

# Sweet.js - Hygienic Macros for JavaScript

Tim Disney  
UC Santa Cruz

Nate Faubion

David Herman  
Mozilla

Cormac Flanagan  
UC Santa Cruz

## Abstract

Sweet.js is a hygienic macro system for JavaScript. It's pretty sweet.

## 1. Introduction

Expressive macros systems have a long history in the design of extensible programming languages going back to Lisp and Scheme [3, 7] 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 pipeline 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:

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

As a first step to designing a Scheme-like macro system for JavaScript, it is necessary to introduce a read step to 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 divide operator.

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

It is necessary to separate the parser and lexer in order to implement a macro system for JavaScript. Our system, `sweet.js`<sup>1</sup>, implements such a separated reader that converts a sequence of tokens into a sequence of token trees (a little analogous to s-expressions) without feedback from the parser.

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

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 always correctly decide whether to parse a sequence of tokens `/x/g` as a regular expression or as division operators and operands (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

[Copyright notice will appear here once 'preprint' option is removed.]

<sup>1</sup> <http://sweetjs.org>

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 the partial AST. For example, an `if` macro could indicate that the then and else branches must be single expressions and then manipulate them appropriately.

This approach, called *enforestation*, was pioneered by Honu [8, 9], which we adapt here for JavaScript<sup>2</sup>. In addition, we make two extensions to the Honu technique. 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.

#### **todo: more about the project like users**

We make 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 [4], work over just the token stream. However, as mentioned in the introduction, 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.

Our macro system includes a reader that transforms a stream of flat tokens into a stream of more structured *token trees*, by matching delimiters.

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. Unfortunately, the syntax of JavaScript presents a challenge to correctly implement the `read` function. In particular, one of the syntax features of JavaScript is support for regular expression literals, which are formed with the forward slash `/`. Delimiters are valid inside of a regular expression literal so the reader must not match delimiters inside of a regular expression literal otherwise it would break the true delimiter structure of the program.

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

Since the reader must ignore delimiters inside regular expression literals, it must know when it is reading between the literal's slashes. However, slash is also used as the divide operator in JavaScript. So, in order for the reader to ignore delimiters inside of regex literals it must first decide if a slash is the beginning of a literal or a divide operator. Unfortunately this is somewhat complicated for JavaScript.

Traditionally, JavaScript parsers are not separated from the lexer; the parser calls the lexer with a flag indicating if the grammatical context accepts either a regular expression literal or a divide token. For JavaScript parsers this technique is sufficient but a macro expander requires the separation of the lexer and reader from the parser.

Note that this problem is not present in Honu since their language does not include regular expression literals.

A key novelty in `sweet.js` is the design and implementation of a correct reader for full ES5 JavaScript<sup>3</sup>. For clarity of presentation, this paper describes the implementation of `read` for the subset of JavaScript shown in Figure 3, which retains the essential complications of a correct version of `read` but elides the full details of the entire JavaScript language.

`Read` takes a sequence of tokens and produces a sequence of token trees. Tokens are punctuators (`+`, `/`, `:`, `;`, `=`, `.`), keywords (`return`, `function`, `if`), variables (`x`), and delimiters (`{`, `(`, `}`, `)`): see Figure 3. We write a sequence of tokens separated by a dot so the source string `"foo(/x/)"` is lexed into the sequence of six tokens `foo · ( · / · x · / · )`.

A token tree is similar to a token with the addition of the regular expression literal (`/x/`) and instead of individual opening and closing delimiter tokens a delimiter nests a sequence of inner token trees, so that the sequence `foo · ( · / · x · / · )` would be the token tree sequence `foo · (/x/)`.

The key idea of `read` is to maintain a 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. Note that this may correspond

<sup>2</sup> Honu does not support regular expression literals which simplifies their reader.

<sup>3</sup> Our implementation also has initial support for the upcoming ES6 version of JavaScript

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. `100 + /x/g`) the slash will always be a regular expression literal.

