrow nil(), \\ N_5 \rightarrow cons(N_0,N_0), \ N_4 \rightarrow L, \ N_3 \rightarrow L, \ N_2 \rightarrow cons(N_0,N_{10}), \ N_2 \rightarrow cons(N_{13},N_0), \ N_2 \rightarrow nil(), \\ N_2 \rightarrow cons(N_0,N_0), \ N_1 \rightarrow cons(N_0,N_2) \end{array}$

Suppose we added $N_{18} \rightarrow nil(), N_{18} \rightarrow cons(D, D)$; we obtain the productions

 $\begin{array}{l} N_{19} \rightarrow cons(N_0,\,N_0), \ N_{19} \rightarrow nil(), \ N_{18} \rightarrow nil(), \ N_{18} \rightarrow cons(N_0,\,N_0), \ N_{17} \rightarrow L, \ N_{16} \rightarrow nil(), \\ N_{16} \rightarrow cons(N_0,\,N_0), \ N_{15} \rightarrow cons(N_0,\,N_0), \ N_{15} \rightarrow nil(), \ N_{14} \rightarrow cons(N_0,\,N_0), \ N_{14} \rightarrow nil() \end{array}$

In each case, other nonterminals derive only D.

The resulting grammars can be further simplified by minimization [16], but minimization is not needed for identifying dead code, since minimization does not affect whether a nonterminal derives only D.

6 Dead code elimination, implementation, and extensions

Consider all function parameters and subexpressions whose associated liveness patterns are D. These parts of the program are dead, so we eliminate them by replacing them with $_$.

Example 3. Suppose we need to know whether the result of *odd* is *nil* or *cons* and therefore added $N_{18} \rightarrow nil(), N_{18} \rightarrow cons(D, D)$; eliminating dead code based on the simplified grammar in Example 2, where N_1 to N_{13} all have only D on the right hand sides, yields:

$$dd(x) \stackrel{\Delta}{=} \operatorname{if} null(x) \operatorname{then} nil$$

else $cons(-, -)$ (18)

Suppose we need the result of function *getmin*, given in Figure 1; analyzing and eliminating dead code yields the following function along with functions *listsecond* and *getmin*:

```
 \min(x) \stackrel{\Delta}{=} \quad \text{if } null(x) \text{ then } nil \\ \text{else if } null(cdr(x)) \text{ then } cons(triple(\_, car(x), \_), nil) \\ \text{else let } v = \min(car(x), 2nd(car(v))), -), v) 
 (19)
```

As another example, if the result of *minmax* is used as argument to *len* instead of *listsecond*, then our algorithm finds that the entire *triple* constructions in *minmax* are dead. However, if the result of *minmax* is used as argument to *odd*, then none of the subexpressions in *minmax* is dead, since the triple is used in every odd recursive call.

Our dead code elimination preserves semantics in the sense that, if the original program terminates with a value, then the new program terminates with the same value.

Two further optimizations are possible but need further study. First, minmax in (19) can be further optimized by removing the *triple* constructions and selectors. Second, when the result of minmax is used as argument to odd, there is no dead code in minmax, but the triple in every even call is indeed dead. One needs to unfold the definition of minmax to remove such dead computations.

Eliminating dead code may improve efficiency in many ways. First, the resulting programs can run faster and use less space. Additionally, compilation of the optimized programs takes less time and also less space, which is especially desirable when using libraries. Furthermore, smaller programs are easier to understand and maintain, yielding higher software productivity.

Implementation. We have implemented the analysis in a prototype system. The implementation uses the Synthesizer Generator [44]. The algorithm for simplifying the grammars is written in the Synthesizer Generator Scripting Language, STk, a dialect of Scheme, and consists of about 300 lines of code. Other parts of the system support editing of programs, display of nonterminals at program points, construction of grammars, highlighting of dead code, etc., and consist of about 3000 lines of SSL, the Synthesizer Generator Specification Language. All the grammars for the examples in this paper are generated automatically using the system.

