Code optimization

Advanced Compiler Construction Michel Schinz – 2020-03-26

Optimization

Goal: rewrite the program to a new one that is:

- behaviorally equivalent to the original one,
- better in some respect e.g. faster, smaller, more energy-efficient, etc. Optimizations can be broadly split in two classes:
- **machine-independent optimizations** are high-level and do not depend on the target architecture,
- **machine-dependent optimizations** are low-level and depend on the target architecture.

This lesson: machine-independent, rewriting optimizations.

The importance of IRs

IRs and optimizations

Intermediate representations (IRs) have a dramatic impact on optimizations, which generally work in two steps:

1. the program is analyzed to find optimization opportunities,

2. the program is rewritten based on the analysis. The IR should make both steps as easy as possible.

Case 1: constant propagation

Consider the following program fragment in some imaginary IR: x $\, \leftarrow \, 7$

••• Outoctic

Question: can all occurrences of x be replaced by 7? Answer: it depends on the IR:

- if it allows multiple assignments, no (further data-flow analyses are required),
- if it disallows multiple assignment, yes!

Case 2: inlining

Inlining replaces a call to a function by a copy of the body of that function, with parameters replaced by the actual arguments.

The IR used also has a dramatic impact on it, as we can see if we try to do inlining on the AST – which might look sensible at first.

Other simple optimizations

Multiple assignments make most simple optimizations hard:

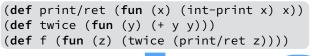
- common subexpression elimination, which consists in avoiding the repeated evaluation of expressions,
- (simple) *dead code elimination*, which consists in removing assignments to variables whose value is not used later,

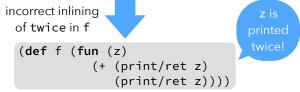
- etc.

Common problem: analyses are required to distinguish the various "versions" of a variable that appear in the program.

Conclusion: a good IR should not allow multiple assignments to a variable!

Naïve inlining: problem #1





Possible solution: bind actual parameters to variables (using a let) to ensure that they are evaluated *at most* once.

Naïve inlining: problem #2



(**def** print/ret (**fun** (z) z))

Possible solution: bind actual parameters to variables (using a let) to ensure that they are evaluated *at least* once.

Easy inlining

Common solution:

bind actual arguments to variables before using them in the body of the inlined function.

However:

the IR can also avoid the problem by ensuring that actual parameters are always atoms (variables/constants).

Conclusion:

a good IR should only allow atomic arguments to functions.

IR comparison

Conclusion:

- standard RTL/CFG is:

- bad as its variables are mutable, but
- good as it allows only atoms as function arguments,
- RTL/CFG in SSA form and CPS/L3 are:
- good as their variables are immutable,
- good as they only allow atoms as function arguments.

Simple CPS/L₃ optimizations

Rewriting optimizations

The rewriting optimizations for CPS/L₃ are specified as a set of rewriting rules of the form T \rightarrow_{opt} T'.

These rules rewrite a CPS/L3 term T to an equivalent – but hopefully more efficient – term T'.

Optimization contexts

Rewriting rules can only be applied in specific locations. For example, it would be incorrect to try to rewrite the parameter list of a function.

We express this constraint by specifying all the **contexts** in which it is valid to perform a rewrite, where a context is a term with a single **hole** denoted by \Box .

The hole of a context C can be plugged with a term T, an operation written as C[T].

For example, if C is (if \Box ct cf), then C[(= x y)] is

(if (= x y) ct cf).

(Non-)shrinking rules

We can distinguish two classes of rewriting rules:

- 1. **shrinking rules** rewrite a term to an equivalent but smaller one, and can be applied at will,
- 2. **non-shrinking rules** rewrite a term to an equivalent but potentially larger one, and must be applied carefully.

Except for inlining, all optimizations we will see are shrinking.

