

# Laziness by Need

Stephen Chang

Northeastern University  
stchang@ccs.neu.edu

**Abstract.** Lazy functional programming has many benefits that strict functional languages can simulate via lazy data constructors. In recognition, ML, Scheme, and other strict functional languages have supported lazy stream programming with `delay` and `force` for several decades. Unfortunately, the manual insertion of `delay` and `force` can be tedious and error-prone.

We present a semantics-based refactoring that helps strict programmers manage manual lazy programming. The refactoring uses a static analysis to identify where additional `delays` and `forces` might be needed to achieve the desired simplification and performance benefits, once the programmer has added the initial lazy data constructors. The paper presents a correctness argument for the underlying transformations and some preliminary experiences with a prototype tool implementation.

## 1 Laziness in a Strict World

A lazy functional language naturally supports the construction of reusable components and their composition into reasonably efficient programs [12]. For example, the solution to a puzzle may consist of a generator that produces an easily-constructed stream of all *possible* solutions and a filter that extracts the desired *valid* solutions. Due to laziness, only a portion of the possible solutions are explored. Put differently, lazy composition appears to naturally recover the desired degree of efficiency without imposing a contorted programming style.

Unfortunately, programming in a lazy language comes at a cost. Not only are data constructors lazy, but all functions are as well. This pervasiveness of laziness makes it difficult to predict the behavior and time/space performance of lazy programs. As several researchers noticed [2,6,15,16,23], however, most programs need only a small amount of laziness. In response, people have repeatedly proposed lazy programming in strict functional languages [1,8,20,25,27]. In fact, Scheme [22] and ML [3] have supported manual stream programming with `delay` and `force` for decades. Using `delay` and macros, a programmer can easily turn an eager, Lisp-style list constructor into a lazy one [11], while `force` retrieves the value from a delayed computation.

However, merely switching from eager constructors to lazy ones is often not enough to achieve the performance benefits of laziness. The insertion of one `delay` tends to require additional `delays` elsewhere in the program to achieve the desired lazy behavior. Since these additional `delay` insertions depend on

the value flow of the program, it can be difficult to determine where to insert them, especially in the presence of higher-order functions. In short, manual lazy programming is challenging and error-prone.

In response, we introduce a static analysis-based refactoring that assists programmers with the task of inserting `delays` and accompanying `forces`. We imagine a programmer who wishes to create a lazy generator and starts using lazy constructs in the obvious places. Our transformation then inserts additional `delays` and `forces` to achieve the desired lazy performance benefit.

The paper is organized as follows. The second section introduces some motivating examples. Section 3 presents the analysis-based program transformation, and section 4 argues its correctness. Section 5 sketches a prototype implementation, and section 6 describes real-world applications. Section 7 compares our approach with other attempts at taming laziness. Finally, section 8 lists some ideas for future work.

## 2 Motivating Examples

Nearly every modern strict programming language supports laziness, either via `delay` and `force`, or in the form of a streams or other lazy data structure library. None of these languages offer much help, however, in figuring out the right way to use these forms. To illustrate the problems, this section presents three examples in three distinct languages, typed and untyped. The first one, in Racket [10], shows how conventional program reorganizations can eliminate the performance benefits of laziness without warning. The second, in Scala [19], demonstrates how laziness propagates across function calls. The third example illustrates the difficulties of developing an idiomatic lazy  $n$ -queens algorithm in a strict language like OCaml [14]. That is, the problems of programming lazily in a strict language are universal across many languages.

### 2.1 Reorganizations Interfere with Laziness

Using `delay` and `force` occasionally confuses even the most experienced programmers. This subsection retells a recent story involving a senior Racket developer. A game tree is a data structure representing all possible sequences of moves in a game. It is frequently employed in AI algorithms to calculate an optimal next move, and it is also useful for game developers wishing to experiment with the rules of a game. For anything but the simplest games, however, the multitude of available moves at each game state results in an unwieldy or even infinite game tree. Thus, laziness is frequently utilized to manage such trees.

The Racket code to generate a game tree might roughly look like this:

```
;; A GameTree (short: GT) is one of:
;; -- (GT-Leaf GameState)
;; -- (GT-Node GameState Player [ListOf Move])

;; A Move is a (Move Name Position GameTree)
```

```

;; gen-GT : GameState Player -> GameTree
(define (gen-GT game-state player)
  (if (final-state? game-state)
      (GT-Leaf game-state)
      (GT-Node game-state player (calc-next-moves game-state player))))

;; calc-next-moves : GameState Player -> [ListOf Move]
(define (calc-next-moves game-state player)
  ((for each possible attacker and target in game-state:))
   (define new-state ...)
   (define new-player ...)
   (Move attacker target (gen-GT new-state new-player)))

```

A game tree is created with the `gen-GT` function, which takes a game state and the current active player. If the given state is a final state, then a `GT-Leaf` node is created. Otherwise, a `GT-Node` is created with the current game state, the current player, and a list of moves from the given game state. The `calc-next-moves` function creates a list of `Move` structures, where each move contains a new game tree starting from the game state resulting from the move.

An upcoming, Racket-based programming book utilizes such a game tree. Initially, only a small game is implemented, so `Move` is defined as a strict constructor. As the book progresses, however, the game tree becomes unwieldy as more features are added to the game. In response, the third argument of the `Move` structure is changed to be lazy, meaning the call to the `Move` constructor implicitly wraps the third argument with a `delay`.<sup>1</sup> With the lazy `Move` constructor, the code above generates only the first node of a game tree.

To prepare the book for typesetting, an author reorganized the definition of `calc-next-moves` in a seemingly innocuous fashion to fit it within the margins of a page:

```

;; calc-next-moves : GameState Player -> [ListOf Move]
(define (calc-next-moves game-state player)
  ((for each possible attacker and target in game-state:))
   (define new-state ...)
   (define new-player ...)
   (define new-gt (gen-GT new-state new-player))
   (Move attacker target new-gt))

```

The underlined code above pulls the generation of the game tree into a separate definition. As the astute reader will recognize, the new game tree is no longer created lazily. Even though the `Move` constructor is lazy in the third position, the benefits of laziness are lost. Even worse, such a performance bug is easily unnoticed because the program still passes all unit tests.

In contrast, our laziness transformation recognizes that the `new-gt` variable flows into the lazy position of the `Move` constructor, and in turn, proposes a `delay` around the construction of the new game tree.

---

<sup>1</sup> Specifically, `Move` becomes a macro that expands to a private constructor call where the third argument is delayed. This is a common idiom in Lisp-like languages [11].

## 2.2 Laziness Must Propagate

A 2009 blog post<sup>2</sup> illustrates a related tricky situation in the following Scala [19] example. Scala delays method arguments whose type is marked with `=>`, as in:<sup>3</sup>

```
def foo[A,B](a: A, b: => B): B = ...
```

When `foo` is called, its second argument is not evaluated until its value is needed inside the function body. However, if another function, `bar`, calls `foo`:

```
def bar[C,A,B](c: C, a: A, b: B): B = { ... foo(a, b) }
```

the `b` argument is evaluated when `bar` is called, thus negating the benefit of laziness in `foo`. To recover it, we must delay the third argument to `bar`:

```
def bar[C,A,B](c: C, a: A, b: => B): B = ...
```

If yet another function calls `bar` then that function must delay its argument as well. For programs with complex call graphs, the required delay points may be scattered throughout the program, making programmer errors more likely. Our transformation is designed to help with just such situations.

## 2.3 Idiomatic Lazy Programming in a Strict Language

The *n*-queens problem makes an illustrative playground for advertising lazy programming. An idiomatic lazy solution to such a puzzle may consist of just two parts: a part that places *n* queens at arbitrary positions on an *n* by *n* chess board, and a part for deciding whether a particular placement is a solution to the puzzle. Given these two components, a one-line function calculates a solution:

```
let nqueens n = hd (filter isValid all_placements)
```

The `all_placements` variable stands for a stream of all possible placements of *n* queens; `filter isValid` eliminates placements with conflicting queens; and `hd` picks the first valid one. Lazy evaluation guarantees that `filter isValid` traverses `all_placements` for just enough placements to find the first solution.

The approach cleanly separates two distinct concerns. While `all_placements` may ignore the rules of the puzzle, it is the task of `isValid` to enforce them. If the components were large, two different programmers could tackle them in parallel. All they would have to agree on is the representation of queen placements, for which we choose a list of board coordinates  $(r, c)$ . The rest of the section explains how an OCaml [14] programmer may develop such a lazy algorithm. Here is `all_placements`:

```
let process_row r qss_so_far =
  foldr (fun qs new_qss -> (map (fun c -> (r,c)::qs) (rng n)) @ new_qss)
    [] qss_so_far
```

```
let all_placements = foldl process_row [[]] (rng n)
```

<sup>2</sup> <http://pchiusano.blogspot.com/2009/05/optional-laziness-doesnt-quite-cut-it.html>

<sup>3</sup> The `=>` syntax specifies “by-name” parameter passing for this position but the distinction between “by-name” and “lazy” is inconsequential here.

