

Enabling Optimizations through Demodularization

Blake Johnson

A thesis submitted to the faculty of
Brigham Young University
in partial fulfillment of the requirements for the degree of
Master of Science

Eric Mercer, Chair
Christophe Giraud-Carrier
Quinn Snell

Department of Computer Science
Brigham Young University
December 2015

Copyright © 2015 Blake Johnson
All Rights Reserved

ABSTRACT

Enabling Optimizations through Demodularization

Blake Johnson

Department of Computer Science, BYU

Master of Science

Programmers want to write modular programs to increase maintainability and create abstractions, but modularity hampers optimizations, especially when modules are compiled separately or written in different languages. In languages with syntactic extension capabilities, each module in a program can be written in a separate language, and the module system must ensure that the modules interoperate correctly. In Racket, the module system ensures this by separating module code into phases for runtime and compile-time and allowing phased imports and exports inside modules. We present an algorithm, called demodularization, that combines all executable code from a phased modular program into a single module that can then be optimized as a whole program. The demodularized programs have the same behavior as their modular counterparts but are easier to optimize. We show that programs maintain their meaning through an operational semantics of the demodularization process and verify that performance increases by running the algorithm and existing optimizations on Racket programs.

Keywords: macros, Racket, modules, optimization

ACKNOWLEDGMENTS

Jay McCarthy, Matthew Flatt, Eric Mercer

Contents

List of Figures	vi
List of Listings	vii
1 Introduction	1
2 The Racket Module System	5
2.1 The Racket Programming Language	5
2.2 Macros	5
2.3 Modules	7
2.4 Languages	8
2.5 Separate Compilation	9
2.6 Phases	9
2.7 Program Evaluation	10
3 Intuition	11
3.1 Example Program	11
3.2 Compilation	12
3.3 Demodularization	13
3.4 Optimization	14
4 Model	17
4.1 A Module Language	17
4.2 Compilation	18

4.3	Evaluation	19
4.4	Demodularization	21
4.5	Equivalence	23
5	Implementation	25
5.1	Racket bytecode	25
5.2	Demodularization	26
5.3	Optimization	27
5.4	Decompilation	27
5.5	Limitations	28
6	Evaluation	29
7	Related Work	31
8	Conclusion	33
	References	35

List of Figures

2.1	Module dependencies for the while-test program	9
4.1	<i>mod</i> language grammar	17
4.2	Compiled language grammar	18
4.3	Extensions to compiled language grammar	19
4.4	<i>module require</i> rule for modular evaluation	20
4.5	<i>var ref</i> rule for modular evaluation	20
4.6	<i>add</i> rule for modular evaluation	20
4.7	<i>set!</i> rule for modular evaluation	20
4.8	<i>expression done</i> rule for modular evaluation	21
4.9	<i>module done</i> rule for modular evaluation	21
4.10	<i>module done already</i> rule for modular evaluation	21
4.11	Demodularization algorithm	21
4.12	Demodularization algorithm	22
4.13	Demodularization algorithm	22
4.14	Demodularization algorithm	22
4.15	Bisimulation of a program before and after demodularization	23

List of Listings

1	a macro implementation of a <code>while</code> loop	6
2	use and expansion of a <code>while</code> loop	6
3	<code>counter.rkt</code> : A simple Racket module implementing a counter	7
4	<code>while-test.rkt</code> : A Racket module that uses other modules	7
5	<code>while-lang.rkt</code> : A Racket module implementing a language with <code>while</code> loops	8
6	<code>main.rkt</code> module with queue usage	11
7	<code>queue.rkt</code> module	12
8	<code>long-queue.rkt</code> module	13
9	<code>short-queue.rkt</code> module	13
10	<code>main.rkt</code> module after macro expansion	14
11	<code>main.rkt</code> module after demodularization	15
12	<code>main.rkt</code> module after optimization	15
13	Example program written in <code>kernel</code> language	26
14	Bytecode representation of program from Listing 13	26

Chapter 1

Introduction

Programmers should not have to sacrifice the software engineering goals of modular design and good abstractions for performance. Instead, their tools should make running a well-designed program as efficient as possible.

Many languages provide features for creating modular programs which enable separate compilation and module reuse. Some languages provide expressive macro systems, which enable programmers to extend the compiler in arbitrary ways. Combining module systems with expressive macro systems allow programmers to write modular programs with each module written in its own domain-specific language. A compiler for such a language must ensure that modular programs have the same meaning independent of the order in which the modules are compiled. A phased module system, like the one described by Flatt [5] for Racket, is a way to allow both separately compiled modules and expressive macros in a language.

Modular programs are difficult to optimize because the compiler has little to no information about values that come from other modules when compiling a single module. Existing optimizations have even less information when modules can extend the compiler. Good abstractions are meant to obscure internal implementations so that it is easier for programmers to reason about their programs, but this obscurity also limits information available for optimizations. In contrast, non-modular programs are simpler to optimize because the compiler has information about every value in the program. For example, inlining a function in a non-modular program is trivial because every use of the function is

known at compile-time.

Some languages avoid the problem of optimizing modular programs by not allowing modules, while others do optimizations at link time, and others use inlining. Not allowing modules defeats the benefits of modular design. Languages such as C and Scheme that do not have well developed module systems make it difficult to manage large programs written in them. Link time optimizations can be too low level to do useful optimizations. Linking is usually performed on already compiled modules, so any optimizations are performed on machine code and are missing source-level information. Inlining must be heuristic-based, and good heuristics are hard to develop. Some work has been done on these heuristics for the ML programming language [1].

Our solution for optimizing modular programs, called demodularization, is to transform a modular program into a non-modular program by combining all runtime code and data in the program into a single module. In a phased module system, finding all of the runtime values is not trivial. Phased module systems allow programmers to refer to the same module while writing compiler extensions and while writing normal programs. A demodularized program does not need to include modules that are only needed during compile-time, but whether or not the module is needed only at compile-time is not obvious from just examining the module in isolation.