Other cases require additional context to disambiguate. For example, if the previous token tree was a parentheses (e.g. dividing the result of a call `foo(100)/ 10`) 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) /y/g
```

One of the most complicated cases is a slash following 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 both object literals and blocks with labeled statements can nest:

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

Here, 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.

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.

## 2.1 Proving Read

To show that our read algorithm correctly distinguishes between the divide operator and a regular expression literal, 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 3, and generates ASTs in the abstract syntax shown in Figure 1. A parser for the nonterminal *Program* is a function from a

**Figure 1: AST for Simplified JavaScript**

---


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


---

sequence of tokens to an AST.

$$Program :: Token^* \rightarrow AST$$

We use the notation  $Program_e ::= SourceElements_e$  to mean 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 standard [6] and many of the nonterminal names we use here correspond to shortened versions of nonterminals in the standard.

It is mostly straightforward to extend the algorithm presented here for the simplified language to the sweet.js implementation for full 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 similar except for that all rules with delimiters and regular expression literals will change:

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

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

- The traditional parsing strategy is, given a token stream  $s$ , to parse  $s$  into an AST  $e$  using a traditional parser.
- The second parsing strategy first reads  $s$  into a token tree stream  $t = \text{read}(s, \epsilon, \text{false})$ , and then parses this token tree stream  $t$  into an AST  $e$ .

**Theorem 1 (Parse Equivalence).**

$\forall s.$

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

*Proof.* By showing parse equivalence for each non-terminal in the grammar. Details in the appendix.  $\square$

## 3. Writing Macros

The sweet.js system provides two kinds of macros: *rule* macros (analogous to `syntax-rules` in Scheme [1]) and *case* macros (analogous to `syntax-case` in Scheme [5]). Rule macros are the simpler of the two and work by matching on a pattern and generating a template:

**Figure 2: Simplified Read Algorithm**

<i>Punctuator</i> ::= /   +   :   ;   =   .	
<i>Keyword</i> ::= return   function   if	
<i>Token</i> ::= <i>Punctuator</i>   <i>Keyword</i>	
	<i>x</i>   {   (   }   )
<i>TokenTree</i> ::= <i>Punctuator</i>   <i>Keyword</i>	
	<i>x</i>   / <i>x</i> /   ( <i>t</i> )   { <i>t</i> }
<i>x</i> ∈ <i>Variable</i>	
<i>s</i> ∈ <i>Token</i> *	
<i>t, p</i> ∈ <i>TokenTree</i> *	
	<i>isExprPrefix</i> : <i>TokenTree</i> * → <i>Bool</i> → <i>Int</i>
	<i>isExprPrefix</i> ( $\epsilon$ , <i>true</i> , <i>l</i> ) = <i>true</i>
	<i>isExprPrefix</i> ( <i>p</i> · /, <i>b</i> , <i>l</i> ) = <i>true</i>
	<i>isExprPrefix</i> ( <i>p</i> · +, <i>b</i> , <i>l</i> ) = <i>true</i>
	<i>isExprPrefix</i> ( <i>p</i> · =, <i>b</i> , <i>l</i> ) = <i>true</i>
	<i>isExprPrefix</i> ( <i>p</i> · :, <i>b</i> , <i>l</i> ) = <i>b</i>