Brackets denote lists, `rng n` is `[1..n]`, `::` is infix `cons`, and `@` is infix `append`. All possible placements are generated by adding one coordinate at a time. The `process_row` function, given a row `r` and a list of placements `qss_so_far`, duplicates each placement in `qss_so_far` `n` times, adding to each copy a new coordinate of `r` with a different column `c`, and then appends all these new placements to the final list of all placements. The `process_row` function is called `n` times, once per row. The result of evaluating `all_placements` looks like this:

```
[[ (n,1); (n-1,1); ... ; (1,1) ];
  ... ;
  [ (n,n); (n-1,n); ... ; (1,n) ]]
```

where each line represents one possible placement.

Since OCaml is strict, however, using `all_placements` with the `nqueens` function from earlier generates all possible placements before testing each one of them for validity. This computation is obviously time consuming and performs far more work than necessary. For instance, here is the timing for `n = 8` queens:<sup>4</sup>

```
real 0m52.122s  user 0m51.399s  sys 0m0.468s
```

If the programmer switches to lazy lists to represent `all_placements`, then only a portion of the possible placements should be explored. Specifically, all instances of `cons (::)` are replaced with its lazy variant, represented with `::l` below. In this setting, lazy `cons` is defined using OCaml's `Lazy` module and is `cons` with a delayed rest list. It is also necessary to add `forces` where appropriate.<sup>5</sup> For example, here is `append (@)` and `map` with lazy `cons` (`[]` also represents the empty lazy list):<sup>6</sup>

```
let rec (@) lst1 lst2 =
  match force lst1 with
  | [] -> lst2
  | x::lxs -> x::ldelay (xs @ lst2)

let rec map f lst =
  match force lst with
  | [] -> []
  | x::lxs -> f x::ldelay (map f xs)
```

Running this program, however, surprises our lazy-strict programmer:

```
real 1m3.720s  user 1m3.072s  sys 0m0.352s
```

With lazy `cons` and `force`, the program runs even slower than the strict version. Using lazy `cons` naïvely does not seem to generate the expected performance gains. Additional `delays` and `forces` are required, though it is not immediately obvious where to insert them. This step is precisely where our analysis-based refactoring transformation helps a programmer. In this particular case, our transformation would insert a `delay` in the `foldr` function:

<sup>4</sup> Run on an Intel i7-2600k, 16GB memory machine using the Linux `time` command.

<sup>5</sup> “Appropriate” here means we avoid Wadler et al.’s [27] “odd” errors.

<sup>6</sup> OCaml’s delaying construct is `lazy` but for clarity and consistency with the rest of the paper we continue to use `delay`. Also, in ML languages, the `delay` is explicit.

```

let rec foldr f base lst =
  match force lst with
  | [] -> base
  | x::lzxs -> f x (delay (foldr f base xs))

```

This perhaps unobvious `delay` is needed because `f`'s second argument eventually flows to a lazy `cons` in `append` (`@`). Without this `delay`, the list of all queen placements is evaluated prematurely. With this refactoring, and an appropriate insertion of `forces`, the lazy-strict programmer sees a dramatic improvement:

```

real 0m3.103s   user 0m3.068s   sys 0m0.024s

```

Lazy programmers are already familiar with such benefits, but our refactoring transformation enables strict programmers to reap the same benefits as well.

### 3 Refactoring For Laziness

The heart of our refactoring is a whole-program analysis that calculates where values may flow. Our transformation uses the results of the analysis to insert `delays` and `forces`. Section 3.1 describes the core of our strict language. We then present our analysis in three steps: section 3.2 explains the analysis rules for our language; section 3.3 extends the language and analysis with lazy forms: `delay`, `force`, and lazy `cons` (`lcons`); and section 3.4 extends the analysis again to calculate the potential insertion points for `delay` and `force`. Finally, section 3.5 defines the refactoring transformation function.

#### 3.1 Language Syntax

Our starting point is an untyped<sup>7</sup> functional core language. The language is *strict* and uses a standard expression notation:

$$\begin{aligned}
 e \in \text{Exp} = & n \mid b \mid x \mid \lambda(x \dots).e \mid ee \dots \mid oe \mid \text{zero? } e \mid \text{not } e \mid \text{if } ee \ e \\
 & \mid \text{let } x = e \text{ in } e \mid \text{null} \mid \text{cons } ee \mid \text{first } e \mid \text{rest } e \mid \text{null? } e \\
 n \in \mathbb{Z}, \quad & b \in \text{Bool} = \text{true} \mid \text{false}, \quad x \in \text{Var}, \quad o \in \text{Op} = + \mid - \mid * \mid / \mid < \mid > \mid = \mid \text{or} \mid \text{and}
 \end{aligned}$$

There are integers, booleans, variables,  $\lambda$ s, applications, boolean and arithmetic primitives, conditionals, (non-recursive) lets, and eager lists and list operations. Here are the values, where both components of a non-empty list must be values:

$$v \in \text{Val} = n \mid b \mid \lambda(x \dots).e \mid \text{null} \mid \text{cons } v \ v$$

A program  $p$  consists of two pieces: a series of mutually referential function definitions and an expression that may call the functions:

$$p \in \text{Prog} = d \dots e \qquad d \in \text{Def} = \text{define } f(x \dots) = e$$


---

<sup>7</sup> Standard type systems cannot adequately express the flow of laziness and thus cannot solve the `delay`-insertion problems from section 2. A type error can signal a missing `force`, but a type system will not suggest where to add performance-related `delays`. Thus we omit types for this first step in our research.

### 3.2 Analysis Step 1: 0-CFA

Our initial analysis is based on 0-CFA [13,24,26]. The analysis assumes that each subexpression has a unique label  $\ell$ , also drawn from  $Var$ , but that the set of labels and the set of variables in a program are disjoint. The analysis computes an abstract environment  $\widehat{\rho}$  that maps elements of  $Var$  to sets of abstract values:

$$\widehat{\rho} \in Env = Var \rightarrow \mathcal{P}(\widehat{Val}) \quad \ell \in Var \quad \widehat{v} \in \widehat{Val} = \mathbf{val} \mid \lambda(x \dots).\ell \mid \mathbf{cons} \ell \ell$$

A set  $\widehat{\rho}(x)$  or  $\widehat{\rho}(\ell)$  represents an approximation of all possible values that can be bound to  $x$  or observed at  $\ell$ , respectively, during evaluation of the program.

The analysis uses abstract representations of values,  $\widehat{v}$ , where  $\mathbf{val}$  stands for all literals in the language. In addition,  $\lambda(x \dots).\ell$  are abstract function values where the body is represented with a label, and  $(\mathbf{cons} \ell \ell)$  are abstract list values where the  $\ell$ 's are the labels of the respective pieces. We overload the  $\widehat{\cdot}$  notation to denote a function that converts a concrete value to its abstract counterpart:

$$\begin{array}{ccc} \widehat{n} = \mathbf{val} & \widehat{b} = \mathbf{val} & \widehat{\mathbf{null}} = \mathbf{val} \\ \lambda(\widehat{x \dots}).e^\ell = \lambda(x \dots).\ell & \widehat{\mathbf{cons} v_1^{\ell_1} v_2^{\ell_2}} = \mathbf{cons} \ell_1 \ell_2 & \widehat{\cdot} : Val \rightarrow \widehat{Val} \end{array}$$

We present our analysis with a standard [18], constraints-based specification, where notation  $\widehat{\rho} \models p$  means  $\widehat{\rho}$  is an acceptable approximation of program  $p$ . Figures 1 and 2 show the analysis for programs and expressions, respectively.

The  $[prog]$  rule specifies that environment  $\widehat{\rho}$  satisfies program  $p = d \dots e$  if it satisfies all definitions  $d \dots$  as well as the expression  $e$  in the program. The  $[def]$  rule says that  $\widehat{\rho}$  satisfies a definition if the corresponding abstract  $\lambda$ -value is included for variable  $f$  in  $\widehat{\rho}$ , and if  $\widehat{\rho}$  satisfies the function body as well.

