

Sweet Hygienic Macros for Full ES5 JavaScript

Tim Disney
UC Santa Cruz

Nathan Faubion
Flow Corp.

David Herman
Mozilla Corp.

Cormac Flanagan
UC Santa Cruz

Abstract

Lisp and Scheme have demonstrated the power of macros to enable programmers to evolve and craft languages. In languages with more complex syntax, macros have had less success. In part, this has been due to the difficulty in building expressive hygienic macro systems for such languages. JavaScript in particular presents unique challenges for macro systems due to ambiguities in the lexing stage that force the JavaScript lexer and parser to be intertwined.

In this paper we present a novel solution to the lexing ambiguity of JavaScript that enables us to cleanly separate the JavaScript lexer and parser. This separated lexer and parser allows us to implement an expressive hygienic macro system for JavaScript called `sweet.js`.

1. Introduction

Expressive macros systems have a long history in the design of extensible programming languages going back to Lisp and Scheme [11, 18] as a powerful tool that enables programmers to craft their own languages.

While macro systems have found success in Lisp-derived languages, they have not been widely adopted in languages such as JavaScript. In part, this failure is due to the difficulty in implementing macro systems in languages that are not fully delimited. A key feature of a sufficiently expressive macro system is the ability for macros to manipulate unparsed and unexpanded subexpressions. In a language with parentheses like Scheme, manipulating unparsed subexpressions is simple:

```
(if (> denom 0)
    (/ x denom)
    (error "divide by zero"))
```

The Scheme *reader* converts the source string into nested s-expressions, which macros can easily manipulate. Since each subexpression of the `if` form is a fully delimited s-expression, it is easy to implement `if` as a macro.

Conceptually, the Scheme compiler lexes a source string into a stream of tokens which are then read into s-expressions before being macro expanded and parsed into an abstract syntax tree.

$$\text{lexer} \xrightarrow{\text{Token}^*} \text{reader} \xrightarrow{\text{Sexpr}} \text{expander/parser} \xrightarrow{\text{AST}}$$

As a first step to designing a Scheme-like macro system for JavaScript, it is necessary to introduce a read step into the compiler pipeline. However, the design of a correct reader for full JavaScript turns out to be surprisingly subtle, due to ambiguities in how regular expression literals (such as `/[0-9]*/`) and the divide operator (`/`) should be lexed. In traditional JavaScript compilers, the parser and lexer are intertwined. Rather than run the entire program through the lexer once to get a sequence of tokens, the parser calls out to the lexer from a given grammatical context with a flag to indicate if the lexer should accept a regular expression or a divide operator, and the input character stream is tokenized accordingly. So if the parser is in a context accepts a regular expression, the characters `"/x/"` will be lexed into the single token `/x/` otherwise it will lex into the individual tokens `/`, `x`, and `/`.

$$\text{lexer} \xrightleftharpoons[\text{Token}^*]{\text{feedback}} \text{parser} \xrightarrow{\text{AST}}$$

It is necessary to separate the parser and lexer in order to implement a macro system for JavaScript. Our JavaScript macro system, `sweet.js`, includes a separate reader that converts a sequence of tokens into a sequence of token trees (a little analogous to s-expressions) without feedback from the parser.

$$\text{lexer} \xrightarrow{\text{Token}^*} \text{reader} \xrightarrow{\text{TokenTree}^*} \text{expander/parser} \xrightarrow{\text{AST}}$$

Figure 1: The sweet.js editor



This enables us to finally separate the JavaScript lexer and parser and build a fully hygienic macro system for JavaScript. The reader records sufficient history information in order to correctly decide whether to parse a sequence of tokens `/x/g` as a regular expression or as division operators (as in `4.0/x/g`). Surprisingly, this history information needs to be remembered from arbitrarily far back in the token stream.

Once JavaScript source has been correctly read, there are still a number of challenges to building an expressive macro system. The lack of parentheses in particular make writing declarative macro definitions difficult. For example, the `if` statement in JavaScript allows undelimited then and else branches:

```
if (denom > 0)
  x / denom;
else
  throw "divide by zero";
```

It is necessary to know where the then and else branches end to correctly implement an `if` macro but this is complicated by the lack of delimiters.

The solution to this problem that we take is by progressively building a partial AST during macro expansion. Macros can then match against and manipulate this partial AST. For example, an `if` macro could indicate that the then

and else branches must be single statements and then manipulate them appropriately.

This approach, called *enforestation*, was pioneered by Honu [23, 24], which we adapt here for JavaScript¹. In addition, we make two extensions to the Honu technique that enable more expressive macros to be built. First, as described in Section 4.1 we add support for infix macros, which allow macros to match syntax both before and after the macro identifier. Secondly, we implement the `invoke` pattern class, described in Section 4.2, which allows macro authors to extend the patterns used to match syntax.

Sweet.js is implemented in JavaScript and takes source code written with sweet.js macros and outputs the expanded source that can be run in any JavaScript environment. The project web page² includes an interactive browser-based editor that makes it simple to try out writing macros without requiring any installation. Figure 1 shows the editor in action; a macro implementing classes is being edited in the left pane and the right pane shows the expanded output. There is already an active community using sweet.js to implement sig-

¹ Honu does not support regular expression literals which simplifies their reader.

² <http://sweetjs.org>

Figure 2: AST for Simplified JavaScript

$$\begin{aligned}
 e \in AST &::= x \mid /r/ \mid \{x: e\} \mid (e) \mid e.x \mid e(e) \\
 &\mid e / e \mid e + e \mid e = e \mid \{e\} \mid x:e \mid \text{if } (e) \ e \\
 &\mid \text{return} \mid \text{return } e \\
 &\mid \text{function } x (x) \{e\} \mid e \ e
 \end{aligned}$$

nificant features from the upcoming ES6 version of JavaScript [19] and pattern matching for JavaScript [8].

This paper makes the following contributions:

- a reader implementation that addresses the complexities of JavaScript, including regular expression literals, together with a correctness proof for this reader.
- a macro system with the ability to define infix macros that match on arbitrary surrounding syntax.
- the `invoke` primitive to allow custom parser classes and more powerful matching

2. Reading JavaScript

Parsers give structure to unstructured source code. In parsers without macro systems this is usually accomplished by a lexer (which converts a character stream to a token stream) and a parser (which converts the token stream into an AST according to a context-free grammar). A system with macros must implement a macro *expander* that sits between the lexer and parser. Some macros systems, such as the C preprocessor [15], work over just the token stream. However, to implement truly expressive Scheme-like macros that can manipulate groups of unparsed tokens, it is necessary to structure the token stream via a *reader*, which performs delimiter matching and enables macros to manipulate delimiter grouped tokens.

As mentioned in the introduction, the design of a correct reader for JavaScript is surprisingly subtle due to ambiguities between lexing regular expression literals and the divide operator. This disambiguation is critical to the correct implementation of `read` because delimiters can appear inside of a regular expression literal. If the reader failed to distinguish between a regular expression/divide operator, it could incorrectly match delimiters that appear inside of a regular expression.

```
function makeRegex() {
  return /{/;
}
```

A key novelty in `sweet.js` is the design and implementation of a reader that correctly distinguishes between regular expression literals and the divide operator for full ES5 JavaScript³. For clarity of presentation, this paper describes

³ Our implementation also has initial support for the upcoming ES6 version of JavaScript.

the implementation of `read` for the subset of JavaScript shown in Figure 4, which retains the essential complications of a the full version of `read`.

In our formalism in Figure 3, `read` takes a *Token* sequence. Tokens are the output of a very simple lexer, which we do not define here, that performs the standard grouping of keywords, punctators, the (unmatched) delimiters, and variable identifiers, but does not distinguish between regular expressions and divide (i.e. there is just the ambiguous / token).

$$\begin{aligned}
 Punctuator &::= / \mid + \mid : \mid ; \mid = \mid . \\
 Keyword &::= \text{return} \mid \text{function} \mid \text{if} \\
 Token &::= x \mid Punctuator \mid Keyword \\
 &\mid \{ \mid \} \mid (\mid) \\
 x, y &\in Variable \\
 s &\in Token^*
 \end{aligned}$$

The job of `read` is then to produces a correct *TokenTree* sequence. Token trees include regular expression literals `/r/`, where `r` is the regular expression body. We simplify regular expression bodies to just a variable and the individual delimiters, which captures the essential problems of parsing regular expression literals. Token trees also include fully matched delimiters with nested token tree sequences $\underline{(t)}$ and $\underline{\{t\}}$ rather than individual delimiters (note that we write token tree delimiters with an underline to distinguish them from token delimiters).

$$\begin{aligned}
 k \in TokenTree &::= x \mid Punctuator \mid Keyword \\
 &\mid /r/ \mid \underline{(t)} \mid \underline{\{t\}} \\
 r \in RegexPat &::= x \mid \{ \mid \} \mid (\mid) \\
 t, p &\in TokenTree^*
 \end{aligned}$$

Each token and token tree also carries their line number from the original source string. Line numbers are needed because there are edge cases in the JavaScript grammar where newlines influence parsing. For example, the following function returns the object literal `{x: y}` as expected.

```
function f(y) {
  return { x: y }
}
```

However, adding a newline causes this function to return `undefined`, because the grammar calls for an implicit semi-colon to be inserted after the `return` keyword.

```
function g(y) {
  return
  { x: y }
}
```

For clarity of presentation, we leave token line numbers implicit unless we require them, in which case we use the notation $\{^l$ where l is a line number.