	<i>isExprPrefix</i> ( <i>p</i> · return <sup><i>l</i></sup> , <i>b</i> , <i>l'</i> ) = <i>false</i> if <i>l</i> ≠ <i>l'</i>
	<i>isExprPrefix</i> ( <i>p</i> · return <sup><i>l</i></sup> , <i>b</i> , <i>l'</i> ) = <i>true</i> if <i>l</i> = <i>l'</i>
	<i>isExprPrefix</i> ( <i>p</i> , <i>b</i> , <i>l</i> ) = <i>false</i>
<i>read</i> : <i>Token</i> * → <i>TokenTree</i> * → <i>Bool</i> → <i>TokenTree</i> *	
<i>read</i> (/ · <i>x</i> · / · <i>s</i> , $\epsilon$ , <i>b</i> )	= / <i>x</i> / · <i>read</i> ( <i>s</i> , / <i>x</i> /, <i>b</i> )
<i>read</i> (/ · <i>x</i> · / · <i>s</i> , <i>p</i> · <i>p'</i> , <i>b</i> )	= / <i>x</i> / · <i>read</i> ( <i>s</i> , <i>p</i> · <i>p'</i> · / <i>x</i> /, <i>b</i> )
if <i>p'</i> ∈ <i>Punctuator</i> or <i>Keyword</i>	
<i>read</i> (/ · <i>x</i> · / · <i>s</i> , <i>p</i> · if · ( <i>t</i> ), <i>b</i> )	= / <i>x</i> / · <i>read</i> ( <i>s</i> , <i>p</i> · if · ( <i>t</i> ) · / <i>x</i> /, <i>b</i> )
<i>read</i> (/ · <i>x</i> · / · <i>s</i> , <i>p</i> · <i>p'</i> · function <sup><i>l</i></sup> · <i>x</i> · ( <i>t</i> ) · { <i>t'</i> }, <i>b</i> )	= / <i>x</i> / · <i>read</i> ( <i>s</i> , <i>p</i> · <i>p'</i> · function <sup><i>l</i></sup> · <i>x</i> · ( <i>t</i> ) · { <i>t'</i> } · / <i>x</i> /, <i>b</i> )
if <i>false</i> = <i>isExprPrefix</i> ( <i>p'</i> , <i>b</i> , <i>l</i> )	
<i>read</i> (/ · <i>x</i> · / · <i>s</i> , <i>p</i> · { <i>t</i> } <sup><i>l</i></sup> , <i>b</i> )	= / <i>x</i> / · <i>read</i> ( <i>s</i> , <i>p</i> · { <i>t</i> } <sup><i>l</i></sup> · / <i>x</i> /, <i>b</i> )
where <i>false</i> = <i>isExprPrefix</i> ( <i>p</i> , <i>b</i> , <i>l</i> )	
<i>read</i> (/ · <i>s</i> , <i>p</i> · <i>x</i> , <i>b</i> )	= / · <i>read</i> ( <i>s</i> , <i>p</i> · <i>x</i> · /, <i>b</i> )
<i>read</i> (/ · <i>s</i> , <i>p</i> · / <i>x</i> /, <i>b</i> )	= / · <i>read</i> ( <i>s</i> , <i>p</i> · / <i>x</i> / · /, <i>b</i> )
<i>read</i> (/ · <i>s</i> , <i>p</i> · ( <i>t</i> ), <i>b</i> )	= / · <i>read</i> ( <i>s</i> , <i>p</i> · ( <i>t</i> ) · /, <i>b</i> )
<i>read</i> (/ · <i>s</i> , <i>p</i> · <i>p'</i> · function <sup><i>l</i></sup> · <i>x</i> · ( <i>t</i> ) · { <i>t'</i> }, <i>b</i> )	= / · <i>read</i> ( <i>s</i> , <i>p</i> · <i>p'</i> · function <sup><i>l</i></sup> · <i>x</i> · ( <i>t</i> ) · { <i>t'</i> } · /, <i>b</i> )
if <i>true</i> = <i>isExprPrefix</i> ( <i>p'</i> , <i>b</i> , <i>l</i> )	
<i>read</i> (/ · <i>s</i> , <i>p</i> · { <i>t</i> } <sup><i>l</i></sup> , <i>b</i> )	= / · <i>read</i> ( <i>s</i> , <i>p</i> · { <i>t</i> } <sup><i>l</i></sup> · /, <i>b</i> )
where <i>true</i> = <i>isExprPrefix</i> ( <i>p</i> , <i>b</i> , <i>l</i> )	
<i>read</i> (( · <i>s</i> · ) · <i>s'</i> , <i>p</i> , <i>b</i> )	= ( <i>t</i> ) · <i>read</i> ( <i>s'</i> , <i>p</i> · ( <i>t</i> ), <i>b</i> )
where <i>s</i> contains no unmatched )	where <i>t</i> = <i>read</i> ( <i>s</i> , $\epsilon$ , <i>false</i> )
<i>read</i> ({ <sup><i>l</i></sup> · <i>s</i> · } · <i>s'</i> , <i>p</i> , <i>b</i> )	= { <sup><i>l</i></sup> · <i>t</i> · } · <i>read</i> ( <i>s'</i> , <i>p</i> · { <sup><i>l</i></sup> · <i>t</i> · }, <i>b</i> )
where <i>s</i> contains no unmatched {	where <i>t</i> = <i>read</i> ( <i>s</i> , $\epsilon$ , <i>isExprPrefix</i> ( <i>p</i> , <i>b</i> , <i>l</i> ))
<i>read</i> ( <i>x</i> · <i>s</i> , <i>p</i> , <i>b</i> )	= <i>x</i> · <i>read</i> ( <i>s</i> , <i>p</i> · <i>x</i> , <i>b</i> )
<i>read</i> ( $\epsilon$ , <i>p</i> , <i>b</i> )	= $\epsilon$