In figure 2, the  $[num]$ ,  $[bool]$ , and  $[null]$  rules show that  $\mathbf{val}$  represents these literals in the analysis. The  $[var]$  rule connects variables  $x$  and their labels  $\ell$ , specifying that all values bound to  $x$  should also be observable at  $\ell$ . The  $[lam]$  rule for an  $\ell$ -labeled  $\lambda$  says that its abstract version must be in  $\widehat{\rho}(\ell)$  and that  $\widehat{\rho}$  must satisfy its body. The  $[app]$  rule says that  $\widehat{\rho}$  must satisfy the function and arguments in an application. In addition, for each possible  $\lambda$  in the function position, the arguments must be bound to the corresponding parameters of that  $\lambda$  and the result of evaluating the  $\lambda$ 's body must also be a result for the application itself. The  $[let]$  rule has similar constraints. The  $[op]$ ,  $[zero?]$ ,  $[not]$ , and  $[null?]$  rules require that  $\widehat{\rho}$  satisfy a primitive's operands and uses  $\mathbf{val}$  as the result. The  $[if]$  rule requires that  $\widehat{\rho}$  satisfy the test expression and the two branches, and that any resulting values in the branches also be a result for the entire

$$\begin{array}{l} \widehat{\rho} \models d \dots e \text{ iff} \\ \widehat{\rho} \models_d d \wedge \dots \wedge \widehat{\rho} \models_e e \end{array} \quad [prog] \quad \left| \quad \begin{array}{l} \widehat{\rho} \models_d \mathbf{define} f(x \dots) = e^\ell \text{ iff} \\ \lambda(x \dots).\ell \in \widehat{\rho}(f) \wedge \widehat{\rho} \models_e e^\ell \end{array} \quad [def]$$

Fig. 1. 0-CFA analysis on programs

$\hat{\rho} \models_e n^\ell$ iff $\mathbf{val} \in \hat{\rho}(\ell)$	[num]	$\hat{\rho} \models_e (\mathbf{zero}? e_1^{\ell_1})^\ell$ iff	[zero?]
$\hat{\rho} \models_e b^\ell$ iff $\mathbf{val} \in \hat{\rho}(\ell)$	[bool]	$\hat{\rho} \models_e e_1^{\ell_1} \wedge \mathbf{val} \in \hat{\rho}(\ell)$	
$\hat{\rho} \models_e x^\ell$ iff $\hat{\rho}(x) \subseteq \hat{\rho}(\ell)$	[var]	$\hat{\rho} \models_e (\mathbf{not} e_1^{\ell_1})^\ell$ iff	[not]
$\hat{\rho} \models_e (\lambda(x \dots). e_0^{\ell_0})^\ell$ iff	[lam]	$\hat{\rho} \models_e e_1^{\ell_1} \wedge \mathbf{val} \in \hat{\rho}(\ell)$	
$\lambda(x \dots). \ell_0 \in \hat{\rho}(\ell) \wedge \hat{\rho} \models_e e_0^{\ell_0}$		$\hat{\rho} \models_e (\mathbf{if} e_1^{\ell_1} e_2^{\ell_2} e_3^{\ell_3})^\ell$ iff	[if]
$\hat{\rho} \models_e (e_f^{\ell_f} e_1^{\ell_1} \dots)^\ell$ iff	[app]	$\hat{\rho} \models_e e_1^{\ell_1} \wedge \hat{\rho} \models_e e_2^{\ell_2} \wedge \hat{\rho}(\ell_2) \subseteq \hat{\rho}(\ell)$	
$\hat{\rho} \models_e e_f^{\ell_f} \wedge \hat{\rho} \models_e e_1^{\ell_1} \wedge \dots \wedge$		$\wedge \hat{\rho} \models_e e_3^{\ell_3} \wedge \hat{\rho}(\ell_3) \subseteq \hat{\rho}(\ell)$	
$(\forall \lambda(x_1 \dots). \ell_0 \in \hat{\rho}(\ell_f) :$		$\hat{\rho} \models_e \mathbf{null}^\ell$ iff $\mathbf{val} \in \hat{\rho}(\ell)$	[null]
$\hat{\rho}(\ell_1) \subseteq \hat{\rho}(x_1) \wedge \dots \wedge$		$\hat{\rho} \models_e (\mathbf{null}? e_1^{\ell_1})^\ell$ iff	[null?]
$\hat{\rho}(\ell_0) \subseteq \hat{\rho}(\ell))$		$\hat{\rho} \models_e e_1^{\ell_1} \wedge \mathbf{val} \in \hat{\rho}(\ell)$	
$\hat{\rho} \models_e (\mathbf{let} x = e_1^{\ell_1} \mathbf{in} e_0^{\ell_0})^\ell$ iff	[let]	$\hat{\rho} \models_e (\mathbf{cons} e_1^{\ell_1} e_2^{\ell_2})^\ell$ iff	[cons]
$\hat{\rho} \models_e e_1^{\ell_1} \wedge \hat{\rho}(\ell_1) \subseteq \hat{\rho}(x) \wedge$		$\hat{\rho} \models_e e_1^{\ell_1} \wedge \hat{\rho} \models_e e_2^{\ell_2} \wedge (\mathbf{cons} \ell_1 \ell_2) \in \hat{\rho}(\ell)$	
$\hat{\rho} \models_e e_0^{\ell_0} \wedge \hat{\rho}(\ell_0) \subseteq \hat{\rho}(\ell)$		$\hat{\rho} \models_e (\mathbf{first} e_1^{\ell_1})^\ell$ iff $\hat{\rho} \models_e e_1^{\ell_1} \wedge$	[first]
$\hat{\rho} \models_e (o e_1^{\ell_1} e_2^{\ell_2})^\ell$ iff	[op]	$(\forall (\mathbf{cons} \ell_2 \_ ) \in \hat{\rho}(\ell_1) : \hat{\rho}(\ell_2) \subseteq \hat{\rho}(\ell))$	
$\hat{\rho} \models_e e_1^{\ell_1} \wedge \hat{\rho} \models_e e_2^{\ell_2} \wedge \mathbf{val} \in \hat{\rho}(\ell)$		$\hat{\rho} \models_e (\mathbf{rest} e_1^{\ell_1})^\ell$ iff $\hat{\rho} \models_e e_1^{\ell_1} \wedge$	[rest]
		$(\forall (\mathbf{cons} \_ \ell_2) \in \hat{\rho}(\ell_1) : \hat{\rho}(\ell_2) \subseteq \hat{\rho}(\ell))$	

**Fig. 2.** Step 1: 0-CFA analysis on expressions

expression. The [cons] rule for an  $\ell$ -labeled, eager **cons** requires that  $\hat{\rho}$  satisfy both arguments and that a corresponding abstract **cons** value be in  $\hat{\rho}(\ell)$ . Finally, the [first] and [rest] rules require satisfiability of their arguments and that the appropriate piece of any **cons** arguments be a result of the entire expression.

### 3.3 Analysis Step 2: Adding delay and force

Next we extend our language and analysis with lazy forms:

$$e \in \mathit{Exp} = \dots \mid \mathbf{delay} e \mid \mathbf{force} e \mid \mathbf{lcons} e e$$

$$\text{where } \mathbf{lcons} e_1 e_2 \stackrel{df}{=} \mathbf{cons} e_1 (\mathbf{delay} e_2)$$

The language is still strict but **delay** introduces promises. A **force** term recursively forces all nested **delays**. Lazy **cons** (**lcons**) is only lazy in its rest argument and **first** and **rest** work with both **lcons** and **cons** values so that **rest** (**lcons**  $v e$ ) results in (**delay**  $e$ ).

We add promises and lazy lists to the sets of values and abstract values, and  $\hat{\cdot}$  is similarly extended. The abstract representation of a **delay** replaces the labeled delayed expression with just the label and the abstract **lcons** is similar.

$$\begin{array}{l}
v \in \text{Val} = \dots \mid \mathbf{delay} \ e \mid \mathbf{lcons} \ v \ e \\
\widehat{v} \in \widehat{\text{Val}} = \dots \mid \mathbf{delay} \ \ell \mid \mathbf{lcons} \ \ell \ \ell \\
\dots \quad \widehat{\mathbf{delay}} \ e^\ell = \mathbf{delay} \ \ell \quad \mathbf{lcons} \ \widehat{v_1^{\ell_1}} \ e_2^{\ell_2} = \mathbf{lcons} \ \ell_1 \ \ell_2 \quad \boxed{\widehat{\cdot} : \text{Val} \rightarrow \widehat{\text{Val}}}
\end{array}$$

Figure 3 presents the new and extended analysis rules. The  $[\mathit{delay}]$  rule specifies that for an  $\ell$ -labeled  $\mathbf{delay}$ , the corresponding abstract  $\mathbf{delay}$  must be in  $\widehat{\rho}(\ell)$  and  $\widehat{\rho}$  must satisfy the delayed subexpression. In addition, the values of the delayed subexpression must also be in  $\widehat{\rho}(\ell)$ . This means that the analysis approximates evaluation of a promise with both a promise and the result of forcing that promise. We discuss the rationale for this constraint below. The  $[\mathit{force}]$  rule says that  $\widehat{\rho}$  must satisfy the argument and that non- $\mathbf{delay}$  arguments are propagated to the outer  $\ell$  label. Since the  $[\mathit{delay}]$  rule already approximates evaluation of the delayed expression, the  $[\mathit{force}]$  rule does not have any such constraints.