We write a sequence by separating elements with a dot so the source string `"foo(/)/)"` is lexed into a sequence of six tokens `foo · (· / ·) · / ·)`. The equivalent token tree sequence is `foo · (/)/`.

Figure 3: Simplified Read Algorithm

<i>Punctuator</i>	$::=$	$/ \mid + \mid : \mid ; \mid = \mid .$		
<i>Keyword</i>	$::=$	<code>return</code> \mid <code>function</code> \mid <code>if</code>		
<i>Token</i>	$::=$	$x \mid \textit{Punctuator} \mid \textit{Keyword}$		
		$\mid \{ \mid \} \mid (\mid)$		
$k \in \textit{TokenTree}$	$::=$	$x \mid \textit{Punctuator} \mid \textit{Keyword}$		
		$\mid /r/ \mid (t) \mid \{t\}$		
$r \in \textit{RegexPat}$	$::=$	$x \mid \{ \mid \} \mid (\mid)$		
x	\in	<i>Variable</i>		
s	\in	\textit{Token}^*		
t, p	\in	$\textit{TokenTree}^*$		
			$\text{isExprPrefix} : \textit{TokenTree}^* \rightarrow \textit{Bool} \rightarrow \textit{Int} \rightarrow \textit{Bool}$	
			$\text{isExprPrefix}(\epsilon, \text{true}, l)$	$= \text{true}$
			$\text{isExprPrefix}(p \cdot /, b, l)$	$= \text{true}$
			$\text{isExprPrefix}(p \cdot +, b, l)$	$= \text{true}$
			$\text{isExprPrefix}(p \cdot =, b, l)$	$= \text{true}$
			$\text{isExprPrefix}(p \cdot :, b, l)$	$= b$
			$\text{isExprPrefix}(p \cdot \text{return}^l, b, l')$	$= \text{false} \quad \text{if } l \neq l'$
			$\text{isExprPrefix}(p \cdot \text{return}^l, b, l')$	$= \text{true} \quad \text{if } l = l'$
			$\text{isExprPrefix}(p, b, l)$	$= \text{false}$
$\text{read} : \textit{Token}^* \rightarrow \textit{TokenTree}^* \rightarrow \textit{Bool} \rightarrow \textit{TokenTree}^*$				
$\text{read}(/ \cdot r \cdot / \cdot s, \epsilon, b)$	$=$	$/r/ \cdot \text{read}(s, /r/, b)$		
$\text{read}(/ \cdot r \cdot / \cdot s, p \cdot k, b)$	$=$	$/r/ \cdot \text{read}(s, p \cdot k \cdot /r/, b)$		
$\text{if } k \in \textit{Punctuator} \cup \textit{Keyword}$				
$\text{read}(/ \cdot r \cdot / \cdot s, p \cdot \text{if} \cdot (t), b)$	$=$	$/r/ \cdot \text{read}(s, p \cdot \text{if} \cdot (t) \cdot /r/, b)$		
$\text{read}(/ \cdot r \cdot / \cdot s, p \cdot \text{function}^l \cdot x \cdot (t) \cdot \{t'\}, b)$	$=$	$/r/ \cdot \text{read}(s, p \cdot \text{function}^l \cdot x \cdot (t) \cdot \{t'\} \cdot /r/, b)$		
$\text{if } \text{isExprPrefix}(p, b, l) = \text{false}$				
$\text{read}(/ \cdot r \cdot / \cdot s, p \cdot \{t\}^l, b)$	$=$	$/r/ \cdot \text{read}(s, p \cdot \{t\}^l \cdot /r/, b)$		
$\text{if } \text{isExprPrefix}(p, b, l) = \text{false}$				
$\text{read}(/ \cdot s, p \cdot x, b)$	$=$	$/ \cdot \text{read}(s, p \cdot x \cdot /, b)$		
$\text{read}(/ \cdot s, p \cdot /x/, b)$	$=$	$/ \cdot \text{read}(s, p \cdot /x/ \cdot /, b)$		
$\text{read}(/ \cdot s, p \cdot (t), b)$	$=$	$/ \cdot \text{read}(s, p \cdot (t) \cdot /, b)$		
$\text{read}(/ \cdot s, p \cdot \text{function}^l \cdot x \cdot (t) \cdot \{t'\}, b)$	$=$	$/ \cdot \text{read}(s, p \cdot \text{function}^l \cdot x \cdot (t) \cdot \{t'\} \cdot /, b)$		
$\text{if } \text{isExprPrefix}(p, b, l) = \text{true}$				
$\text{read}(/ \cdot s, p \cdot \{t\}^l, b)$	$=$	$/ \cdot \text{read}(s, p \cdot \{t\}^l \cdot /, b)$		
$\text{if } \text{isExprPrefix}(p, b, l) = \text{true}$				
$\text{read}((\cdot s \cdot) \cdot s', p, b)$	$=$	$(t) \cdot \text{read}(s', p \cdot (t), b)$		
$\text{where } s \text{ contains no unmatched })$		$\text{where } t = \text{read}(s, \epsilon, \text{false})$		
$\text{read}(\{^l \cdot s \cdot \} \cdot s', p, b)$	$=$	$\{t\}^l \cdot \text{read}(s', p \cdot \{t\}^l, b)$		
$\text{where } s \text{ contains no unmatched } \}$		$\text{where } t = \text{read}(s, \epsilon, \text{isExprPrefix}(p, b, l))$		
$\text{read}(x \cdot s, p, b)$	$=$	$x \cdot \text{read}(s, p \cdot x, b)$		
$\text{read}(\epsilon, p, b)$	$=$	ϵ		

The key idea of `read` is to maintain of prefix of already read token trees. When the reader comes to a slash and needs to decide if it should read the slash as a divide token or the start of a regular expression literal, it consults the prefix. Looking back at most five tokens trees in the prefix is sufficient to disambiguate the slash token. Not that this may correspond to looking back an unbounded distance in the original token stream.

Some of the cases of `read` are relatively obvious. For example, if the token just read was one of the binary operators (e.g. the `+` in `f + . / . }` `/`) the slash will always be the start of a regular expression literal.

Other cases require additional context to disambiguate. For example, if the previous token tree was a parentheses (e.g. `foo \cdot (\cdot x \cdot) \cdot / \cdot y \cdot`) then slash will be the divide operator, *unless* the token tree before the parentheses was the keyword `if`, in which case it is actually the start of a regular expression (since single statement `if` bodies do not require braces).

```
if (x) /}/ // regex
```

One of the most complicated cases is a slash after curly braces. Part of the complication here is that curly braces can be either an object literal (in which case the slash should be a divide) or it could be a block (in which case the slash should be a regular expression), but even more problematic is that