A program with a single module is effectively a non-modular program. After demodularization, a program becomes a single module, so existing optimizers have more information. Also, demodularization enables new optimizations that need whole program information. Demodularizing programs in a language with an expressive macro system is feasible and useful.

We explain the Racket module system in Chapter 2. In Chapter 3, we explain demodularization at a high level with a detailed example. Next, we use the operational semantics model of the demodularization process to explain why demodularization is correct (Chapter 4), then describe an actual implementation for Racket (Chapter 5), followed by exper-

imental results of demodularizing and optimizing real-world Racket programs (Chapter 6). The operational semantics model removes the unnecessary details of the full implementation so the demodularization process is easier to understand and verify. The actual implementation presents interesting difficulties that the model does not. The experimental results show that demodularization improves performance, especially when a program is highly modular.

Chapter 2

The Racket Module System

2.1 The Racket Programming Language

The Racket Programming Language is a platform for creating powerful abstractions. These abstractions are created through the use of Racket’s macro and module systems. The macro system, a heritage from LISP [1] and Scheme [2], allows programmers to add new features to their programs in an integrated way. The module system allows programmers to separate their programs into logical parts and hide implementation details. Together, they allow programmers to create Languages as Libraries [3] that are suitable for specific tasks.

2.2 Macros

Macros are the main way programmers add new features to Racket. Essentially, macros are functions (known as transformers) whose domain and range are syntax objects. Syntax objects are data structures that contain the raw syntax of a program, along with lexical information and other properties associated with the syntax. If a programmer needs the power, they can write transformers using all of the features of Racket, as long as the function takes in and produces syntax objects. Racket provides pattern-matching for syntax objects that makes writing macros simpler when that is all that a programmer needs.

The **define-syntax** form is used to identify a macro definition that is a function from syntax to syntax. The helper function, **syntax-case** matches syntax patterns so that it is easy to use them in the output of the macro.

Listing 1 shows how to use one of these macros to write a **while** loop. The macro

```

(define-syntax (while stx)
  (syntax-case stx ()
    [(while test body ...)
     #'(begin
         (define (while-loop)
           (when test
             body ...
             (while-loop)))
         (while-loop))]))

```

Listing 1: a macro implementation of a `while` loop

turns a `while` loop into a recursive function and a call to that recursive function. Whatever syntax used by the programmer for the `test` and `body ...` forms will be substituted into the output of the macro. Listing 2 shows a use of the `while` macro and what it expands into after running the macro. Racket's macros are hygienic [], which means that identifiers created from

```

(while (< x 10)
  (printf "x is ~a\n" x)
  (set! x (add1 x)))

=>

(begin
  (define (while-loop)
    (when (< x 10)
      (printf "x is ~a\n" x)
      (set! x (add1 x))
      (while-loop))))
(while-loop))

```

Listing 2: use and expansion of a `while` loop

programs will not clash with identifiers created from macros. In this example that means that the Racket program using a `while` loop will not be able to call the `while-loop` function that the macro creates, and that if the program has another definition for `while-loop`, it will be independent from the definition introduced by the macro. Also, each use of the macro will have distinct definitions for `while-loop` so that the macro can be nested and used multiple times.

2.3 Modules

Racket modules are a way of grouping definitions, expressions, and macros into groups. Modules also allow control over imports and exports, using the racket forms **require** and **provide** respectively. Listing 3 shows a simple Racket module and some of the features of the Racket module system. The `counter` module contains a variable definition and functions

```
#lang racket/base
(provide get-counter-val count-up)

(define x 0)
(define (get-counter-val) x)
(define (count-up . v) (set! x (add1 x)))
```

Listing 3: `counter.rkt`: A simple Racket module implementing a counter

for getting and incrementing that variable. The module only exports the functions, so the module encapsulates the variable. If a programmer wanted to use this module, they would import it using **require** as shown in Listing 4.

```
#lang racket/base
(require "while-lang.rkt")
(require "counter.rkt")

(while (< (get-counter-val) 4)
  (printf "loop runtime val is ~a\n" (get-counter-val))
  (count-up))
```

Listing 4: `while-test.rkt`: A Racket module that uses other modules

It is also possible to use modules to help write macros. For example, if a programmer wanted to change the `while` macro so that it reported how many expressions were in each while loop, they could use the `counter` module inside the `while` definition as seen in Listing 5. It is possible to add code to the `while` macro because it is just a function from syntax to syntax that runs at compile-time. Also, the **require** form includes a **for-syntax** form to indicate that the `counter` module is needed at compile-time (along with the `racket/base`

```

#lang racket/base
(require (for-syntax racket/base
                    "counter.rkt"))

(provide while)

(define-syntax (while stx)
  (syntax-case stx ()
    [(while test body ...)
     (begin
       (for-each count-up
                  (syntax->list #'(body ...)))
       (printf "while loop contains ~a expressions\n"
               (get-counter-val))
       #'(begin
            (define (while-loop)
              (when test
                body ...
                (while-loop)))
            (while-loop)))))]))

```

Listing 5: while-lang.rkt: A Racket module implementing a language with `while` loops

module for `printf`). Now this module can be used to extend the language of other programs by adding a `while` loop, but it is just a library that is included like any other library.

2.4 Languages

The ability to use modules as languages allows programmers to write specialized domain-specific languages (DSLs) to better express solutions to domain problems. Using this ability, programmers have added languages for Typed Racket [\[1\]](#), Object-Oriented Racket [\[2\]](#), Logic Programming [\[3\]](#), and more [\[4\]](#). Also, because they are modules, it is possible to use them all in the same program and use the right language for each specific problem. The language creating facilities of Racket extend to more than just collections of functions and macros; they also allow changing the parser and changing the meaning of things like function application syntax.