We also add a rule for  $\mathbf{lcons}$  and extend the  $[\mathit{first}]$  and  $[\mathit{rest}]$  rules to handle  $\mathbf{lcons}$  values. The  $[\mathit{lcons}]$  rule requires that  $\widehat{\rho}$  satisfy the arguments and requires a corresponding abstract  $\mathbf{lcons}$  at the expressions's  $\ell$  label. The  $[\mathit{first}]$  rule handles  $\mathbf{lcons}$  values just like  $\mathbf{cons}$  values. For the  $[\mathit{rest}]$  rule, a  $\mathbf{delay}$  with the  $\mathbf{lcons}$ 's second component is a possible result of the expression. Just like the  $[\mathit{delay}]$  rule, the  $[\mathit{rest}]$  rule assumes that the lazy component of the  $\mathbf{lcons}$  is both forced and unforced, and thus there is another constraint that propagates the values of the (undelayed) second component to the outer label.

**Implicit Forcing.** In our analysis,  $\mathbf{delays}$  are both evaluated and unevaluated. We assume that during evaluation, a programmer does not want an unforced  $\mathbf{delay}$  to appear in a strict position. For example, if the analysis discovers an unforced  $\mathbf{delay}$  as the function in an application, we assume that the programmer forgot a  $\mathbf{force}$  and analyze that function call anyway. This makes our analysis quite conservative but minimizes the effect of any laziness-related errors in the computed control flow. On the technical side, implicit forcing also facilitates the proof of a safety theorem for the transformation (see subsection 4.3).

$$\begin{array}{l|l}
\widehat{\rho} \models_e (\mathbf{delay} \ e_1^{\ell_1})^\ell \text{ iff} & [\mathit{delay}] \\
(\mathbf{delay} \ \ell_1) \in \widehat{\rho}(\ell) \wedge \widehat{\rho} \models_e e_1^{\ell_1} \wedge \widehat{\rho}(\ell_1) \subseteq \widehat{\rho}(\ell) & \\
\widehat{\rho} \models_e (\mathbf{force} \ e_1^{\ell_1})^\ell \text{ iff} & [\mathit{force}] \\
\widehat{\rho} \models_e e_1^{\ell_1} \wedge (\forall \widehat{v} \in \widehat{\rho}(\ell_1), \widehat{v} \notin \mathbf{delay} : \widehat{v} \in \widehat{\rho}(\ell)) & \\
\widehat{\rho} \models_e (\mathbf{lcons} \ e_1^{\ell_1} \ e_2^{\ell_2})^\ell \text{ iff} & [\mathit{lcons}] \\
\widehat{\rho} \models_e e_1^{\ell_1} \wedge \widehat{\rho} \models_e e_2^{\ell_2} \wedge (\mathbf{lcons} \ \ell_1 \ \ell_2) \in \widehat{\rho}(\ell) & \\
\widehat{\rho} \models_e (\mathbf{first} \ e_1^{\ell_1})^\ell \text{ iff} \dots \wedge [\mathit{first}] & \\
(\forall (\mathbf{lcons} \ \ell_2 \ \_)) \subseteq \widehat{\rho}(\ell_1) : & \\
\widehat{\rho}(\ell_2) \subseteq \widehat{\rho}(\ell) & \\
\widehat{\rho} \models_e (\mathbf{rest} \ e_1^{\ell_1})^\ell \text{ iff} \dots \wedge [\mathit{rest}] & \\
(\forall (\mathbf{lcons} \ \_ \ell_2)) \in \widehat{\rho}(\ell_1) : & \\
(\mathbf{delay} \ \ell_2) \in \widehat{\rho}(\ell) \wedge & \\
\widehat{\rho}(\ell_2) \subseteq \widehat{\rho}(\ell) &
\end{array}$$

**Fig. 3.** Step 2: Analysis with lazy forms

$ \begin{aligned} & (\hat{\rho}, \hat{\mathcal{D}}) \models_e (e_f^{\ell_f} e_1^{\ell_1} \dots)^\ell \text{ iff} & [app] \\ & (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_f^{\ell_f} \wedge (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_1^{\ell_1} \wedge \dots \wedge \\ & (\forall \lambda(x_1 \dots). \ell_0 \in \hat{\rho}(\ell_f) : \\ & \quad \hat{\rho}(\ell_1) \subseteq \hat{\rho}(x_1) \wedge \dots \wedge \\ & \quad \boxed{(\arg \ell_1) \in \hat{\rho}(x_1) \wedge \dots}_1 \wedge \\ & \quad \boxed{(\forall \hat{v} \in \hat{\rho}(\ell_0), \hat{v} \notin \arg : \hat{v} \in \hat{\rho}(\ell))}_2) \end{aligned} $	$ \begin{aligned} & (\hat{\rho}, \hat{\mathcal{D}}) \models_e (\mathbf{delay} e_1^{\ell_1})^\ell \text{ iff} & [delay] \\ & (\mathbf{delay} \ell_1) \in \hat{\rho}(\ell) \wedge \\ & (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_1^{\ell_1} \wedge \hat{\rho}(\ell_1) \subseteq \hat{\rho}(\ell) \wedge \\ & (\forall x \in fv(e_1) : (\forall(\arg \ell_2) \in \hat{\rho}(x) : \\ & \quad \boxed{\ell_2 \in \hat{\mathcal{D}}}_3 \wedge \boxed{(\mathbf{darg} \ell_2) \in \hat{\rho}(x)}_4)) \end{aligned} $
$ \begin{aligned} & (\hat{\rho}, \hat{\mathcal{D}}) \models_e (\mathbf{let} x = e_1^{\ell_1} \text{ in } e_0^{\ell_0})^\ell \text{ iff} & [let] \\ & (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_1^{\ell_1} \wedge \hat{\rho}(\ell_1) \subseteq \hat{\rho}(x) \wedge \\ & \boxed{(\arg \ell_1) \in \hat{\rho}(x)}_1 \wedge (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_0^{\ell_0} \wedge \\ & \boxed{(\forall \hat{v} \in \hat{\rho}(\ell_0), \hat{v} \notin \arg : \hat{v} \in \hat{\rho}(\ell))}_2) \end{aligned} $	$ \begin{aligned} & (\hat{\rho}, \hat{\mathcal{D}}) \models_e (\mathbf{lcons} e_1^{\ell_1} e_2^{\ell_2})^\ell \text{ iff} & [lcons] \\ & (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_1^{\ell_1} \wedge (\hat{\rho}, \hat{\mathcal{D}}) \models_e e_2^{\ell_2} \wedge \\ & (\mathbf{lcons} \ell_1 \ell_2) \in \hat{\rho}(\ell) \wedge \\ & (\forall x \in fv(e_2) : (\forall(\arg \ell_3) \in \hat{\rho}(x) : \\ & \quad \boxed{\ell_3 \in \hat{\mathcal{D}}}_3 \wedge \boxed{(\mathbf{darg} \ell_3) \in \hat{\rho}(x)}_4)) \end{aligned} $

Fig. 4. Step 3a: Calculating flow to lazy positions

### 3.4 Analysis Step 3: Laziness Analysis

Our final refinement revises the analysis to calculate three additional sets, which are used to insert additional **delays** and **forces** in the program:

$$\hat{\mathcal{D}} \in DPos = \mathcal{P}(Var), \quad \hat{\mathcal{S}} \in SPos = \mathcal{P}(Var), \quad \hat{\mathcal{F}} \in FPos = \mathcal{P}(Var \cup (Var \times Var))$$

Intuitively,  $\hat{\mathcal{D}}$  is a set of labels representing function arguments that flow to lazy positions and  $\hat{\mathcal{S}}$  is a set of labels representing arguments that flow to strict positions. Our transformation then delays arguments that reach a lazy position but not a strict position. Additionally,  $\hat{\mathcal{F}}$  collects the labels where a delayed value may appear—both those manually inserted by the programmer and those suggested by the analysis—and is used by the transformation to insert **forces**.

We first describe how the analysis computes  $\hat{\mathcal{D}}$ . The key is to track the flow of arguments from an application into a function body and for this, we introduce a special abstract value ( $\arg \ell$ ), where  $\ell$  labels an argument in a function call.

$$\hat{v} \in \widehat{Val} = \dots \mid \arg \ell$$

Figure 4 presents revised analysis rules related to  $\hat{\mathcal{D}}$ . To reduce clutter, we express the analysis result as  $(\hat{\rho}, \hat{\mathcal{D}})$ , temporarily omitting  $\hat{\mathcal{S}}$  and  $\hat{\mathcal{F}}$ . In the new  $[app]$  and  $[let]$  rules, additional constraints (box 1) specify that for each labeled argument, an **arg** abstract value with a matching label must be in  $\hat{\rho}$  for the corresponding parameter. We are only interested in the flow of arguments within a function's body, so the result-propagating constraint filters out **arg** values (box 2).