Figure 4: Simplified ES5 Grammar

$PrimaryExpr_x$	$::= x$
$PrimaryExpr_{/r/}$	$::= / \cdot x \cdot /$
$PrimaryExpr_{\{\{x:e\}\}}$	$::= \{ \cdot x \cdot : \cdot AssignExpr_e \cdot \}$
$PrimaryExpr_{(e)}$	$::= (\cdot AssignExpr_e \cdot)$
$MemberExpr_e$	$::= PrimaryExpr_e$
$MemberExpr_e$	$::= Function_e$
$MemberExpr_{e.x}$	$::= MemberExpr_e \cdot \cdot x$
$CallExpr_e (e')$	$::= MemberExpr_e \cdot (\cdot AssignExpr_{e'} \cdot)$
$CallExpr_e (e')$	$::= CallExpr_e \cdot (\cdot AssignExpr_{e'} \cdot)$
$CallExpr_{e.x}$	$::= CallExpr_e \cdot x$
$BinaryExpr_e$	$::= CallExpr_e$
$BinaryExpr_e / e'$	$::= BinaryExpr_e \cdot / \cdot BinaryExpr_{e'}$
$BinaryExpr_e + e'$	$::= BinaryExpr_e \cdot + \cdot BinaryExpr_{e'}$
$AssignExpr_e$	$::= BinaryExpr_e$
$AssignExpr_{e = e'}$	$::= CallExpr_e \cdot = \cdot AssignExpr_{e'}$
$StmtList_e$	$::= Stmt_e$
$StmtList_{e e'}$	$::= StmtList_e \cdot Stmt_{e'}$
$Stmt_{\{e\}}$	$::= \{ \cdot StmtList_e \cdot \}$
$Stmt_x: e$	$::= x \cdot : \cdot Stmt_e$
$Stmt_e$	$::= AssignExpr_e ; \quad \text{where lookahead} \neq \{ \text{or function} \}$
$Stmt_{if (e) e'}$	$::= if \cdot (\cdot AssignExpr_e \cdot) \cdot Stmt_{e'}$
$Stmt_{return}$	$::= return$
$Stmt_{return e}$	$::= return \cdot [\text{no line terminator here}] AssignExpr_e \cdot ;$
$Function_{function x (x') \{e\}}$	$::= function \cdot x \cdot (\cdot x' \cdot) \cdot \{ \cdot SourceElements_e \cdot \}$
$SourceElement_e$	$::= Stmt_e$
$SourceElement_e$	$::= Function_e$
$SourceElements_e$	$::= SourceElement_e$
$SourceElements_{e e'}$	$::= SourceElements_e \cdot SourceElement_{e'}$
$Program_e$	$::= SourceElements_e$
$Program$	$::= \epsilon$

both object literals and blocks with labeled statements can nest. For example, in the following code snippet the outer curly brace is a block with a labeled statement `x`, which is another block with a labeled statement `y` followed by a regular expression literal.

```
{
  x:{y: z} /}/ // regex
}
```

But if we change the code slightly, the outer curly braces become an object literal and `x` is a property so the inner curly braces are also an object literal and thus the slash is a divide operator.

```
o = {
  x:{y: z} /x/g // divide
}
```

While it is unlikely that a programmer would attempt to intentionally perform division on an object literal, it is not a parse error. In fact, this is not even a runtime error since JavaScript will implicitly convert the object to a number (technically `NaN`) and then perform the division (yielding `NaN`).

The reader handles these cases by checking if the prefix of a curly brace block forces the curly to be an object literal or a statement block and then setting a boolean flag to be used while reading the tokens inside of the braces.

Based on this discussion, our reader is implemented as a function that takes a sequence of tokens, a prefix of previously read token trees, a boolean indicating if the token stream currently being read is inside an object literal, and returns a sequence of token trees.

$$\text{read} : \text{Token}^* \rightarrow \text{TokenTree}^* \rightarrow \text{Bool} \rightarrow \text{TokenTree}^*$$

The implementation of `read` shown in Figure 3 includes an auxiliary function `isExprPrefix` used to determine if the prefix for a curly brace indicates that the braces should be part of an expression (i.e. the braces are an object literal) or if they should be a block statement.

Interestingly, the `isExprPrefix` function must also be used when the prefix before a slash contains a function definition. This is because there are two kinds of function definitions in JavaScript, function expressions and function declarations, and these also affect how slash is read. For example, a slash following a function declaration is always the start of a regular expression:

```
function f() {}
/}/ // regex
```

However, a slash following a function expression is a divide operator:

```
x = function f() { }
/y/g // divide
```

As in the object literal case, it is unlikely that a programmer would attempt to intentionally divide a function expression but it is not an error to do so.

2.1 Proving Read

To show that our `read` algorithm correctly distinguishes divide operations from regular expression literals, we show that a parser defined over normal tokens produces the same AST as a parser defined over the token trees produced from `read`.

The parser for normal tokens is defined in Figure 4, and generates ASTs in the abstract syntax shown in Figure 2. A parser for the nonterminal *Program* is a function from a sequence of tokens to an AST.

$$\text{Program} :: \text{Token}^* \rightarrow \text{AST}$$

We use notation whereby the grammar production $\text{Program}_e ::= \text{SourceElements}_e$ means to match the input sequence with SourceElements_e and produce the resulting AST e .

Note that the grammar we present here is a simplified version of the grammar specified in the ECMAScript 5.1 standard [17] and many nonterminal names are shortened versions of nonterminals in the standard. It is mostly straightforward to extend the algorithm presented here to the full `sweet.js` implementation for ES5 JavaScript.

In addition to the *Program* parser just described, we also define a parser *Program'* that works over token trees. The rules of the two parsers are identical, except that all rules with delimiters and regular expression literals changes as follows:

$$\begin{aligned} \text{PrimaryExpr}_{/r/} &::= / \cdot x \cdot / \\ \text{PrimaryExpr}'_{/r/} &::= /r/ \\ \text{PrimaryExpr}_{(e)} &::= (\cdot \text{AssignExpr}_e \cdot) \\ \text{PrimaryExpr}'_{(e)} &::= (\text{AssignExpr}'_e) \end{aligned}$$

To prove that `read` is correct, we show that the following two parsing strategies give identical behavior:

- The traditional parsing strategy takes a token sequence s and parses s into an AST e using the traditional parser Program_e .
- The second parsing strategy first reads s into a token tree sequence $\text{read}(s, \epsilon, \text{false})$, and then parses this token tree sequence into an AST e via $\text{Program}'_e$.

Theorem 1 (Parse Equivalence).

$$\forall s.$$

$$s \in \text{Program}_e \Leftrightarrow \text{read}(s, \epsilon, \text{false}) \in \text{Program}'_e$$

Proof. The proof proceeds by induction on ASTs to show that parse equivalence holds between all corresponding nonterminals in the two grammars. We present the details of this proof in Appendix A. \square

3. Writing Macros

The `sweet.js` system provides two kinds of macros: *rule* macros (analogous to `syntax-rules` in Scheme [5]) and *case* macros (analogous to `syntax-case` in Scheme [16]). Rule macros are the simpler of the two and work by matching on a pattern and generating a template:

```
macro <name> {
  rule { <pattern> } => { <template> }
}
```

For example, the following macro introduces a new function definition form:

```
macro def {
  rule {
    $name ($params (,) ...) { $body ... }
  } => {
    function $name ($params ...) {
      $body ...
    }
  }
}
def id (x) { return x; }
// expands to:
// function id (x) { return x; }
```

The pattern is matched against the syntax following the macro name. Identifiers in a pattern that begin with \$ are *pattern variables* and bind the syntax they match in the template (identifiers that do not begin with \$ are matched literally). The ellipses (\$params (,)...) mean match zero or more tokens separated by commas.

The above example show the power of matching delimited groups of syntax (i.e. matching all the tokens inside the function body). In order for macros to be convenient in a language like JavaScript it is necessary to have the ability to match logical groupings of syntax that are not fully delimited. For example, consider the `let` macro:

```
macro let {
  rule { $id = $init:expr } => {
    var $id = $init
  }
}
let x = 40 + 2;
// expands to:
// var x = 40 + 2;
```

The initialization of a `let` can be an arbitrary expression so we use the *pattern class* `:expr` to match the largest possible expression, so in this example the entire expression `40 + 2` is bound to `$init`.

Along with the template-based rule macros, `sweet.js` also provides the more powerful procedural *case* macros. While rule macros just provide term rewriting templates, case macros allow macro authors to use JavaScript code to create and manipulate syntax. For example, the following macro creates a string from the contents of a file.

```
macro fromFile {
  case {_ ($path) } => {
    var fname = unwrapSyntax("#{ $path}");
    var f = readFile(fname);
    letstx $content = [makeValue(f, #{here})];
    return "#{ $content}"
  }
}
var s = fromFile ("./file.txt")
// expands to:
// s = "contents of file";
```

Syntax bound in the macro pattern can be referenced inside of templates using the notation `#{ ... }`. Here we take the file name token matched by the macro (`#{ $path }`), unwrap it, and `read4` the file into a variable. The `makeValue` function is used to create a string syntax object (`#{here}` is just a placeholder for the lexical context used for hygiene) from the context of the file and `letstx` (analogous to `with-syntax` in Scheme) binds that new syntax object to a pattern variable `$content` that can be used inside of a template.

