

CSC_52064 Compilation

Mini Go

version 1 — January 9, 2026

The goal is to build a compiler for a tiny fragment of **the Go programming language**, called **Mini Go** in the following, to x86-64 assembly. This fragment contains integers, Booleans, strings, structures, and pointers. It is compatible with Go. This means that Go can be used as a reference when needed.

The syntax of **Mini Go** is described in Sec. 1. A parser is provided, together with abstract syntax. You have to implement static type checking (Sec. 2) and code generation (Sec. 3).

1 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 Go symbols. Similarly, do not confuse grammar parentheses with terminal symbols (and).

1.1 Lexical Conventions

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

- delimited by `/*` and `*/` (and not nested);
- starting from `//` and extending to the end of line.

Identifiers follow the regular expression $\langle ident \rangle$:

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

The following identifiers are keywords:

```

else    false  for      func    if
import  nil    package  return  struct
true    type   var

```

Integer literals follow the regular expression $\langle integer \rangle$:

$$\begin{aligned}
\langle hexa \rangle &::= 0-9 \mid a-f \mid A-F \\
\langle integer \rangle &::= \langle digit \rangle^+ \\
&\quad \mid (0x \mid 0X) \langle hexa \rangle^+
\end{aligned}$$

Integer literals must be in the range -2^{63} to $2^{63} - 1$. A string literal $\langle string \rangle$ is enclosed with quotes ("). There are four escape sequences: `\"` (for the character `"`), `\n` (for a newline character), `\t` (for a tabulation character), and `\\` (for the character `\`).

Automatic insertion of semi-colons. To save the programmer the trouble of writing semicolons at the end of lines that contain instructions, the lexical analyzer automatically inserts a semicolon when it encounters a carriage return and the previously emitted token was part of the following set:

$$\langle ident \rangle \mid \langle integer \rangle \mid \langle string \rangle \mid \text{true} \mid \text{false} \mid \text{nil} \mid \text{return} \mid ++ \mid - \mid) \mid \}$$

1.2 Syntax

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

| operation | associativity |
|---|---------------|
| <code> </code> | left |
| <code>&&</code> | left |
| <code>==, !=, >, >=, <, <=</code> | left |
| <code>+, -</code> | left |
| <code>*, /, %</code> | left |
| <code>- (unary), * (unary), &, !</code> | — |
| <code>.</code> | left |

Syntactic sugar. We have the following equivalences:

- instruction `for b` stands for `for true b` (infinite loop).
- instruction `for i1; e; i2 b` stands for `{ i1; for e { b i2 }}`.
- instruction `if e b` stands for `if e b else {}`.
- instruction `x1, ..., xn := e1, ..., em` stands for `var x1, ..., xn = e1, ..., em`.

Syntactic sugar is eliminated at parsing time (in the provided parser).

| | | |
|--------------------------------|-------|---|
| $\langle file \rangle$ | $::=$ | <code>package main ; (import "fmt" ;)? $\langle decl \rangle^*$ EOF</code> |
| $\langle decl \rangle$ | $::=$ | $\langle structure \rangle$ $\langle function \rangle$ |
| $\langle structure \rangle$ | $::=$ | <code>type $\langle ident \rangle$ struct { ($\langle vars \rangle^+ ; ?$)? } ;</code> |
| $\langle function \rangle$ | $::=$ | <code>func $\langle ident \rangle$ (($\langle vars \rangle^*$, ?)?) $\langle return_type \rangle$? $\langle bloc \rangle$;</code> |
| $\langle vars \rangle$ | $::=$ | $\langle ident \rangle^+ \langle type \rangle$ |
| $\langle return_type \rangle$ | $::=$ | $\langle type \rangle$ $(\langle type \rangle^+ , ?)$ |
| $\langle type \rangle$ | $::=$ | $\langle ident \rangle$ $* \langle type \rangle$ |
| $\langle expr \rangle$ | $::=$ | $\langle integer \rangle$ $\langle string \rangle$ <code>true</code> <code>false</code> <code>nil</code> $(\langle expr \rangle)$ $\langle ident \rangle$ $\langle expr \rangle . \langle ident \rangle$ $\langle ident \rangle (\langle expr \rangle^* ,)$ <code>fmt . Print ($\langle expr \rangle^* ,)$</code> <code>! $\langle expr \rangle$ - $\langle expr \rangle$ & $\langle expr \rangle$ * $\langle expr \rangle$</code> $\langle expr \rangle \langle operator \rangle \langle expr \rangle$ |
| $\langle operator \rangle$ | $::=$ | <code>== != < <= > >=</code> <code>+ - * / % && </code> |
| $\langle bloc \rangle$ | $::=$ | <code>{ ($\langle stmt \rangle^+ ; ?$)? }</code> |
| $\langle stmt \rangle$ | $::=$ | $\langle simple_stmt \rangle$ $\langle bloc \rangle$ $\langle stmt_if \rangle$ <code>var $\langle ident \rangle^+ \langle type \rangle$? (= $\langle expr \rangle^+$)?</code> <code>return $\langle expr \rangle^*$,</code> <code>for $\langle bloc \rangle$</code> <code>for $\langle expr \rangle \langle bloc \rangle$</code> <code>for $\langle simple_stmt \rangle$? ; $\langle expr \rangle$; $\langle simple_stmt \rangle$? $\langle bloc \rangle$</code> |
| $\langle simple_stmt \rangle$ | $::=$ | $\langle expr \rangle$ $\langle expr \rangle (++ -)$ $\langle expr \rangle^+ = \langle expr \rangle^+$ $\langle ident \rangle^+ := \langle expr \rangle^+$ |
| $\langle stmt_if \rangle$ | $::=$ | <code>if $\langle expr \rangle \langle bloc \rangle$ (else ($\langle bloc \rangle$ $\langle stmt_if \rangle$))?</code> |