Recall that  $\hat{\mathcal{D}}$  is to contain labels of arguments that reach lazy positions. Specifically, if an ( $\arg \ell$ ) value flows to a **delay** or the second position of an

$$\begin{array}{l|l}
(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}}) \models_e (\mathbf{force} \ e_1^\ell)^\ell \text{ iff } [force] & (\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}}) \models_e S[e^\ell] \text{ iff } \dots \wedge [strict] \\
(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}}) \models_e e_1^{\ell_1} \wedge & \boxed{(\forall(\mathbf{arg} \ \ell_1) \in \widehat{\rho}(\ell) : \ell_1 \in \widehat{\mathcal{S}})}_5 \wedge \\
(\forall \widehat{v} \in \widehat{\rho}(\ell_1), \widehat{v} \notin \mathbf{delay} : \widehat{v} \in \widehat{\rho}(\ell)) \wedge & \boxed{(\exists \mathbf{delay} \in \widehat{\rho}(\ell) \Rightarrow \ell \in \widehat{\mathcal{F}})}_6 \wedge \\
\boxed{(\forall(\mathbf{darg} \ \ell_2) \in \widehat{\rho}(\ell_1) : \ell_2 \in \widehat{\mathcal{S}})}_5 & \boxed{(\forall(\mathbf{darg} \ \ell_2) \in \widehat{\rho}(\ell) : (\ell, \ell_2) \in \widehat{\mathcal{F}})}_7
\end{array}$$

where  $S \in SCtx = [] \ e \dots \mid o [] \ e \mid o \ v [] \mid \mathbf{if} [] \ e_1 \ e_2$   
 $\mid \mathbf{zero}? [] \mid \mathbf{not} [] \mid \mathbf{null}? [] \mid \mathbf{first} [] \mid \mathbf{rest} []$

**Fig. 5.** Step 3b: Calculating flow to strict positions

$\mathbf{lcons}$ , then  $\ell$  must be in  $\widehat{\mathcal{D}}$  (box 3) ( $fv(e)$  calculates free variables in  $e$ ). If an  $\ell$ -labeled argument reaches a lazy position, the transformation *may* decide to delay that argument, so the analysis must additionally track it for the purposes of inserting **forces**. To this end, we introduce another abstract value ( $\mathbf{darg} \ \ell$ ),

$$\widehat{v} \in \widehat{Val} = \dots \mid \mathbf{darg} \ \ell$$

and insert it when needed (box 4). While  $(\mathbf{arg} \ \ell)$  can represent any argument,  $(\mathbf{darg} \ \ell)$  only represents arguments that reach a lazy position (i.e.,  $\ell \in \widehat{\mathcal{D}}$ ).

Figure 5 presents revised analysis rules involving  $\widehat{\mathcal{S}}$  and  $\widehat{\mathcal{F}}$ . These rules use the full analysis result  $(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}})$ . Here,  $\widehat{\mathcal{S}}$  represents arguments that reach a strict position so the new  $[force]$  rule dictates that if an  $(\mathbf{arg} \ \ell)$  is the argument of a **force**, then  $\ell$  must be in  $\widehat{\mathcal{S}}$  (box 5). However, a **force** is not the only expression that requires the value of a promise. There are several other contexts where a **delay** should not appear and the  $[strict]$  rule deals with these strict contexts  $S$ : the operator in an application, the operands in the primitive operations, and the test in an **if** expression. Expressions involving these strict positions have three additional constraints. The first specifies that if an  $(\mathbf{arg} \ \ell_1)$  appears in any of these positions, then  $\ell_1$  should also be in  $\widehat{\mathcal{S}}$  (box 5). The second and third additional constraints show how  $\widehat{\mathcal{F}}$  is computed. Recall that  $\widehat{\mathcal{F}}$  determines where to insert **forces** in the program. The second  $[strict]$  constraint says that if any **delay** flows to a strict position  $\ell$ , then  $\ell$  is added to  $\widehat{\mathcal{F}}$  (box 6). This indicates that a programmer-inserted **delay** has reached a strict position and should be forced. Finally, the third constraint dictates that if a  $(\mathbf{darg} \ \ell_2)$  value flows to a strict label  $\ell$ , then a pair  $(\ell, \ell_2)$  is required to be in  $\widehat{\mathcal{F}}$  (box 7), indicating that the analysis *may* insert a **delay** at  $\ell_2$ , thus requiring a **force** at  $\ell$ .

### 3.5 The Refactoring Transformation

Figure 6 specifies our refactoring as a function  $\varphi$  that transforms a program  $p$  using analysis result  $(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}})$ . The  $\varphi_e$  function wraps expression  $e^\ell$  with

$$\begin{array}{c}
\boxed{\varphi : Prog \times Env \times DPos \times SPos \times FPos \rightarrow Prog} \\
\varphi[\langle\langle \mathbf{define} \ f(x \dots) = e_1 \dots e \rangle\rangle_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}} = \langle\langle \mathbf{define} \ f(x \dots) = \varphi_e[e_1]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}} \dots \varphi_e[e]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}} \rangle\rangle \\
\\
\boxed{\varphi_e : Exp \times Env \times DPos \times SPos \times FPos \rightarrow Exp} \\
\varphi_e[e^\ell]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}} = \langle\langle \mathbf{delay}^* \ (\varphi_e[e]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}})^\ell \rangle\rangle^{\ell_1}, \quad \text{if } \ell \in \widehat{\mathcal{D}}, \ell \notin \widehat{\mathcal{S}}, \ell_1 \notin \text{dom}(\widehat{\rho}) \quad (\dagger) \\
\varphi_e[e^\ell]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}} = \langle\langle \mathbf{force} \ (\varphi_e[e]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}})^\ell \rangle\rangle^{\ell_1}, \quad \text{if } \ell \in \widehat{\mathcal{F}}, \ell_1 \notin \text{dom}(\widehat{\rho}), \quad (\ddagger) \\
\text{or } \exists \ell_2. (\ell, \ell_2) \in \widehat{\mathcal{F}}, \ell_2 \in \widehat{\mathcal{D}}, \ell_2 \notin \widehat{\mathcal{S}}, \ell_1 \notin \text{dom}(\widehat{\rho}) \\
\dots
\end{array}$$

**Fig. 6.** Transformation function  $\varphi$

$\mathbf{delay}^*$  if  $\ell$  is in  $\widehat{\mathcal{D}}$  and not in  $\widehat{\mathcal{S}}$ . In other words,  $e$  is delayed if it flows to a lazy position but not a strict position. With the following correctness section in mind, we extend the set of expressions with  $\mathbf{delay}^*$ , which is exactly like  $\mathbf{delay}$  and merely distinguishes programmer-inserted  $\mathbf{delays}$  from those inserted by the our transformation. The new  $\mathbf{delay}^*$  expression is given a fresh label  $\ell_1$ . In two cases,  $\varphi_e$  inserts a  $\mathbf{force}$  around an expression  $e^\ell$ . First, if  $\ell$  is in  $\widehat{\mathcal{F}}$ , it means  $\ell$  is a strict position and a programmer-inserted  $\mathbf{delay}$  reaches this strict position and must be forced. Second, an expression  $e^\ell$  is also wrapped with  $\mathbf{force}$  if there is some  $\ell_2$  such that  $(\ell, \ell_2)$  is in  $\widehat{\mathcal{F}}$  and the analysis says to delay the expression at  $\ell_2$ , i.e.,  $\ell_2 \in \widehat{\mathcal{D}}$  and  $\ell_2 \notin \widehat{\mathcal{S}}$ . This ensures that transformation-inserted  $\mathbf{delay}^*$ s are also properly forced. All remaining clauses in the definition of  $\varphi_e$ , represented with ellipses, traverse the structure of  $e$  in a homomorphic manner.

## 4 Correctness

Our refactoring for laziness is not semantics-preserving. For example, non-terminating programs may be transformed into terminating ones or exceptions may be delayed indefinitely. Nevertheless, we can prove our analysis sound and the  $\varphi$  transformation safe, meaning that unforced promises cannot cause exceptions.

