

Gopher Academy Blog



Robin Eklind

☆ Dec 19, 2018

🕒 14 min read

LLVM IR and Go

In this post, we'll look at how to build Go programs – such as compilers and static analysis tools – that interact with the LLVM compiler framework using the LLVM IR assembly language.

TL;DR we wrote a library for interacting with LLVM IR in pure Go, see links to [code](#) and [example projects](#).

1. [Quick primer on LLVM IR](#)
2. [LLVM IR library in pure Go](#)
3. [Closing notes](#)
4. [Further resources](#)

Quick primer on LLVM IR

(For those already familiar with LLVM IR, feel free to [jump to the next section](#)).

[LLVM IR](#) is a low-level intermediate representation used by the [LLVM compiler framework](#). You can think of LLVM IR as a platform-independent assembly language with an infinite number of function local registers.

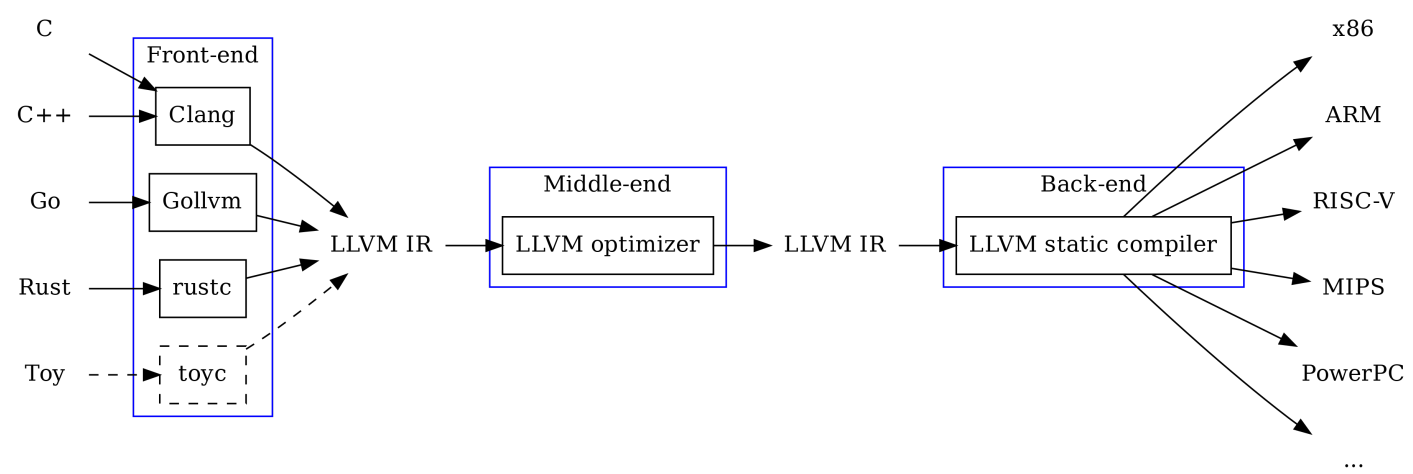
When developing compilers there are huge benefits with compiling your source language to an intermediate representation (IR)¹ instead of compiling directly to a target architecture (e.g. x86). As many optimization techniques are general (e.g. dead code elimination, constant propagation), these optimization passes may be performed directly on the IR level and thus shared between all targets².

Compilers are therefore often split into three components, the front-end, middle-end and back-end; each with a specific task that takes IR as input and/or produces IR as output.

Front-end: compiles source language to IR.

Middle-end: optimizes IR.

Back-end: compiles IR to machine code.



Example program in LLVM IR assembly

To get a glimpse of what LLVM IR assembly may look like, let's consider the following C program.

```

1 int f(int a, int b) {
2     return a + 2*b;
3 }
4
5 int main() {
6     return f(10, 20);
7 }
  
```

Using [Clang³](#), the above C code compiles to the following LLVM IR assembly.

```

1 define i32 @f(i32 %a, i32 %b) {
2     ; <label>:0
3     %1 = mul i32 2, %b
4     %2 = add i32 %a, %1
5     ret i32 %2
6 }
7
8 define i32 @main() {
9     ; <label>:0
10    %1 = call i32 @f(i32 10, i32 20)
11    ret i32 %1
12 }
  
```

By looking at the LLVM IR assembly above, we may observe a few noteworthy details about LLVM IR, namely:

LLVM IR is statically typed (i.e. 32-bit integer values are denoted with the `i32` type).

