

DRAFT: 2020-04-08 22:22:09-04:00
YOCTO-CFA

by
Leandro Facchinetti

A dissertation submitted to Johns Hopkins University
in conformity with the requirements for the degree of Doctor of Philosophy

Baltimore, Maryland
August 2020

Contents

Contents	ii
1 Developing an Analyzer	1
1.1 The Analyzed Language: Yocto-JavaScript	1
1.1.1 Values in Yocto-JavaScript	2
1.1.2 Operations in Yocto-JavaScript	3
1.1.3 The Computational Power of Yocto-JavaScript	4
1.1.4 A Formal Grammar for Yocto-JavaScript	5
1.2 The Analyzer Language: TypeScript	6
1.3 Step 0: Substitution-Based Interpreter	6
1.3.1 Architecture	7
1.3.2 Data Structures to Represent Yocto-JavaScript Programs . . .	9
1.3.3 An Expression That Already Is a Value	11
1.3.4 A Call Involving Immediate Functions	13
1.3.5 Substitution in Function Definitions	15
1.3.6 Name Mismatch	16
1.3.7 Name Reuse	17
1.3.8 Substitution in Function Calls	18
1.3.9 An Argument That Is Not Immediate	18
1.3.10 A Function That Is Not Immediate	20

1.3.11	Continuing to Run After a Function Call	21
1.3.12	A Reference to an Undefined Variable	21
1.3.13	The Entire Runner	22
1.3.14	An Operational Semantics for the Interpreter	23
1.3.15	Parser	24
1.3.16	Stringifier	27
1.3.17	Programs That Do Not Terminate	28
1.4	Step 1: Environment-Based Interpreter	31
1.4.1	Environments and Closures	31
1.4.2	Adding an Environment to the Runner	33
1.4.3	A Function Definition	34
1.4.4	A Function Call	35
1.4.5	Name Reuse	36
1.4.6	A Variable Reference	36
	Bibliography	38

Chapter 1

Developing an Analyzer

1.1 The Analyzed Language: Yocto-JavaScript

Our first decision when developing an analyzer is which language it should analyze. In this dissertation we are interested in analysis techniques for higher-order functions, a feature which is supported by most languages, including JavaScript, Java, Python, Ruby, and so forth.

From all these options, we would like to choose JavaScript because it is the most popular language among programmers [24, 16], but JavaScript has many features besides higher-order functions that would complicate our analyzer, so we support only a *subset* of JavaScript features that are related to higher-order functions, resulting in a language that we call *Yocto-JavaScript* ($\text{JavaScript} \times 10^{-24}$). By design, every Yocto-JavaScript program is also a JavaScript program, but the converse does not hold.

Advanced

On the surface the choice of analyzed language is important because it determines how difficult the analyzer is to develop, but the analyzed language may also influence the analyzer's precision and running time. For example, there

is an analysis technique called k -CFA [27] that may be slower when applied to a language with higher-order functions than when applied to a language with objects, because the algorithmic complexity of the former is exponential and of the latter is polynomial [19].

Technical Terms

Yocto-JavaScript is a representation of something called the λ -calculus [28, § 6].

1.1.1 Values in Yocto-JavaScript

JavaScript has many kinds of values: strings (for example, "Leandro"), numbers (for example, 29), arrays (for example, ["Leandro", 29]), objects (for example, { name: "Leandro", age: 29 }), and so forth. From all these kinds of values, Yocto-JavaScript supports only one: functions.

An Yocto-JavaScript function is written as `<parameter> => <body>`, for example, `x => x`, in which the `<parameter>` is called `x` and the `<body>` is a reference to the variable `x` (see § 1.1.2 for more on variable references). An Yocto-JavaScript function must have exactly one parameter. Because an Yocto-JavaScript function is a value, it may be passed as argument in a function call or returned as the result of a function call (see § 1.1.2 for more on function calls).

Technical Terms

The notation we use for writing functions is something called *arrow function expressions* [20]. The function given as example is called the *identity* function. The ability of acting as values is what characterizes these functions as *higher-order*.

1.1.2 Operations in Yocto-JavaScript

JavaScript has many operations: strings may have its characters accessed (for example, `"Leandro"[2]`, which results in `"a"`), numbers may be added together (for example, `29 + 1`, which results in `30`), and so forth. From all these operations, Yocto-JavaScript supports only two: functions may be called and variables may be referenced.

A function call is written as `<function>(<argument>)`, for example, `f(a)`, in which the `<function>` is a hypothetical function `f` and the `<argument>` is a hypothetical argument `a`. An Yocto-JavaScript function call must have exactly one argument (because an Yocto-JavaScript function must have exactly one parameter; see § 1.1.1). A variable reference is written as a bare identifier, for example, `x`.

The following is a complete Yocto-JavaScript program that exemplifies all the supported operations:

```
(x => x)(y => y)
```

This program is a function call in which the `<function>` is `x => x` and the `<argument>` is `y => y`. When called, an Yocto-JavaScript function returns the result of computing its `<body>` and the `<body>` of `x => x` is a reference to the variable `x`, so `x => x` is a function that returns its argument unchanged and the final result of the example above is `y => y`.

In general, all kinds of Yocto-JavaScript expressions (function definitions, function calls, and variable references) may appear in the `<body>` of a function definition, or as the `<function>` or `<argument>` of a call; for example, in the program `(f(a))(b)` the function call `f(a)` appears as the `<function>` of a call.

We use parentheses to resolve ambiguities on where function definitions start and end, and in which order operations are computed. For example, given hypothetical functions `f`, `g`, and `h`, in `(f(g))(h)` the call `f(g)` happens first and the

result is a function that is called with `h`, and in `f(g(h))` the call `g(h)` happens first and the result is passed as argument to `f`. If there are no parentheses, then nested function definitions are read right-to-left and a sequence of function calls are read left-to-right; for example, `x => y => x` is equivalent to `x => (y => x)` and `f(a)(b)` is equivalent to `(f(a))(b)`.

Technical Terms

The order in which operations are computed is something called their *precedence*, and operations that happen first are said to have *higher precedence*. Because of the order in which they are read, function definitions are said to be *right-associative* and function calls are said to be *left-associative*.

Advanced

1.1.3 The Computational Power of Yocto-JavaScript

Yocto-JavaScript has only a few features, which makes it the ideal language for discussing the analysis of higher-order functions, but is it *too* simple? In other words, in the process of pairing down JavaScript to define Yocto-JavaScript, have we removed features that make the language incapable of some computations? Perhaps surprisingly, the answer is negative: Yocto-JavaScript is equivalent to JavaScript (and Java, Python, Ruby, and so forth) in the sense that, with some effort, any program in any one of these languages may be translated into an equivalent program in any other of these languages [28, § 6].

As an example of how to carry out this translation, consider a JavaScript function of two parameters: `(x, y) => x`. This function is not supported by Yocto-JavaScript because it does not have exactly one parameter (see § 1.1.1), but we may encode it as a function that receives the first parameter and returns another function that receives the second parameter: `x => (y => x)`.

Similarly, we may encode a call with multiple arguments as a sequence of calls that passes one argument at a time; for example, $f(a, b)$ may be encoded as $(f(a))(b)$.