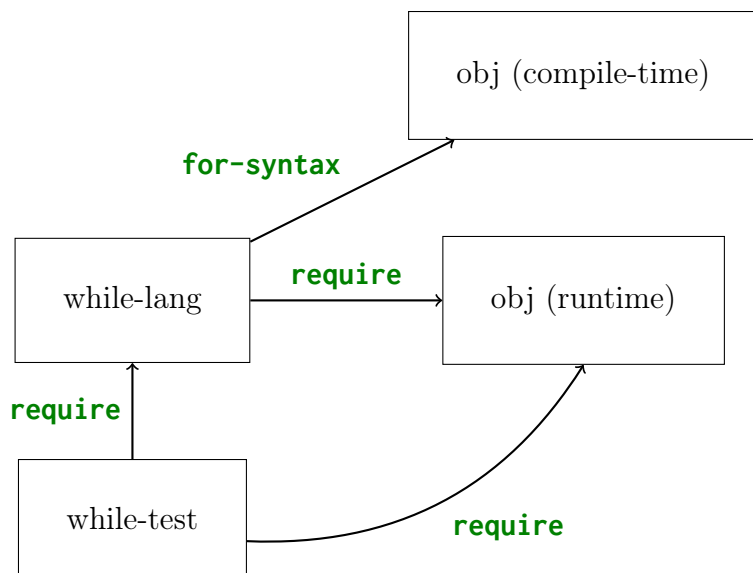


Figure 2.1: Module dependencies for the `while-test` program

2.5 Separate Compilation

In the presence of macros that can run arbitrary code, compilation of a single module becomes complicated. Compilation and execution of code has to be interleaved in order to produce the correct compiled program. All of the listings in the previous section form a program with the `while-test` module as the main module. Figure 2.1 shows the relationship between the different modules in the program. To compile the program, both the `while-lang` and `counter` need to be compiled before `while-test`, but also the `counter` module needs to be compiled and ready to use while compiling `while-test` because the counter functions are used inside the `while` macro. While compiling `while-test`, the `while` macro will be expanded and code from the `counter` module will run. The value of the counter is not shared between compiling `while-test` and running it.

2.6 Phases

The way the Racket module system allows for reliable separate compilation is by separating compilation and execution into phases. Phase 0 is the execution phase of the main module of a program, what could be considered running the program. Phase 1 is the compilation

phase of the main module, which could include execution of code inside macros. Phase 2 is the compilation phase of any modules that need to execute during phase 1, and higher phases are used as needed to compile and run code used in lower phases. It is possible for modules to reference modules at higher phases by using the **for-syntax** form as seen above, and it is also possible to reference modules at lower phases by using the **for-template** form.

In the example above, the **while-test** module is compiled at phase 1, which triggers compilation of the **counter** module at phase 2, so that the **counter** module can be instantiated at phase 1 to count how many expressions are in the **while** loop. The phase 1 instantiation of the **counter** module is separate from the phase 0 instantiation of the module. This is a guarantee of the Racket module system, that “Any effects of the instantiation of the module’s phase 1 due to compilation on the Racket runtime system are discarded,” which is known as the "Separate Compilation Guarantee" [].

2.7 Program Evaluation

Running a separately compiled Racket program involves following all of the **require** forms and running all phase 0 code in the order they were imported. It is possible that phase 0 code will come through a module that is required by using **for-syntax** if somewhere along the line code is then required using **for-template**.

By understanding how Racket programs are compiled and evaluated, it is apparent that only phase 0 code is necessary to run the program. This is the basis for how demodularization can recover whole programs from separately compiled Racket modules.

Chapter 3

Intuition

Demodularization is the process of collecting all phase 0 code required by a program into a single module. This is done by tracing through the **require** graph starting at a program's main module. The following example program illustrates the need for demodularization and how it is done.

3.1 Example Program

A programmer wants to use a queue library where the library uses different backing structures depending on the length of the queue. Listing 6 shows an example of using such a library. The library is implemented as a macro, show in Listing 7, that switches between two imple-

```
#lang racket/base
(require "queue.rkt")

(with-queue (1 2 3 4 5 6)
  (enqueue 4)
  (displayln (dequeue))
  (displayln (dequeue)))
```

Listing 6: `main.rkt` module with queue usage

mentations depending on the length of the initial queue. Listing 8 and Listing 9 show the two queue implementations, with the `long-queue` using vectors and the `short-queue` using lists. The `....` signifies elided code (but `...` has meaning in racket).

```

#lang racket/base
(require (for-syntax racket/base
                    racket/syntax)
         "short-queue.rkt"
         "long-queue.rkt")

(define-syntax (with-queue stx)
  (syntax-case stx ()
    [(with-queue (v ...) e ...)
     (begin
      (define type
        (if (> (length (syntax->list #'(v ...))) 5)
            'long
            'short))
      (define make-queue
        (format-id #'stx "make-~a-queue" type))
      (define enqueue (format-id #'stx "~a-enqueue" type))
      (define dequeue (format-id #'stx "~a-dequeue" type))
      #`(let ([q (#,make-queue v ...)])
          (define (#,(datum->syntax stx 'dequeue))
            (#,dequeue q))
          (define (#,(datum->syntax stx 'enqueue) x)
            (#,enqueue q x))
          e ...)))]))

(provide with-queue)

```

Listing 7: queue.rkt module

3.2 Compilation

Compiling this program involves expanding the `with-queue` macro, which will turn the program into Listing 10.

The expanded program contains references to `long-queue` operations because the length of the initial queue was over 5. While the expanded program is not the same as a fully compiled program, it can stand in for one with regard to optimization and demodularization.