Figure 1: Grammar of Mini Go.

2 Static Typing

Static types τ are given by the following abstract syntax:

$$\tau ::= \text{int} \mid \text{bool} \mid \text{string} \mid S \mid * \tau$$

where S is a structure name. A typing context Γ contains structures (written S), variables (written $x : \tau$), and functions (written $f(\tau_1, \dots, \tau_n) \Rightarrow \tau'_1, \dots, \tau'_m$ with $n \geq 0$ and $m \geq 0$).

Well-formed types. The judgment $\Gamma \vdash \tau \text{ bf}$ means “the type τ is well-formed in environment Γ ”. It is defined as follows:

$$\frac{}{\Gamma \vdash \text{int} \text{ bf}} \quad \frac{}{\Gamma \vdash \text{bool} \text{ bf}} \quad \frac{}{\Gamma \vdash \text{string} \text{ bf}} \quad \frac{S \in \Gamma}{\Gamma \vdash S \text{ bf}} \quad \frac{\Gamma \vdash \tau \text{ bf}}{\Gamma \vdash * \tau \text{ bf}}$$

Structure fields. We write $S\{x : \tau\}$ the fact that structure S has a field x of type τ .

Typing expressions. The judgment $\Gamma \vdash e : \tau$ means “in the context Γ , the expression e is well typed with type τ ”. The judgment $\Gamma \vdash_l e : \tau$ additionally means that e is a left value. These two judgments are defined by the following set of rules.

$$\begin{array}{c} \frac{c \text{ constant of type } \tau}{\Gamma \vdash c : \tau} \quad \frac{\Gamma \vdash \tau \text{ bf}}{\Gamma \vdash \text{nil} : * \tau} \quad \frac{x : \tau \in \Gamma}{\Gamma \vdash_l x : \tau} \quad \frac{\Gamma \vdash_l e : \tau \quad e \neq _}{\Gamma \vdash e : \tau} \\[10pt] \frac{\Gamma \vdash e : S \quad S\{x : \tau\}}{\Gamma \vdash e.x : \tau} \quad \frac{\Gamma \vdash e : *S \quad e \neq \text{nil} \quad S\{x : \tau\}}{\Gamma \vdash e.x : \tau} \\[10pt] \frac{\Gamma \vdash_l e : \tau' \quad \Gamma \vdash e.x : \tau}{\Gamma \vdash_l e.x : \tau} \quad \frac{\Gamma \vdash e : \text{int}}{\Gamma \vdash - e : \text{int}} \quad \frac{\Gamma \vdash e : \text{bool}}{\Gamma \vdash !e : \text{bool}} \\[10pt] \frac{\Gamma \vdash e : * \tau \quad e \neq \text{nil}}{\Gamma \vdash_l *e : \tau} \quad \frac{\Gamma \vdash_l e : \tau}{\Gamma \vdash \&e : * \tau} \quad \frac{S \in \Gamma}{\Gamma \vdash \text{new}(S) : *S} \\[10pt] \frac{\Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau \quad op \in \{=, !=\} \quad e_1 \neq \text{nil} \vee e_2 \neq \text{nil}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{bool}} \\[10pt] \frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int} \quad op \in \{<, <=, >, >=\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{bool}} \\[10pt] \frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int} \quad op \in \{+, -, *, /, \%\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{int}} \\[10pt] \frac{\Gamma \vdash e_1 : \text{bool} \quad \Gamma \vdash e_2 : \text{bool} \quad op \in \{\&\&, ||\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{bool}} \\[10pt] \frac{\Gamma \vdash f(e_1, \dots, e_n) \Rightarrow \tau_1}{\Gamma \vdash f(e_1, \dots, e_n) : \tau_1} \end{array}$$

