

Exercises on Compilers

Jordi Cortadella

February 7, 2022

Semantic Analysis

1. We look at a simple language with an exception mechanism:

$$\begin{aligned} S &\rightarrow \text{throw id} \\ S &\rightarrow S \text{ catch id} \Rightarrow S \\ S &\rightarrow S \text{ or } S \\ S &\rightarrow \text{other} \end{aligned}$$

A **throw** statement throws a named exception. This is caught by the nearest enclosing **catch** statement (i.e., where the throw statement is in the left sub-statement of the catch statement) using the same name, whereby the statement after the arrow in the catch statement is executed. An **or** statement is a non-deterministic choice between the two statements, so either one can be executed. *other* is a statement that does not throw any exceptions. We want the type checker to ensure that all possible exceptions are caught and that no catch statement is superfluous, i.e., that the exception it catches can, in fact, be thrown by its left sub-statement. Write type-check functions that implement these checks.

Hint: Let the type of a statement be the set of possible exceptions it can throw.

Solution:

We define a new synthesized attribute for S called $S.\text{except}$ that represents the set of uncaught exceptions of the statements executed by S . The attribute is calculated as follows:

$$\begin{aligned} S \rightarrow \text{throw id} & \quad \{S.\text{except} = \{\text{id}\}\} \\ S \rightarrow S_1 \text{ catch } \{\text{id}\} \Rightarrow S_2 & \quad \{S.\text{except} = (S_1.\text{except} \setminus \{\text{id}\}) \cup S_2.\text{except}\} \\ S \rightarrow S_1 \text{ or } S_2 & \quad \{S.\text{except} = S_1.\text{except} \cup S_2.\text{except}\} \\ S \rightarrow \text{other} & \quad \{S.\text{except} = \emptyset\} \end{aligned}$$

An error exists if $S.\text{except} \neq \emptyset$ at the topmost statement of a function.

A **catch** statement is redundant if $\text{id} \notin S_1$ in the second production, meaning that the statement is prepared to catch **id** but S_1 never throws it.

2. [From Aho-Sethi-Ullman's book.]

Let us consider the grammar of a simple language, represented by the nonterminal P , consisting of a sequence of declarations D followed by a single expression E .

$$\begin{aligned} P &\rightarrow D ; E \\ D &\rightarrow D ; D \mid \text{id} : T \\ T &\rightarrow \text{char} \mid \text{integer} \mid \text{array} [\text{num}] \text{ of } T \mid *T \\ E &\rightarrow \text{literal} \mid \text{num} \mid \text{id} \mid E \text{ mod } E \mid E [E] \mid E* \end{aligned}$$

One program generated by the grammar is:

```
key: integer;
key mod 1999
```

The language has two basic types, *char* and *integer*; a third basic type *type_error* is used to signal errors. For arrays, the number inside the square brackets represents its size. For example,

```
array [256] of char
```

leads to the type expression $\text{array}(256, \text{char})$. The prefix operator $*$ in declarations builds a pointer type, so $*\text{integer}$ leads to the type expression $\text{pointer}(\text{integer})$.

The action associated with the production $D \rightarrow \text{id} : T$ saves a type in the symbol-table entry for an identifier. The action $\text{addtype}(\text{id.entry}, T.\text{type})$ is applied to the synthesized attribute *entry* pointing to the symbol-table entry for *id* and a type expression represented by the synthesized attribute *entry* of the nonterminal T .

If T generates *char* or *integer*, then $T.\text{type}$ is defined to be *char* or *integer*, respectively. The upper bound of an array is obtained from the attribute *val* of token *num* that gives the integer represented by *num*. Since the declarations appear before the expression, we can be sure that the types of all declared identifiers will be saved in the symbol table before the expression is checked.

Next, part of the semantic actions that save the type of an identifier in the symbol table is shown:

$$\begin{aligned} P &\rightarrow D ; E \\ D &\rightarrow D ; D \\ D &\rightarrow \text{id} : T && \{ \text{addtype}(\text{id.entry}, T.\text{type}) \} \\ T &\rightarrow \text{char} && \{ T.\text{type} := \text{char} \} \\ T &\rightarrow \text{integer} && \{ T.\text{type} := \text{integer} \} \\ T &\rightarrow *T_1 && \{ T.\text{type} := \text{pointer}(T_1.\text{type}) \} \\ T &\rightarrow \text{array} [\text{num}] \text{ of } T_1 && \{ T.\text{type} := \text{array}(\text{num.val}, T_1.\text{type}) \} \end{aligned}$$

The tokens *literal* and *num* have type *char* and *integer*, respectively. The expressions formed by applying the *mod* operator to two subexpressions of type *integer* has type *integer*; otherwise, its type is *type_error*. In an array reference $E_1[E_2]$, the index expression E_2 must have type *integer*. The type of $E*$ is the type t of the object pointed to by the pointer E .

Add the type checking actions of the grammar to synthesize the attribute $E.\text{type}$ for the expression. We can use the function $\text{lookup}(\text{id.entry})$ to fetch the type saved in the symbol table for the entry corresponding to the identifier *id*.