With the expanded program, it is possible to see how cross-module optimizations would be difficult. The program has calls to functions in the `long-queue` module, so the optimizer does not have access to the definitions of the functions while finishing the compilation

```

#lang racket/base
(define (make-long-queue . vs)
  .... make-vector ....)

(define (long-enqueue q v)
  .... vector-set! ....)

(define (long-dequeue q)
  .... vector-ref ....)

(provide (all-defined-out))

```

Listing 8: long-queue.rkt module

```

#lang racket/base
(define (make-short-queue . vs)
  .... list ....)

(define (short-enqueue q v)
  .... cons ....)

(define (short-dequeue q)
  .... list-ref ....)

(provide (all-defined-out))

```

Listing 9: short-queue.rkt module

of the main module.

3.3 Demodularization

Demodularization of the program proceeds by following the **require** graph and including phase 0 code. The `main` module requires the `queue` module, but the `queue` module only has phase 1 code, so it is not included in the final module. The `queue` module requires the `long-queue` and `short-queue` modules, which have various phase 0 definitions that get included in the order they appear in their original modules. Listing 11 shows what the demodularized example program would look like.

The demodularized program should have the exact same behavior at runtime as the

```

(module main racket/base
  (%module-begin
    (require "queue.rkt")
    (let ((q (make-long-queue 1 2 3 4 5 6)))
      (define (dequeue) (long-dequeue q))
      (define (enqueue x) (long-enqueue q x))
      (enqueue 4)
      (displayln (dequeue))
      (displayln (dequeue))))))

```

Listing 10: `main.rkt` module after macro expansion

modular program. At this point, the only thing that changes is that the Racket runtime does not need to follow the `require` graph because all phase 0 definitions and expressions are in a single module. Also, both queue implementations appear in the demodularized program even though only one of them is used.

3.4 Optimization

After a program has been demodularized, the optimizer now has all the information it needs to do optimizations on the whole program. Listing 12 shows what the example program would look like after optimization.

The optimizer is able to inline the various functions because it has access to all of the definitions within the module. The optimizer also removes the extra queue implementation because it is dead code. Therefore, demodularization enables optimization by providing the optimizer access to the whole program in one place.

```

(module main racket/base
  (%module-begin
    (define (make-short-queue . vs)
      .... list ....)

    (define (short-enqueue q v)
      .... cons ....)

    (define (short-dequeue q)
      .... list-ref ....)

    (define (make-long-queue . vs)
      .... make-vector ....)

    (define (long-enqueue q v)
      .... vector-set! ....)

    (define (long-dequeue q)
      .... vector-ref ....)

    (let ((q (make-long-queue 1 2 3 4 5 6)))
      (define (dequeue) (long-dequeue q))
      (define (enqueue x) (long-enqueue q x))
      (enqueue 4)
      (displayln (dequeue))
      (displayln (dequeue))))))

```

Listing 11: main.rkt module after demodularization

```

(module main racket/base
  (%module-begin
    (let ((q (.... make-vector .... 1 2 3 4 5 6 ....)))
      (.... vector-set! .... 4 ....)
      (displayln (.... vector-ref ....))
      (displayln (.... vector-ref ....)))))

```

Listing 12: main.rkt module after optimization

Chapter 4

Model

We can understand the specifics of demodularization by describing it as an algorithm for a simple language with a well defined semantics.

4.1 A Module Language

The *mod* language (Figure 4.1) contains only the features necessary to write modular programs where it is possible to observe the effects of module evaluation order.

```
program ::= (mod ...)
mod ::= (module id (req ...) (def ...) expr)
req ::= (require id @ phase)
expr ::= var
        | val
        | (+ expr expr)
        | (set! var expr)
        | (begin expr expr)
val ::= number
        | (quote id @ phase)
var ::= id
        | (ref id @ phase)
id ::= variable-not-otherwise-mentioned
phase ::= number
def ::= (define id @ phase as expr)
```

Figure 4.1: *mod* language grammar

A program in *mod* consists of a list of modules that can refer to each other. Each module has a name, any number of imports, any number of definitions, and sequenced code expressions. All definitions in a module are exposed as exports to other modules, but to use definitions from another module, the program must import it through a **require** expression. Both **require** and **define** expressions have phase annotations; this simulates the interactions

between modules in a language with macros and a language tower without requiring a model of macro expansion. The language includes variable references, numbers, addition, and mutation. Mutation makes module evaluation order observable, and addition represents the work that a module does. In addition to numbers and variables, there are two special forms of values and references that model the interaction of macros with the module system. A **quote** expression is like a reference to syntax at runtime. A **ref** expression is like a macro that can only do one thing: refer to a variable at a phase.

4.2 Compilation

We have to compile *mod* programs before demodularizing them, just like in the Racket implementation. In Racket, compiling expands all macros in a program and changes definitions and variable references to refer to memory locations. In *mod*, compiling eliminates **ref** expressions, turns definitions into **set!** expressions, changes variable references to include module information, and sorts code into phases. Compilation in both cases still leaves behind a relatively high-level language, but the language is free of syntactic extensions. This is important for demodularization because otherwise macro expansion would have to be part of the algorithm, which would complicate it and possibly duplicate work. The grammar in Figure 4.2 specifies the compiled language for *mod*.

```

program ::= (mod ...)
mod ::= (module id (req ...) (code ...))
code ::= (phase expr)
req ::= (require id @ phase)
expr ::= var
        | val
        | (+ expr expr)
        | (set! var expr)
        | (begin expr expr)
val ::= number
        | (quote id @ phase)
var ::= (id id)
id ::= variable-not-otherwise-mentioned
phase ::= number

```

Figure 4.2: Compiled language grammar

The grammar no longer has definitions, all variables now include module references, and code is sorted into phases. The actual compilation function is not relevant to demodularization.