We have used the system to analyze dozens of examples. The lengths of those programs range from dozens of lines to over a thousand lines. The analysis, although written in STk, is very efficient. Our original motivation for studying this general problem was for identifying appropriate intermediate results to cache and use for incremental computation [33]. There, we propose a method, called cacheand-prune, that first transforms a program to cache all intermediate results, then reuses them in a computation on incremented input, and finally prunes out cached values that are not used. Reusing cached values often produces asymptotic speedup, but leaving in unused values can be extremely inefficient. The analysis method studied in this paper, when adopted for pruning, is extremely effective. The pruned programs consistently run faster, use less space, and are smaller in code size. We also used the analysis for eliminating dead code in deriving incremental programs [34]. There, the speedup is often asymptotic. For example, dead code elimination enables incremental selection sort to improve from $O(n^2)$ time to O(n) time. Figure 6 summarizes the experimental results for a number of examples. Program minmax is as in Figure 1. Programs incsort and incout [34] are derived incremental programs for selection sort and outer product, respectively, where dead code after incrementalization is to be eliminated. Programs cachebin and cachelcs [31] are dynamic-programming programs transformed from straightforward exponential-time programs for computing binomial coefficients and longest common subsequences, respectively, with intermediate results cached, reused, and to be pruned. Program calend is a collection of calendrical calculation functions [12], and program takr is a 100-function version of TAK that tries to defeat cache memory effects [47].

| program name | functions of interest | total program points | dead program points | | number of initial productions | | time (ms) | analysis time (ms) (excl. GC) |
|--------------------|-----------------------------------|----------------------------|---------------------------|----------------|-------------------------------------|------------------------|---|---|
| minmax minmax | getlen getmin | 81 81 | $50 \\ 32$ | | 111 111 | $90 \\ 150$ | $\begin{array}{c} 0.007 \\ 0.016 \end{array}$ | $\begin{array}{c} 0.005 \\ 0.011 \end{array}$ |
| incsort incout | sort' out' | 108 117 | 84 62 | | 143 151 | 34 78 | 0.006 | 0.005 |
| cachebin | bin | 74 | 7 | 67 | 90 | 114 | 0.014 | 0.010 |
| cachelcs calend | lcs gregorian2absolute | 101 1551 | 12 1359 | 89 192 | | 206 229 | 0.038 | 0.033 |
| calend | islamic-date | 1551 | 1205 | 346 | 1839 | 419 | 0.083 | 0.069 |
| calend calend | eastern-orthodox-Xmas yahrzeit | $1551 \\ 1551$ | $1176 \\ 1123$ | $375 \\ 428$ | 1839 1839 | $ 461 \\ 485 $ | $\begin{array}{c} 0.086\\ 0.086\end{array}$ | |
| takr takr | run-takr tak99 | $2804 \\ 2804$ | 0 4 | $2804 \\ 2800$ | | $2805 \\ 2801$ | $0.403 \\ 0.419$ | $0.304 \\ 0.310$ |

Fig. 6. Experimental results.

When using dead code analysis for incrementalization and for pruning unused intermediate results, there is always a particular function of interest, shown in Figure 6. For general programs, especially libraries, such as the calend example, there may not be a single function that is of interest, so we have applied the analysis on several different functions of interested.

The size of a program is precisely captured by the total number of program points, which for most programs is about twice the number of lines of code. The number of dead program points depends on both the program and the function of interest. For example, for libraries, such as the **calend** program, much dead code is found, whereas for **takr**, all 100 functions other than the driver function **run-takr**, are involved in calling each other. Our highlighting allows us to easily see the resulting live or dead slices. For example, for several functions in the **calend** program, only the slice for date, not year or month, is needed. We can see the number of initial productions is roughly linear in the size of the given program, and the number of resulting productions is roughly linear in the number of live program points.