Technical Terms

All the languages we are considering are said to be equivalent in terms of *computational power*: they are all *Turing complete* [28, § 7]. The translation technique for functions with multiple arguments is called *currying* [28, page 163].

For our goal of exploring analysis techniques, we are concerned only with computational power, but it is worth noting that programmers are more interested in other language properties: Does the language promote writing programs of higher quality? (It most probably does not [10].) Does the language improve productivity? Does the language work well for the domain of the problem? (For example, we would probably write an operating system in C and a web application in JavaScript, not the other way around.) Is the language more expressive than others? (Perhaps surprisingly, it is possible to make formal arguments about expressiveness without resorting to personal preference and anecdotal evidence [13].) Despite having the same computational power as other languages, Yocto-JavaScript fares badly in these other aspects: it is remarkably unproductive and inexpressive.

1.1.4 A Formal Grammar for Yocto-JavaScript

The description of Yocto-JavaScript given so far has been informal; the following is a grammar in *Backus–Naur Form* (BNF) [18] [9, § 4.2] that formalizes it:

$e ::= x \Rightarrow e \mid e(e) \mid x$	Expressions
$x ::= \text{<A Valid JavaScript Identifier>}$	Variables

1.2 The Analyzer Language: TypeScript

After choosing our analyzed language (Yocto-JavaScript; see § 1.1), we must decide in which language to develop the analyzer itself. Despite our analyzed language being based on JavaScript, we may choose to develop the analyzer in any language (for example, JavaScript, Java, Python, Ruby, and so forth), because the analyzer treats the analyzed program as data. Still, from all these options, JavaScript does offer some advantages: it is the most popular [24, 16], and it includes convenient tools to manipulate JavaScript programs (and therefore Yocto-JavaScript programs as well; see § 1.3.15 and § 1.3.16). But JavaScript lacks a way to express the *types* of data structures, functions, and so forth, which we will need (for example, see § 1.3.2), so we choose to implement our analyzer in a JavaScript extension with support for types called *TypeScript* [7, 29, 11].

1.3 Step 0: Substitution-Based Interpreter

Having chosen the analyzed language (Yocto-JavaScript; see § 1.1) and the language in which to develop the analyzer itself (TypeScript; see § 1.2), we are ready to start the series of Steps in the development of the analyzer. The first Step is an interpreter that executes Yocto-JavaScript programs and produces the same outputs that would be produced by a regular JavaScript interpreter. This is a good starting point for two reasons: first, this interpreter is the basis upon which we will build the analyzer; and second, the outputs of this interpreter are the ground truth against which we will validate the outputs of the analyzer.

1.3.1 Architecture

Our interpreter is defined as a function called `evaluate()`, which receives as parameter an Yocto-JavaScript program represented as a string and returns the result of running it.

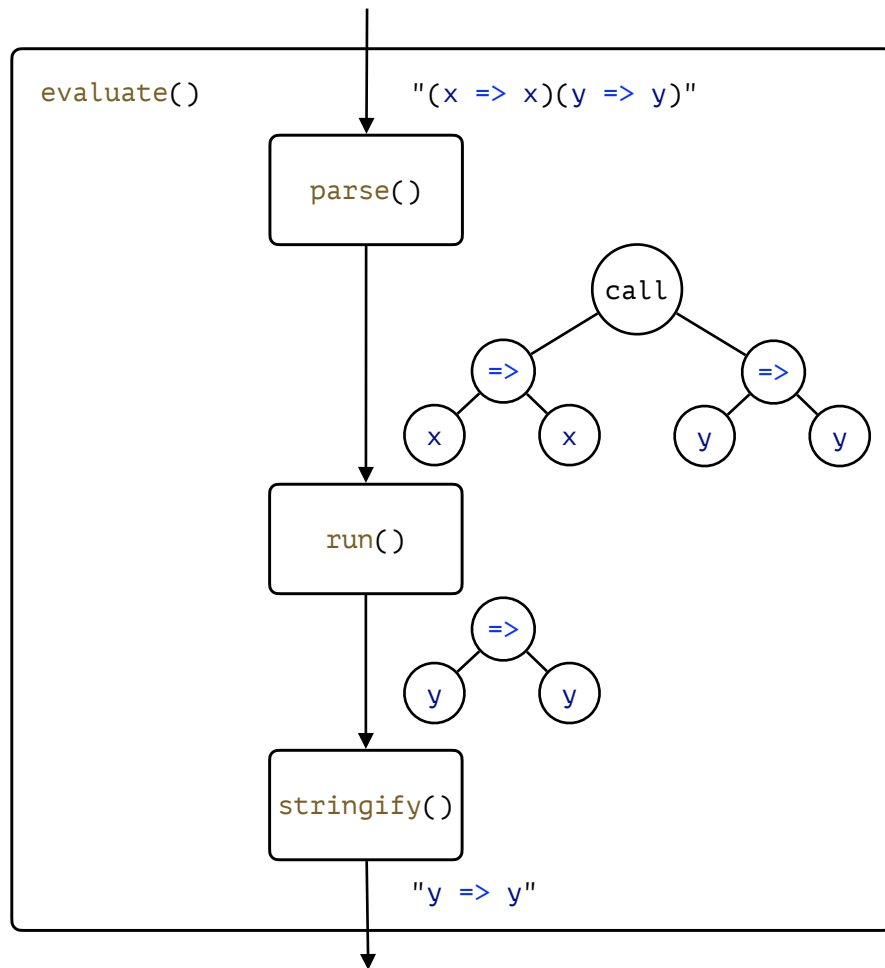
Advanced

The `evaluate()` function is named after a native JavaScript function called `eval()` [22], which is similar to `evaluate()` but for JavaScript programs instead of Yocto-JavaScript.

The following are two examples of how we will be able to use `evaluate()` by the end of Step 0 (the `>` represents the console):

```
> evaluate("x => x")
"x => x"
> evaluate("(x => x)(y => y)")
"y => y"
```

The implementation of `evaluate()` is separated into three parts called `parse()`, `run()`, and `stringify()`:



```
export function evaluate(input: string): string {  
  return stringify(run(parse(input)));  
}
```

The `parse()` function prepares the `input` for interpretation, converting it from a string into more convenient data structures (see § 1.3.2 for more on these data structures). The `run()` function is responsible for the interpretation itself. The `stringify()` function converts the outputs of `run()` into a human-readable format. In the following sections (§ 1.3.2–§ 1.3.14) we address the implementation of `run()`, deferring `parse()` to § 1.3.15 and `stringify()` to § 1.3.16.

In later Steps the implementations of `run()` and `stringify()` will change, but the architecture and therefore the implementations of `evaluate()` and `parse()` will remain the same.

1.3.2 Data Structures to Represent Yocto-JavaScript Programs

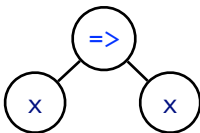
The `evaluate()` function receives an Yocto-JavaScript program represented as a string (see § 1.3.1), which is convenient for humans to write and read, but inconvenient for `run()` to manipulate directly, because `run()` is concerned with the *structure* of the program instead of the *text*: from `run()`'s perspective it does not matter, for example, whether a function is written as `x => x` or as `x=>x`. So before `run()` starts interpreting the program, `parse()` transforms it from a string into more convenient data structures (see § 1.3.15 for `parse()`'s implementation).