Optimization contexts

C_{opt} ::= □

 $\begin{array}{l} (\texttt{let}_p \ ((\texttt{n} \ (\texttt{p} \ \texttt{n}_1 \ldots))) \ C_{opt}) \\ (\texttt{let}_c \ ((\texttt{c}_1 \ \texttt{e}_1) \ \ldots \ (\texttt{c}_i \ (\texttt{cnt} \ (\texttt{n}_{i,1} \ \ldots) \ C_{opt})) \ \ldots \ (\texttt{c}_k \ \texttt{e}_k)) \ \texttt{e}) \\ (\texttt{let}_c \ ((\texttt{c}_1 \ \texttt{e}_1) \ \ldots) \ C_{opt}) \\ (\texttt{let}_f \ ((\texttt{f}_1 \ \texttt{e}_1) \ \ldots \ (\texttt{f}_i \ (\texttt{fun} \ (\texttt{n}_{i,1} \ \ldots) \ C_{opt})) \ \ldots \ (\texttt{f}_k \ \texttt{e}_k)) \ \texttt{e}) \\ (\texttt{let}_f \ ((\texttt{f}_1 \ \texttt{e}_1) \ \ldots) \ C_{opt}) \end{array}$

Optimization relation

By combining the optimization rewriting rules – presented later – and the optimization contexts, it is possible to specify the optimization relation \Rightarrow_{opt} that rewrites a term to an optimized version: $C_{opt}[T] \Rightarrow_{opt} C_{opt}[T']$ where $T \Rightarrow_{opt} T'$

Dead code elimination

 $(let_p ((n (p v_1 ...))) e)$

**opt e
[when n is not free in e and p ∉ { byte-read, byte-write, block-set! }]

 $\begin{array}{l} (\texttt{let}_{f} \ ((n_{1} \ f_{1}) \ \dots \ (n_{i} \ f_{i}) \ \dots \ (n_{k} \ f_{k})) \ e) \\ \twoheadrightarrow_{opt} \ (\texttt{let}_{f} \ ((n_{1} \ f_{1}) \ \dots \ (n_{k} \ f_{k})) \ e) \\ [when \ n_{i} \ is \ not \ free \ in \ \{f_{1}, \ \dots, \ f_{i-1}, \ f_{i+1}, \ \dots \ f_{k}, \ e\}] \end{array}$

The rule for continuations is similar to the one for functions.

Dead code elimination

Limitation:

Does not eliminate dead, mutually-recursive functions. Solution:

- start from the main expression of the program, and

- identify all functions transitively reachable from it. All unreachable functions are dead.

CSE

 $\begin{array}{l} (\texttt{let}_p \ ((n_1 \ (+ \ v_1 \ v_2))) \\ C_{opt}[(\texttt{let}_p \ ((n_2 \ (+ \ v_1 \ v_2))) \ e)]) \\ \twoheadrightarrow_{opt} \ (\texttt{let}_p \ ((n_1 \ (+ \ v_1 \ v_2))) \ C_{opt}[e\{n_2 \rightarrow n_1\}]) \end{array}$

 $\begin{array}{l} (let_{p} ((n_{1} (-v_{1} v_{2}))) \\ C_{opt}[(let_{p} ((n_{2} (-v_{1} v_{2}))) e)]) \\ \twoheadrightarrow_{opt} (let_{p} ((n_{1} (-v_{1} v_{2}))) C_{opt}[e\{n_{2} \rightarrow n_{1}\}]) \end{array}$

etc.

CSE

Limitation:

...)

η-reduction

Constant folding (1)

 $(let_p ((n (+ |_1 |_2))) e)$ $\Rightarrow_{opt} e\{n \rightarrow (|_1 + |_2)\}$ $[when |_1 and |_2 are integer literals]$

 $(let_p ((n(-l_1 l_2)))e)$

 $\Rightarrow_{opt} e\{n \rightarrow (|_1 - |_2)\}$ [when $|_1$ and $|_2$ are integer literals]

(let_p ((n (* $I_1 I_2$))) e)

 $\Rightarrow_{opt} e\{n \rightarrow (|_1 \times |_2)\}$ [when $|_1$ and $|_2$ are integer literals]

etc.

Constant folding (2)

(if $(= v v) c_t c_f$) $\Rightarrow_{opt} (app_c c_t)$

 $\begin{array}{l} (\text{if } (= l_1 \mid_2) \ c_t \ c_f) \\ \twoheadrightarrow_{opt} \ (app_c \ c_t) \\ [when \mid_1 \ and \mid_2 \ are \ literals \ and \mid_1 = l_2] \end{array}$

etc.

Neutral/absorbing elements

 $\begin{array}{l} (let_{p} \ ((n \ (\star \ 0 \ v))) e) \\ \begin{array}{l} \bullet_{opt} \ e\{n \rightarrow 0\} \\ (let_{p} \ ((n \ (\star \ v \ 0))) e) \\ \begin{array}{l} \bullet_{opt} \ e\{n \rightarrow 0\} \end{array} \end{array}$

etc.