The analysis time for simplifying the grammars, in milliseconds, is measured on an Ultra 10 with 299MHz CPU and 124 MB main memory. We can see that the analysis time is roughly linear in the number of live program points. This is important, especially for analyzing libraries, where being linear in the size of the entire program is clearly not good. We achieved this high efficiency by a careful but simple optimization in our simplification algorithm: after adding a new production, we consider only productions in extended forms whose righthand sides use the left-hand side symbol of the new production. This makes the analysis proceed in an incremental fashion, and only program points that are not dead are followed.

To summarize, our method produces precise analysis results as desired. The analysis is also very fast compared with other reported analyses using constraints. For example, Heintze's analysis takes on the order of seconds for programs of 100 lines to over 1000 lines [19].

Figure 7 is a screen dump of the system on a small example of four functions and constructors nil and cons. Program points are annotated with nonterminals highlighted in red. The shaded region contains the function of interest. The two sets of productions are the original set and resulting set. Dead code (function bigfun as well as the first argument of cons in function f) is highlighted in green.



Fig. 7. A prototype implementation.

Extensions. We believe that our method for dead code analysis can be extended to handle side effects. The extension is to use graph grammars instead of tree grammars. The ideas of including L and D as terminals, constructing grammars based on program points as well as the semantics of program constructs connecting these points, and doing grammar simplifications are the same. Recent work by Sagiv, Reps, and Wilhelm [46] uses graph grammars for shape analysis. We believe we can make similar use of graph grammars for dead code analysis in the presence of destructive updates.

Our method can also be extended to handle higher-order functions in two ways, and we have worked out this extension in the second way. First, we can simply apply a control-flow analysis [50] before we do dead code analysis. This allows our method to handle complete programs that contain higher-order functions. Second, we can directly construct productions corresponding to function abstraction and application and add rules for simplifying them. This is similar to how Henglein [20] addresses higher-order binding-time analysis and how Heintze [19] handles higher-order functions for analyzing sets of values for ML programs. Similar use of constraints has been studied for stopping deforestation for higherorder programs [49]. Our extension adds two constraints/productions for each lambda expression and uses two additional rules for simplification; it is not yet implemented. Handling higher-order functions does not increase the time complexity of our algorithms. In fact, for a language with higher-order functions but not recursive data construction, the constraints may be simplified in worst-case almost linear time [20].

Our method is described here for an untyped language, but the analysis results provide an important kind of type information; the analysis may also be adopted to enhance soft typing; and the analysis applies to typed languages as well. For example, consider the third set of productions in Example 2. The grammar at each program point gives its liveness together with the shape of data. Dead code should be reported to the programmer before, or at least at the same time as, type errors such as 3rd(cons(1,2)) in the dead code. Live code may have its type inferred by small refinements of our rules. For example, if we replace L by Boolean for the condition in rule (5) of Figure 2, we have $N_{17} \rightarrow$ Boolean in the third set of productions in Example 2, and thus everything there is precisely typed. Finally, for a typed language, possible values are restricted also by type information, so the overall analysis results can be more precise, e.g., type information about the value of an expression e can help restrict the grammar at e when e is the argument of a primitive function c?.

7 Related work and conclusion

Our backward dependence analysis uses liveness patterns, which are domain projections, to specify sufficient information. Wadler and Hughes use projections for strictness analysis [56]. Their analysis is also backward but seeks necessary rather than sufficient information, and it uses a fixed finite abstract domain for all programs. Launchbury uses projections for binding-time analysis of partially static data structures in partial evaluation [27]. It is a forward analysis equivalent to strictness analysis and uses a fixed finite abstract domain as well [28]. Mogensen [39], De Niel, and others [11] also use projections, based on grammars in particular, for binding-time analysis and program bifurcation, but they use only a restricted class of regular tree grammars. Another kind of analysis is escape analysis [42, 13, 4], but existing methods can not express as precise information as we do.