Technical Terms

The process of converting a program represented as a string into more convenient data structures is known as *parsing*, and the data structures are called the *Abstract Syntax Tree* (AST) of the program [9, § 4].

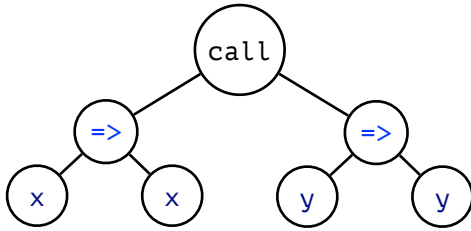
The following are two examples of Yocto-JavaScript programs followed by the data structures used to represent them, first in a high-level graphical representation and then in an equivalent low-level textual representation:

```
> parse("x => x")
```



```
{
  "type": "ArrowFunctionExpression",
  "params": [
    {
      "type": "Identifier",
      "name": "x"
    }
  ],
  "body": {
    "type": "Identifier",
    "name": "x"
  }
}
```

```
> parse("(x => x)(y => y)")
```



```
{
  "type": "CallExpression",
  "callee": {
    "type": "ArrowFunctionExpression",
    "params": [
      {
        "type": "Identifier",
        "name": "x"
      }
    ],
    "body": {
      "type": "Identifier",
      "name": "x"
    }
  },
  "arguments": [
    {
      "type": "ArrowFunctionExpression",
      "params": [
        {
          "type": "Identifier",
          "name": "y"
        }
      ],
      "body": {
        "type": "Identifier",
        "name": "y"
      }
    }
  ]
}
```

We choose to represent Yocto-JavaScript programs with the data structures above because they follow a specification called ESTree [5], and by adhering to this specification we may reuse tools from the JavaScript ecosystem (see § 1.3.15 and

§ 1.3.16).

In general, the data structures used to represent Yocto-JavaScript programs are of the following types (written as TypeScript types adapted from the ESTree types [8] to include only the features supported by Yocto-JavaScript):

```
type Expression = ArrowFunctionExpression | CallExpression | Identifier;

type ArrowFunctionExpression = {
  type: "ArrowFunctionExpression";
  params: [Identifier];
  body: Expression;
};

type CallExpression = {
  type: "CallExpression";
  callee: Expression;
  arguments: [Expression];
};

type Identifier = {
  type: "Identifier";
  name: string;
};
```

Advanced

The definitions above correspond to elements of the Yocto-JavaScript grammar (see § 1.1.4); for example, `Expression` corresponds to e .

In later Steps almost everything about the interpreter will change, but the data structures used to represent Yocto-JavaScript programs will remain the same.

1.3.3 An Expression That Already Is a Value

Example Program: `x => x`

Current Output: —

Expected Output: `x => x`

We start the definition of `run()` by considering the example above. As mentioned in § 1.3.2, the `run()` function receives as parameter an Yocto-JavaScript program represented as an `Expression`. The `run()` function is then responsible for interpreting the program and producing a value. In Yocto-JavaScript, the only kind of value is a function (see § 1.1.1), so we start the implementation of `run()` with the following (we use `throw` as a placeholder for code that has not be written yet to prevent the TypeScript compiler from signaling type errors):

```
type Value = ArrowFunctionExpression;

function run(expression: Expression): Value {
  throw new Error("NOT IMPLEMENTED YET");
}
```

The first thing that `run()` has to do is determine which type of `expression` it is given:

```
function run(expression: Expression): Value {
  switch (expression.type) {
    case "ArrowFunctionExpression":
      throw new Error("NOT IMPLEMENTED YET");
    case "CallExpression":
      throw new Error("NOT IMPLEMENTED YET");
    case "Identifier":
      throw new Error("NOT IMPLEMENTED YET");
  }
}
```

In our current example, the `expression` already is a `Value`, so it may be returned unchanged:

```
// run()
case "ArrowFunctionExpression":
  return expression;
```

1.3.4 A Call Involving Immediate Functions

Example Program: `(x => x)(y => y)`

Current Output: NOT IMPLEMENTED YET

Expected Output: `y => y`

Interpreting function calls is the main responsibility of our interpreter. There are several techniques to do this and in Step 0 we use one of the simplest: when the interpreter encounters a function call, it substitutes the variable references in the body of the function that is called with the argument that is passed. This is similar to how we reason about functions in mathematics; for example, given the function $f(x) = x + 1$, we calculate $f(29)$ by substituting the references to x in f with the argument 29: $f(29) = 29 + 1$. The implementation of this substitution technique starts in this section and will only be complete in § 1.3.8.

In the example we are considering both the function that is called `(x => x)` and the argument `(y => y)` are immediate functions, as opposed to being the result of other operations, so for now we may limit the interpreter to handle only this case:

```
// run()
case "CallExpression":
  if (
    expression.callee.type !== "ArrowFunctionExpression" ||
    expression.arguments[0].type !== "ArrowFunctionExpression"
  )
    throw new Error("NOT IMPLEMENTED YET");
  throw new Error("NOT IMPLEMENTED YET");
```

Next, we unpack the called function (using something called *destructuring assignment* [21]) and the argument:

```
// run()
case "CallExpression":
  if (
    expression.callee.type !== "ArrowFunctionExpression" ||
    expression.arguments[0].type !== "ArrowFunctionExpression"
  )
```

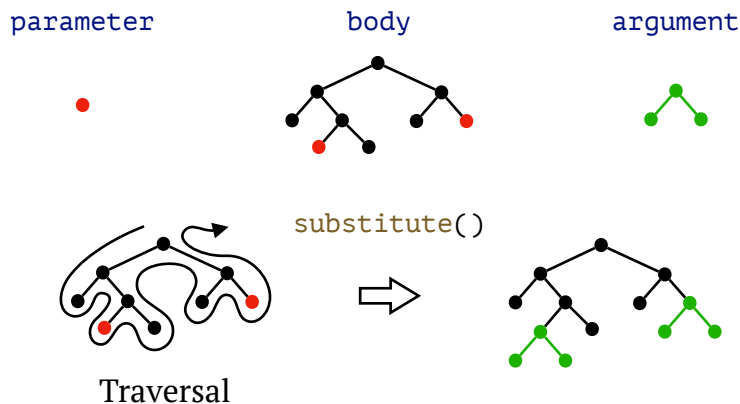


```

    throw new Error("NOT IMPLEMENTED YET");
  const {
    params: [parameter],
    body,
  } = expression.callee;
  const argument = expression.arguments[0];
  throw new Error("NOT IMPLEMENTED YET");