4. Enforestation

The token tree structure produced by the reader is sufficient to implement an expressive macro system for a fully delimited language like Scheme. However since most of the syntax forms in a language like JavaScript are only partially delimited, it is necessary to provide additional structure during expansion that allows macros to manipulate un delimited or partially delimited groups of tokens. As an example, consider the `let` macro described earlier:

```
macro let {
  rule { $id = $init:expr } => {
    var $id = $init
  }
}
let x = 40 + 2;
// expands to:
// var x = 40 + 2;
```

Like many syntactic forms in JavaScript, the variable initialization expression is an un delimited group of tokens. Building an expressive macro system requires that the macro can match and manipulate patterns such as an expression.

`Sweet.js` groups tokens by transforming a token tree into a *term tree* through *enforestation* [23]. Enforestation works by progressively recognizing syntax forms (e.g. literals, identifiers, expressions, and statements) during expansion.

A term tree is a kind of proto-AST that represents a partial parse of the program. As the expander passes through the token trees, it creates term trees that contain unexpanded sub trees that will be expanded once all macro definitions have been discovered in the current scope (as discussed in Section 5.1).

For an example of how enforestation progresses, consider the following use of the `let` macro:

```
macro let {
  rule { $id = $init:expr } => {
    var $id = $init
  }
}
function foo(x) {
  let y = 40 + 2;
  return x + y;
}
foo(100);
```

⁴Note that `readFile` will depend on the particular JavaScript environment in which `sweet.js` is being run (e.g. in `node.js` it might be `fs.readFile` and in the browser it might an `XMLHttpRequest`).

Enforestation begins by loading the `let` macro into the macro environment and converting the function declaration into a term tree (we use angle brackets to denote a term tree). Notice that the body of the function has not yet been enforested in the first pass.

```
<fn: foo,
  args: (x),
  body: {
    let y = 40 + 2;
    return x + y;
  }>
foo(100);
```

Next, a term tree is created for the function call.

```
<fn: foo,
  params: (x),
  body: {
    let y = 40 + 2;
    return x + y;
  }>
<call: foo, args: (100)>
```

On the second pass through the top level scope the expander moves into the function body. The use of `let` is expanded away and the `var` and `return` term trees are created.

```
<fn: foo,
  params: (x),
  body: {
    <var: x, init: <op: +, left: 40, right: 2>
    <return: <op: +, left: x right: y>
  }>
<call: foo, args: (100)>
```

The additional structure provided by the term trees allow macros to match undelimited groups of tokens like binary expressions, such as `<op: +, left: 40, right: 2>`.

The enforestation technique described here was first proposed for the Honu language [23]. We take their technique and extend it with two features described in the following sections: infix macros and the invoke pattern class.

4.1 Infix Macros

The macros we have described so far must all be prefixed by the macro identifier and syntax after the macro name is matched. This is sufficient for many kinds of macros but some syntax forms require the macro identifier to sit between patterns.

For example, the upcoming ES6 version of JavaScript includes a shorthand syntax for defining functions with arrow notation:

```
id = (x) => x
// equivalent to:
// id = (function(x) { return x; }).bind(this);
```

One of the primary goals of `sweet.js` is to enable the kinds of syntax extension previously only done by the standardization committee. But, prefix macros are not capable of implementing syntax extensions like arrow functions since the macro name (`=>`) is in the middle of its syntax arguments.

Honu addressed this need for more flexible syntax transformation forms in a limited fashion via user definable unary and binary operators⁵. Definable operators provide the ability to set custom precedence and associativity along with a syntax transformation. For example, the exponentiation operator could be defined as:

```
operator ^ 10 left {$lhs, $rhs} => {
  Math.pow($lhs, $rhs)
}
2 ^ 100;
// expands to:
// Math.pow(2, 100);
```

Unfortunately, definable operators are limited in that their operands must be expressions (this limitation allows setting custom precedence and associativity rules). This restriction means it is impossible to define syntax forms such as arrow functions via definable operators since the parameter list is not an expression.

To address this limitation, `sweet.js` introduces *infix macros*, which allow a macro to match syntax that comes before the macro identifier. Infix macros are defined by adding the keyword `infix` after the `rule` or `case` keyword and placing a pipe (`|`) in the pattern where the macro name would appear (the pipe can be thought of as the cursor in the token stream).

For example, the following infix macro implements arrow functions:

```
macro => {
  rule infix {
    ($params ...) | { $body:expr }
  } => {
    function ($params ...) {
      return $body;
    }
  }
}

var id = (x) => x
// expands to:
// var id = function (x) { return x; }
```

Standard prefix macro transformers are invoked with the sequence of tokens following the macro identifier and then return a modified sequence of tokens that are then expanded:

$$transformer : TokenTree^* \rightarrow TokenTree^*$$

To implement infix macros, we modify the type of a transformer to take two arguments, one for the tokens that precede the macro identifier and the other with the tokens that follow. The transformer may then consume from either ends as needed, yielding new preceding and following tokens:

$$transformer_{infix} : (TokenTree^*, TokenTree^*) \rightarrow (TokenTree^*, TokenTree^*)$$

A naive implementation of this transformer type will lead to brittle edge cases. For example:

⁵Sweet.js will also provide support for definable operators though at the time of this writing they have not yet been fully implemented.


```
bar(x) => x
```

Here the `=>` macro is juxtaposed next to a function call, which we did not intend to be valid syntax. The naive expansion results in unparsable code:

```
bar function(x) { return x; }
```

In more subtle cases, a naive expansion might result in code that actually parses but has incorrect semantics, leading to a debugging nightmare.

To avoid this problem we provide the macro transformer with both the previous tokens and their term tree representation.

```
transformerinfix : ((TokenTree*, TermTree*), TokenTree*)
                  → (TokenTree*, TokenTree*)
```

We verify that an infix macro only matches previous tokens within boundaries delimited by the term trees. In our running example:

```
bar(x) => x
```

is first enforested to:

```
<call: bar, args: (x)> => x
```

before invoking the `=>` transformer with both the tokens `bar(x)` and the term tree `<call: bar, args: (x)>`. When the macro attempts to match just `(x)`, it detects that the parentheses splits the term tree `<call: bar, args: (x)>` and fails the match.

While infix macros fill the gap left by definable operators and allow us to write previously undefinable macros such as arrow functions, they also come with two limitations. First, there is no way to set custom precedence or associativity as with operators. It is unclear if this is a fundamental limitation or if there is some technique that might allow custom precedence and associativity for infix macros. The second limitation is that the preceding tokens to an infix macro must be first expanded. This means that any macros that occur before an infix macro will be invoked first.

This behavior introduces an asymmetry between the kinds of syntax an infix macro can match before its identifier and the kinds of syntax it can match after since the syntax following the identifier can contain unexpanded macros.

Even with these limitations infix macros work as a powerful complement to definable operators and greatly extend the kinds of syntax forms that can be implemented with macros.

4.2 Invoke and Pattern Classes

Extensibility is the guiding design principle of any expressive macro system. Since the entire point of macros is to extend the expressive power of the language, so too should the macro system itself be extensible by users. To that end `sweet.js` provides a mechanism to extend the patterns used by macros to match their arguments.

As motivation, consider the following macro that matches against color options:

```
macro colors_options {
  rule { (red) } => { ["#FF0000"] }
  rule { (green) } => { ["#00FF00"] }
  rule { (blue) } => { ["#0000FF"] }
}
r = colors_options (red)
g = colors_options (green)
// expands to:
// r = ["#FF0000"]
// g = ["#00FF00"]
```

Macros can have multiple rules and the first rule to match (from top to bottom) is used. While this macro seems to work, attempting to generalize it quickly leads to a mess:

```
macro colors_options {
  rule { (red, red) } => {
    ["#FF0000", "#FF0000"]
  }
  rule { (red, green) } => {
    ["#FF0000", "#00FF00"]
  }
  rule { (red, blue) } => {
    ["#FF0000", "#0000FF"]
  }
  // ... etc.
}
r = colors_options (red, green)
g = colors_options (green, blue)
// expands to:
// r = ["#FF0000", "#00FF00"]
// g = ["#00FF00", "#0000FF"]
```

While it is possible to solve this problem through the use of case macros, the declarative intent is quickly lost in a mess of token manipulation code.