### 4.1 Language Semantics

To establish soundness, we use Flanagan and Felleisen's [9] technique, which relies on a reduction semantics. The semantics is based on evaluation contexts, which are expressions with a hole in place of one subexpression:

$$\begin{array}{l}
E \in Ctx = [] \mid v \dots E e \dots \mid o E e \mid o v E \mid \mathbf{let} \ x = E \ \mathbf{in} \ e \mid \mathbf{if} \ E \ e \ e \mid \mathbf{zero?} \ E \\
\mid \mathbf{not} \ E \mid \mathbf{null?} \ E \mid \mathbf{force} \ E \mid \mathbf{cons} \ E \ e \mid \mathbf{cons} \ v \ E \mid \mathbf{lcons} \ E \ e \mid \mathbf{first} \ E \mid \mathbf{rest} \ E
\end{array}$$

A reduction step  $\mapsto$  is defined as follows, where  $\rightarrow$  is specified in figure 7:

$$E[e] \mapsto E[e'] \quad \text{iff} \quad e \rightarrow e'$$

A conventional  $\delta$  function evaluates primitives and is elided. We again assume that subexpressions are uniquely labeled but since labels do not affect evaluation, they are implicit in the reduction rules, though we do mention them explicitly in the theorems. Since our analysis does not distinguish memoizing promises from non-memoizing ones, neither does our semantics. To evaluate complete programs, we parameterize  $\mapsto$  over definitions  $d \dots$ , and add a look-up rule:

$$E[f] \mapsto_{d \dots} E[\lambda(x \dots).e], \quad \text{if } (\mathbf{define } f(x \dots) = e) \in d \dots$$

Thus, the result of evaluating a program  $p = d \dots e$  is the result of reducing  $e$  with  $\mapsto_{d \dots}$ . We often drop the  $d \dots$  subscript to reduce clutter.

## Exceptions

Our  $\rightarrow$  reduction thus far is partial, as is the (elided)  $\delta$  function. If certain expressions show up in the hole of the evaluation context, e.g., **first null** or division by 0, we consider the evaluation stuck. To handle stuck expressions, we add an exception **exn** to our semantics. We assume that  $\delta$  returns **exn** for invalid operands of primitives and we extend  $\rightarrow$  with the exception-producing reductions in figure 8.

The (apx) rule says that application of non- $\lambda$ s results in an exception. The (fstx) and (rstx) rules state that reducing **first** or **rest** with anything but a non-empty list is an exception as well. The (strictx) and (strictx\*) reductions partially override some reductions from figure 7 and specify that an exception occurs when an unforced promise appears in a context where the value of that promise is required. These contexts are exactly the strict contexts  $S$  from figure 5. We introduce **dexn** and **dexn\*** to indicate when a **delay** or **delay\*** causes an exception; otherwise these tokens behave just like **exn**. We also extend  $\mapsto$ :

$$E[\mathbf{exn}] \mapsto \mathbf{exn}$$

A conventional well-definedness theorem summarizes the language's semantics.

$(\lambda(x \dots).e) v \dots \rightarrow e\{x := v, \dots\}$	(ap)	$\mathbf{null? null} \rightarrow \mathbf{true}$	(nuln)
$o v_1 v_2 \rightarrow \delta o v_1 v_2$	(op)	$\mathbf{null? } v \rightarrow \mathbf{false}, v \neq \mathbf{null}$	(nul)
$\mathbf{let } x = v \mathbf{ in } e \rightarrow e\{x := v\}$	(let)	$\mathbf{first} (\mathbf{cons } v_1 v_2) \rightarrow v_1$	(fstc)
$\mathbf{if false } e_1 e_2 \rightarrow e_2$	(iff)	$\mathbf{first} (\mathbf{lcons } v e) \rightarrow v$	(fstlc)
$\mathbf{if } v e_1 e_2 \rightarrow e_1, v \neq \mathbf{false}$	(if)	$\mathbf{rest} (\mathbf{cons } v_1 v_2) \rightarrow v_2$	(rstc)
$\mathbf{zero? } 0 \rightarrow \mathbf{true}$	(z0)	$\mathbf{rest} (\mathbf{lcons } v e) \rightarrow \mathbf{delay } e$	(rstlc)
$\mathbf{zero? } v \rightarrow \mathbf{false}, v \neq 0$	(z)	$\mathbf{force} (\mathbf{delay } e) \rightarrow \mathbf{force } e$	(ford)
$\mathbf{not false} \rightarrow \mathbf{true}$	(notf)	$\mathbf{force } v \rightarrow v, v \neq \mathbf{delay } e$	(forv)
$\mathbf{not } v \rightarrow \mathbf{false}, v \neq \mathbf{false}$	(not)		

**Fig. 7.** Call-by-value reduction semantics

$v v_1 \dots \rightarrow \text{exn}$ , if $v \neq \lambda(x \dots).e$	(apx)	$S[\text{delay } e] \rightarrow \text{dexn}$	(strictx)
<b>first</b> $v \rightarrow \text{exn}$ , if $v \notin \text{cons}$ or $\text{lcons}$	(fstx)	$S[\text{delay}^* e] \rightarrow \text{dexn}^*$	(strictx*)
<b>rest</b> $v \rightarrow \text{exn}$ , if $v \notin \text{cons}$ or $\text{lcons}$	(rstx)		

**Fig. 8.** Exception producing reductions

**Theorem 1 (Well-Definedness).** *A program  $p$  either reduces to a value  $v$ ; starts an infinitely long chain of reductions; or reduces to  $\text{exn}$ .*

## 4.2 Soundness of the Analysis

Before stating the soundness theorem, we first extend our analysis for exceptions:

$$(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}}) \models_e \text{exn}^\ell \quad [\text{exn}]$$

Lemma 1 states that  $\mapsto$  preserves  $\models_e$ . We use notation  $\widehat{\rho} \models_e e$  when we are not interested in  $\widehat{\mathcal{D}}$ ,  $\widehat{\mathcal{S}}$ , and  $\widehat{\mathcal{F}}$ , which are only used for transformation. This means  $\widehat{\rho}$  satisfies only the constraints from sections 3.2 and 3.3.

**Lemma 1 (Preservation).** *If  $\widehat{\rho} \models_e e$  and  $e \mapsto e'$ , then  $\widehat{\rho} \models_e e'$ .*

We now state our soundness theorem, where  $\mapsto$  is the reflexive-transitive closure of  $\mapsto$ . The theorem says that if an expression in a program reduces to an  $\ell$ -labeled value, then any acceptable analysis result  $\widehat{\rho}$  correctly predicts that value.

**Theorem 2 (Soundness).** *For all  $\widehat{\rho} \models p$ ,  $p = d \dots e$ , if  $e \mapsto_{d \dots} E[v^\ell]$ ,  $\widehat{v} \in \widehat{\rho}(\ell)$ .*

## 4.3 Safety of Refactoring

We show that refactoring for laziness cannot raise an exception due to a **delay** or **delay**<sup>\*</sup> reaching a strict position. To start, we define a function  $\xi$  that derives a satisfactory abstract environment for a  $\varphi$ -transformed program:

$$\xi[\widehat{\rho}]_p = \widehat{\rho}', \text{ where}$$

$$\boxed{\xi : Env \times Prog \rightarrow Env}$$

$$\forall \ell, x \in \text{dom}(\widehat{\rho}) : \widehat{\rho}'(\ell) = \widehat{\rho}(\ell) \cup \{(\text{delay}^* \ell_1) \mid (\text{darg } \ell_1) \in \widehat{\rho}(\ell), (\text{delay}^* e_1^{\ell_1}) \in p\} \quad (1)$$

$$\widehat{\rho}'(x) = \widehat{\rho}(x) \cup \{(\text{delay}^* \ell_1) \mid (\text{darg } \ell_1) \in \widehat{\rho}(x), (\text{delay}^* e_1^{\ell_1}) \in p\}$$

$$\forall (\text{delay}^* e_1^{\ell_1})^\ell \in p, \ell \notin \text{dom}(\widehat{\rho}) : \quad (2)$$

$$\widehat{\rho}'(\ell) = \widehat{\rho}(\ell_1) \cup \{(\text{delay}^* \ell_1)\} \cup \{(\text{delay}^* \ell_2) \mid (\text{darg } \ell_2) \in \widehat{\rho}(\ell_1), (\text{delay}^* e_2^{\ell_2}) \in p\}$$

$$\forall (\text{force } e_1^{\ell_1})^\ell \in p, \ell \notin \text{dom}(\widehat{\rho}) : \widehat{\rho}'(\ell) = \{\widehat{v} \mid \widehat{v} \in \widehat{\rho}(\ell_1), \widehat{v} \notin \text{delay}\} \quad (3)$$

The  $\xi$  function takes environment  $\widehat{\rho}$  and a program  $p$  and returns a new environment  $\widehat{\rho}'$ . Part 1 of the definition copies  $\widehat{\rho}$  entries to  $\widehat{\rho}'$ , except **darg** values are replaced with **delay**<sup>\*</sup>s when there is a corresponding **delay**<sup>\*</sup> in  $p$ . Parts 2 and 3 add new  $\widehat{\rho}'$  entries for **delay**<sup>\*</sup>s and **forces** not accounted for in  $\widehat{\rho}$ . When the given  $p$  is a  $\varphi$ -transformed program, then the resulting  $\widehat{\rho}'$  satisfies that program.

**Lemma 2.** *If  $(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}}) \models p$ , then  $\xi[[\widehat{\rho}]]_{\varphi[[p]]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}}} \models \varphi[[p]]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}}$ .*

Finally, theorem 3 states the safety property. It says that evaluating a transformed program cannot generate an exception due to `delays` or `delay*`s.