```

Finally, we setup an auxiliary function called `substitute()` that implements the traversal of the `body` looking for references to `parameter` and substituting them:



```

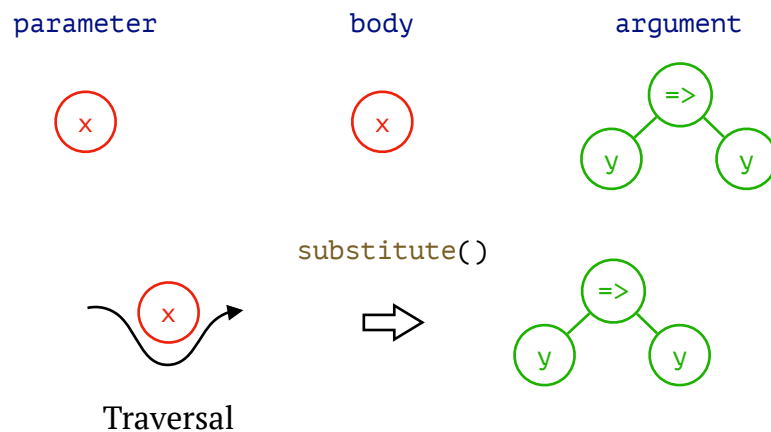
// run()
case "CallExpression":
  if (
    expression.callee.type !== "ArrowFunctionExpression" ||
    expression.arguments[0].type !== "ArrowFunctionExpression"
  )
    throw new Error("NOT IMPLEMENTED YET");
  const {
    params: [parameter],
    body,
  } = expression.callee;
  const argument = expression.arguments[0];
  const substitutedBody = substitute(body);
  if (substitutedBody.type !== "ArrowFunctionExpression")
    throw new Error("NOT IMPLEMENTED YET");
  return substitutedBody;
function substitute(expression: Expression): Expression {
  throw new Error("NOT IMPLEMENTED YET");
}

```

Similar to `run()` itself, `substitute()` starts by determining which type of `expression` is passed to it:

```
function substitute(expression: Expression): Expression {
  switch (expression.type) {
    case "ArrowFunctionExpression":
      throw new Error("NOT IMPLEMENTED YET");
    case "CallExpression":
      throw new Error("NOT IMPLEMENTED YET");
    case "Identifier":
      throw new Error("NOT IMPLEMENTED YET");
  }
}
```

In our current example the `expression` is `x`, an `Identifier`, and it must be substituted with the `argument`:



```
// substitute()
case "Identifier":
  return argument;
```

1.3.5 Substitution in Function Definitions

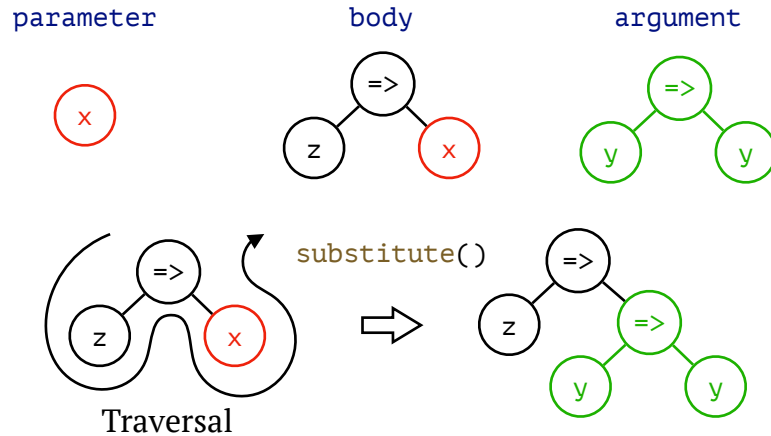
Example Program: `(x => z => x)(y => y)`

Current Output: NOT IMPLEMENTED YET

Expected Output: `z => y => y`

When `substitute()` (see § 1.3.4) starts traversing the `body` of the example above, the `expression` is an `ArrowFunctionExpression` (`z => x`), and

we want substitution to proceed deeper to find and substitute `x`, so we call `substitute()` recursively (we use a feature called *spread syntax* [23] to build an `expression` based on the existing one with a new `body`):



```
// substitute()
case "ArrowFunctionExpression":
  return {
    ...expression,
    body: substitute(expression.body),
  };
```

1.3.6 Name Mismatch

Example Program: `(x => z => z)(y => y)`

Current Output: `z => y => y`

Expected Output: `z => z`

The implementation of `substitute()` in the case of `Identifier` introduced in § 1.3.4 *always* substitutes variable references, regardless of whether they refer to the `parameter` of the called function. For example, in the program above `substitute()` is substituting the `z` even though the `parameter` is `x`. To fix this, we check whether the variable reference matches the `parameter`, and if it does not then we prevent the substitution by returning the variable reference unchanged:

```
// substitute()
case "Identifier":
  if (expression.name !== parameter.name) return expression;
  return argument;
```

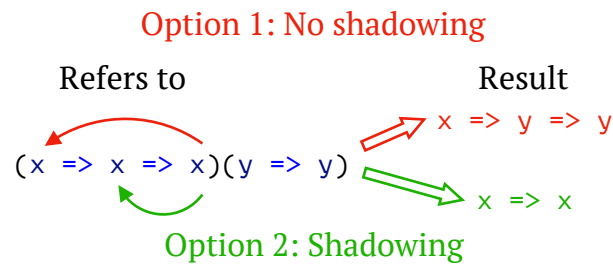
1.3.7 Name Reuse

Example Program: $(x \Rightarrow x \Rightarrow x)(y \Rightarrow y)$

Current Output: $x \Rightarrow y \Rightarrow y$

Expected Output: $x \Rightarrow x$

In the program above, there are two options for the variable reference x on the right of the second \Rightarrow : it may refer to the first (outer) x on the left of the first \Rightarrow , in which case the output of the program would be $x \Rightarrow y \Rightarrow y$; or it may refer to the second (inner) x on the left of the second \Rightarrow , in which case the output of the program would be $x \Rightarrow x$:



Currently `substitute()` is implementing Option 1, but this leads to a serious issue: we are not able to reason about the inner function $x \Rightarrow x$ independently; we must know where it appears and whether a variable called x is already defined there.

Technical Terms

We say that the problem with Option 1 is that it defeats something called *local reasoning*. We say that Option 2 exhibits a behavior called *shadowing*, and that the outer x is *shadowed* by the inner x , because there is no way to refer to the outer x from the body of the inner function.