Our solution to this problem is the `:invoke` pattern class. This pattern class takes as a parameter a macro name which is inserted into the token tree stream before matching. If the inserted macro successfully matches its arguments, the result of its expansion is bound to the pattern variable. This makes declarative options simple to write:

```
macro color {
  rule { red } => { "#FF0000" }
  rule { green } => { "#00FF00" }
  rule { blue } => { "#0000FF" }
}
macro colors_options {
  rule { ($opt:invoke(color) (,) ...) } => {
    [$opt (,) ...]
  }
}
colors_options (red, green, blue, blue)
// expands to:
// ["#FF0000", "#00FF00", "#0000FF", "#0000FF"]
```

Here the `color` macro plays the role of enumerating the valid colors that can be matched and bound to `$opt`. If the token is not in one of the rules for `color` (e.g. `orange`), the `colors_options` macro will fail to match.

As a sweet added bit of sugar, the use of `:invoke` can be inferred by just using the name of a macro in scope (i.e. `$opt:color` is equivalent to `$opt:invoke(color)`).

```
macro color {
```

```

rule { red }    => { "#FF0000" }
rule { green } => { "#00FF00" }
rule { blue }  => { "#0000FF" }
}
macro colors_options {
  rule { ($opt:color (,) ...) } => {
    [$opt (,) ...]
  }
}
colors_options (red, green, blue, blue)
// expands to:
// ["#FF0000", "#00FF00", "#0000FF", "#0000FF"]

```

By using `:invoke`, macro writers can encode patterns like alternates, optional tokens, keyword classes, numeric classes, and more in a declarative style.

5. Hygiene

Maintaining hygiene during macro expansion is perhaps the single most critical feature of an expressive macro system. The hygiene condition enables macros to be true syntactic *abstractions* by removing the burden of reasoning about a macro's implementation details from the user of a macro.

Our implementation of hygiene for `sweet.js` follows the Scheme approach [9, 16] of tracking the lexical context in each syntax object.

5.1 Macro Binding Limitation

Unfortunately, the syntax of JavaScript does place a limitations on our system that is not present in Scheme. In Scheme, definitions and uses of macros can be freely mixed in a given lexical scope. In particular, a macro can be used in an internal definition before its definition:

```

(define (foo)
  (define y (id 100))
  (define-syntax-rule (id x) x)
  (+ y y))
;; expands to:
;; (define (foo)
;;   (define y 100)
;;   (+ y y))

```

In order to support mixing use and definition, the Scheme expander must make multiple passes through a scope. The first pass registers any macro definitions and the second pass expands any uses of macros discovered during the first pass. Critically, during the first pass the expander does not expand inside internal definitions.

For example, after the first pass of expansion the above example will become:

```

(define (foo)
  (define y (id 100))
  (+ y y))

```

where the `id` macro has been registered in the macro environment. During the second pass, the expander will expand macros inside of internal definitions and so the example becomes:

```

(define (foo)

```

```

  (define y 100)
  (+ y y))

```

Scheme is able to defer expansion of macros inside of internal definitions because internal definitions are fully delimited. The expander can skip over all of the syntax inside of the delimiter because there actually is an *inside* to skip over. However, this is not true for JavaScript since `var` statements (JavaScript's equivalent of Scheme's internal definitions) are not delimited and so it is not possible to use a macro that has not yet been defined in a `var` statement in `sweet.js`.

For example, this will fail:

```

function foo() {
  var y = id 100;
  macro id { rule { $x } => { $x } }
  return y + y;
}
// fails!

```

When the expander reaches the `var` statement during the first pass, it has not yet loaded the `id` macro into the macro environment and so `id 100` will be left unexpanded. Since `id 100` is not a valid expression, this will fail to parse. Without delimiters there is no way to defer expansion of the right-hand side of a `var` statement.

Note that at first glance it might appear that the semicolon could serve to delimit a `var` statement but this will not work for two reasons. First, a semicolon is just a token which might be consumed by a macro. Second, the JavaScript specification calls for missing semicolons to be automatically inserted by the parser, which means it is not guaranteed that a semicolon will close every `var` statement.

Though the expander cannot defer expansion of `var` statements, it still does two passes so that the second pass can expand inside of delimiters. For example, a macro can be used inside of a function body that appears before the macro definition:

```

function foo() {
  return id 100;
}
macro id { rule { $x } => { $x } }
// expands to:
// function foo() {
//   return 100;
// }

```

While it is unfortunate that the syntax of JavaScript prevents fully general mixing of macro use and definition, the primary need for flexible macro definition locations is when writing mutually recursive macros, which is fully supported with our approach.

6. Related Work

Macros have been extensively used and studied in Lisp [11, 22] and related languages for many years. Scheme in particular has embraced macros, pioneering the development of declarative definitions [18] and working out the hygiene conditions for term rewriting macros (rule macros) [5] and

procedural macros (case macros) [16] that enable true composability. In addition there has been work to integrate procedural macros and module systems [10, 13]. Racket takes it even further by extending the Scheme macro system with deep hooks into the compilation process [9, 29].

In addition, there are a number of macro systems for languages with more traditional syntax that are not fully delimited. As mentioned before, `sweet.js` is most closely related to Honu [23, 24]. In contrast with Honu, which does not include regular expression literals, we solve the reader ambiguity problem for JavaScript and introduce infix macros along with the `invoke` pattern class.

Macro systems that use a similar technique as Honu include Fortress [2] and Dylan [3] however they only provide support for term rewriting macros (our `rule` macros). Dylan’s auxiliary rules are similar to our `invoke` pattern class. Nemerle [27] also uses a similar technique but does not allow local definitions of macros.

C++ templates [1] are a powerful compile time meta programming facility for C++. In contrast to C++ templates `sweet.js` provides more syntactic flexibility in the definable syntactic forms along with the ability to define transformations in the host language rather than the just the template language.

Template Haskell [26] makes a treadoff by forcing the macro call sites to always be demarcated. The means that macros are always a second class citizen; macros in Haskell cannot seamlessly build a language on top of Haskell in the same way that Scheme and Racket can.

Some systems such as SugarJ [7], ometa [31], Xtc [14], Xoc [6], and Polyglot [21] provide extensible grammars but require the programmer to reason about parser details. Multi stage systems such as mython [25] and MetaML [20, 28] can also be used to create macros systems like MacroML [12]. Systems like Stratego [30] transforms syntax using its own language, separate from the host language. Metaborg [4] and SugarJ [7] are syntax extension facilities built on top of Stratego.

7. Conclusion

We have presented `sweet.js`, a hygienic macro system for JavaScript. Our macro system follows in the footsteps of Scheme by providing declarative and procedural macro definition functionality.

In addition we presented an algorithm for read that correctly handles key ambiguities in the JavaScript grammar that allows the lexer and parser to be separated for the first time in JavaScript.

References

- [1] Andrei Alexandrescu. *Modern C++ design: generic programming and design patterns applied*. 2001.
- [2] E Allen, R Culpepper, and JD Nielsen. Growing a syntax. *Proceedings of Workshop on Foundations of Object-Oriented Languages*, 2009.
- [3] Jonathan Bachrach, Keith Playford, and Chandler Street. D-Expressions : Lisp Power , Dylan Style. *Style DeKalb IL*, 1999.
- [4] M Bravenboer and E Visser. Concrete syntax for objects: domain-specific language embedding and assimilation without restrictions. *ACM SIGPLAN Notices*, 2004.
- [5] William Clinger. Macros that work. In *Proceedings of the 18th ACM SIGPLAN-SIGACT symposium on Principles of programming languages - POPL ’91*, pages 155–162, New York, New York, USA, January 1991. ACM Press.
- [6] R Cox, T Bergan, and AT Clements. Xoc, an extension-oriented compiler for systems programming. *ACM SIGARCH Computer Architecture News*, 2008.
- [7] S Erdweg, T Rendel, C Kästner, and K Ostermann. SugarJ: library-based syntactic language extensibility. *ACM SIGPLAN Notices*, 2011.
- [8] Nate Faubion. The Sparkler project. <https://github.com/natefaubion/sparkler>.
- [9] M Flatt and R Culpepper. Macros that Work Together. *Journal of Functional Programming*, 2012.
- [10] Matthew Flatt. Composable and compilable macros: You Want it When? *ACM SIGPLAN Notices*, 37(9):72–83, September 2002.
- [11] JK Foderaro, KL Sklower, and K Layer. *The FRANZ Lisp Manual*. 1983.
- [12] SE Ganz, A Sabry, and W Taha. Macros as multi-stage computations: type-safe, generative, binding macros in MacroML. *ACM SIGPLAN Notices*, 2001.
- [13] Abdulaziz Ghuloum and R. Kent Dybvig. Implicit phasing for R6RS libraries. *ACM SIGPLAN Notices*, 42(9):303, October 2007.
- [14] R Grimm. Better extensibility through modular syntax. *ACM SIGPLAN Notices*, 2006.
- [15] S P Harbison and Jr. Steele, G. L. *C: A Reference Manual*. Prentice-Hall, 1984.
- [16] R Hieb, RK Dybvig, and C Bruggeman. Syntactic abstraction in scheme. *Lisp and symbolic computation*, 5(4):295–326, 1992.
- [17] E. C. M. A. International. *ECMA-262 ECMAScript Language Specification*. Number June. ECMA (European Association for Standardizing Information and Communication Systems), 5.1 edition, 2011.
- [18] E. E. Kohlbecker and M. Wand. Macro-by-example: Deriving syntactic transformations from their specifications. In *Proceedings of the 14th ACM SIGACT-SIGPLAN symposium on Principles of programming languages - POPL ’87*, pages 77–84, New York, New York, USA, October 1987. ACM Press.
- [19] James Long. The es6-macros project. <https://github.com/jlongster/es6-macros>.
- [20] M Martel and T Sheard. Introduction to multi-stage programming using metaml. 1997.
- [21] N Nystrom, MR Clarkson, and AC Myers. Polyglot: An extensible compiler framework for Java. *Compiler Construction*,