Local variables are scoped to each function (i.e. `%1` in the `@main` function is different from `%1` in the `@f` function).

Unnamed (temporary) registers are assigned local IDs (e.g. `%1`, `%2`) from an incrementing counter in each function.

Each function may use an infinite number of registers (i.e. we are not limited to 32 general purpose registers).

Global identifiers (e.g. `@f`) and local identifiers (e.g. `%a`, `%1`) are distinguished by their prefix (`@` and `%`, respectively).

Most instructions do what you'd think, `mul` performs multiplication, `add` addition, etc.

Line comments are prefixed with `;` as is quite common for assembly languages.

The structure of LLVM IR assembly

The contents of an LLVM IR assembly file denotes a [module](#). A module contains zero or more top-level entities, such as [global variables](#) and [functions](#).

A function declaration contains zero basic blocks and a function definition contains one or more basic blocks (i.e. the body of the function).

A more detailed example of an LLVM IR module is given below, including the global definition `@foo` and the function definition `@f` containing three basic blocks (`%entry`, `%block_1` and `%block_2`).

```
1  ; Global variable initialized to the 32-bit integer value 21.
2  @foo = global i32 21
3
4  ; f returns 42 if the condition cond is true, and 0 otherwise.
5  define i32 @f(i1 %cond) {
6      ; Entry basic block of function containing zero non-branching instructions and
7      ; a
8      ; conditional branching terminator instruction.
9      entry:
10         ; The conditional br terminator transfers control flow to block_1 if %cond
11         ; is true, and to block_2 otherwise.
12         br i1 %cond, label %block_1, label %block_2
13
14     ; Basic block containing two non-branching instructions and a return
15     ; terminator.
16     block_1:
17         %tmp = load i32, i32* @foo
18         %result = mul i32 %tmp, 2
19         ret i32 %result
20
21     ; Basic block with zero non-branching instructions and a return terminator.
22     block_2:
23         ret i32 0
24 }
```

Basic block

A [basic block](#) is a sequence of zero or more non-branching instructions followed by a branching instruction (referred to as the terminator instruction). The key idea behind a basic block is that if a single instruction of the basic block is executed, then all instructions of the basic block are executed. This notion simplifies control flow analysis.

Instruction

An instruction is a non-branching LLVM IR instruction, usually performing a computation or accessing memory (e.g. [add](#), [load](#)), but not changing the control flow of the program.

Terminator instruction

A [terminator instruction](#) is at the end of each basic block, and determines where to transfer control flow once the basic block finishes executing. For instance [ret](#) terminators returns control flow back to the caller function, and [br](#) terminators branches control flow either conditionally or unconditionally.

Static Single Assignment form

One very important property of LLVM IR is that it is in [SSA](#)-form (Static Single Assignment), which essentially means that each register is assigned exactly once. This property simplifies data flow analysis.

To handle variables that are assigned more than once in the original source code, a notion of [phi](#) instructions are used in LLVM IR. A [phi](#) instruction essentially returns one value from a set of incoming values, based on the control flow path taken during execution to reach the phi instruction. Each incoming value is therefore associated with a predecessor basic block.