4.3 Evaluation

We evaluate the compiled language using a small-step reduction semantics. Because the reduction rules are syntactic, we extend the compiled language further with evaluation contexts, a heap representation, and a stack representation to keep track of the order to instantiate modules. These extensions are in Figure 4.3. An expression of the form:

$$(\sigma / \text{program} / ((id \text{ phase}) \dots) / ((id \text{ phase}) \dots))$$

represents the state of the machine during evaluation. σ represents the heap of the program, and when evaluation finishes represents the output of the program. The list of modules is the code of program in the compiled language. The first list of $(id \text{ phase})$ pairs is the list of modules to evaluate, and the second list is the modules that have already been evaluated.

$$\begin{aligned} E &::= [] \\ &\quad | (+ E \text{ expr}) \\ &\quad | (+ \text{ val } E) \\ &\quad | (\text{set! } \text{var } E) \\ &\quad | (\text{begin } E \text{ expr}) \\ \sigma &::= ([\text{var val}] \dots) \\ \text{state} &::= (\text{program} / (\text{inst } \dots) / (\text{inst } \dots)) \\ \text{inst} &::= (id \text{ phase}) \\ st &::= (\sigma / \text{state}) \end{aligned}$$

Figure 4.3: Extensions to compiled language grammar

The reduction rules in Figures 4.4–4.10 evaluate a compiled program that starts with an empty heap, the program code, a stack that contains the identifier of the main module at phase 0, and an empty completed module list.

The *module require* rule in Figure 4.4 matches a program with a **require** expression in the module at the top of the evaluation stack and evaluates it by removing the **require**

$$\begin{array}{l}
(\sigma / ((mod_0 \dots (\text{module } id_0 ((\text{require } id_{new} @ phase_{new}) req \dots) (code \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots)) \longrightarrow [\text{module require}] \\
(\sigma / ((mod_0 \dots (\text{module } id_0 (req \dots) (code \dots)) mod_n \dots) / ((id_{new} (+ (- phase_{new}) phase_0)) (id_0 phase_0) inst_n \dots) / (inst_d \dots)))
\end{array}$$

Figure 4.4: *module require* rule for modular evaluation

expression from the module and pushing the required module onto the evaluation stack with the phase shifted appropriately. The current module is still on the stack and will continue evaluating after the required module is done evaluating. The subsequent rules all apply only when the phase relative to the main module is zero.

$$\begin{array}{l}
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[var]) \longrightarrow [\text{var ref}] \\
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[val]) \\
\text{where } val = \text{lookup } [\sigma, var]
\end{array}$$

Figure 4.5: *var ref* rule for modular evaluation

The *var ref* rule in Figure 4.5 looks up a variable in the heap and replaces the variable with its current value. The lookup function is a simple list lookup function that matches the variable name with its occurrence in the heap.

$$\begin{array}{l}
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[(+ number_0 number_1)]) \longrightarrow [\text{add}] \\
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[(+ number_0 number_1)])
\end{array}$$

Figure 4.6: *add* rule for modular evaluation

The *add* rule in Figure 4.6 replaces an addition expression of numbers with the result of computing their sum. The rule uses the standard Racket implementation of plus to compute the sum.

$$\begin{array}{l}
(\sigma_0 / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[(\text{set! } var \text{ val})]) \longrightarrow [\text{set!}] \\
(\sigma_1 / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[val]) \\
\text{where } \sigma_1 = \text{assign } [\sigma_0, var, val]
\end{array}$$

Figure 4.7: *set!* rule for modular evaluation

The *set!* rule in Figure 4.7 installs a value for a variable into the heap and reduces to the value. The assign function either replaces the existing value for a variable in the heap or installs a new entry in the heap if it does not exist already.

When an expression is a value, the *expression done* rule in Figure 4.8 matches and removes the expression from the module.

$$\begin{array}{l}
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[(\text{begin val } expr)] \longrightarrow [\text{expression done}] \\
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 E) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))[expr]
\end{array}$$

Figure 4.8: *expression done* rule for modular evaluation

$$\begin{array}{l}
(\sigma / ((mod_0 \dots (\text{module } id_0 () (code_0 \dots (phase_0 val) code_n \dots)) mod_n \dots) / ((id_0 phase_0) inst_n \dots) / (inst_d \dots))) \longrightarrow [\text{module done}] \\
(\sigma / ((mod_0 \dots mod_n \dots) / (inst_n \dots) / ((id_0 phase_0) inst_d \dots)))
\end{array}$$

Figure 4.9: *module done* rule for modular evaluation

When there are no more expressions left in a module, the *module done* rule in Figure 4.9 applies by removing the module from the program and placing a reference to it in the list of finished modules.

$$\begin{array}{l}
(\sigma / ((mod \dots) / (inst_0 inst_n \dots) / (inst_{d0} \dots inst_0 inst_{dn} \dots))) \longrightarrow [\text{module done already}] \\
(\sigma / ((mod \dots) / (inst_n \dots) / (inst_{d0} \dots inst_0 inst_{dn} \dots)))
\end{array}$$

Figure 4.10: *module done already* rule for modular evaluation

The *module done already* rule in Figure 4.10 applies when the current module on the stack is in the finished list, so that modules are not evaluated multiple times.