Exercise

CPS/L₃ contains the following block primitives:

- block-alloc-n size
- block-tag block
- block-size block
- block-get block index
- block-set! block index value

Informally describe three rewriting optimizations that could be performed on these primitives, and in which conditions they could be performed.

Block primitives

(letp ((b (block-alloc-ks))) C_{opt}[(letp ((t (block-set! biv))) C'_{opt}[(letp ((n (block-get bi))) e)])]) *•_{opt} (letp ((b (block-alloc-ks))) C_{opt}[(letp ((t (block-set! biv))) C'_{opt}[e{n→v}])]) [when tag k identifies a block that is not modified after initialization, e.g. a closure block]

CPS/L₃ inlining

(Non-)shrinking inlining

We can distinguish two kinds of inlining:

- 1. **shrinking inlining**, for functions/continuations that are applied exactly once,
- 2. non-shrinking inlining, for other functions/continuations.

Shrinking inlining can be applied at will, non-shrinking cannot.

Inlining heuristics (1)

Heuristics must be used to decide when to perform non-shriking inlining. They typically combine several factors, like:

- the size of the candidate function smaller ones should be inlined more eagerly than bigger ones,
- the number of times the candidate is called in the whole program a function called only a few times should be inlined,

(continued on next slide)

Inlining

1. Shrinking: (let_f ((f₁ e₁) ... (f_i (fun (c_i n_{i,1} ...) e_i)) ... (f_k e_k)) $C_{opt}[(app_{f} f_{i} c m_{1} ...)])$ \Rightarrow_{opt} (let_f ((f₁ e₁) ... (f_k e_k)) $C_{opt}[e_{i}\{c_{i}\rightarrow c\}\{n_{i,1}\rightarrow m_{1}\}...])$ [when f_i is not free in C_{opt}, e₁, ..., e_n] 2. Non-shrinking (fresh versions of bound names should be created to preserve their global uniqueness): (let_f (... (f_i (fun (c_i n_{i,1} ...) e_i)) ...) $C_{opt}[(app_{f} f_{i} c m_{1} ...)])$ \Rightarrow_{opt} (let_f (... (f_i (fun (c_i n_{i,1} -..) e_i)) ...) $C_{opt}[e_{i}\{c_{i}\rightarrow c\}\{n_{i,1}\rightarrow m_{1}\}...])$

Inlining heuristics (2)

- the nature of the candidate not much gain can be expected from the inlining of a recursive function,
- the kind of arguments passed to the candidate, and/or the way these are used in the candidate – constant arguments could lead to further reductions in the inlined candidate, especially if it combines them with other constants,
- etc.

Exercise

Imagine an imperative intermediate language equipped with a return statement to return from the current function to its caller.

- 1. Describe the problem that would appear when inlining a function containing such a **return** statement.
- 2. Explain how a return statement could be encoded in CPS/L_3 and why such an encoding would not suffer from the above problem.

CPS/L₃ "contification"

Contification

Contification: transforms functions into continuations. Interesting optimization as it transforms functions, which are expensive (closures) into continuations, which are cheap.

Contification example

Example: the loop function in the L_3 example below can be contified, leading to efficient compiled code.

Contifiability

A CPS/ L_3 function is contifiable if and only if it always returns to the same location – because then it does not need a return continuation.

- Non-recursive case: true iff that function is only used in appf nodes, in function position, and always passed the same return continuation.
- Recursive case: slightly more involved see later.

Non-recursive contification

The contification of the non-recursive function f is given by: $\begin{array}{l} (\texttt{let}_f \ ((\texttt{f}(\texttt{fun} \ (\texttt{ca}_1 \ldots) \texttt{e}))) \\ C_{\texttt{opt}}[C'_{\texttt{opt}}[(\texttt{app}_f \ \texttt{fc}_0 \ \texttt{n}_{1,1} \ldots), (\texttt{app}_f \ \texttt{fc}_0 \ \texttt{n}_{2,1}, \ldots), \ldots]]) \\ \twoheadrightarrow_{\texttt{opt}} C_{\texttt{opt}}[(\texttt{let}_c \ ((\texttt{m} \ (\texttt{cnt} \ (\texttt{a}_1 \ldots) \ \texttt{e}[\texttt{c} \rightarrow \texttt{c}_0\}))) \\ C'_{\texttt{opt}}[(\texttt{app}_c \ \texttt{m} \ \texttt{n}_{1,1} \ldots), (\texttt{app}_c \ \texttt{m} \ \texttt{n}_{2,1} \ldots), \ldots])] \end{array}$

where:

- f does not appear free in C_{opt} or C'_{opt} ,

- $C^\prime_{\,\rm opt}$ is the smallest (multi-hole) context enclosing all applications of f,
- c_0 is the (single) return continuation that is passed to function f.