```
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

Figure 3: Simplified ES5 Grammar

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

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`.

#### (some better examples than `def` and `var`?)

Along with the template-based rule macros, `sweet.js` also provides the more powerful case macros. Instead of a template, the body of a case macro is JavaScript that is run when the macro is invoked.

```
macro log {
  case { _ ($msg) } => {
    // need a good example here...
  }
}
```

## 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 undelimited 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 undelimited 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* [8]. 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);
```

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 [8]. 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 operators<sup>4</sup>. 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:

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

<sup>4</sup>Sweet.js will also provide support for definable operators though at the time of this writing they have not yet been fully implemented.

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:

```
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 unparseable 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.

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

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:

```
macro id { rule { $x } => { $x } }
macro m {
  rule infix { $first $second | } => {
    $first + $second
  }
}
foo("sweet")
id 100 m
// expands to:
// foo
// ("sweet") + 100
```

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.

**todo: talk about honu-style custom pattern classes? requires explaining the pattern environment which might be too deep.**

## 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 [2, 5] of tracking the lexical context in each syntax object.

**todo: do we need examples of hygiene?**

### 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 off a `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. Implementation

Sweet.js is written in JavaScript and runs in the major JS environments (i.e. the browser and node.js). This is in contrast to Honu which translates its code to Racket code and reuses the hygienic expansion machinery already built in Racket. While this simplifies the implementation of Honu it also requires an installation of Racket which in some cases is not feasible (e.g. `sweet.js` is able to run in mobile device browsers).

## 7. Related Work

Earlier treatments of macro-by-example in [7].

- Scheme/Racket
- Honu
- Template Haskell
- Nemerle [10]
- Scala
- Closure

## References

- [1] 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.
- [2] M Flatt and R Culpepper. Macros that Work Together. *Journal of Functional Programming*, 2012.
- [3] JK Foderaro, KL Sklower, and K Layer. *The FRANZ Lisp Manual*. 1983.
- [4] S P Harbison and Jr. Steele, G. L. *C: A Reference Manual*. Prentice-Hall, 1984.
- [5] R Hieb, RK Dybvig, and C Bruggeman. Syntactic abstraction in scheme. *Lisp and symbolic computation*, 5(4):295–326, 1992.

- [6] 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.
- [7] 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.
- [8] Jon Rafkind. *Syntactic extension for languages with implicitly delimited and infix syntax*. PhD thesis, 2013.
- [9] 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.
- [10] K Skalski, M Moskal, and P Olszta. Meta-programming in Nemerle. *Proceedings Generative Programming and Component Engineering*, 2004.

## A. Read Proof

**todo: define RegexPrefix/DivPrefix**

**Theorem 2** (Parse Equivalence for Program).

$\forall s.$

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

*Proof.* For the left-to-right direction, there are two production rules for  $\text{Program}_e$ .

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

A similar argument holds for the right-to-left direction.  $\square$

**Lemma 1** (Parse Equivalence for SourceElements).

$\forall s.$

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

*Proof.* For the left-to-right direction there are two production rules for  $\text{SourceElements}_e$ .

- $s \in \text{SourceElement}_e$ . This holds by Lemma 2.
- $s \in \text{SourceElements}_e \text{ SourceElement}_{e'}$ . 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 12,  $\text{read}(s', \epsilon, \text{false}) \in \text{RegexPrefix}$ ). Thus  $t \in \text{SourceElements}'_{e \cdot e'}$ .

The argument is similar for the right-to-left direction.  $\square$

**Lemma 2** (Parse Equivalence for SourceElement).

$\forall s, p \in \text{RegexPrefix}$ .

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

*Proof.* For the left-to-right direction there are two production rules for  $\text{SourceElement}_e$ .

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

The argument is similar for the right-to-left direction.  $\square$

**Lemma 3** (Parse Equivalence for FunctionDecl).

$\forall s, p, b.$

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

*Proof.* For the left-to-right direction there is one production rule for  $\text{FunctionDecl}_e$ :

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

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

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

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 right-to-left direction.  $\square$

**Lemma 4** (Parse Equivalence for Stmt).

$\forall s, p \in \text{RegexPrefix}$ .

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

*Proof.* For the left-to-right direction we have several production rules for  $\text{Stmt}_e$ .

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

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

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

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

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

where

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

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

- $s \in \text{return}$ . Since  $s = \text{return}$  and  $\text{read}(s, p, \text{false}) = \text{return} \in \text{Stmt}'_{\text{return}}$  we have our result directly.
- $s \in \text{return AssignExpr}_e ;$ . 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  $t \in \text{Stmt}'_{\text{return } e}$ .

- $s \in x : Stmt_e$ . We have  $s = x : : s'$  where  $s' \in Stmt_e$ . Then:

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

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

The argument for the right-to-left direction is similar.  $\square$

**Lemma 5** (Parse Equivalence for StmtList).

$\forall s, p \in \text{RegexPrefix}$ .

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

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

- $s \in Stmt_e$ . This follows by Lemma 4 since  $p \in \text{RegexPrefix}$ .
- $s \in StmtList_e Stmt_{e'}$ . So  $s = s' \cdot s''$  where  $s' \in StmtList_e$  and  $s'' \in Stmt_{e'}$ . 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  $\text{read}(s', p, \text{false}) \in StmtList'_e$  and by Lemma 4,  $\text{read}(s'', p \cdot \text{read}(s', p, \text{false}), \text{false}) \in Stmt'_{e'}$  (since by Lemma 13,  $p \cdot \text{read}(s', p, \text{false}) \in \text{RegexPrefix}$ ). Thus,  $t \in StmtList'_{e \cdot e'}$ .