Several analyses are in the same spirit as ours. The necessity interpretation by Jones and Le Métayer [24] uses necessity patterns that correspond to a restricted class of liveness patterns. Necessity patterns specify only heads and tails of list values. The absence analysis by Hughes [21] uses contexts that correspond to a restricted class of liveness patterns. Even if it is extended for recursive data types, it handles only a finite domain of list contexts where every head context and every tail context is the same. The analysis for pruning by Liu, Stoller, and Teitelbaum [33] uses projections to specify specific components of tuple values and thereby provide more accurate information. However, methods used there for handling unbounded growth of such projections are crude. Wand and Siveroni's recent work [58] discusses safe elimination of dead variables but does not handle data constructions. Our method of replacing all dead code (including dead variables) by a dummy constant _ is simple, direct, and more general than their method; in particular, it is safe to simply remove dead function parameters.

The idea of using regular tree grammars for program flow analysis is due to Jones and Muchnick [22], where it is used mainly for shape analysis and hence for improving storage allocation. It is later used to describe other data flow information such as types and binding times [38, 39, 2, 11, 59, 51, 45]. In particular, the analysis for backward slicing by Reps and Turnidge [45] explicitly adopts regular tree grammars to represent projections. It is closest in goal and scope to our analysis. However, it uses only a limited class of regular tree grammars, in which each nonterminal appears on the left side of one production, and each right side is one of five forms, corresponding to L, D, atom, pair, and atom | pair. It forces grammars to be deterministic in a most approximate way, and it gives no algorithms for computing the least fixed point from the set of equations. Our work uses general regular tree grammars extended with L and D. We also use productions of extended forms to make the framework more flexible. We give efficient algorithms for constructing and simplifying the grammars. Compared with [45], we also handle more program constructs, namely, binding expressions and user-defined constructors of arbitrary arity.

Our treatment is rigorous, since we have adopted the view that regular-treegrammar-based program analysis is also abstract interpretation and approximations can be built into the grammar transformers as a set of constraints [9]. We extend the grammars and handle L and D specially in grammar manipulations. The result can also be viewed as using program-based finite grammar domains for yielding precise and efficient analysis methods. Another standard way to obtain the analysis result is to do a fixed point computation using general grammar transformers on potentially infinite grammar domains and use approximation operations to guarantee termination. Approximation operations provide a more general solution and make the analysis framework more modular and flexible [9]. In a separate paper [30], we describe three approximation operations that together produce significantly more precise analysis results than previous methods. Each operation is efficient, but due to their generality and interaction, that work does not have an exact characterization of the total number of iterations needed. The finite domains described in this work make a complete analysis easy, and it yields a most precise liveness pattern for the data at each program point.

Regular-tree-grammar-based program analysis can be reformulated as setconstraint-based analysis [18, 19, 9], but we do not know any work that treats precise and efficient dead code analysis for recursive data as we do. Melski and Reps [35, 36] show the interconvertibility of a class of set constraints and context-free-language reachability and, at the end of [35], they show how general CFL-reachability can be applied nicely to program slicing. That essentially addresses the same problem we do with a similar framework, but their description is sketchy, with little discussion about correctness and with no results from implementation, experiments, or applications.

The method and algorithms for dead code elimination studied here have many applications: program slicing and specialization [60, 45], strength reduction, finite differencing, and incrementalization [7, 41, 34, 32], caching intermediate results

for program improvement [33], deforestation and fusion [55, 6], as well as compiletime garbage collection [24, 21, 42, 57]. The analysis results also provide a kind of type information.

The overall goal of this work is to analyze dead data and eliminate computations of them across recursions and loops, possibly interleaved with wrappers such as classes in object-oriented programs. This paper discusses techniques for recursion. The basic ideas should extend to loops. Pugh and Rosser's work has started this direction; it extends slicing to symbolically capture particular iterations in a loop [43]. Object-oriented programming is used widely, but crossclass optimization heavily depends on inlining, which often causes code blow-up. Grammar-based analysis and transformation can be applied to methods across classes without inlining. A direct application would be to improve techniques for eliminating dead data members, as noted by Sweeney and Tip [53]

Even though this paper focuses on dead code analysis and dead code elimination for recursive data, the framework for representing recursive substructures using general regular tree grammars and the algorithms for computing them applies to other analyses and optimizations on recursive data as well, e.g., binding-time analysis for partial evaluation [27, 39]. We have recently developed a binding-time analysis using the same framework.