Typing function calls. The judgment $\Gamma \vdash f(e_1, \dots, e_n) \Rightarrow \tau_1, \dots, \tau_m$ means “in the context Γ , the function call $f(e_1, \dots, e_n)$ is well typed and returns m values of types τ_1, \dots, τ_m ”. It is defined as follows:

$$\frac{f(\tau_1, \dots, \tau_n) \Rightarrow \tau'_1, \dots, \tau'_m \in \Gamma \quad \forall i, \Gamma \vdash e_i : \tau_i}{\Gamma \vdash f(e_1, \dots, e_n) \Rightarrow \tau'_1, \dots, \tau'_m}$$

$$\frac{n \geq 2 \quad f(\tau_1, \dots, \tau_n) \Rightarrow \tau'_1, \dots, \tau'_m \in \Gamma \quad \Gamma \vdash g(e_1, \dots, e_k) \Rightarrow \tau_1, \dots, \tau_n}{\Gamma \vdash f(g(e_1, \dots, e_k)) \Rightarrow \tau'_1, \dots, \tau'_m}$$

The second rule allows us to pass the n results of function g to a function f expecting n parameters.

Typing statements. The judgment $\Gamma \vdash s$ means “in the context Γ , the statement s is well typed”. It is defined as follows:

$$\frac{\Gamma \vdash_l e : \text{int}}{\Gamma \vdash e++} \quad \frac{\Gamma \vdash_l e : \text{int}}{\Gamma \vdash e-}$$

$$\frac{\forall i, \Gamma \vdash e_i : \tau_i}{\Gamma \vdash \text{fmt.Print}(e_1, \dots, e_n)} \quad \frac{n \geq 2 \quad \Gamma \vdash f(e_1, \dots, e_k) \Rightarrow \tau_1, \dots, \tau_n}{\Gamma \vdash \text{fmt.Print}(f(e_1, \dots, e_k))}$$

$$\frac{\forall i, \Gamma \vdash_l e_i : \tau_i \quad \forall i, \Gamma \vdash e'_i : \tau_i}{\Gamma \vdash e_1, \dots, e_n = e'_1, \dots, e'_n} \quad \frac{\forall i, \Gamma \vdash_l e_i : \tau_i \quad \Gamma \vdash f(e'_1, \dots, e'_m) \Rightarrow \tau_1, \dots, \tau_n}{\Gamma \vdash e_1, \dots, e_n = f(e'_1, \dots, e'_m)}$$

$$\frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash b_1 \quad \Gamma \vdash b_2}{\Gamma \vdash \text{if}(e) b_1 \text{ else } b_2} \quad \frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash b}{\Gamma \vdash \text{for } e \text{ b}}$$

$$\frac{\forall i, \Gamma \vdash e_i : \tau_i}{\Gamma \vdash \text{return } e_1, \dots, e_n} \quad \frac{\Gamma \vdash f(e_1, \dots, e_k) \Rightarrow \tau_1, \dots, \tau_n}{\Gamma \vdash \text{return } f(e_1, \dots, e_k)}$$

$$\frac{\Gamma \vdash \tau \text{ bf} \quad \Gamma + x_1 : \tau, \dots, x_n : \tau \vdash \{s_2; \dots; s_m\}}{\Gamma \vdash \{\text{var } x_1, \dots, x_n \tau; s_2; \dots; s_m\}}$$

$$\frac{\Gamma \vdash \tau \text{ bf} \quad \forall i, \Gamma \vdash e_i : \tau \quad \Gamma + x_1 : \tau, \dots, x_n : \tau \vdash \{s_2; \dots; s_n\}}{\Gamma \vdash \{\text{var } x_1, \dots, x_n \tau = e_1, \dots, e_n; s_2; \dots; s_n\}}$$

