

## Static Typing of a Fragment of C Language

*Based on material from INF564 course given at  
École Polytechnique by Jean-Christophe Filliâtre*

### 1 Introduction

The goal is to build a typechecker for a tiny fragment of the C language, called Mini C in the following. it contains integers and pointers to structures. It is fully compatible with C. This means a C compiler such as gcc can be used as a reference.

**Differences wrt C.** Any Mini C program is a legal C program. Yet, Mini C has limitations wrt C. Here are some of them:

- There is no variable initialization. To initialize a variable, one has to use an assignment;
- the only types are integers (Basic signed 32 integer type `int`), pointers to the structures (`struct id *`), and void pointer type, `void *`, (e.g. used for the return type of `malloc`);
- There is no pointer arithmetic (and no memory deallocation);
- Mini C has fewer keywords than C.

**Predefined Functions.** The following functions are predefined:

```
int putchar(int c);  
void *malloc(int n);
```

(But there is no need for `#include` in Mini C for testing.)

### 2 Syntax

We use the following notations in grammars:

$\langle rule \rangle^*$	repeats $\langle rule \rangle$ an arbitrary number of times (including zero)
$\langle rule \rangle_t^*$	repeats $\langle rule \rangle$ an arbitrary number of times (including zero), with separator $t$
$\langle rule \rangle^+$	repeats $\langle rule \rangle$ at least once
$\langle rule \rangle_t^+$	repeats $\langle rule \rangle$ at least once, with separator $t$
$\langle rule \rangle?$	use $\langle rule \rangle$ optionally
$( \langle rule \rangle )$	grouping

Be careful not to confuse “\*” and “+” with “\*” and “+” that are C symbols. Similarly, do not confuse grammar parentheses with terminal symbols ( and ).

## 2.1 Lexical Conventions

Spaces, tabs, and newlines are blanks. Comments are of two kinds:

- from `/*` to `*/` and not nested;
- from `//` to the end of the line.

Identifiers follow the regular expression  $\langle \textit{ident} \rangle$  :

$$\begin{aligned} \langle \textit{digit} \rangle &::= 0-9 \\ \langle \textit{alpha} \rangle &::= \text{a-z} \mid \text{A-Z} \\ \langle \textit{ident} \rangle &::= (\langle \textit{alpha} \rangle \mid \_)(\langle \textit{alpha} \rangle \mid \langle \textit{digit} \rangle \mid \_)^* \end{aligned}$$

The following identifiers are keywords:

`int struct if else while return sizeof`

Last, integer literals follow the regular expression  $\langle \textit{integer} \rangle$  :

$$\begin{aligned} \langle \textit{integer} \rangle &::= 0 \\ &\mid 1-9 \langle \textit{digit} \rangle^* \\ &\mid 0 \langle \textit{digit-octal} \rangle^+ \\ &\mid 0\text{x} \langle \textit{digit-hexa} \rangle^+ \\ &\mid ' \langle \textit{character} \rangle ' \\ \langle \textit{digit-octal} \rangle &::= 0-7 \\ \langle \textit{digit-hexa} \rangle &::= 0-9 \mid \text{a-f} \mid \text{A-F} \\ \langle \textit{character} \rangle &::= \text{any ASCII character with a code in } [32, 127], \\ &\text{other than } \backslash, ', \text{ and } " \\ &\mid \backslash \backslash \mid \backslash ' \mid \backslash " \\ &\mid \backslash \text{x} \langle \textit{digit-hexa} \rangle \langle \textit{digit-hexa} \rangle \end{aligned}$$

## 2.2 Syntax

The grammar of source files is given in Fig. 1. The entry point is  $\langle \textit{file} \rangle$ . Associativity and priorities are given below, from lowest to strongest priority.

operation	associativity	priority
<code>=</code>	right	lowest
<code>  </code>	left	
<code>&amp;&amp;</code>	left	
<code>==</code> <code>!=</code>	left	
<code>&lt;</code> <code>&lt;=</code> <code>&gt;</code> <code>&gt;=</code>	left	↓
<code>+</code> <code>-</code>	left	
<code>*</code> <code>/</code>	left	
<code>!</code> <code>-</code> (unary)	right	
<code>-&gt;</code>	left	strongest