References

- 1. A. V. Aho, R. Sethi, and J. D. Ullman. Compilers, Principles, Techniques, and Tools. Addison-Wesley, Reading, Mass., 1986.
 2. A. Aiken and B. R. Murphy. Static type inference in a dynamically typed lan-
- guage. In Conference Record of the 18th Annual ACM Symposium on Principles of Programming Languages. ACM, New York, Jan. 1991.
 3. K. Arnold and J. Golsing. The Java Programming Language. Addison-Wesley,
- Reading, Mass., 1996.
- 4. B. Blanchet. Escape analysis: correctness proof, implementation and experimental results. In Conference Record of the 25th Annual ACM Symposium on Principles
- of Programming Languages, pages 25-37. ACM, New York, Jan. 1998.
 5. R. Bodík and R. Gupta. Partial dead code elimination using slicing transformations. In Proceedings of the ACM SIGPLAN '97 Conference on Programming Language Design and Implementation, pages 159-170. ACM, New York, June 1997.
- 6. W.-N. Chin. Safe fusion of functional expressions. In LFP 1992 [29], pages 11-20. 7. J. Cocke and K. Kennedy. An algorithm for reduction of operator strength. Commun. ACM, 20(11):850-856, Nov. 1977.
- 8. J. Cocke and J. T. Schwartz. Programming languages and their compilers; preliminary notes. Technical report, Courant Institute of Mathematical Sciences, New York University, 1970.
- 9. P. Cousot and R. Cousot. Formal language, grammar and set-constraint-based program analysis by abstract interpretation. In Proceedings of the 7th International Conference on Functional Programming Languages and Computer Architecture, pages 170–181. ACM, New York, June 1995.
- 10. R. Cytron, J. Ferrante, B. K. Rosen, M. M. Wegman, and F. K. Zadeck. Efficiently computing static single assignment form and the control dependence graph. ACM Trans. Program. Lang. Syst., 13(4):451-490, Oct. 1991.
- 11. A. De Niel, E. Bevers, and K. De Vlaminck. Program bifurcation for a polymorphically typed functional language. In Proceedings of the Symposium on Partial Evaluation and Semantics-Based Program Manipulation, pages 142–153. ACM, New York, June 1991.
- 12. N. Dershowitz and E. M. Reingold. Calendrical calculations. Software-Practice and Experience, 20(9):899-928, Sept. 1990.
- 13. A. Deutsch. On the complexity of escape analysis. In Conference Record of the 24th Annual ACM Symposium on Principles of Programming Languages, pages 358-371. ACM, New York, Jan. 1997.

