

Implementation and comparison of register allocation to translate LLVM-- to x86

William Welle Tange

2023

Contents

1	Introduction	2
2	Intermediate Representation	2
2.1	Static Single-Assignment form	3
3	Control Flow Analysis	5
3.1	Building a CFG	5
4	Liveness Analysis	5
4.1	Dataflow Analysis	6
4.2	Interference Graph	8
5	Graph Coloring	9
5.1	Coloring by simplification	9
5.2	Coalescing	9
6	Linear Scan	9
7	Instruction Selection	9
7.1	LLVM-- instruction set	9
7.2	Translating LLVM-- to x86	10
7.2.1	Implementing call instructions	10
7.2.2	Implementing phi instructions	10
7.3	Assessing Correctness	11
7.3.1	Debugger Harness	11
8	Evaluation	12
8.1	Benchmarking	12
8.1.1	benches/fib.ll	12
8.1.2	benches/ackermann.ll	13
8.1.3	benches/factori32.ll	13
8.1.4	benches/factori64.ll	14
8.1.5	benches/sieven.ll	14
8.1.6	benches/subset.ll	14
8.1.7	benches/sha256.ll	15
8.1.8	benches/fannkuch-redux.ll	15
8.2	Comparison to other work and ideas for future work	15
9	Conclusion	15

1 Introduction

Compilation refers to the process of translating from one language to another, most often from a high-level programming language intended for humans to work with, to machine- or bytecode intended to be executed on a target architecture. This process can be divided into several distinct phases, which are grouped into one of two stages colloquially referred to as the *frontend* and *backend* (see Figure 1). The former is translating a high-level programming language to an *intermediate representation* (IR) and the latter is translating IR to executable machine code of a target architecture or bytecode of a target *virtual machine* (VM).



Figure 1: Compiler phases, backend highlighted

Most operations of a general-purpose programming language are translated to a set of control, logic, and arithmetic instructions to be executed sequentially on a computer processor: a single circuit/chip, referred to as the *central processing unit* (CPU), the design of which has varied and evolved over time.

Most CPUs are *register machines*, in that they use a limited set of *general-purpose registers* (GPRs) to store working values in combination with *random access memory* (RAM) for mid-term, and other I/O peripherals for long-term storage. This can largely be attributed to performance, as register machines routinely outperform *stack machines* [1] that are often used in VMs. Although register machines generally also have a stack available, as being limited to mere bytes of storage is simply infeasible for most large scale applications, using it compared to GPRs is orders of magnitudes slower [2]. Because of this, a crucial part of the backend stage for an optimizing compiler is assigning each variable of the source program to a GPR in such a way that maximizes performance without sacrificing correctness.

The process of assigning each variable to a GPR is referred to as *register allocation*, and can be approached in several different ways. This paper will seek to implement graph coloring and linear scan and evaluate them in terms of runtime performance after compilation. The primary sources will be *Modern Compiler Implementation in ML* [3] and *Compilers: Principles, techniques, and tools* [4], in addition to publications concerning the linear scan approach.

The implementation is written in OCaml for the most part, although there is some marginal use of Python to interact with the LLDB debugger.

2 Intermediate Representation

Although compilers can feasibly translate from the source language to machine code of the target architecture directly, which could even be more efficient, doing so hinders the portability as machine code targeting architecture *A* isn't necessarily useful for targeting architecture *B* further down the line.

Instead of translating directly from a source language to target machine code, most compiler frontends emit IR which is intended as an abstraction over particular CPU implementations. It is by no means executable on either *A* or *B*, but it is much closer to the instruction set executed on a CPU than the source language originally fed to the frontend. This helps portability immensely, as a frontend emitting an IR no longer needs to specialize to a particular architecture, instead it can target as many as the backend supports. Likewise, a compiler backend will support any frontend provided that they emit correct IR.

Because of this, compilation of high-level programming languages will often go no further than emitting IR, then leave the rest for a subsequent backend toolchain. This naturally leads to the LLVM toolchain, which is by far the most used compiler backend in practice. Most languages, if targeting native execution on a CPU, will have an LLVM implementation. This is true for C/C++ (Clang compiler [5]), D (LDC compiler [6]), Swift (official Swift compiler [7]) and Rust (official `rustc` compiler [8]). The use of LLVM allows for them to be compiled to completely unforeseen architectures, like web browsers with WebAssembly or even graphics cards with compute kernels [9].

Conversely, managed languages such as Java and C# target their own respective VMs, so they typically don't have much to gain from the LLVM toolchain. That said, some projects have attempted to implement LLVM in their translation of bytecode to native architectures, such as LLILC which promised both *just-in-time* (JIT) as well as *ahead-of-time* (AOT) compilation [10]. Similarly, Kotlin (the spiritual successor of Java) that ordinarily targets the JVM has a native backend that compiles directly to machine code using LLVM, circumventing the need for a VM entirely [11].