Recursive contifiability

A set of mutually-recursive functions $F = \{ f_1, ..., f_n \}$ is contifiable – which we write Cnt(F) – if and only if every function in F is always used in one of the following two ways:

1. applied to a common return continuation, or

2. called in tail position by a function in F.

Intuitively, this ensures that all functions in F eventually return through the common continuation.

Example

As an example, functions even and odd in the CPS/L3 translation of the following L3 term are contifiable:

(letrec

((even (fun (x) (if (= 0 x) #t (odd (- x 1))))) (odd (fun (x) (if (= 0 x) #f (even (- x 1)))))) (even 12))

Cnt(F = {even, odd}) is satisfied since:

- the single use of odd is a tail call from $even \in F$,

- even is tail-called from odd \in F and called with the continuation of the

letrec statement – the common return continuation c_0 for this example.

Recursive contification

Given a set of mutually-recursive functions

 $(let_f ((f_1 e_1) (f_2 e_2) ... (f_n e_n))$

e)

the condition Cnt(F) for some $F \subseteq \{ f_1, ..., f_n \}$ can only be true if all the non tail calls to functions in F appear either:

- in the term e, or

- in the body of exactly one function $f_i \not\in F.$

Therefore, two separate rewriting rules must be defined, one per case.

Recursive contification #1

Case 1: all non tail calls to functions in F = { $f_1, ..., f_i$ } appear in the body of the let $_f$, and Cnt(F) holds:

 $\begin{array}{c} (\texttt{let}_f \ ((f_1 \ (\texttt{fun} \ (c_1 \ a_{1,1} \ ...) \ e_1)) \ ... \ (f_n \ ...)) \\ C_{opt}[e]) \\ \twoheadrightarrow_{opt} \ (\texttt{let}_f \ ((f_{i+1} \ (\texttt{fun} \ (c_{i+1} \ a_{i+1,1} \ ...) \ e_{i+1})) \ ... \ (f_n \ ...)) \\ C_{opt}[(\texttt{let}_c \ ((m_1 \ (\texttt{cnt} \ (a_{1,1} \ ...) \ e_1^* \{c_1 \rightarrow c_0\})) \ ...) \\ e^*)]) \end{array}$

*)])

where $f_1, ..., f_i$ do not appear free in C_{opt} and e is minimal. Note: the term t* is t with all applications of contified functions transformed to continuation applications.

Recursive contification #2

Case 2: all non tail calls to functions in F = { $f_1,\,...,\,f_i$ } appear in the body of the function $f_n,$ and Cnt(F) holds:

 $\begin{array}{c} (\mathsf{let}_{f} \ ((f_{1} \ (\mathsf{fun} \ (c_{1} \ a_{1,1} \ ...) \ e_{1})) \ ... \\ (f_{n} \ (\mathsf{fun} \ (c_{n} \ a_{n,1} \ ...) \ C_{opt}[e_{n}]))) \ e) \\ \stackrel{\textbf{\rightsquigarrow_{opt}}}{} (\mathsf{let}_{f} \ ((f_{i+1} \ (\mathsf{fun} \ (c_{i+1} \ a_{i+1,1} \ ...) \ e_{i+1})) \ ... \\ (f_{n} \ (\mathsf{fun} \ (c_{n} \ a_{n,1} \ ...) \\ C_{opt}[(\mathsf{let}_{c} \ ((\mathsf{m}_{1} \ (\mathsf{cnt} \ (a_{1,1} \ ...) \ e_{1}*[c_{1} \ \rightarrow c_{0}])) \\ \ ...) \\ e_{n}^{\star})])) \ e) \end{array}$

where $f_1, ..., f_i$ do not appear free in C_{opt} and e_n is minimal.

Contifiable subsets

Given a let f term defining a set of functions $F = \{ f_1, ..., f_n \}$, the subsets of F of potentially contifiable functions are obtained by:

1. building the tail-call graph of its functions, identifying the functions that call each-other in tail position,

2. extracting the strongly-connected components of that graph. A given set of strongly-connected functions $F_i \subseteq F$ is then either contifiable together, i.e. $Cnt(F_i)$, or not contifiable at all – i.e. none of its subsets of functions are contifiable.