4.4 Demodularization

Figures 4.11–4.14 shows the demodularization algorithm for the compiled language. The algorithm takes as input an *id* specifying the main module of the program, a working list of phased requires, and a list of modules that make up the program.

$$\text{demod} \llbracket id, (), (mod_0 \dots (\text{module } id () (code_0 \dots (0 \text{ expr}) code_n \dots)) mod_n \dots) \rrbracket = ((\text{module } id () ((0 \text{ expr}))))$$

Figure 4.11: Demodularization algorithm

The first rule in Figure 4.11 applies when the main module has no requires left and there are no requires in the working list, meaning the algorithm can terminate with just the phase 0 code of the main module remaining. The second and third rules in Figure 4.12 apply when the main module requires the next module in the module list. Both rules add the required module's *id* to the working list of required modules so that the algorithm will follow the complete require graph. The second rule handles the case where the required

$$\begin{aligned}
&\text{demod } [id, (), ((\text{module } id ((\text{require } id_r @ \text{phase}_r) req_m \dots) (code_{m0} \dots (0 \text{ expr}_m) code_{mn} \dots)) = \text{demod } [id, ((id, (- \text{phase}_r))), \\
&\quad (\text{module } id_r (req_r \dots) (code_{r0} \dots (\text{phase}_r \text{ expr}_r) code_{rn} \dots)) \quad ((\text{module } id (req_m \dots) (code_{m0} \dots (0 (\text{begin } \text{expr}_r \text{ expr}_m)) code_{mn} \dots)) \\
&\quad \text{mod}_n \dots)] \quad (\text{module } id_r (req_r \dots) (code_{r0} \dots (\text{phase}_r \text{ expr}_r) code_{rn} \dots)) \\
&\quad \text{mod}_n \dots)] \\
&\text{where } \text{phase}_{nr} = (- \text{phase}_r) \\
&\text{demod } [id, (), ((\text{module } id ((\text{require } id_r @ \text{phase}_r) req_m \dots) (code_m \dots)) = \text{demod } [id, ((id, (- \text{phase}_r))), \\
&\quad (\text{module } id_r (req_r \dots) (code_r \dots)) \quad ((\text{module } id (req_m \dots) (code_m \dots)) \\
&\quad \text{mod}_n \dots)] \quad (\text{module } id_r (req_r \dots) (code_r \dots)) \\
&\quad \text{mod}_n \dots)]
\end{aligned}$$

Figure 4.12: Demodularization algorithm

module has code that will be at phase 0 for the main module and puts that code before the existing phase 0 code of the main module. The next three rules in Figure 4.13 apply when

$$\begin{aligned}
&\text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), &= \text{demod } [id, ((id, (- \text{phase}_c \text{ phase}_c)) (id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), \\
&\quad ((\text{module } id (req_m \dots) (code_{m0} \dots (0 \text{ expr}_m) code_{mn} \dots)) &\quad ((\text{module } id (req_m \dots) (code_{m0} \dots (0 (\text{begin } \text{expr}_r \text{ expr}_m)) code_{mn} \dots)) \\
&\quad (\text{module } id_c ((\text{require } id_r @ \text{phase}_r) req_c \dots) (code_c \dots)) &\quad (\text{module } id_c (req_c \dots) (code_c \dots)) \\
&\quad (\text{module } id_r (req_r \dots) (code_{r0} \dots (\text{phase}_c \text{ expr}_r) code_{rn} \dots)) &\quad (\text{module } id_r (req_r \dots) (code_{r0} \dots (\text{phase}_c \text{ expr}_r) code_{rn} \dots)) \\
&\quad \text{mod}_n \dots)] &\quad \text{mod}_n \dots)] \\
&\text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), &= \text{demod } [id, ((id, (- \text{phase}_c \text{ phase}_c)) (id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), \\
&\quad ((\text{module } id (req_m \dots) (code_m \dots)) &\quad ((\text{module } id (req_m \dots) (code_m \dots)) \\
&\quad (\text{module } id_c ((\text{require } id_r @ \text{phase}_r) req_c \dots) (code_c \dots)) &\quad (\text{module } id_c (req_c \dots) (code_c \dots)) \\
&\quad (\text{module } id_r (req_r \dots) (code_r \dots)) &\quad (\text{module } id_r (req_r \dots) (code_r \dots)) \\
&\quad \text{mod}_n \dots)] &\quad \text{mod}_n \dots)] \\
&\text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), &= \text{demod } [id, ((id_n \text{ phase}_n) \dots), \\
&\quad ((\text{module } id (req_m \dots) (code_m \dots)) &\quad ((\text{module } id (req_m \dots) (code_m \dots)) \\
&\quad (\text{module } id_c () (code_c \dots)) &\quad (\text{module } id_c () (code_c \dots)) \\
&\quad \text{mod}_n \dots)] &\quad \text{mod}_n \dots)]
\end{aligned}$$

Figure 4.13: Demodularization algorithm

handling the working list of required modules. The first of the three is similar to the second rule because it extracts code that will be in phase 0 of the main module and inserts it into the main module. The second rule of the three working list rules handles the case where there is no matching code for phase 0. The third rule removes an entry from the working list when the module has no more requires.

$$\begin{aligned}
&\text{demod } [id, (), (\text{mod}_0 \dots (\text{module } id (req \dots) (code \dots)) \text{mod}_n \dots)] &= \text{demod } [id, (), ((\text{module } id (req \dots) (code \dots)) \text{mod}_0 \dots \text{mod}_n \dots)] \\
&\text{demod } [id, (), ((\text{module } id ((\text{require } id_r @ \text{phase}_r) req_m \dots) (code_m \dots)) = \text{demod } [id, (), ((\text{module } id ((\text{require } id_r @ \text{phase}_r) req_m \dots) (code_m \dots)) \\
&\quad \text{mod}_i \dots (\text{module } id_r (req_r \dots) (code_r \dots)) \text{mod}_n \dots)] &\quad (\text{module } id_r (req_r \dots) (code_r \dots)) \text{mod}_i \dots \text{mod}_n \dots)] \\
&\text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), &= \text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), \\
&\quad ((\text{module } id (req_m \dots) (code_m \dots)) &\quad ((\text{module } id (req_m \dots) (code_m \dots)) \\
&\quad \text{mod}_i \dots &\quad (\text{module } id_c (req_c \dots) (code_c \dots)) \\
&\quad (\text{module } id_c (req_c \dots) (code_c \dots)) &\quad \text{mod}_i \dots \\
&\quad \text{mod}_n \dots)] &\quad \text{mod}_n \dots)] \\
&\text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), &= \text{demod } [id, ((id_c \text{ phase}_c) (id_n \text{ phase}_n) \dots), \\
&\quad ((\text{module } id (req_m \dots) (code_m \dots)) &\quad ((\text{module } id (req_m \dots) (code_m \dots)) \\
&\quad (\text{module } id_c ((\text{require } id_r @ \text{phase}_r) req_c \dots) (code_c \dots)) &\quad (\text{module } id_c ((\text{require } id_r @ \text{phase}_r) req_c \dots) (code_c \dots)) \\
&\quad \text{mod}_i \dots &\quad (\text{module } id_r (req_r \dots) (code_r \dots)) \\
&\quad (\text{module } id_r (req_r \dots) (code_r \dots)) &\quad \text{mod}_i \dots \\
&\quad \text{mod}_n \dots)] &\quad \text{mod}_n \dots)]
\end{aligned}$$