We avoid this issue by modifying `substitute()` to implement Option 2, which is also the choice of JavaScript and every other popular programming language. We change `substitute()`'s behavior when encountering a function definition

so that if the parameter of the function definition matches the parameter that `substitute()` is looking for, then `substitute()` returns the function unchanged, preventing further substitution (there is no recursive call to `substitute()` in this case):

```
// substitute()
case "ArrowFunctionExpression":
  if (expression.params[0].name === parameter.name) return expression;
  return {
    ...expression,
    body: substitute(expression.body),
  };
```

1.3.8 Substitution in Function Calls

Example Program: $(x \Rightarrow z \Rightarrow x(x))(y \Rightarrow y)$

Current Output: NOT IMPLEMENTED YET

Expected Output: $z \Rightarrow (y \Rightarrow y)(y \Rightarrow y)$

This case is similar to § 1.3.5: all `substitute()` has to do is continue traversing the function call recursively:

```
// substitute()
case "CallExpression":
  return {
    ...expression,
    callee: substitute(expression.callee),
    arguments: [substitute(expression.arguments[0])],
  };
```

1.3.9 An Argument That Is Not Immediate

Example Program: $(x \Rightarrow z \Rightarrow x)((a \Rightarrow a)(y \Rightarrow y))$

Current Output: NOT IMPLEMENTED YET

Expected Output: $z \Rightarrow y \Rightarrow y$

In all example programs we considered so far the argument to a function call was an immediate function definition, but in general arguments may be the result of function calls themselves. We fix this by calling `run()` recursively on the argument (we also remove the check that the argument is an immediate function definition; if it is, then the recursive call to `run()` returns the immediate function unchanged; see § 1.3.3):

```
// run()
case "CallExpression":
  if (expression.callee.type !== "ArrowFunctionExpression")
    throw new Error("NOT IMPLEMENTED YET");
  const {
    params: [parameter],
    body,
  } = expression.callee;
  const argument = run(expression.arguments[0]);
  const substitutedBody = substitute(body);
  if (substitutedBody.type !== "ArrowFunctionExpression")
    throw new Error("NOT IMPLEMENTED YET");
  return substitutedBody;
  function substitute(expression: Expression): Expression {
    // ...
  }
```

Technical Terms

This technique of calling `run()` recursively to produce an immediate function for the argument characterizes the interpreter as *big-step*.

Advanced

The notion that the argument is interpreted to produce a value as soon as the function call is encountered characterizes Yocto-JavaScript as a *call-by-value* language [25]. JavaScript itself and most other popular programming languages are call-by-value as well, but there is a notable exception, Haskell, which is a *call-by-need* language. In a call-by-need language the argument is interpreted only if it is *needed*, for example, if it is used in the function po-

sition of another call (see § 1.3.10), or if it is the result of the program (see § 1.3.11). In a call-by-need language the result of the program above would be $z \Rightarrow ((a \Rightarrow a)(y \Rightarrow y))$. And there is yet another policy for when to interpret arguments called *call-by-name*: the difference between call-by-name and call-by-need is that in a call-by-name language the an argument may be computed multiple times if it is used multiple times, but in a call-by-need language an argument is guaranteed to be computed at most once.

1.3.10 A Function That Is Not Immediate

Example Program: $((z \Rightarrow z)(x \Rightarrow x))(y \Rightarrow y)$

Current Output: NOT IMPLEMENTED YET

Expected Output: $y \Rightarrow y$

This is the dual of § 1.3.9 for the called function, and the solution is the same: to call `run()` recursively (we also remove the check of whether the function is immediate):

```
// run()
case "CallExpression":
  const {
    params: [parameter],
    body,
  } = run(expression.callee);
  const argument = run(expression.arguments[0]);
  const substitutedBody = substitute(body);
  if (substitutedBody.type !== "ArrowFunctionExpression")
    throw new Error("NOT IMPLEMENTED YET");
  return substitutedBody;
function substitute(expression: Expression): Expression {
  // ...
}
```

1.3.11 Continuing to Run After a Function Call

Example Program: $(x \Rightarrow (z \Rightarrow z)(x))(y \Rightarrow y)$

Current Output: NOT IMPLEMENTED YET

Expected Output: $y \Rightarrow y$

This is similar to § 1.3.9 and § 1.3.10: the result of substitution may be not an immediate function but another call, and more work may be necessary to interpret it. We solve this with yet another recursive call to `evaluate()` (we also remove yet another check and inline the `substitutedBody` variable):

```
// run()
case "CallExpression":
  const {
    params: [parameter],
    body,
  } = run(expression.callee);
  const argument = run(expression.arguments[0]);
  return run(substitute(body));
  function substitute(expression: Expression): Expression {
    // ...
  }
```

1.3.12 A Reference to an Undefined Variable

Example Program: $(x \Rightarrow y)(y \Rightarrow y)$

Current Output: NOT IMPLEMENTED YET

Expected Output: Reference to undefined variable: `y`

The only case in which `run()` may encounter a variable reference directly is if the referenced variable is undefined, otherwise `substitute()` would have already substituted it (see § 1.3.4–§ 1.3.11). In this case, we throw an exception:

```
// run()
case "Identifier":
  throw new Error(Reference to undefined variable: ${expression.name});
```


If the reference to an undefined variable occurs in the body of a function that is not called, then we do not reach the case addressed in this section and an exception is not thrown. For example, the result of running the program `x => y` is `x => y`, not an exception. This is consistent with JavaScript's behavior.

1.3.13 The Entire Runner

The implementation of the `run()` function is complete:

```
1  type Value = ArrowFunctionExpression;
2
3  function run(expression: Expression): Value {
4      switch (expression.type) {
5          case "ArrowFunctionExpression":
6              return expression;
7          case "CallExpression":
8              const {
9                  params: [parameter],
10                 body,
11             } = run(expression.callee);
12             const argument = run(expression.arguments[0]);
13             return run(substitute(body));
14             function substitute(expression: Expression): Expression {
15                 switch (expression.type) {
16                     case "ArrowFunctionExpression":
17                         if (expression.params[0].name === parameter.name) return expression;
18                         return {
19                             ...expression,
20                             body: substitute(expression.body),
21                         };
22                     case "CallExpression":
23                         return {
24                             ...expression,
25                             callee: substitute(expression.callee),
26                             arguments: [substitute(expression.arguments[0])],
27                         };
28                     case "Identifier":
29                         if (expression.name !== parameter.name) return expression;
30                         return argument;
31                 }
32             }
33             case "Identifier":
```

```

34         throw new Error(Reference to undefined variable: ${expression.name});
35     }
36 }

```

Advanced

1.3.14 An Operational Semantics for the Interpreter

What we accomplished so far in this section is more than defining an interpreter for Yocto-JavaScript; we also defined formally the *meaning* of Yocto-JavaScript programs: an Yocto-JavaScript program means what the interpreter produces for it. The definition of the meaning of programs in a language is something called the *semantics* of the language, and there are several techniques to specify semantics; the one we are using so far is known as a *definitional interpreter* [26].

A definitional interpreter has some advantages over other techniques for specifying semantics: it is easier to understand for most programmers, and it is executable. But a definitional interpreter also has one disadvantage: to understand the meaning of an Yocto-JavaScript program we have to understand an interpreter written in TypeScript. To address this, there are other techniques for defining semantics that do not depend on other programming languages, and in this section we introduce one of them: *operational semantics* [17, 14, 15].