2003.

- [22] Kent M. Pitman. Special forms in Lisp. In *Proceedings of the 1980 ACM conference on LISP and functional programming - LFP '80*, pages 179–187, New York, New York, USA, August 1980. ACM Press.
- [23] Jon Rafkind. *Syntactic extension for languages with implicitly delimited and infix syntax*. PhD thesis, 2013.
- [24] Jon Rafkind and Matthew Flatt. Honu: Syntactic Extension for Algebraic Notation through Enforestation. *Proceedings of the 11th International Conference on Generative Programming and Component Engineering*, 2012.
- [25] J Riehl. Language embedding and optimization in mython. *ACM Sigplan Notices*, 2009.
- [26] T Sheard and SP Jones. Template meta-programming for Haskell. *Proceedings of the 2002 ACM SIGPLAN workshop on Haskell*, 2002.
- [27] K Skalski, M Moskal, and P Olszta. Meta-programming in Nemerle. *Proceedings Generative Programming and Component Engineering*, 2004.
- [28] W Taha and T Sheard. Multi-stage programming with explicit annotations. *ACM SIGPLAN Notices*, 1997.
- [29] S Tobin-Hochstadt and V St-Amour. Languages as libraries. *PLDI '11 Proceedings of the 32nd ACM SIGPLAN conference on Programming language design and implementation*, 2011.
- [30] E Visser. Program transformation with Stratego/XT. *Domain-Specific Program Generation*, 2004.
- [31] A Warth and I Piumarta. OMeta: an object-oriented language for pattern matching. *Proceedings of the 2007 symposium on Dynamic languages*, 2007.

A. Read Proof

To help reason about prefixes, we use two disjoint sets, *RegexPrefix* and *DividePrefix*, that contain the prefixes that determine if a slash should be either a divide or the start of a regular expression. These sets are parametrized by a boolean indicating if the prefix is inside an object literal or a block statement:

$$\begin{aligned} \text{RegexPrefix}_b &::= \epsilon \\ &| p \cdot k \quad \text{if } k \in \text{Punctuator} \cup \text{Keyword} \\ &| p \cdot \text{if} \cdot \underline{(t)} \\ &| p \cdot \text{function}^l \cdot x \cdot \underline{(t)} \cdot \{t'\} \\ &\quad \text{if isExprPrefix}(p, b, l) = \text{false} \\ &| p \cdot \{t\}^l \\ &\quad \text{if isExprPrefix}(p, b, l) = \text{false} \end{aligned}$$

$$\begin{aligned} \text{DividePrefix}_b &::= p \cdot x \\ &| p \cdot /r/ \\ &| p \cdot k \cdot \underline{(t)} \quad \text{if } k \neq \text{if} \\ &\quad \text{if isExprPrefix}(p, b, l) = \text{true} \\ &| p \cdot \{t\}^l \\ &\quad \text{if isExprPrefix}(p, b, l) = \text{true} \end{aligned}$$

These sets correspond to the prefixes in the read function that determine if a slash should be regular expression or a divide operator.

Theorem 2 (Parse Equivalence for Program).

$$\forall s \in \text{Token}^*.$$

$$\begin{aligned} &s \in \text{Program}_e \\ \Leftrightarrow &\text{read}(s, \epsilon, \text{false}) \in \text{Program}'_e \end{aligned}$$

Proof. For the left-to-right direction, there are two production rules for *Program_e*.

- $s \in \epsilon$. The result is immediate.
- $s \in \text{SourceElements}_e$. Holds by Lemma 1.

A similar argument holds for the other direction. \square

Lemma 1 (Parse Equivalence for SourceElements).

$$\forall s \in \text{Token}^*.$$

$$\begin{aligned} &s \in \text{SourceElements}_e \\ \Leftrightarrow &\text{read}(s, \epsilon, \text{false}) \in \text{SourceElements}'_e \end{aligned}$$

Proof. For the left-to-right direction there are two production rules for *SourceElements_e*.

- $s \in \text{SourceElement}_e$. This holds by Lemma 2.
- $s \in \text{SourceElements}_e \cdot \text{SourceElement}_e$. We have $s = s' \cdot s''$ where $s' \in \text{SourceElements}_e$ and $s'' \in \text{SourceElement}_e$.

$$\begin{aligned} t &= \text{read}(s' \cdot s'', \epsilon, \text{false}) \\ &= \text{read}(s', \epsilon, \text{false}) \cdot \text{read}(s'', \text{read}(s', \epsilon, \text{false}), \text{false}) \end{aligned}$$

By induction $\text{read}(s', \epsilon, \text{false}) \in \text{SourceElements}'_e$ and by Lemma 2, $\text{read}(s'', \text{read}(s', \epsilon, \text{false}), \text{false}) \in$

$\text{SourceElement}'_{e'}$ (since by Lemma 11, $\text{read}(s', \epsilon, \text{false}) \in \text{RegexPrefix}_b$). Thus $t \in \text{SourceElements}'_{e'}$.

The argument is similar for the other direction. \square

Lemma 2 (Parse Equivalence for SourceElement).

$$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b.$$

$$\begin{aligned} &s \in \text{SourceElement}_e \\ \Leftrightarrow &\text{read}(s, p, \text{false}) \in \text{SourceElement}'_e \end{aligned}$$

Proof. For the left-to-right direction there are two production rules for *SourceElement_e*.

- $s \in \text{Stmt}_e$. This holds by Lemma 4 since $p \in \text{RegexPrefix}_b$.
- $s \in \text{FunctionDecl}_e$. This holds by Lemma 3.

The argument is similar for the other direction. \square

Lemma 3 (Parse Equivalence for Function).

$$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{TokenTree}^*.$$

$$\begin{aligned} &s \in \text{Function}_e \\ \Leftrightarrow &\text{read}(s, p, b) \in \text{Function}'_e \end{aligned}$$

Proof. For the left-to-right direction there is one production rule for *Function_e*:

$$s \in \text{function} \cdot x \cdot (\cdot x' \cdot) \cdot \{ \cdot \text{SourceElements}_e \cdot \}$$

We have $s = \text{function} \cdot x \cdot (\cdot x \cdot) \cdot \{ \cdot s' \cdot \}$ where $s' \in \text{SourceElements}_e$ so:

$$\begin{aligned} t &= \text{read}(\text{function} \cdot x \cdot (\cdot x \cdot) \cdot \{ \cdot s' \cdot \}, p, b) \\ &= \text{function} \cdot x \cdot \underline{(x)} \cdot \{t'\} \end{aligned}$$

where $t' = \text{read}(s', \epsilon, \text{false})$. Since by Lemma 1, $t' \in \text{SourceElements}'_e$ we have $t \in \text{FunctionDecl}'_{\text{function } x(x) \{e\}}$.

The argument is similar for the other direction. \square

Lemma 4 (Parse Equivalence for Stmt).

$$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b.$$

$$\begin{aligned} &s \in \text{Stmt}_e \\ \Leftrightarrow &\text{read}(s, p, \text{false}) \in \text{Stmt}'_e \end{aligned}$$

Proof. For the left-to-right direction we have several production rules for *Stmt_e*.

