

Token Expression Determinizer

...

“TED”

Team and Responsibilities



Konstantin Itskov
Manager



Theodore Ahlfeld
Coding Guru



Matthew Haigh
Testing



Gideon Mendels
Architect

Overview

The TED programming language is a web-parsing language designed to simplify web scraping and serve as a bridge between complex high-level web-scraping languages like Javascript and imperative programming languages like C.

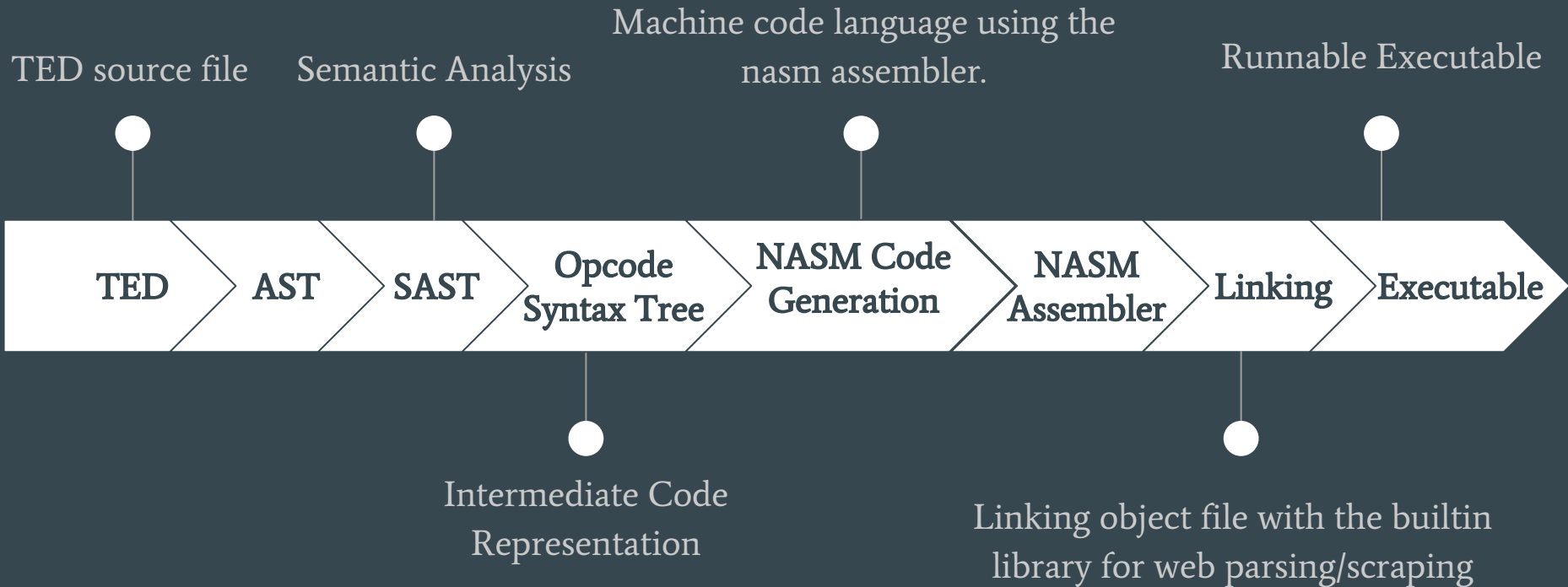
History

- Several group members have used web scraping professionally.
- Per discussion decided to compile to x86 assembly and add web manipulation.
- Original project design turned out to be java-to-java compiler.

Why is it interesting?

TED is the first language designed for web scraping with the power of a low level programming language.

Code Generation



Language Syntax

Code Generation

Similar to GCC code generation

Highlights

- Specialized Data Types designed for web parsing (Page, Element)
- Simple C-like syntax without pointers
- No memory management
- Web parsing with built-in CSS selection

```
4 SECTION .text
5 main:
6   push rbp
7   mov rbp, rsp
8   sub rsp, 40H      TED
9   mov rdi, Str0
10  call pageFetch
11  mov qword [rbp-8H], rax
12  mov rsi, Str1
13  mov rdi, [rbp-8H]
14  call pageFind
15  mov qword [rbp-18H], rax
16  mov rdi, [rbp-18H]
17  call listHead
18  mov qword [rbp-20H], rax
19  mov rsi, Str8
20  mov rdi, [rbp-20H]
21  call elementAttr
22  mov qword [rbp-28H], rax
23  mov rdx, [rbp-28H]
24  mov rsi, Str3
25  mov rdi, [stdout]
26
27 main:      ; Function begin
28   push    rbp
29   mov     rbp, rsp
30   sub     rsp, 48      GCC
31   mov     dword [rbp-24H], edi
32   mov     qword [rbp-30H], rsi
33   mov     edi, ?_001
34   call   pageFetch
35   mov     qword [rbp-20H], rax
36   mov     rax, qword [rbp-20H]
37   mov     esi, ?_002
38   mov     rdi, rax
39   call   pageFind
40   mov     qword [rbp-18H], rax
41   mov     rax, qword [rbp-18H]
42   mov     rdi, rax
43   call   listHead
44   mov     esi, ?_003
45   mov     rdi, rax
46   call   elementAttr
47   mov     rsi, rax
48   mov     edi, ?_004
```