Figure 4.14: Demodularization algorithm

The final four rules in Figure 4.14 rearrange the program’s module list so that modules that require each other are adjacent in the list and the other rules can apply.

4.5 Equivalence

We claim that the programs will evaluate to the same answers before and after demodularization.

Theorem. *Evaluating a program P a number of steps n and the demodularized program P' a number of steps m will be bisimilar with respect to the stores of the programs as shown in Figure 4.15.*

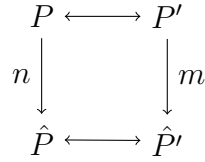


Figure 4.15: Bisimulation of a program before and after demodularization

Proof. The proof is by induction on the DAG of required modules in a program. For each shape of possible requires in a program, we would show the evaluation before and after demodularization end up running the same module code in the same order and result in the same stores. □

The full proof of the theorem would not be instructive because of the complexities of the implementation of the model and the demodularization algorithm.

Chapter 5

Implementation

The implementation of demodularization for the Racket module system is included in the Racket distribution as part of the build tool `raco`. Users can pass in a module they would like to demodularize and optimize to the tool and it will compile all the necessary modules, run the demodularization algorithm on them, decompile the resulting module, and run optimizations on it to produce the final module. The demodularization algorithm runs on compiled modules, which are stored in Racket bytecode files.

5.1 Racket bytecode

Racket bytecode is different from most other forms of bytecode in that it mostly maintains the structure of the abstract syntax tree from the original program. The main aspect that changes between fully-expanded Racket programs and Racket bytecode is the use of stack positions instead of variables and the addition of stack-manipulating instructions. A full explanation of Racket bytecode and what it means can be found in ??.

A table is created in the bytecode for the top level bindings in a module, which is known as the prefix. All top level bindings have an entry in the prefix, including bindings that come from other modules, and any use of a top level binding in a program is replaced with a numeric reference to the entry. Listing 13 shows a simple program written in the most basic language that Racket supports: the `kernel` language. When this code is compiled, it turns into the bytecode in Listing ??.

In the bytecode, all references to the variable `y` have been replaced with references

```

(module hello '%kernel
  (%require racket/private/misc)
  (define-values (y) (random))
  (let-values ([ (x) #f ])
    (set! x 5)
    (displayln x)
    (displayln (+ x y))))

```

Listing 13: Example program written in kernel language

```

(module hello
  (prefix (y (module-variable racket/private/misc displayln 7)))
  (requires '%kernel racket/private/misc)
  (def-values (toplevel 0)
    (app (primval 273)))
  (let-one
    (box-env
      (install-value 0 5
        (seq
          (app (toplevel 1) (local 1))
          (app (toplevel 1)
            (app (primval 247) (local 3) (toplevel 0)))))))

```

Listing 14: Bytecode representation of program from Listing 13

to the prefix in the form of `(toplevel 0)`. All references to `displayln` have also been replaced with `(toplevel 1)` references, which in turn references element 7 of the prefix of the `racket/private/misc` module. The references to `x` are now all different because the stack changes between each usage of `x`. The references to `random` and `+` are replaced with references to primitive implementations in the Racket runtime.

5.2 Demodularization

The implementation of demodularization for the Racket module system is written in Racket and consumes and produces Racket bytecode. The library for reading and writing Racket bytecode in Racket was mainly used for debugging purposes and was incomplete, so the first task was to fully implement the Racket bytecode library.

The actual algorithm for demodularization is similar to the model in the previous chapter in that it traces requires and includes all phase 0 code in the final module, but it also has to deal with the module prefix and references to the prefix. When the algorithm includes a required module, it also includes that module’s prefix appended to the end of the main module’s prefix. Then, it must adjust all references in that module’s code to point at the new combined prefix. Also, the algorithm tracks cross-module references and rewrites them to point to the new prefix as well. In the example in Listing 14, when the algorithm includes the module `racket/private/misc` it must rewrite the reference to `(toplevel 1)` to whatever the new position for `displayln` is in the combined prefix.

5.3 Optimization

For demodularization to be useful, the program needs to be optimized after demodularizing it. To optimize the demodularized bytecode, we wanted the existing Racket optimizer built into the Racket compiler so that we get all existing optimizations “for free” and any new optimizations in the future will work on both regularly compiled and demodularized programs. The existing optimizer works on an intermediate form between fully-expanded Racket code and Racket bytecode. This intermediate form only exists as C data structures in the implementation of Racket. Therefore, to use the existing optimizer, Racket bytecode needed to be decompiled into this intermediate form.

5.4 Decompilation

The decompiler was written in C, so that it could interoperate with the optimizer. There are three main differences between Racket bytecode and the intermediate form: stack positions, cyclic closures, and reference arguments. In Racket bytecode and in the intermediate form, all references to bindings are numeric, but in bytecode the references are stack positions and in the intermediate form the references are lexical. For example, the bytecode in Listing ?? has three references to `x`, in the form of `(install-value 0 5)`, `(local 1)`, and `(local 3)`.