- $s \in \{ \cdot \text{StmtList}_e \cdot \}$. We have $s = \{ \cdot s' \cdot \}$ where $s' \in \text{StmtList}_e$. Then:

$$\begin{aligned} t &= \text{read}(\{ \cdot s' \cdot \}, p, \text{false}) \\ &= \underline{\{t'\}} \end{aligned}$$

where $t' = \text{read}(s', \epsilon, \text{false})$. By Lemma 5, $t' \in \text{StmtList}'_e$ so $t \in \text{Stmt}'_{\{e\}}$.

- $s \in \text{AssignExpr}_e$;. We have $s = s' ;$; where $s' \in \text{AssignExpr}_e$. Then:

$$\begin{aligned} t &= \text{read}(s' ;, p, \text{false}) \\ &= \text{read}(s', p, \text{false}) ; \end{aligned}$$

Since $p \in \text{RegexPrefix}_b$ by Lemma 6, $\text{read}(s', p, \text{false}) \in \text{AssignExpr}'_e$ we have $t \in \text{Stmt}'_e$.

- $s \in \text{if} \cdot (\cdot \text{AssignExpr}_e \cdot) \cdot \text{Stmt}_{e'}$. We have $s = \text{if} \cdot (\cdot s' \cdot) \cdot s''$ where $s' \in \text{AssignExpr}_e$ and $s'' \in \text{Stmt}_{e'}$.

$$\begin{aligned} t &= \text{read}(\text{if} \cdot (\cdot s' \cdot) \cdot s'', p, \text{false}) \\ &= \text{if} \cdot \underline{(t')} \cdot t'' \end{aligned}$$

where

$$\begin{aligned} t' &= \text{read}(s', \epsilon, \text{false}) \\ t'' &= \text{read}(s'', p \cdot \text{if} \cdot \underline{(t')}, \text{false}) \end{aligned}$$

By Lemma 6 $t' \in \text{AssignExpr}'_e$. By induction, $t'' \in \text{Stmt}'_{e'}$ (since $p \cdot \underline{(t')} \in \text{RegexPrefix}_b$). Thus $t \in \text{Stmt}'_{\text{if} \cdot (\epsilon) \cdot e'}$.

- $s \in \text{return}$. Since $s = \text{return}$ and $\text{read}(s, p, \text{false}) = \text{return} \in \text{Stmt}'_{\text{return}}$ we have our result directly.
- $s \in \text{return} \cdot \text{AssignExpr}_e$;. We have $s = \text{return} \cdot s' ;$; where $s' \in \text{AssignExpr}_e$. Then:

$$\begin{aligned} t &= \text{read}(\text{return} \cdot s' ;, p, \text{false}) \\ &= \text{return} \cdot \text{read}(s', p \cdot \text{return}, \text{false}) ; \end{aligned}$$

By Lemma 6, $\text{read}(s', p \cdot \text{return}, \text{false}) \in \text{AssignExpr}'_e$ thus $t \in \text{Stmt}'_{\text{return} \cdot e}$.

- $s \in x : : \cdot \text{Stmt}_e$. We have $s = x : : \cdot s'$ where $s' \in \text{Stmt}_e$. Then:

$$\begin{aligned} t &= \text{read}(x : : \cdot s', p, \text{false}) \\ &= x : : \cdot \text{read}(s', p \cdot x : :, \text{false}) \end{aligned}$$

By induction, $\text{read}(s', p \cdot x : :, \text{false}) \in \text{Stmt}'_e$ so $t \in \text{Stmt}'_{x : : e}$.

The argument for the other direction is similar. \square

Lemma 5 (Parse Equivalence for StmtList).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$.

$$\begin{aligned} s &\in \text{StmtList}_e \\ \Leftrightarrow \text{read}(s, p, \text{false}) &\in \text{StmtList}'_e \end{aligned}$$

Proof. For the left-to-right direction we have two production rules for StmtList_e .

- $s \in \text{Stmt}_e$. This follows by Lemma 4 since $p \in \text{RegexPrefix}_b$.
- $s \in \text{StmtList}_e \cdot \text{Stmt}_{e'}$. So $s = s' \cdot s''$ where $s' \in \text{StmtList}_e$ and $s'' \in \text{Stmt}_{e'}$. Then,

$$\begin{aligned} t &= \text{read}(s' \cdot s'', p, \text{false}) \\ &= \text{read}(s', p, \text{false}) \cdot \text{read}(s'', p \cdot \text{read}(s', p, \text{false}), \text{false}) \end{aligned}$$

By induction $\text{read}(s', p, \text{false}) \in \text{StmtList}'_e$ and by Lemma 4, $\text{read}(s'', p \cdot \text{read}(s', p, \text{false}), \text{false}) \in \text{Stmt}'_{e'}$ (since by Lemma 12, $p \cdot \text{read}(s', p, \text{false}) \in \text{RegexPrefix}_b$). Thus, $t \in \text{StmtList}'_{e \cdot e'}$.

The argument for the other direction is similar. \square

Lemma 6 (Parse Equivalence for AssignExpr).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$. Where if $b = \text{false}$ then $s \neq \{ \cdot s' \}$, we have that

$$\begin{aligned} s &\in \text{AssignExpr}_e \\ \Leftrightarrow \text{read}(s, p, b) &\in \text{AssignExpr}'_e \end{aligned}$$

Proof. The constraint on s when b is false is due to the lookahead check in the production $\text{Stmt}_e ::= \text{AssignExpr}_e$;. Lemma 10 will make use of this constraint.

For the left-to-right direction we have two production rules for AssignExpr_e .

- $s \in \text{BinaryExpr}_e$. This holds by Lemma 7.
- $s \in \text{CallExpr}_e \cdot = \cdot \text{AssignExpr}_{e'}$. Then $s = s' \cdot = \cdot s''$ where $s' \in \text{CallExpr}_e$ and $s'' \in \text{AssignExpr}_{e'}$. Then:

$$\begin{aligned} t &= \text{read}(s' \cdot = \cdot s'', p, b) \\ &= \text{read}(s', p, b) \cdot \text{read}(= \cdot s'', p \cdot \text{read}(s', p, b), b) \\ &= \text{read}(s', p, b) \cdot = \cdot \text{read}(s'', p \cdot \text{read}(s', p, b) \cdot =, b) \end{aligned}$$

Since $p \in \text{RegexPrefix}_b$ by Lemma 8, $\text{read}(s', p, b) \in \text{CallExpr}'_e$ and by induction $\text{read}(s'', p \cdot \text{read}(s', p, b) \cdot =, b) \in \text{AssignExpr}'_{e'}$ (since $p \cdot \text{read}(s', p, b) \cdot = \in \text{RegexPrefix}_b$) so $t \in \text{AssignExpr}'_{e \cdot e'}$.

The argument for the other direction is similar. \square

Lemma 7 (Parse Equivalence for BinaryExpr).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$. Where if $b = \text{false}$ then $s \neq \{ \cdot s' \}$, we have that

$$\begin{aligned} s &\in \text{BinaryExpr}_e \\ \Leftrightarrow \text{read}(s, p, b) &\in \text{BinaryExpr}'_e \end{aligned}$$

Proof. For the left-to-right direction there are three productions for BinaryExpr_e .

- $s \in \text{CallExpr}_e$. This holds by Lemma 8 since $p \in \text{RegexPrefix}_b$.
- $s \in \text{BinaryExpr}_e \cdot / \cdot \text{BinaryExpr}_{e'}$. We have $s = s' \cdot / \cdot s''$ where $s' \in \text{BinaryExpr}_e$ and $s'' \in \text{BinaryExpr}_{e'}$. Then:

$$\begin{aligned} t &= \text{read}(s' \cdot / \cdot s'', p, b) \\ &= \text{read}(s', p, b) \cdot \text{read}(/ \cdot s'', p \cdot \text{read}(s', p, b), b) \\ &= \text{read}(s', p, b) \cdot / \cdot \text{read}(s'', p \cdot \text{read}(s', p, b) \cdot /, b) \\ &\quad (\text{since } p \cdot \text{read}(s', p, b) \in \text{DividePrefix}_b) \end{aligned}$$

By induction $\text{read}(s', p, b) \in \text{BinaryExpr}'_e$ and $\text{read}(s'', p \cdot \text{read}(s', p, b) \cdot /, b) \in \text{BinaryExpr}'_{e'}$ (since $p \cdot \text{read}(s', p, b) \cdot / \in \text{RegexPrefix}_b$) thus

- $s \in \text{BinaryExpr}_e \cdot + \cdot \text{BinaryExpr}_{e'}$. We have $s = s' \cdot + \cdot s''$ where $s' \in \text{BinaryExpr}_e$ and $s'' \in \text{BinaryExpr}_{e'}$. Then:

$$\begin{aligned} t &= \text{read}(s' \cdot + \cdot s'', p, b) \\ &= \text{read}(s', p, b) \cdot + \cdot \text{read}(s'', p \cdot \text{read}(s', p, b) \cdot +, b) \end{aligned}$$

By induction $\text{read}(s', p, b) \in \text{BinaryExpr}'_e$ and $\text{read}(s'', p \cdot \text{read}(s', p, b) \cdot +, b) \in \text{BinaryExpr}'_{e'}$ (since $p \cdot \text{read}(s', p, b) \cdot + \in \text{RegexPrefix}_b$) thus $t \in \text{BinaryExpr}'_{e+e'}$.