The Infamous GCD

```
1 int gcd(int x, int y) {
2     int a = x;
3     int b = y;
4     while (a!=b) {
5         if (a > b) {
6             a = a-b;
7         } else {
8             b = b - a;
9         }
10    }
11    return a;
12 }
```

```
23
24 SECTION .text
25 gcd:
26 > push rbp
27 > mov rbp, rsp
28 > sub rsp, 40H
29 > mov rax, [rbp-18H]
30 > mov qword [rbp-8H], rax
31 > mov rax, [rbp-20H]
32 > mov qword [rbp-10H], rax
33 > mov [rbp-18H], rdi
34 > mov [rbp-20H], rsi
35 gcd0w:
36 > mov rcx, [rbp-10H]
37 > mov rax, [rbp-8H]
38 > cmp rax, rcx
39 > setne dl, 1
40 > cmp dl, 1
41 > jnz gcd0wend
42 > mov rcx, [rbp-10H]
43 > mov rax, [rbp-8H]
44 > cmp rax, rcx
45 > setg dl, 1
46 > cmp dl, 1
47 > jz gcd0w0t
48 > mov rcx, [rbp-8H]
49 > mov rax, [rbp-10H]
50 > sub rax, rcx
51 > mov qword [rbp-10H], rax
52 > jmp gcd0wbend
53 gcd0w0t:
54 > mov rcx, [rbp-10H]
55 > mov rax, [rbp-8H]
56 > sub rax, rcx
57 > mov qword [rbp-8H], rax
58 gcd0wbend:
59 > jmp gcd0w
60 gcd0wend:
61 > mov rax, [rbp-8H]
62 > leave
63 > ret
64
```

Demo

Built In Functions

List

- `listNew();`
- `listHead(list);`
- `listTail(list);`
- `listSet(list, data);`
- `listAddAfter(list, data);`
- `listRemove(list, index);`
- `listConcat (list1, list2);`
- `listAddLast(list, data);`

Page

- `pageFetch("http://www.sample.com/");`
- `pageFind(page, "#sample_id");`
- `pageURL(page);`
- `pageHTML(page);`
- `pageRoot(page);`

Element

- `elementText(element);`
- `elementType(element);`
- `elementAttr(element, "sample");`
- `elementChildren(page, element);`

Functionality

The built in functions work by communicating over underlying integrational layer with the PhantomJS interpreter that serves as the functionality for all of the easy to use functionalities of web scraping and parsing. That library opens the entire library of css-selector language that allows for easy selection of information to be collected from the parsed web-page that the developer is visiting.

CSS-Selectors Overview

This is a subset selection language similar to the way regex functioning. Below is just a tiny sample of the array of data selection available.

- “*” - Selects all elements.
 - “.class” - Selects all elements with the given class.
 - [name=”value”] - Selects elements that have the specified attribute with a value exactly equal to a certain value.
 - “parent > child” - Selects all direct child elements specified by “child” of elements specified by “parent”.
-

Regression Testing

TEST 1

The first series of tests are basic parsing and variable declarations for syntax only.

TEST 2

Next we introduced string, Page, List and Element and a series of tests were developed for declaration and syntax only.

TEST 3

Next we needed to implement *print* so tests were designed to print integers and strings.

TEST 4

For loop and while loops

TEST 5

The library was becoming functional so this round of tests included Lists, file, and web data

TEST 6

The final rounds of testing involved modifying existing tests to match TED as the language evolved. While we maintained the original language design, syntax and implementation changed as TED became more complex

Future Work

1. Improve language syntax by introducing nested function definitions and better function invocation methods.
2. Introduce syntax for formatting the web-scraped data to shape it in a meaningfully presentable format such as csv and ascii tables as well as mysql insert queries.
3. Remove the dependency on the PhantomJS layer and build the functionality directly into the language.
4. Improve syntax compromises that were made due to implementation such as declaring all variables prior to function calls, improving readability of built-in functions, etc.

Questions?