$\langle file \rangle$	$::=$	$\langle decl \rangle^* EOF$
$\langle decl \rangle$	$::=$	$\langle decl\_typ \rangle \mid \langle decl\_fct \rangle$
$\langle decl\_vars \rangle$	$::=$	$int \langle ident \rangle^+ ;$ $\mid$ $struct \langle ident \rangle (* \langle ident \rangle)^+ ;$
$\langle decl\_typ \rangle$	$::=$	$struct \langle ident \rangle \{ \langle decl\_vars \rangle^* \} ;$
$\langle decl\_fct \rangle$	$::=$	$int \langle ident \rangle ( \langle param \rangle^* ) \langle bloc \rangle$ $\mid$ $struct \langle ident \rangle * \langle ident \rangle ( \langle param \rangle^* ) \langle bloc \rangle$
$\langle param \rangle$	$::=$	$int \langle ident \rangle \mid struct \langle ident \rangle * \langle ident \rangle$
$\langle expr \rangle$	$::=$	$\langle integer \rangle$ $\mid$ $\langle ident \rangle$ $\mid$ $\langle expr \rangle \rightarrow \langle ident \rangle$ $\mid$ $\langle ident \rangle ( \langle expr \rangle^* )$ $\mid$ $! \langle expr \rangle \mid - \langle expr \rangle$ $\mid$ $\langle expr \rangle \langle binop \rangle \langle expr \rangle$ $\mid$ $\langle ident \rangle = \langle expr \rangle$ $\mid$ $\langle expr \rangle \rightarrow \langle ident \rangle = \langle expr \rangle$ $\mid$ $sizeof ( struct \langle ident \rangle )$ $\mid$ $malloc ( struct \langle ident \rangle )$ $\mid$ $( \langle expr \rangle )$
$\langle binop \rangle$	$::=$	$== \mid != \mid < \mid <= \mid > \mid >= \mid + \mid - \mid * \mid / \mid \&\& \mid   $
$\langle stmt \rangle$	$::=$	$;$ $\mid$ $\langle expr \rangle ;$ $\mid$ $if ( \langle expr \rangle ) \langle stmt \rangle$ $\mid$ $if ( \langle expr \rangle ) \langle stmt \rangle else \langle stmt \rangle$ $\mid$ $while ( \langle expr \rangle ) \langle stmt \rangle$ $\mid$ $\langle bloc \rangle$ $\mid$ $return \langle expr \rangle ;$
$\langle bloc \rangle$	$::=$	$\{ \langle decl\_vars \rangle^* \langle stmt \rangle^* \}$

Figure 1: Grammar of Mini C.

### 3 Static Typing

Once parsing phase is completed (provided in the lab assignment), we explain how to perform static typing of Mini C.

**Types and Typing Environments.** Expressions have types  $\tau$  with the following abstract syntax

$$\tau ::= \text{int} \mid \text{struct } id * \mid \text{void}^*$$

where  $id$  stands for a structure name. We introduce the relation  $\equiv$  over types as the smallest reflexive and symmetric relation that additionally satisfies the equation  $\text{void}^* \equiv \text{struct } id *$ .

A typing environment  $\Gamma$  is a sequence of variable declarations  $\tau x$ , structure declarations  $\text{struct } S \{\tau_1 x_1 \cdots \tau_n x_n\}$  and function declarations  $\tau f(\tau_1, \dots, \tau_n)$ . We write  $\text{struct } S \{\tau x\}$  to indicate that structure  $S$  has a field  $x$  with type  $\tau$ .

We say that a type  $\tau$  is *well-formed* in environment  $\Gamma$ , and we write  $\Gamma \vdash \tau \text{ bf}$ , if all structure names in  $\tau$  correspond to structures declared in  $\Gamma$ .

**Uniqueness Rules.** In addition to the typing rules below (for structure declarations, expressions, statements and function declarations), we have to check for uniqueness

- of structure names over the whole file;
- of structure fields inside a *single* structure;
- of function parameters;
- of local variables inside a *single* block;
- of function names over the whole file.

#### 3.1 Adding Structure Declarations to the typing environment

A file is a list of structure and function declarations (there are no global variables in Mini C). We first add structure declarations to the typing environment. To this end, we introduce the judgment  $\Gamma \vdash d \rightarrow \Gamma'$  meaning “in environment  $\Gamma$ , declaration  $d$  is well-formed and outputs environment  $\Gamma'$ ”. It is defined as follows.

$$\frac{\forall i, \Gamma, \text{struct } id \{\tau_1 x_1 \cdots \tau_n x_n\} \vdash \tau_i \text{ bf}}{\Gamma \vdash \text{struct } id \{\tau_1 x_1; \cdots \tau_n x_n\} \rightarrow \{\text{struct } id \{\tau_1 x_1 \cdots \tau_n x_n\}\} \cup \Gamma}$$

Note that types  $\tau_i$  may only refer to the structure  $id$  via pointer types (including the case where structure definition is recursive).