First, we extend the grammar from § 1.1.4 with the notion of values that is equivalent to the type `Value` (see § 1.3.13, line 1):

$$v ::= x \Rightarrow e \quad \text{Values}$$

Next, we define a *relation* $e \Rightarrow v$ using *inference rules* that are equivalent to the behavior of `run()` (see § 1.3.13, lines 3–36):

$$\frac{}{v \Rightarrow v} \quad \frac{e_f \Rightarrow x_p \Rightarrow e_b \quad e_a \Rightarrow v_a \quad e_b[x_p \setminus v_a] \Rightarrow v}{e_f(e_a) \Rightarrow v}$$

Finally, we define a *metafunction* $e[x \setminus v] = e$ that is equivalent to the behavior of `substitute()` (see § 1.3.13, lines 14–32):

$$\begin{aligned} (x \Rightarrow e)[x_p \setminus v_a] &= x \Rightarrow (e[x_p \setminus v_a]) && \text{if } x \neq x_p \\ (x_p \Rightarrow e)[x_p \setminus v_a] &= x_p \Rightarrow e \\ (e_f(e_a))[x_p \setminus v_a] &= (e_f[x_p \setminus v_a])(e_a[x_p \setminus v_a]) \\ x[x_p \setminus v_a] &= x && \text{if } x \neq x_p \\ x_p[x_p \setminus v_a] &= v_a \end{aligned}$$

1.3.15 Parser

The parser is responsible for converting an Yocto-JavaScript program written as a string into data structures that are more convenient for the runner to manipulate (see § 1.3.1 for a high-level view of the architecture and § 1.3.2 for the definition of the data structures). We choose to represent Yocto-JavaScript programs with data structures that are compatible with a specification for representing JavaScript programs called ESTree [5, 8] because it allows us to reuse tools from the JavaScript ecosystem, including a parser called Esprima [3], and the Esprima Interactive Online Demonstration [4], which shows the data structures used to represent a given program.

Our strategy to implement the Yocto-JavaScript parser is to delegate most of the work to Esprima and check that the program is using only features supported by Yocto-JavaScript. The following is the full implementation of the parser:

```
1 function parse(input: string): Expression {
2   const program = esprima.parseScript(input, {}, checkFeatures);
```

```

3  const expression = (program as any).body[0].expression as Expression;
4  return expression;
5  function checkFeatures(node: estree.Node): void {
6      switch (node.type) {
7          case "Program":
8              if (node.body.length !== 1)
9                  throw new Error(
10                     "Unsupported Yocto-JavaScript feature: Program with multiple statements"
11                 );
12             break;
13         case "ExpressionStatement":
14             break;
15         case "ArrowFunctionExpression":
16             break;
17         case "CallExpression":
18             if (node.arguments.length !== 1)
19                 throw new Error(
20                     "Unsupported Yocto-JavaScript feature: CallExpression with multiple argu
21                 );
22             break;
23         case "Identifier":
24             break;
25         default:
26             throw new Error(Unsupported Yocto-JavaScript feature: ${node.type});
27     }
28 }
29 }

```

Line 1: The parser is defined as a function called `parse()`, which receives the program `input` represented as a `string` and returns an `Expression` (see § 1.3.2).

Line 2: Call `esprima.parseScript()`, which parses the `input` as if it were a JavaScript program and produces a data structure following the ESTree specification. The `esprima.parseScript()` function also detects syntax errors, for example, in the program `x =>`, which is missing the function body. We pass as argument to `esprima.parseScript()` a function called `checkFeatures()` which is called with every fragment of data structure that represents a part of the program. The purpose of `checkFeatures()`

is to check that the program uses only the features that are supported by Yocto-JavaScript.

Line 3: Extract the single `Expression` from within the `Program` returned by `esprima.parseScript()`. The `as <something>` forms sidestep the TypeScript type checker and assert that the `expression` is of the correct type. This is safe to do because of `checkFeatures()`.

Line 5: The `checkFeatures()` function, which is passed to `esprima.parseScript()` is called with every fragment of data structure used to represent the program. These fragments are called *nodes*, because the data structure as a whole forms a *tree*, also known as the *Abstract Syntax Tree* (AST) of the program (see § 1.3.2). The `checkFeatures()` does not return anything (`void`); its purpose is only to throw an exception in case the program uses a feature that is not supported by Yocto-JavaScript.

Lines 6, 7, 13, 15, 17, 23, 25: Similar to `run()` and `substitute()` (see § 1.3.13), `checkFeatures()` starts by determining which type of `estree.Node` it is given.

Lines 8–11: Check that the `Program` contains a single statement. This prevents programs such as `x => x; y => y`.

Lines 13, 15: `ExpressionStatements` and `ArrowFunctionExpressions` are supported in Yocto-JavaScript unconditionally. We could check that the `ArrowFunctionExpression` includes only one parameter and that this parameter is a variable (as opposed to being a pattern such as `[x, y]`, for example), but this would be redundant because Esprima already calls `checkFeatures()` with other unsupported `nodes` that subsume these cases. For example, given the program `(x, y) => x`, which is a function of

multiple parameters, Esprima calls `checkFeatures()` with a `node` of type `SequenceExpression`. Similarly, given the program `([x, y]) => x`, which is a function in which the parameter is a pattern, Esprima calls `checkFeatures()` with a `node` of type `ArrayExpression`.

Lines 18–21: Check that the `CallExpression` contains a single argument. This prevents programs such as `f(a, b)`.

Line 23: `Identifiers` are supported in Yocto-JavaScript unconditionally. An `Identifier` may be an expression or the parameter of an `ArrowFunctionExpression`.

Line 26: All other types of `estree.Node` are not supported by Yocto-JavaScript. This includes programs such as `29` (`estree.Literal`) and `const f = x => x` (`estree.VariableDeclarator`).

In later Steps almost everything about the interpreter will change, but the parser will remain the same.

1.3.16 Stringifier

The stringifier transforms a `Value` produced by `run()` into a human-readable format (see § 1.3.1 for a high-level view of the architecture). Similar to what happened in the parser (see § 1.3.15), we may implement the stringifier by reusing existing tools from the JavaScript ecosystem, because we are representing Yocto-JavaScript programs and values with data structures that follow the ESTree specification. In particular, we use a library called Escodegen [2] to generate a string representation of an ESTree data structure, and a library called Prettier [6] to format that string. The following is the full implementation of the stringifier:

```
1 function stringify(value: Value): string {  
2   return prettier
```

```

3     .format(escodegen.generate(value), {
4         parser: "babel",
5         semi: false,
6         arrowParens: "avoid",
7     })
8     .trim();
9 }

```

Line 4: Prettier needs to parse and regenerate the string representing the value, in an architecture similar to that of the interpreter (see § 1.3.1), and it may use different parsers. We choose Babel [1], which is the default parser for JavaScript (Prettier also supports formatting other languages, for example, TypeScript and Markdown).

Line 5: Instruct Prettier not to produce semicolons at the end of the line, for example, `x => y` instead of `x => y;`.

Line 6: Instruct Prettier not to wrap parameters in parentheses, for example, `x => y` instead of `(x) => y`.

Line 8: Remove the newline at the end of the output produced by Prettier.

1.3.17 Programs That Do Not Terminate

Example Program: `(f => f(f))(f => f(f))`

Current Output: DOES NOT TERMINATE

Expected Output: DOES NOT TERMINATE

Example Program: `(f => (f(f))(f(f)))(f => (f(f))(f(f)))`

Current Output: DOES NOT TERMINATE

Expected Output: DOES NOT TERMINATE

Technical Terms

The first example program $((f \Rightarrow f(f))(f \Rightarrow f(f)))$ is also known as the Ω -combinator. The function $f \Rightarrow f(f)$ that is part of this program is also known as the U -combinator ($\Omega = (U)(U)$).