- 14. R. K. Dybvig. The Scheme Programming Language. Prentice-Hall, Englewood Cliffs, N.J., 1987.
- 15. Proceedings of the 4th International Conference on Functional Programming Languages and Computer Architecture. ACM, New York, Sept. 1989.
- F. Gecseg and M. Steinb. Tree Automata. Akademiai Kiado, Budapest, 1984.
 C. A. Gunter. Semantics of Programming Languages. The MIT Press, Cambridge,
- Mass., 1992.
- 18. N. Heintze. Set-Based Program Analysis. PhD thesis, Department of Computer Science, Carnegie Mellon University, Pittsburgh, Pennsylvania, Oct. 1992.
- 19. N. Heintze. Set-based analysis of ML programs. In Proceedings of the 1994 ACM Conference on LISP and Functional Programming, pages 306-317. ACM, New York, June 1994.
- 20. F. Henglein. Efficient type inference for higher-order binding-time analysis. In Proceedings of the 5th International Conference on Functional Programming Languages and Computer Architecture, volume 523 of Lecture Notes in Computer Science, pages 448-472. Springer-Verlag, Berlin, Aug. 1991.
- J. Hughes. Compile-time analysis of functional programs. In D. Turner, edi-tor, Research Topics in Functional Programming, pages 117-153. Addison-Wesley, Reading, Mass., 1990.
- 22. N. D. Jones and S. S. Muchnick. Flow analysis and optimization of LISP-like structures. In Conference Record of the 6th Annual ACM Symposium on Principles
- of Programming Languages, pages 244-256. ACM, New York, Jan. 1979.
 23. N. D. Jones and S. S. Muchnick. Flow analysis and optimization of LISP-like structures. In S. S. Muchnick and N. D. Jones, editors, Program Flow Analysis, pages 102-131. Prentice-Hall, Englewood Cliffs, N.J., 1981.
 24. S. B. Jones and D. Le Métayer. Compile-time garbage collection by sharing anal-
- ysis. In FPCA 1989 [15], pages 54-74.
- 25. K. Kennedy. Use-definition chains with applications. J. Comput. Lang., 3(3):163-179, 1978.
- 26. J. Knoop, O. Rüthing, and B. Steffen. Partial dead code elimination. In Proceed-ings of the ACM SIGPLAN '94 Conference on Programming Language Design and Implementation, pages 147-158. ACM, New York, June 1994.
- 27. J. Launchbury. Projection Factorisations in Partial Evaluation. PhD thesis, Department of Computing, University of Glasgow, Glasgow, Scotland, 1989.
- I. Launch of Computing, Conversion of Chargew, Bension, Bession, 1990.
 J. Launch bury. Strictness and binding-time analysis: Two for the price of one. In Proceedings of the ACM SIGPLAN '91 Conference on Programming Language Design and Implementation, pages 80–91. ACM, New York, June 1991.
 Proceedings of the 1992 ACM Conference on LISP and Functional Programming.
- ACM, New York, June 1992.
- 30. Y. A. Liu. Dependence analysis for recursive data. In Proceedings of the IEEE 1998 International Conference on Computer Languages, pages 206–215. IEEE CS Press, Los Alamitos, Calif., May 1998. 31. Y. A. Liu and S. D. Stoller. Dynamic programming via static incrementalization.
- In Proceedings of the 8th European Symposium on Programming, volume 1576 of Lecture Notes in Computer Science, pages 288–305. Springer-Verlag, Berlin, Mar.
- 1999. 32. Y. A. Liu, S. D. Stoller, and T. Teitelbaum. Discovering auxiliary information for incremental computation. In Conference Record of the 23rd Annual ACM Symposium on Principles of Programming Languages, pages 157-170. ACM, New York, Jan. 1996. 33. Y. A. Liu, S. D. Stoller, and T. Teitelbaum. Static caching for incremental com-
- putation. ACM Trans. Program. Lang. Syst., 20(3):546-585, May 1998.
- 34. Y. A. Liu and T. Teitelbaum. Systematic derivation of incremental programs. Sci. Comput. Program., 24(1):1-39, Feb. 1995.
- 35. D. Melski and T. Reps. Interconvertibility of set constraints and conext-free language reachability. In Proceedings of the 1997 ACM SIGPLAN Symposium on Partial Evaluation and Semantics-Based Program Manipulation. ACM, New York, June 1997. 36. D. Melski and T. Reps. Interconvertibility of a class of set constraints and context-
- free language reachability. Theoret. Comput. Sci., 1999. Accepted. 37. R. Milner, M. Tofte, and R. Harper. The definition of Standard ML. The MIT
- Press, Cambridge, Mass., 1990.