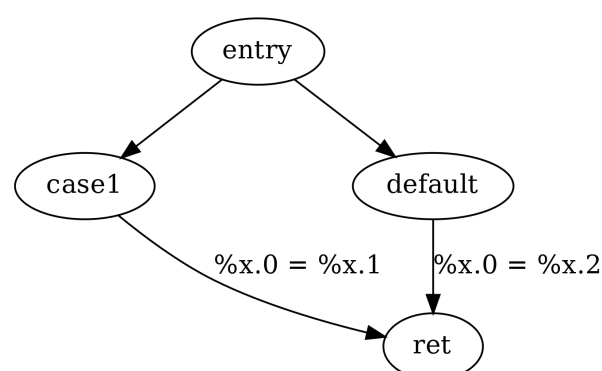
For a concrete example, consider the following LLVM IR function.

```

1  define i32 @f(i32 %a) {
2  ; <label>:0
3      switch i32 %a, label %default [
4          i32 42, label %case1
5      ]
6
7  case1:
8      %x.1 = mul i32 %a, 2
9      br label %ret
10
11 default:
12     %x.2 = mul i32 %a, 3
13     br label %ret
14
15 ret:
16     %x.0 = phi i32 [ %x.2, %default ], [ %x.1, %case1 ]
17     ret i32 %x.0
18 }

```

The [phi](#) instruction (sometimes referred to as [phi](#) nodes) in the above example essentially models the set of possible incoming values as distinct assignment statements, exactly one of which is executed based on the control flow path taken to reach the basic block of the [phi](#) instruction during execution. One way to illustrate the corresponding data flow is as follows:



In general, when developing compilers which translates source code into LLVM IR, all local variables of the source code may be transformed into SSA-form, with the exception of variables of which the address is taken.

To simplify the implementation of LLVM front-ends, one recommendation is to model local variables in the source language as memory allocated variables (using [alloca](#)), model assignments to local variables as [store](#) to memory, and uses of local variables as [load](#) from memory. The reason for this is that it may be non-trivial to directly translate a source language into LLVM IR in SSA-form. As long as the memory accesses follows certain patterns, we may then rely on the `mem2reg` LLVM optimization pass to translate memory allocated local variables to registers in SSA-form (using `phi` nodes where necessary).

LLVM IR library in pure Go

The two main libraries for working with LLVM IR in Go are:

llvm.org/llvm/bindings/go/llvm: the official LLVM bindings for the Go programming language.

github.com/llir/llvm: a pure Go library for interacting with LLVM IR.

The official LLVM bindings for Go uses Cgo to provide access to the rich and powerful API of the LLVM compiler framework, while the `llir/llvm` project is entirely written in Go and relies on LLVM IR to interact with the LLVM compiler framework.

This post focuses on `llir/llvm`, but should generalize to working with other libraries as well.

Why write a new library?

The primary motivation for developing a pure Go library for interacting with LLVM IR was to make it more fun to code compilers and static analysis tools that rely on and interact with the LLVM compiler framework. In part because the compile time of projects relying on the official LLVM bindings for Go could be quite substantial (Thanks to [@aykevl](#), the author of [TinyGo](#), there are now ways to speed up the compile time by dynamically linking against a system-installed version of LLVM⁴).

Another leading motivation was to try and design an idiomatic Go API from the ground up. The main difference between the API of the LLVM bindings for Go and `llir/llvm` is how LLVM values are modelled. In the LLVM bindings for Go, LLVM values are modelled as [a concrete struct type](#), which essentially contains every possible method of every possible LLVM value. My personal experience with using this API is that it was difficult to know what subsets of methods you were allowed to invoke for a given value. For instance, to retrieve the Opcode of an instruction, you'd invoke the [InstructionOpcode](#) method – which is quite intuitive. However, if

Share



you happen to invoke the [Opcode](#) method instead (which is used to retrieve the Opcode of constant expressions), you'd get the runtime errors “*cast<Ty>() argument of incompatible type!*”.

The `llir/llvm` library was therefore designed to provide compile time guarantees by further relying on the Go type system. LLVM values in `llir/llvm` are modelled as [an interface type](#). This approach only exposes the minimum set of methods shared by all values, and if you want to access more specific methods or fields, you'd use a type switch (as illustrated in the [analysis example](#) below).

Usage examples

Now, let's consider a few concrete usage examples. Given that we have a library to work with, what may we wish to do with LLVM IR?

Firstly, we may want to *parse* LLVM IR produced by other tools, such as Clang and the LLVM optimizer `opt` (see the [input example](#) below).