$$\frac{\forall i, e_i \neq \text{nil} \quad \forall i, \Gamma \vdash e_i : \tau_i \quad \Gamma + x_1 : \tau_1, \dots, x_n : \tau_n \vdash \{s_2; \dots; s_m\}}{\Gamma \vdash \{\text{var } x_1, \dots, x_n = e_1, \dots, e_n; s_2; \dots; s_m\}}$$

$$\frac{\Gamma \vdash \tau \text{ bf} \quad \Gamma \vdash f(e_1, \dots, e_k) \Rightarrow \tau^n \quad \Gamma + x_1 : \tau, \dots, x_n : \tau \vdash \{s_2; \dots; s_m\}}{\Gamma \vdash \{\text{var } x_1, \dots, x_n \tau = f(e_1, \dots, e_k); s_2; \dots; s_m\}}$$

$$\frac{\Gamma \vdash f(e_1, \dots, e_k) \Rightarrow \tau_1, \dots, \tau_n \quad \Gamma + x_1 : \tau_1, \dots, x_n : \tau_n \vdash \{s_2; \dots; s_m\}}{\Gamma \vdash \{\text{var } x_1, \dots, x_n = f(e_1, \dots, e_k); s_2; \dots; s_m\}}$$

$$\frac{}{\Gamma \vdash \{\}} \quad \frac{\Gamma \vdash s_1 \quad \Gamma \vdash \{s_2; \dots; s_n\}}{\Gamma \vdash \{s_1; \dots; s_n\}}$$

In addition, all the variables introduced in the *same* block must have distinct names, with the exception of variables named `_`.

Typing a file. The declarations in a file can appear in any order. In particular, functions and structures are mutually recursive. It is suggested to proceed in three steps:

1. We add all the structures (but not their fields) to the environment, checking the uniqueness of the structure names.
2. (a) We add all the functions to the environment, checking the uniqueness of the function names. For a function declaration of the form

$$\text{func } f(x_1 : \tau_1, \dots, x_n : \tau_n) (\tau'_1, \dots, \tau'_m) \{b\}$$

we check that the x_i are pairwise distinct and that all types τ_i and τ'_j are well-formed.

- (b) We check and add all the structure fields to the environment. For a structure declaration of the form

$$\text{type } S \text{ struct } \{ x_1 : \tau_1, \dots, x_n : \tau_n \}$$

we check that the x_i are pairwise distinct and that all types τ_i are well-formed.

3. (a) For each function declaration

$$\text{func } f(x_1 : \tau_1, \dots, x_n : \tau_n) (\tau'_1, \dots, \tau'_m) \{b\}$$

we construct a new environment Γ by adding all the variables $x_i : \tau_i$ to the environment containing the structures and functions, and we type the block b in Γ , *i.e.*, we check $\Gamma \vdash b$. Additionally, we check

- that every **return** statement in b returns a list of m values of the expected types τ'_1, \dots, τ'_m ;
 - if $m > 0$, that any control-flow branch in b reaches a **return** statement;
 - that any local variable introduced in b , other than $_$, is further used.
- (b) We check that there is no “recursive” structure, that is, no structure S with a field (which contains a field, which contains a field, etc.) of type S , without passing through a pointer.

Finally, we check that there exists a function **main** without parameters and without return type and that the file contains **import "fmt"** if and only if there is at least one **fmt.Print** instruction somewhere in the file.

3 Code Generation

The aim is to produce a simple but correct compiler. In particular, we do not attempt to do any kind of register allocation, but simply use the stack to store any intermediate calculations. Of course, it is possible, and even desirable, to use some x86-64 registers locally. Memory is allocated using `malloc` and no attempt will be made to free memory.

Simplification. To make things simpler, you can assume that your project will not be tested on programs passing structures by value, returning structures from functions, or assigning structures. Beside, some program transformations are performed on the typed syntax trees, between static typing and code generation, to

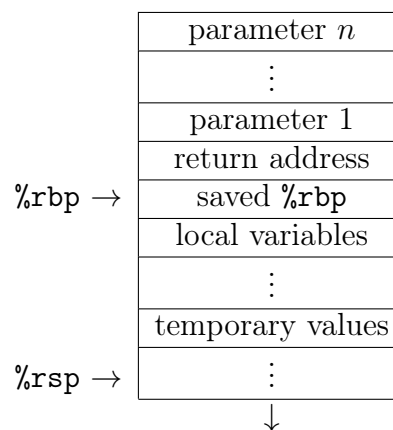
- get rid of multiple assignments;
- get rid of multiple returns;
- allocate all structures on the heap;
- allocate on the heap any variable whose address is taken with `&`.

(This is the purpose of module `Rewrite`, already implemented and used in `minigo.ml`.)