Yocto-JavaScript may express any program that a computer may run (see § 1.1.3), including some programs that do not terminate. For example, consider the first example program above, which is the shortest non-terminating program in Yocto-JavaScript. The following is a trace of the first call to `run()` (see § 1.3.13):

Line

```
3      expression = (f => f(f))(f => f(f))
9      parameter  = f
10     body       = f(f)
12     argument   = f => f(f)
13 substitute(body) = (f => f(f))(f => f(f))
```

The result of substitution is the same as the initial expression, so when it is passed as argument to the recursive call to `run()` in line 13, it causes `run()` to go into an infinite loop.

There are also programs for which interpretation does not terminate that never produce the same expression twice. For example, consider the second example program, which is a variation of the first in which every variable reference to `f` has been replaced with `f(f)`. The following is a trace of the first call to `run()`:

Line		$F = f \Rightarrow (f(f))(f(f))$
3	<code>expression</code>	<code>= F(F)</code>
9	<code>parameter</code>	<code>= f</code>
10	<code>body</code>	<code>= (f(f))(f(f))</code>
12	<code>argument</code>	<code>= F</code>
13	<code>substitute(body)</code>	<code>= (F(F))(F(F))</code>

The result of substitution is an expression that contains the initial expression in addition to some extra work (the second $F(F)$), so when it is passed as argument to the recursive call to `run()` in line 13, it causes `run()` to go into an infinite loop. Unlike what happened in the first example, when interpreting this program the expressions that are passed to `run()` never repeat themselves:

```
(F(F))
(F(F))(F(F))
(F(F))(F(F))(F(F))
(F(F))(F(F))(F(F))(F(F))
//      ⋮
```

Non-termination is what we expect from an interpreter, but not from an analyzer, and as the second example demonstrates, preventing non-termination is not as simple as checking whether `run()` is being called with the same expression multiple times. As we move forward from an interpreter to an analyzer in the next Steps one of the main issues we address is termination: even if it takes a long time, an analyzer must eventually finish its work regardless of the program it is given.

Advanced

Detecting non-termination in an interpreter is a problem that cannot be solved, regardless of the sophistication of the detector and the computational power available to it. The problem, which is known as the *halting problem*, is

said to be *uncomputable* [28, § 8], and is a direct consequence of the Turing completeness of Yocto-JavaScript (see § 1.1.3).

1.4 Step 1: Environment-Based Interpreter

As mentioned in § 1.3.17, the interpreter from Step 0 may not terminate for some programs, and preventing non-termination is one of the main issues we must address when developing an analyzer. The source of non-termination in Step 0 is substitution, which may produce infinitely many new expressions. In Step 1, we modify the interpreter so that it does not perform substitution, and as a consequence the only expressions it considers are the finitely many expressions found in the input program. The interpreter in Step 1 may still not terminate, but that is due to other sources of non-termination that will be addressed in subsequent Steps.

1.4.1 Environments and Closures

In Step 0, when a function was applied we substituted all references to the parameter in the body with the argument. In Step 1, we do not perform the substitution, but keep the body intact along with a map from the parameter to the argument, for example:

	Program	Evaluates to
Step 0	$(x \Rightarrow z \Rightarrow x)(y \Rightarrow y)$	$z \Rightarrow (y \Rightarrow y)$
Step 1	$(x \Rightarrow z \Rightarrow x)(y \Rightarrow y)$	$\langle (z \Rightarrow x), [x \mapsto (y \Rightarrow y)] \rangle$

Step 0 produces the expression $z \Rightarrow (y \Rightarrow y)$ (see § 1.3.5), which does not appear in the original program, but Step 1 produces only expressions that exist in the original program, for example, $z \Rightarrow x$ and $y \Rightarrow y$.

Technical Terms

A map from parameters to arguments (for example, $[x \mapsto (y \Rightarrow y)]$) is something called an *environment*. A function definition along with an environment (for example, $\langle (z \Rightarrow x), [x \mapsto (y \Rightarrow y)] \rangle$) is something called a *closure*.

In Step 1 we have a different notion of *value*: while in Step 0 the interpreter produced a *function*, now it produces a *closure*. It is possible to use the closure produced in Step 1 to recreate the function that would have been produced in Step 0 by substituting the variable references in the function body with the corresponding values from the environment found in the closure. We may do this to check that the outputs of the interpreters are equivalent. But in Step 1 we do not perform substitution in the regular course of interpretation; we add more mappings to the environment and look up variable references in the environment.

Alternative Argument

Another way to reason about an environment-based interpreter is that it is a substitution-based interpreter in which substitutions are delayed until they are needed.

The following are the data structures for Step 1:

```
type Value = Closure;

type Closure = {
  function: ArrowFunctionExpression;
  environment: Environment;
};

type Environment = MapDeepEqual<Identifier["name"], Value>;
```

The `MapDeepEqual` data structure is provided by a JavaScript package developed by the author called Collections Deep Equal [12]. A `MapDeepEqual` is similar to a `Map`, except that the keys are compared by value, not by reference, for example:

```

> new Map([[{ name: "Leandro" }, 29]]).get({ name: "Leandro" });
undefined
> new MapDeepEqual([[{ name: "Leandro" }, 29]]).get({ name: "Leandro" });
29

```

The occurrences of `{ name: "Leandro" }` are objects with the same keys and values, but they are not the same object.

1.4.2 Adding an Environment to the Runner

The runner needs to maintain an environment, so we introduce an auxiliary function called `step()` that receives an `environment` as an extra parameter:

```

1 function run(expression: Expression): Value {
2   return step(expression, new MapDeepEqual());
3   function step(expression: Expression, environment: Environment): Value {
4     switch (expression.type) {
5       case "ArrowFunctionExpression":
6         return expression;
7       case "CallExpression":
8         const {
9           params: [parameter],
10          body,
11        } = step(expression.callee, environment);
12        const argument = step(expression.arguments[0], environment);
13        return step(substitute(body), environment);
14        function substitute(expression: Expression): Expression {
15          switch (expression.type) {
16            case "ArrowFunctionExpression":
17              if (expression.params[0].name === parameter.name)
18                return expression;
19              return {
20                ...expression,
21                body: substitute(expression.body),
22              };
23            case "CallExpression":
24              return {
25                ...expression,
26                callee: substitute(expression.callee),
27                arguments: [substitute(expression.arguments[0])],
28              };
29            case "Identifier":

```

```

30         if (expression.name !== parameter.name) return expression;
31         return argument;
32     }
33 }
34 case "Identifier":
35     throw new Error(Reference to undefined variable: ${expression.name});
36 }
37 }
38 }

```

Line 3: We define the `step()` auxiliary function that receives an `environment` as an extra parameter.

Line 2: The `environment` starts empty.

Lines 11–13: The recursive calls to `run()` are changed to recursive calls to `step()` and the `environment` is propagated.

The listing above does not compile yet because we are not producing closures. In the following sections we fix this by considering how to manage the `environment` for each type of `expression`.