The argument for the other direction is similar. \square

Lemma 8 (Parse Equivalence for CallExpr).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$. Where if $b = \text{false}$ then $s \neq \{ \cdot s' \}$, we have that

$$\begin{aligned} s &\in \text{CallExpr}_e \\ \Leftrightarrow \text{read}(s, p, b) &\in \text{CallExpr}'_e \end{aligned}$$

Proof. For the left-to-right direction there are two production rules for CallExpr_e .

- $s \in \text{MemberExpr}_e \cdot (\cdot \text{AssignExpr}'_{e'} \cdot)$. We have $s = s' \cdot (\cdot s'' \cdot)$ where $s' \in \text{MemberExpr}_e$ and $s'' \in \text{AssignExpr}'_{e'}$. Then

$$\begin{aligned} t &= \text{read}(s' \cdot (\cdot s'' \cdot), p, b) \\ &= \text{read}(s', p, b) \cdot \text{read}(\cdot s'' \cdot, p \cdot \text{read}(s', p, b), b) \end{aligned}$$

Since $p \in \text{RegexPrefix}_b$, by Lemma 9 we have $\text{read}(s', p, b) \in \text{MemberExpr}'_e$ and my Lemma 6 we have $\text{read}(s'', \epsilon, \text{false}) \in \text{AssignExpr}'_{e'}$ thus $t \in \text{CallExpr}'_{e(e')}$.

- $s \in \text{CallExpr}_e \cdot \dots \cdot x$. Then $s = s' \cdot \dots \cdot x$ where $s' \in \text{CallExpr}_e$. Then

$$\begin{aligned} t &= \text{read}(s' \cdot \dots \cdot x, p, b) \\ &= \text{read}(s', p, b) \cdot \dots \cdot x \end{aligned}$$

By induction $\text{read}(s', p, b) \in \text{CallExpr}'_e$. Thus $t \in \text{CallExpr}'_{e \cdot x}$.

The argument for the other direction is similar. \square

Lemma 9 (Parse Equivalence for MemberExpr).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$. Where if $b = \text{false}$ then $s \neq \{ \cdot s' \}$, we have that

$$\begin{aligned} s &\in \text{MemberExpr}_e \\ \Leftrightarrow \text{read}(s, p, b) &\in \text{MemberExpr}'_e \end{aligned}$$

Proof. For the left-to-right direction there are two three production rules for MemberExpr_e .

- $s \in \text{PrimaryExpr}_e$. This follows from Lemma 10 since $p \in \text{RegexPrefix}_b$.
- $s \in \text{Function}_e$. This follows from Lemma 3.
- $s \in \text{MemberExpr}_e \cdot \dots \cdot x$. We have $s = s' \cdot \dots \cdot x$ where $s' \in \text{MemberExpr}_e$. Then

$$\begin{aligned} t &= \text{read}(s' \cdot \dots \cdot x, p, b) \\ &= \text{read}(s', p, b) \cdot \dots \cdot x \end{aligned}$$

By induction $\text{read}(s', p, b) \in \text{MemberExpr}'_e$ thus $t \in \text{MemberExpr}'_{e \cdot x}$.

The argument for the other direction is similar. \square

Lemma 10 (Parse Equivalence for PrimaryExpr).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$. Where if $b = \text{false}$ then $s \neq \{ \cdot s' \}$, we have that

$$\begin{aligned} s &\in \text{PrimaryExpr}_e \\ \Leftrightarrow \text{read}(s, p, b) &\in \text{PrimaryExpr}'_e \end{aligned}$$

Proof. For the left-to-right direction there are several production rules for PrimaryExpr_e .

- $s \in x$. Then $s = x$ and $\text{read}(x, p, b) \in \text{PrimaryExpr}'_x$ directly.
- $s \in / \cdot x \cdot /$. Then $s = / \cdot x \cdot /$ and $\text{read}(/ \cdot x \cdot / , p, b) = /r/ \in \text{PrimaryExpr}'_{/r/}$ since $p \in \text{RegexPrefix}_b$.
- $s \in \{ \cdot x \cdot : \cdot \text{AssignExpr}'_e \cdot \}$. Then $s = \{ \cdot x \cdot : \cdot s' \cdot \}$ where $s' \in \text{AssignExpr}'_e$. Then:

$$\begin{aligned} t &= \text{read}(\{ \cdot x \cdot : \cdot s' \cdot \}, p, \text{true}) \\ &= \underline{\{x \cdot : \cdot t'\}} \end{aligned}$$

where $t' = \text{read}(s', x \cdot :, \text{true})$. By Lemma 6, $t' \in \text{AssignExpr}'_e$ and thus $t \in \text{PrimaryExpr}'_{\{x:e\}}$.

- $s \in (\cdot \text{AssignExpr}'_e \cdot)$. Then $s = (\cdot s' \cdot)$ where $s' \in \text{AssignExpr}'_e$. So,

$$\begin{aligned} t &= \text{read}(\cdot s' \cdot, p, b) \\ &= \underline{(t')} \end{aligned}$$

where $t' = \text{read}(s', \epsilon, \text{false})$. By Lemma 6 $t' \in \text{AssignExpr}'_e$. So, $t \in \text{PrimaryExpr}'_{(e)}$.

For the other direction the argument is similar. \square

Lemma 11 (SourceElement Prefix).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$.

$$\begin{aligned} s &\in \text{SourceElement}_e \\ \Rightarrow \text{read}(s, p, \text{false}) &\in \text{RegexPrefix}_b \end{aligned}$$

Proof. There are two cases:

- $s \in \text{Stmt}_e$. This holds by Lemma 12.
- $s \in \text{Function}_e$. Follows directly. \square

Lemma 12 (Stmt Prefix).

$\forall s \in \text{Token}^*, b \in \text{Bool}, p \in \text{RegexPrefix}_b$.

$$\begin{aligned} \text{read}(s, p, \text{false}) &\in \text{Stmt}'_e \\ \Rightarrow p \cdot \text{read}(s, p, \text{false}) &\in \text{RegexPrefix}_b \end{aligned}$$

Proof. We have several cases:

- $\text{read}(s, p, \text{false}) \in \underline{\{\text{StmtList}'_e\}}$. Follows since $p \in \text{RegexPrefix}_b$.
- $\text{read}(s, p, \text{false}) \in \text{AssignExpr}'_e \cdot ;$. Follows since $t \cdot ; \in \text{RegexPrefix}_b$ for any t .

- $t = \text{read}(s, p, \text{false}) \in \text{if} \cdot (\text{AssignExpr}'_e) \text{ Stmt}'_e$.
Since $t = \text{if} \cdot \overline{(t')} \cdot t''$ where $t'' \in \text{Stmt}'_e$ by induction $t'' \in \text{RegexPrefix}_b$ and thus $p \cdot t \in \text{RegexPrefix}_b$.
- $\text{read}(s, p, \text{false}) \in \text{return}$. Follows since $p \cdot \text{return} \in \text{RegexPrefix}_b$.
- $t = \text{read}(s, p, \text{false}) \in \text{return} \cdot \text{AssignExpr}'_e \cdot ;$. Since $;$ $\in \text{RegexPrefix}_b$ then $p \cdot t \in \text{RegexPrefix}_b$.
- $t = \text{read}(s, p, \text{false}) \in x \cdot : \cdot \text{Stmt}'_e$. Since $t = x \cdot : \cdot t'$ where $t' \in \text{Stmt}'_e$. Since $p \in \text{RegexPrefix}_b$ and by induction $t' \in \text{RegexPrefix}_b$ we have $p \cdot t \in \text{RegexPrefix}_b$.

□