**Theorem 3 (Safety).** *For all  $p$  and  $(\widehat{\rho}, \widehat{\mathcal{D}}, \widehat{\mathcal{S}}, \widehat{\mathcal{F}}) \models p$ , if  $\varphi[[p]]_{\widehat{\rho}\widehat{\mathcal{D}}\widehat{\mathcal{S}}\widehat{\mathcal{F}}} = d \dots e$ , then  $e \not\mapsto_{d \dots} \text{dexc}$ , and  $e \not\mapsto_{d \dots} \text{dexc}^*$ .*

*Proof.* (Sketch) Using Soundness, the analysis rules in figure 5, and Lemma 2.

#### 4.4 Idempotency

Our transformation is not idempotent. Indeed, it may be necessary to refactor a program multiple times to get the “right” amount of laziness. For example:

```
let x = ⟨long computation⟩ in let y = ⟨short computation involving x⟩
in (delay y)
```

The long computation should be delayed but applying our transformation once only delays the short computation. To delay the long computation, a second transformation round is required. In practice, we have observed that one round of laziness refactoring suffices to handle the majority of cases. However, section 6 presents a real-world example requiring multiple transformations so our tool currently allows the programmer to decide how often to apply the refactoring.

## 5 A Prototype Implementation

We have implemented refactoring for laziness as a tool for Racket [10], in the form of a plugin for the DrRacket IDE. It uses laziness analysis to automatically insert `delay` and `force` expressions as needed, with graphical justification.

### 5.1 Constraint Solving Algorithm

Computing our laziness analysis requires two stages: (1) generate a set of constraints from a program, and (2) solve for the least solution using a conventional worklist algorithm [18]. The graph nodes are the variables and labels in the program, plus one node each for  $\widehat{\mathcal{D}}$ ,  $\widehat{\mathcal{S}}$ , and  $\widehat{\mathcal{F}}$ . Without loss of generality, we use only labels for the nodes and  $\widehat{\rho}$  for the analysis result in our description of the algorithm. There exists an edge from node  $\ell_1$  to  $\ell_2$  if there is a constraint where  $\widehat{\rho}(\ell_2)$  depends on  $\widehat{\rho}(\ell_1)$ ; the edge is labeled with that constraint. Thus one can view a node  $\ell$  as the endpoint for a series of data flow paths. To compute  $\widehat{\rho}(\ell)$ , it suffices to traverse all paths from the leaves to  $\ell$ , accumulating values according to the constraints along the way.

The analysis result is incrementally computed in a breadth-first fashion by processing constraints according a worklist of nodes. Processing a constraint

entails adding values to  $\widehat{\rho}$  so the constraint is satisfied. The algorithm starts by processing all constraints where a node depends on a value, e.g., `val`  $\in$   $\widehat{\rho}(\ell)$ ; the nodes on the right-hand side of these constraints constitute the initial worklist. Nodes are then removed from the worklist, one at a time. When a node is removed, the constraints on the out-edges of that node are processed and a neighbor  $\ell$  of the node is added to the worklist if  $\widehat{\rho}(\ell)$  was updated while processing a constraint. A node may appear in the worklist more than once, but only a finite number of times, as shown by the following termination argument.

### Termination and Complexity of Constraint Solving

Inspecting the constraints from section 3 reveals that an expression requires recursive calls only for subexpressions. Thus, a finite program generates a finite number of constraints. For a finite program with finitely many labels and variables, the set of possible abstract values is also finite. Thus, a node can only appear in the worklist a finite number of times, so algorithm must terminate.

We observe in the constraint-solving algorithm that, (1) a node  $\ell$  is added to the worklist only if  $\widehat{\rho}(\ell)$  is updated due to a node on which it depends being in the worklist, and (2) values are only ever added to  $\widehat{\rho}$ ; they are never removed. For a program of size  $n$ , there are  $O(n)$  nodes in the dependency graph. Each node can appear in the worklist  $O(n)$  times, and a data flow path to reach that node could have  $O(n)$  nodes, so it can take  $O(n^2)$  node visits to compute the solution at a particular node. Multiplying by  $O(n)$  total nodes, means the algorithm may have to visit  $O(n^3)$  nodes to compute the solution for all nodes.

## 5.2 Laziness Refactoring Tool

Our prototype tool uses the result of the analysis and the  $\varphi$  function from section 3.5 to insert additional `delays` and `forces`. In contrast to the mathematical version of  $\varphi$ , its implementation avoids inserting `delays` and `forces` around values and does not insert duplicate `delays` or `forces`.

We evaluated a number of examples with our tool including the  $n$ -queens problem from section 2. Figure 9 (top) shows the program in Racket, including timing information and a graphical depiction of the answer. Despite the use of `lcons`,<sup>8</sup> the program takes as long as an eager version of the same program (not shown) to compute an answer. Figure 9 (bot) shows the program after our tool applies the laziness transformation. When the tool is activated, it: (1) computes an analysis result for the program, (2) uses the result to insert `delays` and `forces`, highlighting the added `delays` in yellow and the added `forces` in blue, and (3) adds arrows originating from each inserted `delay`, pointing to the source of the laziness, thus explaining its decision to the programmer in an intuitive manner. Running the transformed program exhibits the desired performance.

---

<sup>8</sup> Though `lcons` is not available in Racket, to match the syntax of our paper, we simulate it with a macro that wraps a `delay` around the second argument of a `cons`.

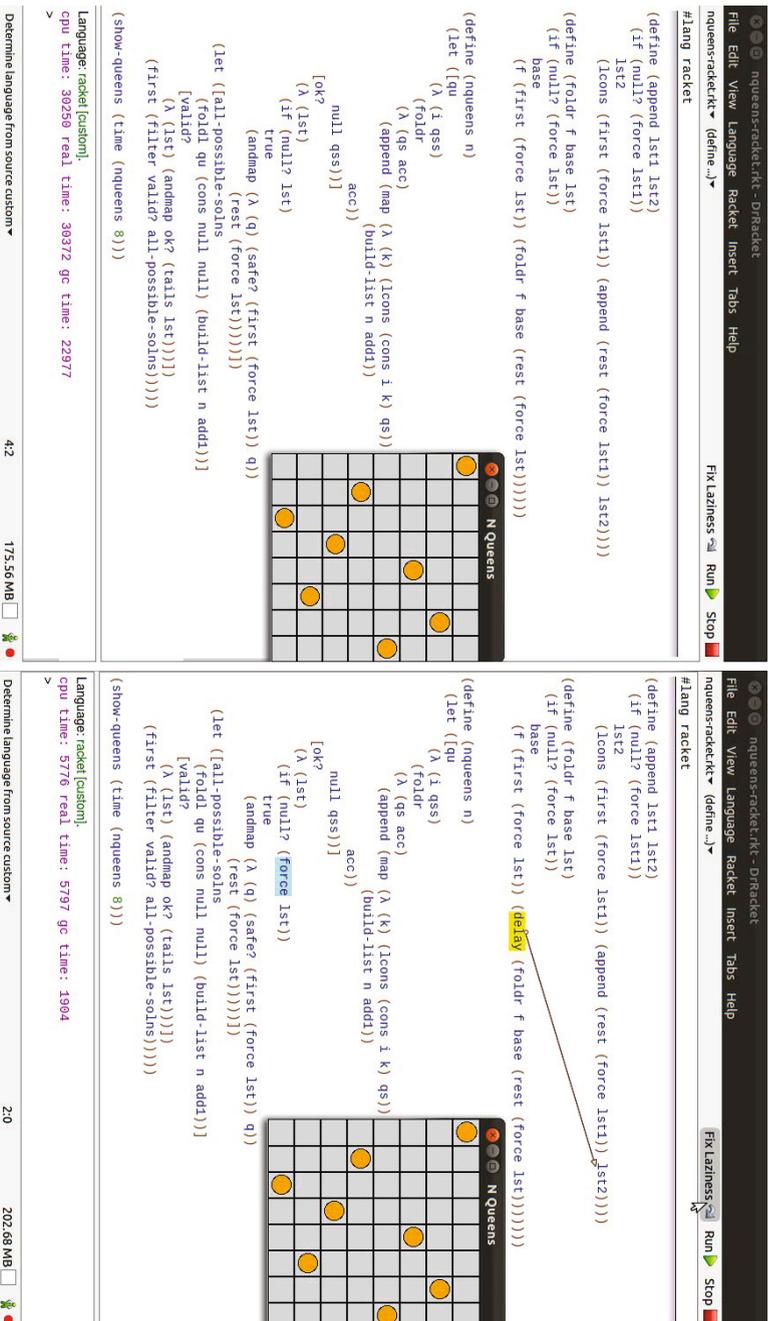


Fig. 9. Evaluating *n*-queens in Racket: only lazy cons (top), after refactoring (bot)

## 6 Laziness in the Large

To further evaluate our idea and our tool, we examined the Racket code base and some user-contributed packages for manual uses of laziness. We found several erroneous attempts at adding laziness and we verified that our tool would have prevented many such errors.<sup>9</sup> We consider this investigation a first confirmation of the usefulness of our tool. The rest of the section describes two of the examples.