Value Representation. We propose a simple compilation scheme but you are free to use any other. A value of type `int`, `bool`, `string`, or `*τ` is a 64-bit word. The value `nil` is the integer 0. The values `false` and `true` are the integers 0 and 1, respectively. A string is a pointer to a 0-terminated string allocated on the data segment (since it is not possible to build strings dynamically in our fragment).

A structure is a heap-allocated block of a size large enough to accommodate all fields. Your compiler maintains, for each field, its position within the structure, relative to the start of the structure. Note that structures can be nested.

Stack Layout. We suggest a compilation scheme where all parameters are passed on the stack (each of them being a 64-bit word, since structures are never passed), and where the return value is in register `%rax` (again, no structure is returned). The stack frame is depicted beside. Local variables are allocated on the stack. (But note that no structure is allocated on the stack, thanks to the rewriting performed prior to code generation; see above.) The top of the stack is used to store intermediate computations, such as the value of the first operand during the evaluation of a binary operation.



The assembly code generated by your compiler should look like:

```
        .text
        .globl main
main:   call of Mini Go function main
        xorq %rax, %rax
        ret
        ... code of Mini Go functions
        ... assembly functions, if needed
        .data
        ... string literals
```

Stack alignment. With recent versions of the `libc`, it is important to have a 16-byte stack alignment when calling library functions such as `malloc` or `printf` (this is required by the System V Application Binary Interface). Since it is not always easy to ensure stack alignment when calling library functions (because of intermediate computations temporarily stored on the stack), it may be convenient to introduce wrappers around library functions, as follows:

```
malloc_
    pushq    %rbp
    movq     %rsp, %rbp
    andq     $-16, %rsp # 16-byte stack alignment
    call     malloc
    movq     %rbp, %rsp
    popq     %rbp
    ret
```

These wrappers are simply concatenated to the generated assembly code — and of course any call to `malloc` is replaced with a call to `malloc_`.

Here is a list of functions from the C standard library that you may want to use (feel free to use any other):

```
void *malloc(size_t size);
    malloc(n) returns a pointer to a freshly heap-allocated block of size n.
    You don't have to free memory.

void *calloc(size_t nmemb, size_t size);
    calloc(n, s) returns a pointer to a freshly heap-allocated block of size  $n \times s$ .
    The memory is set to zero. You don't have to free memory.

int printf(const char *format, ...);
    printf(f,...) writes to standard output according to the format string
    (ignore the return value). Register %rax must be set to zero before calling printf.
```

Important Notice. Grading involves (for one part only) some automated tests using small Go programs with `fmt.Print` commands. They are compiled with your compiler, and the output is compared to the expected output. This means you should be careful in compiling calls to `fmt.Print`.

4 Project Assignment (due March 15, 6pm)

The project must be done **alone or in pair**. It must be delivered **on Moodle**, as a compressed archive containing a directory with your name(s) (*e.g.* `dupont-durand`). Inside this directory, source files of the compiler must be provided (no need to include compiled files). The command `make` must create the compiler, named `minigo`. The compilation may involve any tool (such as `dune` for OCaml) and the `Makefile` can be as simple as a call to such a tool. The command `minigo` may be a script to run the compiler, for instance if the compiler is not implemented in OCaml.

The archive must also contain **a short report** explaining the technical choices and, if any, the issues with the project and the list of whatever is not delivered. The report can be in format ASCII, Markdown, or PDF.

The command line of `minigo` accepts an option (among `--debug`, `--parse-only` and `--type-only`) and exactly one file with extension `.go`. If the file is parsed successfully, the compiler must terminate with code 0 if option `--parse-only` is on the command line. Otherwise, the compiler moves to static type checking. Any type error must be reported as follows:

```
file.go:4:6:
bad arity for function f
```

The location indicates the filename name, the line number, and the column number. Feel free to design your own error messages. The exit code must be 1.

If the file is type-checked successfully, the compiler must exit with code 0 if option `--type-only` is on the command line. Otherwise, the compiler generates x86-64 assembly code in file `file.s` (same name as the input file, but with extension `.s` instead of extension `.go`). The x86-64 file will be compiled and run as follows

```
gcc file.s -o file
./file
```

possibly with option `-no-pie` on the `gcc` command line. The output must be identical to that of the command

```
go run file.go
```

For your convenience, the expected output is provided in `file.out` for each test file.

Tips. It is strongly recommended to proceed construction by construction, whether for typing or code generation, in this order: output, arithmetic, local variables, functions, structures.