Solution: $E \rightarrow \text{literal} \quad \{ E.type := \text{char} \}$ $E \rightarrow \text{num} \quad \{ E.type := \text{integer} \}$ $E \rightarrow \text{id} \quad \{ E.type := \text{lookup}(\text{id.entry}) \}$ $E \rightarrow E_1 \text{ mod } E_2 \quad \{ E.type := (E_1.type = \text{integer} \text{ and } E_2.type = \text{integer}) ? \text{integer} : \text{type_error} \}$ $E \rightarrow E_1[E_2] \quad \{ E.type := (E_2.type = \text{integer} \text{ and } E_1.type = \text{array}(s,t)) ? t : \text{type_error} \}$ $E \rightarrow E_1^* \quad \{ E.type := E_1.type = \text{pointer}(t) ? t : \text{type_error} \}$

3. Let us consider the following fragment of C code:

```
struct {
    int a;
    char b;
} c;

(int *) d[10];

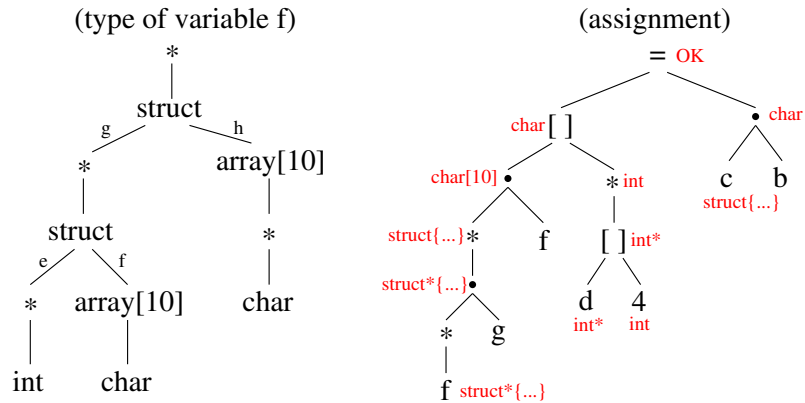
struct {
    struct {
        int *e;
        char f[10];
    } *g;
    (char *) h[10];
} e, *f;
```

```
1: ((*f).g).f[*d[4]] = c.b;
2: e.h[11] = &(c.b);
3: *(e.g + int) = ((*f).g);
4: d[12] = &(c.a+1);
5: c.a = e.2[h];
6: (*f)..h[1];
```

- Give an AST for the type of variable `f`. Use the label `struct` to represent the root node of a struct's AST and the label `array[n]` to represent the root node of an array of n elements. When representing a `struct`, annotate the label of each field at the tree edge associated to the field.
- Give an AST representing the assignment in line 1 and infer the type associated to each node of the AST. Represent the access to a `struct (x.y)` as a binary tree $(.,x,y)$ and the access to an array $(a[i])$ as a binary tree $([],a,i)$. Is it a correct assignment?
- Detect all possible errors in the lines of code 1-6. Indicate when the error is detected (lexical, syntax, semantic, runtime).
- Indicate when the following errors are detected (lexical, syntax, semantic, runtime):
 - `&` can only be used for an addressable memory location.
 - Missing semicolon at the end of a statement.
 - The LHS and RHS types of an assignment are not compatible.
 - A character does not belong to the alphabet of the language.
 - An array index is out of bounds.
 - An array index must be integer.
 - Unclosed parenthesis in an expression.
 - Division by zero.
 - Unclosed quotes `"` of a string.
 - Only integer values can be added to pointers.

Solution:

Below are the ASTs of the type of variable `f` and the assignment. In red, the type of the nodes of the assignment AST.



Errors in lines 1-6:

- Line 1: Only runtime errors could be produced, e.g. access using a NULL pointer, or access to array `f` with an index out of bounds.
- Line 2: Access to array with index out of bounds. This error is typically produced at runtime, although a compiler with powerful static analysis could also detect it during semantic analysis.
- Line 3: `int` is a keyword and the presence in the expression would produce a syntax error. Obtaining a point from an expression `(+)` is a semantic error.
- Line 4: Access to array with index out of bounds. This error is typically produced at runtime, although a compiler with powerful static analysis could also detect it during semantic analysis. Moreover, the expression `c.a+1` is not addressable, i.e., no address can be obtained from it.
- Line 5: syntax error. An identifier must follow the dot that precedes the field of a structure.
- Line 6: syntax error. An identifier must follow the dot that precedes the field of a structure.

Errors in the code:

<code>&</code> can only be used for an addressable memory location	Semantic
Missing semicolon at the end of a statement	Syntax
The LHS and RHS types of an assignment are not compatible	Semantic
A character does not belong to the alphabet of the language	Lexical
An array index is out of bounds	Runtime
An array index must be integer	Semantic
Unclosed parenthesis in an expression	Syntax
Division by zero	Runtime
Unclosed quotes (<code>"</code>) in a string	Lexical
Only integer values can be added to pointers	Semantic