In the intermediate form, all of these references to `x` should be `0`.

Racket bytecode allows for cyclic closures, or closures which contain a reference to themselves. Nothing like this exists in the intermediate form, so to decompile cyclic closures, the decompiler creates a new top level definition for the closure and replaces references to the closure (including the self-references) with references to the top level definition.

Finally, Racket bytecode allows reference arguments (like C++ reference arguments) in functions, but the intermediate form doesn't allow them. The decompiler turns reference arguments into `case-lambda` closures over the reference arguments with one case for getting the argument value and one case for setting the argument value.

5.5 Limitations

Racket provides features that treat modules as first-class objects during runtime. For example, programs can load and evaluate modules at runtime through `dynamic-require`. Because the demodularizer cannot know ahead of time what modules might be loaded at runtime, it disallows programs that use `dynamic-require`. If the modules loaded through `dynamic-require` were completely separate (meaning they do not share any required modules) from the main program, it would be possible to demodularize the program, but in practice most modules required at runtime will share with the main program.

Chapter 6

Evaluation

We tested our implementation of demodularization by selecting existing Racket programs and measuring their execution time before and after demodularization. We also measured the memory usage and compiled bytecode size of the programs. We ran the benchmarks on an Intel Core 2 Quad machine running Ubuntu and ran each program X times. We expect programs to perform better based on how modular the program is, which we measure by counting the number of modules in a program's require graph and how many cross module references occur in the program.

Figure XXX shows the results of running this experiment on XXX Racket programs. On one end of the spectrum, there are programs like XXX which are already basically single module programs, so demodularization does little besides rerun the optimizer on the program. Running the optimizer again may have positive or negative effects on performance, it may unroll loops and inline definitions more aggressively the second time, but some of these "optimizations" may hurt performance. On the other end of the spectrum, highly modular programs like XXX perform much better after demodularization. We expect performance to increase at a linear or even superlinear pace as modularity increases because of the extra information available to the optimizer.

This experiment uses only the existing Racket optimizations, which are intra-module optimizations. Certain optimizations that are not worthwhile to do at the intra-module level have larger payoffs when applied to whole programs. With demodularization, we anticipate that new whole-program optimizations enabled by demodularization will increase

performance even more.

Chapter 7

Related Work

Prior work on whole-program optimization has come in two flavors, depending on how much access to the source code the optimizer has. The first approach assumes full access to the source code and is based on inlining. The second approach only has access to compiled modules and is based on combining modules.

The first approach is based on selectively inlining code across module boundaries because it has full access to the source code of the program [1, 2]. Most of the focus of this approach is finding appropriate heuristics to inline certain functions without ballooning the size of the program and making sure the program still produces the same results. Resulting programs are not completely demodularized; they still have some calls to other modules. Specifically, Chambers et al. [2] show how this approach applies to object-oriented languages like C++ and Java, where they are able to exploit properties of the class systems to choose what to inline. Blume and Appel [1] showed how to deal with inlining in the presence of higher order functions, to make sure the semantics of the program didn't change due to inlining. Their approach led to performance increases of around 8%.

The second approach is taking already compiled modules, combining them into a single module, and optimizing the single module at link time [3, 4]. Most of the work done with this approach optimized at the assembly code level, but because they were able to view the whole program, the performance increases were still valuable. The link-time optimization system by Sutter et al. [3] achieves a 19% speedup on C programs. One of the reasons for starting with compiled modules is so that programs using multiple languages

can be optimized in a common language, like the work done by Debray et al. [4] to combine a program written in both Scheme and Fortran. The main problem with this approach is that the common language has less information for optimization than the source code had. These approaches are similar to demodularization, but they operate at a lower level and work on languages without phased module systems.

Chapter 8

Conclusion

Demodularization is a useful optimization for deploying modular programs. A programmer can write a modular program and get the benefits of separate compilation while developing the program, and then get additional speedups by running the demodularizer on the completed program. Demodularization also enables new optimizations that are not feasible to implement for modular programs. Without module boundaries, inter-procedural analysis is much easier and worthwhile. Also, dead code elimination works much better because the whole program is visible, while in a modular program, only dead code that is private to the module can be eliminated.

In the future, we would like to implement an aggressive dead code elimination algorithm for Racket. We implemented a naive one that does not respect side effects, but shows the potential gains from this optimization; it is able to shrink Racket binaries down from about 2MB to about 100KB. This promising result implies that other low-hanging optimizations should be possible on demodularized programs that can increase performance.

References

- [1] Matthias Blume and Andrew W. Appel. Lambda-splitting: a higher-order approach to cross-module optimizations. In *Proceedings of the second ACM SIGPLAN international conference on Functional programming*, ICFP '97, pages 112–124, New York, NY, USA, 1997. ACM.
- [2] Craig Chambers, Jeffrey Dean, and David Grove. Whole-program optimization of object-oriented languages. Technical report, 1996.
- [3] Bjorn De Sutter, Bruno De Bus, and Koen De Bosschere. Link-time binary rewriting techniques for program compaction. *ACM Trans. Program. Lang. Syst.*, 27(5):882–945, September 2005.
- [4] Saumya K. Debray, Robert Muth, and Scott A. Watterson. Link-time improvement of scheme programs. In *Proceedings of the 8th International Conference on Compiler Construction, Held as Part of the European Joint Conferences on the Theory and Practice of Software, ETAPS'99, CC '99*, pages 76–90, London, UK, UK, 1999. Springer-Verlag.
- [5] Matthew Flatt. Composable and compilable macros:: you want it when? In *Proceedings of the seventh ACM SIGPLAN international conference on Functional programming*, ICFP '02, pages 72–83, New York, NY, USA, 2002. ACM.