#### 3.2 Type-Checking Expressions

We introduce the typing judgment  $\Gamma \vdash e : \tau$  meaning “in environment  $\Gamma$ , expression  $e$  is well-typed, with type  $\tau$ ”. This judgment is defined as follows:

$$\frac{}{\Gamma \vdash 0 : \text{void}^*} \quad \frac{c \text{ integer constant}}{\Gamma \vdash c : \text{int}} \quad \frac{\tau x \in \Gamma}{\Gamma \vdash x : \tau}$$

$$\frac{\Gamma \vdash e : \text{struct } S * \quad \text{struct } S \{\tau x\} \in \Gamma}{\Gamma \vdash e \rightarrow x : \tau} \quad \frac{\text{struct } S \in \Gamma}{\Gamma \vdash \text{sizeof}(\text{struct } S) : \text{int}}$$

$$\begin{array}{c}
\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \equiv \tau_2}{\Gamma \vdash e_1 = e_2 : \tau_1} \\
\frac{\Gamma \vdash e : \tau \quad \tau \equiv \text{int}}{\Gamma \vdash - e : \text{int}} \quad \frac{\Gamma \vdash e : \tau}{\Gamma \vdash ! e : \text{int}} \\
\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \equiv \tau_2 \quad op \in \{==, !=, <, <=, >, >=\}}{\Gamma \vdash e_1 \text{ } op \text{ } e_2 : \text{int}} \\
\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad op \in \{||, \&\&\}}{\Gamma \vdash e_1 \text{ } op \text{ } e_2 : \text{int}} \\
\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \equiv \text{int} \quad \tau_2 \equiv \text{int} \quad op \in \{+, -, *, /\}}{\Gamma \vdash e_1 \text{ } op \text{ } e_2 : \text{int}} \\
\frac{\tau \ f(\tau'_1, \dots, \tau'_n) \in \Gamma \quad \forall i, \Gamma \vdash e_i : \tau_i \quad \tau_i \equiv \tau'_i}{\Gamma \vdash f(e_1, \dots, e_n) : \tau}
\end{array}$$

### 3.3 Type-Checking Statements

We introduce the judgment  $\Gamma \vdash^{\tau_0} s$  meaning “in environment  $\Gamma$ , statement  $s$  is well-typed, for a return type  $\tau_0$ ”. Type  $\tau_0$  stands for the return type of the function in which statement  $s$  occurs. This judgment is defined as follows:

$$\begin{array}{c}
\frac{}{\Gamma \vdash^{\tau_0} ;} \quad \frac{\Gamma \vdash e : \tau}{\Gamma \vdash^{\tau_0} e;} \quad \frac{\Gamma \vdash e : \tau \quad \tau \equiv \tau_0}{\Gamma \vdash^{\tau_0} \text{return } e;} \\
\frac{\Gamma \vdash e : \tau \quad \Gamma \vdash^{\tau_0} s_1 \quad \Gamma \vdash^{\tau_0} s_2}{\Gamma \vdash^{\tau_0} \text{if } (e) \text{ } s_1 \text{ else } s_2} \\
\frac{\Gamma \vdash e : \tau \quad \Gamma \vdash^{\tau_0} s}{\Gamma \vdash^{\tau_0} \text{while}(e) \text{ } s} \\
\frac{\forall j, \Gamma \vdash \tau_j \text{ bf} \quad \forall j, \Gamma + \{\tau_1 \ x_1, \dots, \tau_k \ x_k\} \vdash^{\tau_0} s_j}{\Gamma \vdash^{\tau_0} \{\tau_1 \ x_1 \dots \tau_k \ x_k; s_1 \dots s_n\}}
\end{array}$$

The last rule means that, to type a block with  $k$  local variables and  $n$  statements, we first check that the variable declarations are well-formed and then we type-check each statement in the environment that is augmented with the new declarations.

### 3.4 Type-Checking Function Declarations and Files

Finally, we explain how to type check functions declarations and files.

**Function Declarations.**

$$\frac{\forall i, \Gamma \vdash \tau_i \text{ bf} \quad \{\tau_0 \ f(\tau_1, \dots, \tau_n), \tau_1 \ x_1, \dots, \tau_n \ x_n\} \cup \Gamma \vdash^{\tau_0} b}{\Gamma \vdash \tau_0 \ f(\tau_1 \ x_1, \dots, \tau_n \ x_n) \ b \rightarrow \{\tau_0 \ f(\tau_1, \dots, \tau_n)\} \cup \Gamma}$$

Note that the prototype of function  $f$  is added to the environment before we type-check its body  $b$ , so that recursive functions are allowed.

**Files.** Finally, we introduce the judgment  $\Gamma \vdash_f d_1 \dots d_n$  meaning “in environment  $\Gamma$ , the file made of declarations  $d_1, \dots, d_n$  is well-formed”. Type-checking a file consists in type-checking its declarations in sequence, the environment being augmented with each new declaration.

$$\frac{}{\Gamma \vdash_f \emptyset} \quad \frac{\Gamma \vdash d_1 \rightarrow \Gamma' \quad \Gamma' \vdash_f d_2 \dots d_n}{\Gamma \vdash_f d_1 \ d_2 \dots d_n}$$

**Entry Point.** Finally, we have to check for the existence of a `main` function with type

```
int main();
```