The argument for the right-to-left direction is similar.  $\square$

**Lemma 6** (Parse Equivalence for AssignExpr).

$\forall s, p \in \text{RegexPrefix}$ . If  $b = \text{false}$  then  $s \neq \{ \cdot s' \}$ . If  $b = \text{true}$  then  $s$  is unconstrained.

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

*Proof.* The constraint on  $s$  when  $b$  is false is due to the lookahead check in the production  $Stmt_e ::= \text{AssignExpr}_e ;$ . Lemma 11 will make use of this constraint.

For the left-to-right direction we have two production rules for  $\text{AssignExpr}_e$ .

- $s \in \text{BinaryExpr}_e$ . This holds by Lemma 7.
- $s \in \text{CallExpr}_e = \text{AssignExpr}_{e'}$ . 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  $p \in \text{RegexPrefix}$  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  $p \cdot \text{read}(s', p, b) \cdot = \in \text{RegexPrefix}$ ) so  $t \in \text{AssignExpr}'_{e = e'}$ .

The argument for the right-to-left direction is similar.  $\square$

**Lemma 7** (Parse Equivalence for BinaryExpr).

$\forall s, p \in \text{RegexPrefix}$ . If  $b = \text{false}$  then  $s \neq \{ \cdot s' \}$ . If  $b = \text{true}$  then  $s$  is unconstrained.

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

*Proof.* For the left-to-right direction there are three productions for  $\text{BinaryExpr}_e$ .

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

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

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

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

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

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  $p \cdot \text{read}(s', p, b) \cdot + \in \text{RegexPrefix}$ ) thus  $t \in \text{BinaryExpr}'_{e + e'}$ .

The argument for the right-to-left direction is similar.  $\square$

**Lemma 8** (Parse Equivalence for CallExpr).

$\forall s, p \in \text{RegexPrefix}$ . If  $b = \text{false}$  then  $s \neq \{ \cdot s' \}$ . If  $b = \text{true}$  then  $s$  is unconstrained.

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

*Proof.* For the left-to-right direction there are two production rules for  $\text{CallExpr}_e$ .

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

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

Since  $p \in \text{RegexPrefix}$ , 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  $t \in \text{CallExpr}'_{e (e')}$ .

- $s \in \text{CallExpr}_e \cdot x$ . 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 right-to-left direction is similar.  $\square$

**Lemma 9** (Parse Equivalence for MemberExpr).

$\forall s, p \in \text{RegexPrefix}$ . If  $b = \text{false}$  then  $s \neq \{ \cdot s' \}$ . If  $b = \text{true}$  then  $s$  is unconstrained.

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

*Proof.* For the left-to-right direction there are two three production rules for  $\text{MemberExpr}_e$ .

- $s \in \text{PrimaryExpr}_e$ . This follows from Lemma 11 since  $p \in \text{RegexPrefix}$ .
- $s \in \text{FunctionExpr}_e$ . This follows from Lemma 10.
- $s \in \text{MemberExpr}_e \cdot x$ . 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 right-to-left direction is similar.  $\square$

**Lemma 10** (Parse Equivalence for FunctionExpr).

$\forall s, b, p \in \text{RegexPrefix}$ .

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

*Proof.* Since  $s = \text{function } x (x') \{s'\}$  where  $s' \in \text{SourceElements}_e$  and

$$t = \text{read}(s, p, b) = \text{function} \cdot x \cdot (x') \cdot \{t'\}$$

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

The argument for the right-to-left direction is similar.  $\square$

**Lemma 11** (Parse Equivalence for PrimaryExpr).

$\forall s, p \in \text{RegexPrefix}$ . If  $b = \text{false}$  then  $s \neq \{ \cdot s' \}$ . If  $b = \text{true}$  then  $s$  is unconstrained.

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

*Proof.* For the left-to-right direction there are several production rules for  $\text{PrimaryExpr}_e$ .

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

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

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

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

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

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

For the right-to-left direction the argument is similar.  $\square$

**Lemma 12** (SourceElement Prefix).

$\forall s, p \in \text{RegexPrefix}$ .

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

*Proof.* There are two cases:

- $s \in \text{Stmt}_e$ . This holds by Lemma 13.
- $s \in \text{FunctionDecl}_e$ . Follows directly.  $\square$

**Lemma 13** (Stmt Prefix).

$\forall s, p \in \text{RegexPrefix}$ .

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

*Proof.* We have several cases:

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