- P. Mishra and U. Reddy. Declaration-free type checking. In Conference Record of the 12th Annual ACM Symposium on POPL, pages 7-21. ACM, New York, Jan. 1985.
- T. Mogensen. Separating binding times in language specifications. In FPCA 1989 [15], pages 12-25.
- 40. S. S. Muchnick and N. D. Jones, editors. *Program Flow Analysis: Theory and Applications*. Prentice-Hall, Englewood Cliffs, N.J., 1981.
- 41. R. Paige and S. Koenig. Finite differencing of computable expressions. ACM Trans. Program. Lang. Syst., 4(3):402-454, July 1982.
- 42. Y. G. Park and B. Goldberg. Escape analysis on lists. In Proceedings of the ACM SIGPLAN '92 Conference on Programming Language Design and Implementation, pages 116-127. ACM, New York, June 1992.
- 43. W. Pugh and E. Rosser. Iteration space slicing and its application to communication optimization. In International Conference on Supercomputing, Vienna, Austria, July 1997.
- 44. T. Reps and T. Teitelbaum. The Synthesizer Generator: A System for Constructing Language-Based Editors. Springer-Verlag, New York, 1988.
- T. Reps and T. Turnidge. Program specialization via program slicing. In O. Danvy, R. Glück, and P. Thiemann, editors, *Proceedings of the Dagstuhl Seminar on Partial Evaluation*, volume 1110 of *Lecture Notes in Computer Science*, pages 409-429. Springer-Verlag, Berlin, 1996.
 M. Sagiv, T. Reps, and R. Wilhelm. Solving shape-analysis problems in languages
- M. Sagiv, T. Reps, and R. Wilhelm. Solving shape-analysis problems in languages with destructive updating. ACM Trans. Program. Lang. Syst., 20(1):1-50, Jan. 1998.
- 47. The internet scheme repository. http://www.cs.indiana.edu/scheme-repository/.
- D. S. Scott. Lectures on a mathematical theory of computation. In M. Broy and G. Schmidt, editors, *Theoretical Foundations of Programming Methodology*, pages 145-292. D. Reidel Publishing Company, 1982.
- 49. H. Seidl and M. H. Sørensen. Constraints to stop deforestation. Sci. Comput. Program., 32:73-107, 1998.
- O. Shivers. Control flow analysis in scheme. In Proceedings of the ACM SIGPLAN '88 Conference on Programming Language Design and Implementation. ACM, New York, June 1988.
- 51. M. H. Sørensen. A grammar-based data-flow analysis to stop deforestation. In S. Tison, editor, CAAP'94: Proceedings of the 19th International Colloquium on Trees in Algebra and Programming, volume 787 of Lecture Notes in Computer Science, pages 335-351. Springer-Verlag, Berlin, Apr. 1994.
- 52. G. L. Steele. Common Lisp the Language. Digital Press, 1984.
- 53. P. F. Sweeney and F. Tip. A study of dead data members in c++ applications. In Proceedings of the ACM SIGPLAN '98 Conference on Programming Language Design and Implementation, pages 324-332. ACM, New York, June 1998.
- F. Tip. A survey of program slicing techniques. Journal of Programming Languages, 3(3):121-189, Sept. 1995.
- 55. P. Wadler. Deforestation: Transforming programs to eliminate trees. *Theoret. Comput. Sci.*, 73:231-248, 1990. Special issue of selected papers from the 2nd European Symposium on Programming.
- 56. P. Wadler and R. J. M. Hughes. Projections for strictness analysis. In Proceedings of the 3rd International Conference on Functional Programming Languages and Computer Architecture, volume 274 of Lecture Notes in Computer Science, pages 385-407. Springer-Verlag, Berlin, Sept. 1987.
- M. Wand and W. D. Clinger. Set constraints for destructive array update optimization. In Proceedings of the IEEE 1998 International Conference on Computer Languages, pages 184-193. IEEE CS Press, Los Alamitos, Calif., May 1998.
 M. Wand and I. Siveroni. Constraint systems for useless variable elimination. In
- M. Wand and I. Siveroni. Constraint systems for useless variable elimination. In Conference Record of the 26th Annual ACM Symposium on Principles of Programming Languages, pages 291-302. ACM, New York, Jan. 1999.
- E. Wang and P. M. Hilfinger. Analysis of recursive types in lisp-like languages. In LFP 1992 [29], pages 216-225.
- 60. M. Weiser. Program slicing. IEEE Trans. Softw. Eng., SE-10(4):352-357, July 1984.