The DMdA languages [5] allow students to write contracts for some data structures. These contracts are based on Findler et al.’s lazy contracts [8]. The contracts are primarily implemented via a constructor with a few lazy fields. Additionally, several specialized contract constructors for various data structures call the main constructor. However, since the specialized constructors are implemented with ordinary strict functions, to preserve the intended lazy behavior, the programmer must manually propagate the laziness to the appropriate arguments of these functions, similar to the Scala example from section 2. Thus, a small amount of laziness in the main contract constructor requires several more `delays` scattered all throughout the program. Adding these `delays` becomes tedious as the program grows in complexity and unsurprisingly, a few were left out. Our tool identified the missing `delays`, which the author of the code has confirmed and corrected with commits to the code repository.

A second example concerns queues and dequeues [21] based on implicit recursive slowdown [20, Chapter 11], where laziness enables fast amortized operations and simplifies the implementation. The library contained several performance bugs, as illustrated by this code snippet from a deque enqueue function:

```
define enqueue(elem dq) = ...
  let strictprt = ⟨extract strict part of dq⟩
      newstrictprt = ⟨combine elem and strictprt⟩
      lazyprt = force ⟨extract lazy part of dq⟩
      lazyprt1 = ⟨extracted from lazyprt⟩
      lazyprt2 = ⟨extracted from lazyprt⟩
  in Deque newstrictprt (delay ⟨combine lazyprt1 and lazyprt2⟩)
```

The function enqueues `elem` in deque `dq`, which has a lazy part and a strict part. In one execution path, the lazy part is extracted, forced, and separated into two additional pieces. Clearly, the forcing is unnecessary because neither of the pieces are used before they are inserted back into the new deque. Worse, the extra forcing slows the program significantly. For this example, activating our tool *twice* fixes the performance bug. For a reasonably standard benchmark, the fix reduced the running time by an order of magnitude. The authors of the code have acknowledged the bug and have merged our fix into the code repository.

## 7 Related Work

The idea of combining strict and lazy evaluation is old, but most works involve removing laziness from lazy languages. We approach strict-lazy programming

---

<sup>9</sup> The examples were first translated to work with the syntax in this paper.

from the other, relatively unexplored, end of the spectrum, starting with a strict language and then only adding laziness as needed. This seems worthwhile since empirical studies indicate that most promises in a lazy language are unneeded [6,15,16,23]. Starting with a strict language also alleviates many disadvantages of lazy evaluation such as difficulty reasoning about space/time consumption.

The most well-known related work is strictness analysis [4,17], which calculates when to eagerly evaluate arguments without introducing non-termination. With our work, calculating divergence properties is not sufficient since even terminating programs may require additional laziness, as seen in examples from this paper. Hence we take a different, flow-analysis-based approach.<sup>10</sup> Researchers have also explored other static [7] and dynamic [2,6,15] laziness-removal techniques. However, these efforts all strive to preserve the program’s semantics. We focus on the problem of strict programmers trying to use laziness, but doing so *incorrectly*. Thus our transformation necessarily allows the semantics of the program to change (i.e., from non-terminating to terminating), but hopefully in a way that the programmer intended in the first place.

Sheard [25] shares our vision of a strict language that is also practical for programming lazily. While his language does not require explicit `forces`, the programmer must manually insert all required `delay` annotations.

## 8 Future Work

This paper demonstrates the theoretical and practical feasibility of a novel approach to assist programmers with the introduction of laziness into a strict context. We see several directions for future work. The first is developing a modular analysis. Our transformation requires the whole program and is thus unsatisfactory in the presence of libraries. Also, we intend to develop a typed version of our transformation and tool, so typed strict languages can more easily benefit from laziness as well. We conjecture that expressing strictness information via types may also provide a way to enable a modular laziness-by-need analysis.

**Acknowledgements.** Partial support provided by NSF grant CRI-0855140. Thanks to Matthias Felleisen, Eli Barzilay, David Van Horn, and J. Ian Johnson for feedback on earlier drafts.

## References

1. Abelson, H., Sussman, G.J., Sussman, J.: Structure and Interpretation of Computer Programs. MIT Press (1984)
2. Aditya, S., Arvind, Augustsson, L., Maessen, J.W., Nikhil, R.S.: Semantics of pH: A parallel dialect of Haskell. In: Proc. Haskell Workshop, pp. 34–49 (1995)

---

<sup>10</sup> Interestingly, we conjecture that our approach would be useful to lazy programmers trying to insert strictness *annotations*, such as Haskell’s `seq`, to their programs.

3. Appel, A., Blume, M., Gansner, E., George, L., Huelsbergen, L., MacQueen, D., Reppy, J., Shao, Z.: Standard ML of New Jersey User's Guide (1997)
4. Burn, G.L., Hankin, C.L., Abramsky, S.: Strictness analysis for higher-order functions. *Sci. Comput. Program.* 7, 249–278 (1986)
5. Crestani, M., Sperber, M.: Experience report: growing programming languages for beginning students. In: *Proc. 15th ICFP*, pp. 229–234 (2010)
6. Ennals, R., Peyton Jones, S.: Optimistic evaluation: an adaptive evaluation strategy for non-strict programs. In: *Proc. 8th ICFP*, pp. 287–298 (2003)
7. Faxén, K.F.: Cheap eagerness: speculative evaluation in a lazy functional language. In: *Proc. 5th ICFP*, pp. 150–161 (2000)
8. Findler, R.B., Guo, S.-Y., Rogers, A.: Lazy Contract Checking for Immutable Data Structures. In: Chitil, O., Horváth, Z., Zsók, V. (eds.) *IFL 2007*. LNCS, vol. 5083, pp. 111–128. Springer, Heidelberg (2008)
9. Flanagan, C., Felleisen, M.: Modular and polymorphic set-based analysis: Theory and practice. Tech. Rep. TR96-266, Rice Univ. (1996)
10. Flatt, M., PLT: Reference: Racket. Tech. Rep. PLT-TR-2012-1, PLT Inc. (2012), <http://racket-lang.org/tr1/>
11. Friedman, D., Wise, D.: Cons should not evaluate its arguments. In: *Proc. 3rd ICALP*, pp. 257–281 (1976)
12. Hughes, J.: Why functional programming matters. *Comput. J.* 32, 98–107 (1989)
13. Jones, N.D.: Flow analysis of lambda expressions. Tech. rep., Aarhus Univ. (1981)
14. Leroy, X., Doligez, D., Frisch, A., Garrigue, J., Rémy, D., Vouillon, J.: The OCaml system, release 3.12, Documentation and user's manual. INRIA (July 2011)
15. Maessen, J.W.: Eager Haskell: resource-bounded execution yields efficient iteration. In: *Proc. Haskell Workshop*, pp. 38–50 (2002)
16. Morandat, F., Hill, B., Osvald, L., Vitek, J.: Evaluating the Design of the R Language. In: Noble, J. (ed.) *ECOOP 2012*. LNCS, vol. 7313, pp. 104–131. Springer, Heidelberg (2012)
17. Mycroft, A.: Abstract interpretation and optimising transformations for applicative programs. Ph.D. thesis, Univ. Edinburgh (1981)
18. Nielson, F., Nielson, H.R., Hankin, C.: *Principles of Program Analysis*. Springer (2005)
19. Odersky, M.: The Scala Language Specification, Version 2.9. EPFL (May 2011)
20. Okasaki, C.: *Purely Functional Data Structures*. Cambridge Univ. Press (1998)
21. Hari Prashanth, K.R., Tobin-Hochstadt, S.: Functional data structures for Typed Racket. In: *Proc. Scheme Workshop* (2010)
22. Rees, J., Clinger, W. (eds.): Revised<sup>3</sup> Report on the Algorithmic Language Scheme. *ACM SIGPLAN Notices* (December 1986)
23. Schauser, K.E., Goldstein, S.C.: How much non-strictness do lenient programs require? In: *Proc. 7th FPCA* (1995)
24. Sestoft, P.: Replacing function parameters by global variables. Master's thesis, Univ. Copenhagen (1988)
25. Sheard, T.: A pure language with default strict evaluation order and explicit laziness. In: *2003 Haskell Workshop: New Ideas Session* (2003)
26. Shivers, O.: Control-flow analysis in scheme. In: *Proc. PLDI*, pp. 164–174 (1988)
27. Wadler, P., Taha, W., MacQueen, D.: How to add laziness to a strict language, without even being odd. In: *Proc. Standard ML Workshop* (1998)