Secondly, we may want to *process* LLVM IR to perform analysis of our own (e.g. custom optimization passes) or implement interpreters and Just-in-Time compilers (see the [analysis example](#) below).

Thirdly, we may want to *produce* LLVM IR to be consumed by other tools. This is the approach taken when developing a front-end for a new programming language (see the [output example](#) below).

Input example - Parsing LLVM IR

```
1 // This example program parses an LLVM IR assembly file, and prints the parsed
2 // module to standard output.
3 package main
4
5 import (
6     "fmt"
7
8     "github.com/llir/llvm/asm"
9 )
10
11 func main() {
12     // Parse LLVM IR assembly file.
13     m, err := asm.ParseFile("foo.ll")
14     if err != nil {
15         panic(err)
16     }
17     // process, interpret or optimize LLVM IR.
18
19     // Print LLVM IR module.
20     fmt.Println(m)
21 }
```

Analysis example - Processing LLVM IR


```

1 // This example program analyses an LLVM IR module to produce a callgraph in
2 // Graphviz DOT format.
3 package main
4
5 import (
6     "bytes"
7     "fmt"
8     "io/ioutil"
9
10    "github.com/llir/llvm/asm"
11    "github.com/llir/llvm/ir"
12 )
13
14 func main() {
15     // Parse LLVM IR assembly file.
16     m, err := asm.ParseFile("foo.ll")
17     if err != nil {
18         panic(err)
19     }
20     // Produce callgraph of module.
21     callgraph := genCallgraph(m)
22     // Output callgraph in Graphviz DOT format.
23     if err := ioutil.WriteFile("callgraph.dot", callgraph, 0644); err != nil {
24         panic(err)
25     }
26 }
27
28 // genCallgraph returns the callgraph in Graphviz DOT format of the given LLVM
29 // module.
30 func genCallgraph(m *ir.Module) []byte {
31     buf := &bytes.Buffer{}
32     buf.WriteString("digraph {\n")
33     // For each function of the module.
34     for _, f := range m.Funcs {
35         // Add caller node.
36         caller := f.Ident()
37         fmt.Fprintf(buf, "\t%q\n", caller)
38         // For each basic block of the function.
39         for _, block := range f.Blocks {
40             // For each non-branching instruction of the basic block.
41             for _, inst := range block.Insts {
42                 // Type switch on instruction to find call instructions.
43                 switch inst := inst.(type) {
44                     case *ir.InstCall:
45                         callee := inst.Callee.Ident()
46                         // Add edges from caller to callee.
47                         fmt.Fprintf(buf, "\t%q -> %q\n", caller, callee)
48                     }
49             }
50             // Terminator of basic block.
51             switch term := block.Term.(type) {
52                 case *ir.TermRet:
53                     // do something.
54                     _ = term
55                 }
56             }
57         }
58     }
59     buf.WriteString("}")
60     return buf.Bytes()
61 }

```

Output example - Producing LLVM IR

```

1 // This example produces LLVM IR code equivalent to the following C code,
2 which
3 // implements a pseudo-random number generator.
4 //
5 //     int abs(int x);
6 //
7 //     int seed = 0;
8 //
9 //     // ref: https://en.wikipedia.org/wiki/Linear_congruential_generator
10 //     //     a = 0x15A4E35
11 //     //     c = 1
12 //     int rand(void) {
13 //         seed = seed*0x15A4E35 + 1;
14 //         return abs(seed);
15 //     }
16 package main
17
18 import (
19     "fmt"
20
21     "github.com/llir/llvm/ir"
22     "github.com/llir/llvm/ir/constant"
23     "github.com/llir/llvm/ir/types"
24 )
25
26 func main() {
27     // Create convenience types and constants.
28     i32 := types.I32
29     zero := constant.NewInt(i32, 0)
30     a := constant.NewInt(i32, 0x15A4E35) // multiplier of the PRNG.
31     c := constant.NewInt(i32, 1)         // increment of the PRNG.
32
33     // Create a new LLVM IR module.
34     m := ir.NewModule()
35
36     // Create an external function declaration and append it to the module.
37     //
38     //     int abs(int x);
39     abs := m.NewFunc("abs", i32, ir.NewParam("x", i32))
40
41     // Create a global variable definition and append it to the module.
42     //
43     //     int seed = 0;
44     seed := m.NewGlobalDef("seed", zero)
45
46     // Create a function definition and append it to the module.
47     //
48     //     int rand(void) { ... }
49     rand := m.NewFunc("rand", i32)
50
51     // Create an unnamed entry basic block and append it to the `rand`
52     function.
53     entry := rand.NewBlock("")
54
55     // Create instructions and append them to the entry basic block.
56     tmp1 := entry.NewLoad(seed)
57     tmp2 := entry.NewMul(tmp1, a)
58     tmp3 := entry.NewAdd(tmp2, c)
59     entry.NewStore(tmp3, seed)
60     tmp4 := entry.NewCall(abs, tmp3)
61     entry.NewRet(tmp4)
62
63     // Print the LLVM IR assembly of the module.