1.4.3 A Function Definition

Example Program: `x => x`

Current Output: —

Expected Output: `⟨(x => x), []⟩`

When the interpreter encounters a function definition, it captures the current `environment` in a closure:

```

// step()
case "ArrowFunctionExpression":
    return { function: expression, environment };

```

1.4.4 A Function Call

Example Program: $(x \Rightarrow z \Rightarrow x)(y \Rightarrow y)$

Current Output: —

Expected Output: $\langle (z \Rightarrow x), [x \mapsto (y \Rightarrow y)] \rangle$

First, we remove `substitute()`, which is the goal of Step 1:

```
// step()
case "CallExpression":
  const {
    params: [parameter],
    body,
  } = step(expression.callee, environment);
  const argument = step(expression.arguments[0], environment);
  return step(body, environment);
```

Next, we fix the pattern that matches the result of the interpretation of the called function to take in account the closure:

```
// step()
case "CallExpression":
  const {
    function: {
      params: [parameter],
      body,
    },
    environment: functionEnvironment,
  } = step(expression.callee, environment);
  const argument = step(expression.arguments[0], environment);
  return step(body, environment);
```

Finally, we modify the recursive call to `step()` that interprets the function body so that it receives a new augmented `environment` including a mapping from the `parameter` (for example, `x`) to the `argument` (for example, $\langle (y \Rightarrow y), [] \rangle$):

```
1 // step()
2 case "CallExpression":
3   const {
4     function: {
5       params: [parameter],
```

```

6     body,
7   },
8   environment: functionEnvironment,
9 } = step(expression.callee, environment);
10 const argument = step(expression.arguments[0], environment);
11 return step(
12   body,
13   new MapDeepEqual(environment).set(parameter.name, argument)
14 );

```

1.4.5 Name Reuse

Example Program: $(x \Rightarrow x \Rightarrow z \Rightarrow x)(a \Rightarrow a)(y \Rightarrow y)$

Current Output: $\langle (z \Rightarrow x), [x \mapsto (y \Rightarrow y)] \rangle$

Expected Output: $\langle (z \Rightarrow x), [x \mapsto (y \Rightarrow y)] \rangle$

If a name is reused (for example, x in the example program above), then the second time it is encountered by `step()` it is overwritten in the `environment` (see the call to `set()` in line 13 of § 1.4.4, which overwrites an existing map key). This causes the variable reference to x to refer to the second (inner) x , which is the expected behavior (it is what we called Option 2 in § 1.3.7).

1.4.6 A Variable Reference

Example Program: $(x \Rightarrow x)(y \Rightarrow y)$

Current Output: —

Expected Output: $\langle (y \Rightarrow y), [] \rangle$

Example Program: $(x \Rightarrow y)(y \Rightarrow y)$

Current Output: —

Expected Output: Reference to undefined variable: y

Example Program: $x \Rightarrow y$

Current Output: —

Expected Output: $\langle (x \Rightarrow y), [] \rangle$

When we encounter a variable reference, we look it up in the current environment:

```
// step()
case "Identifier":
  const value = environment.get(expression.name);
  if (value === undefined)
    throw new Error(
      Reference to undefined variable: ${expression.name}
    );
  return value;
```


Bibliography

- [1] Babel. <https://babeljs.io>. Accessed 2020-04-06.
- [2] Escodegen. <https://github.com/eslint/escodegen>. Accessed 2020-02-18.
- [3] Esprima. <https://esprima.org>. Accessed 2020-01-21.
- [4] Esprima interactive online demonstration. <https://esprima.org/demo/parse.html>. Accessed 2020-01-21.
- [5] The estree spec. <https://github.com/estree/estree>. Accessed 2020-01-21.
- [6] Prettier. <https://prettier.io>. Accessed 2020-02-18.
- [7] Typescript homepage. <https://www.typescriptlang.org>. Accessed 2020-01-17.
- [8] Typescript types for the estree spec. <https://github.com/DefinitelyTyped/DefinitelyTyped/blob/1911a1fbbe30af03f0f38a915c4bf0620d251fc6/types/estree/index.d.ts>. Accessed 2020-01-27.
- [9] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., USA, 2006.

- [10] Emery D. Berger, Celeste Hollenbeck, Petr Maj, Olga Vitek, and Jan Vitek. On the impact of programming languages on code quality: A reproduction study. *ACM Trans. Program. Lang. Syst.*, 41(4), October 2019.
- [11] Gavin Bierman, Martín Abadi, and Mads Torgersen. Understanding typescript. In Richard Jones, editor, *ECOOP 2014 – Object-Oriented Programming*, pages 257–281, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg.
- [12] Leandro Facchinetti. Collections deep equal. <https://github.com/leafac/collections-deep-equal>. Accessed 2020-04-01.
- [13] Matthias Felleisen. On the expressive power of programming languages. *Science of Computer Programming*, 17(1):35 – 75, 1991.
- [14] Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. *Semantics Engineering with PLT Redex*. The MIT Press, 1st edition, 2009.
- [15] Mike Grant, Zachary Palmer, and Scott Smith. *Principles of Programming Languages*. 2020.
- [16] JetBrains. The state of developer ecosystem 2019. <https://www.jetbrains.com/lp/devecosystem-2019/>. Accessed 2020-01-14.
- [17] Gilles Kahn. Natural semantics. In *Annual symposium on theoretical aspects of computer science*, pages 22–39. Springer, 1987.
- [18] Matthew Might. The language of languages. <http://matt.might.net/articles/grammars-bnf-ebnf/>. Accessed 2020-01-17.
- [19] Matthew Might, Yannis Smaragdakis, and David Van Horn. Resolving and exploiting the k-cfa paradox: Illuminating functional vs. object-oriented program analysis. In *Proceedings of the 31st ACM SIGPLAN Conference on Pro-*

programming Language Design and Implementation, PLDI '10, page 305–315, New York, NY, USA, 2010. Association for Computing Machinery.

- [20] Mozilla. Arrow function expressions. https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Functions/Arrow_functions. Accessed 2020-01-16.
- [21] Mozilla. Destructuring assignment. https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Operators/Destructuring_assignment. Accessed 2020-01-27.
- [22] Mozilla. eval(). https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Global_Objects/eval. Accessed 2020-02-13.
- [23] Mozilla. Spread syntax. https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Operators/Spread_syntax. Accessed 2020-02-03.
- [24] Stack Overflow. Developer survey results 2019. <https://insights.stackoverflow.com/survey/2019>. Accessed 2020-01-14.
- [25] G.D. Plotkin. Call-by-name, call-by-value and the λ -calculus. *Theoretical Computer Science*, 1(2):125 – 159, 1975.
- [26] John C. Reynolds. Definitional interpreters for higher-order programming languages. In *Proceedings of the ACM Annual Conference - Volume 2*, ACM '72, page 717–740, New York, NY, USA, 1972. Association for Computing Machinery.
- [27] Olin Shivers. *Control-flow analysis of higher-order languages*. PhD thesis, PhD thesis, Carnegie Mellon University, 1991.

- [28] Tom Stuart. *Understanding Computation: From Simple Machines to Impossible Programs*. O'Reilly Media, 1 edition, 2013.
- [29] Basarat Ali Syed. Typescript deep dive. <https://basarat.gitbook.io/typescript/>. Accessed 2020-01-17.