```



```
fmt.Println(m)
}
```

Closing notes

The design and implementation of [llir/llvm](#) has been guided by a community of people who have contributed – not only by writing code – but through shared discussions, pair-programming sessions, bug hunting, profiling investigations, and most of all, a curiosity for learning and taking on exciting challenges.

One particularly challenging part of the [llir/llvm](#) project has been to construct [an EBNF grammar for LLVM IR](#) covering the *entire* LLVM IR assembly language as of LLVM v7.0. This was challenging, not because the process itself is difficult, but because there existed no official grammar covering the entire language. Several community projects have attempted to define a formal grammar for LLVM IR assembly, but these have, to the best of our knowledge, only covered subsets of the language.

The exciting part of having a grammar for LLVM IR is that it enables a lot of interesting projects. For instance, generating syntactically valid LLVM IR assembly to be used for fuzzing tools and libraries consuming LLVM IR (the same approach as taken by [GoSmith](#)). This could be used for cross-validation efforts between LLVM projects implemented in different languages, and also help tease out potential security vulnerabilities and bugs in implementations.

The future is bright, happy hacking!

Further resources

There is a very well written [chapter about LLVM](#) by Chris Lattner – who wrote the initial design of LLVM – in the Architecture of Open Source Applications book.

The [Implement a language with LLVM](#) tutorial – often referred to as the *Kaleidoscope* tutorial – provides great detail on how to implement a simple programming language that compiles to LLVM IR. It goes through the main tasks involved in writing a front-end for LLVM, including lexing, parsing and code generation.

For anyone interested in writing compilers targeting LLVM IR, the [Mapping High Level Constructs to LLVM IR](#) gitbook is warmly recommended.

A good set of slides is [LLVM, in Great Detail](#), which provides an overview of important concepts in LLVM IR, gives an introduction to the LLVM C++ API, and in particular describes very useful LLVM optimization passes.

The [official Go bindings for LLVM](#) is a good fit for many projects, as they expose the LLVM C API which is very powerful and also quite stable.

A good complement to this post is the article [An introduction to LLVM in Go](#).

- 1. The idea of using an IR in compilers is wide spread. GCC uses [GIMPLE](#), Roslyn uses [CIL](#), and LLVM uses [LLVM IR](#). [\[return\]](#)
- 2. Using an IR thus reduces the number of compiler combinations required for n source languages (front-ends) and m target architectures (back-ends) from $n * m$ to $n + m$. [\[return\]](#)
- 3. Compile C to LLVM IR using: `clang -S -emit-llvm -o foo.ll foo.c`. [\[return\]](#)
- 4. The [github.com/aykevl/go-llvm](#) project provides Go bindings to a system-installed LLVM. [\[return\]](#)

[« Introducing Glot the plotting library for Golang](#)

[Go and Apache Arrow: building blocks for data science »](#)

Explore →

- [bleve](#)
- [blog](#)
- [community](#)
- [deployment](#)
- [git](#)
- [go](#)
- [gopher-academy](#)
- [gophercon](#)
- [meta](#)
- [releases](#)
- [revel](#)
- [search](#)
- [text](#)