For example, the add function implemented in C as seen in Figure 2 is compiled to the LLVM IR as seen in Figure 3 after stripping optimization attributes with the `strip` utility. At this level of complexity they are largely equivalent, with the primary difference appearing to be syntactic. LLVM IR keeps the typed binary operations as well as the function construct with a return statement, although the variable names are discarded as IR is generally not for human interaction, except the function name, as this is used for linking with other object files.

```
int add(int a, int b) {
    int sum = a + b;
    return sum;
}
```

Figure 2: Arithmetic function implemented in C

```
define i32 @add (i32 %0, i32 %1) {
    %3 = add i32 %1, %0
    ret i32 %3
}
```

Figure 3: Stripped `clang -O1 -S -emit-llvm`

```
.text
.globl      add
add:
    movl    %esi, %eax
    addl    %edi, %eax
    retq

# -- Begin function add
# @add
```

Figure 4: Output of `clang -S` without debug markers

When the IR is translated to x86 as seen in Figure 4, the virtual variables `%0`, `%1` and `%2` are assigned the registers `%esi`, `%edi` and `%eax` respectively. This is presumably because of several optimizations, as `%rdi` and `%rsi` are used to pass the first and second parameters and the `%rax` as the return value according to the System V AMD64 ABI [12], so by assigning the variables to those, several `movs` are saved. Either way, the assignment is sound and the instruction selection is reasonably close to the source code. There are other changes as well, as marking the target section with `.text` and the `.globl` annotation to export `add` in the symbol table which permits linking with other object files, but that isn't strictly related to register allocation.

2.1 Static Single-Assignment form

In compiler backends, *Static Single-Assignment* (SSA) form is a property that applies to some IR, including LLVM. It is a way of representing a program such that each variable is assigned only once in the scope of a function. Each use of the variable after definition then refers to the value of that single assignment, and any operations are applied by assigning the result of said operation to a freshly defined variable.

This property eliminates redefinitions entirely, making it easier to reason about data flow and in turn also speed up analysis. It's important to note that SSA form itself doesn't inherently mandate immutability, despite only allowing for single assignment. Rather, it provides two approaches to mutable variables: the phi nodes and load/store operations.

Building on the earlier example, consider multiplication as implemented in Figure 5. Translated to LLVM IR with no optimization can be seen in Figure 7, which achieves mutability by allocating a variable on the stack (lines 2-5), loading the value on use and storing on reassignment. This means new values are loaded and assigned upon reentering a block, without reassigning the value of the source variable, which points to whatever slot was available at the time of definition.

```
extern int add(int a, int b);

int mul(int a, int b) {
    int product = 0;
    for (int i = 0; i < a; i++) {
        product = add(product, b);
    }
    return product;
}
```

Figure 5: Multiplication function implemented in C

```
declare i32 @add(i32, i32)

define i32 @mul (i32 %0, i32 %1) {
    %3 = icmp sgt i32 %0, 0
    br i1 %3, label %6, label %4

4:
    %5 = phi i32 [0, %2], [%9, %6]
    ret i32 %5

6:
    %7 = phi i32 [%10, %6], [0, %2]
    %8 = phi i32 [%9, %6], [0, %2]
    %9 = call i32 @add (i32 %8, i32 %1)
    %10 = add i32 %7, 1
    %11 = icmp eq i32 %10, %0
    br i1 %11, label %4, label %6
}
```

Figure 6: Stripped clang -O1 -S -emit-llvm

```
define i32 @mul (i32 %0, i32 %1) {
    %3 = alloca i32
    %4 = alloca i32
    %5 = alloca i32
    %6 = alloca i32
    store i32 %0, i32* %3
    store i32 %1, i32* %4
    store i32 0, i32* %5
    store i32 0, i32* %6
    br label %7

7:
    %8 = load i32, i32* %6
    %9 = load i32, i32* %3
    %10 = icmp slt i32 %8, %9
    br i1 %10, label %11, label %18

11:
    %12 = load i32, i32* %5
    %13 = load i32, i32* %4
    %14 = call i32 @add (i32 %12, i32 %13)
    store i32 %14, i32* %5
    br label %15

15:
    %16 = load i32, i32* %6
    %17 = add i32 %16, 1
    store i32 %17, i32* %6
    br label %7

18:
    %19 = load i32, i32* %5
    ret i32 %19
}
```

Figure 7: Stripped clang -O0 -S -emit-llvm

Another approach is that of Figure 6, which uses so-called 'phi nodes' instead. These are a type of instruction that evaluates to a specified operand depending on which predecessor block was executed immediately before. This circumvents the need for the same level of interior mutability as is achieved by reading/writing to memory, without violating the properties of SSA. By copying a value from the end of a predecessor to the beginning of the current block. A Phi node represents a point in the program where the control flow merges, and it selects the appropriate version of a variable based on the path taken. This ensures that the data flow is well-defined and allows for easy analysis across different control flow paths.

3 Control Flow Analysis

The backend of a compiler takes some form of IR as input, usually a linear sequence of instructions for each separate function. This representation is close to the level of an actual processor by design, but it isn't immediately useful for the further analysis steps needed to generate optimized code for the target architecture. The control flow of a given program refers to the order in which instructions are executed. While the flow of most instructions is linear, in the sense that the next instruction executed is located immediately after, some transfer the flow of execution elsewhere or outright terminate it.

A continuous flow of instructions is referred to as a *basic block*, defined as a sequence of instructions with no branches in or out except for the first instruction (referred to as a *leader*, immediately following either the function entry or label within it) and the last (referred to as a *terminator*, as it either terminates or transfers the flow of execution). In order to represent the order in which instructions are executed, a *control flow graph* (CFG) is introduced. It is a directed graph whose edges represent transfer of control. The type of node varies over the source material, with CFGs of the Appel text [3] constructed over individual instructions and the Aho text [4] over basic blocks. Either way, leaders will have a set of predecessors and terminators a set of successors associated with them.

Terminators have one or two successors in the case of branching or none at all in the case of function exit. Unconditional branching always transfers control to the block labelled, meaning only one successor follow, whereas conditional branching could transfer to either of the two, but because control flow analysis is not concerned with data, both are considered as possible successors.

Successors of node n are denoted $\text{succ}[n]$, and while each node also has a set of predecessors associated, which consists of an unbounded number of nodes from which control may be transferred, it isn't particularly useful in the following analysis.

3.1 Building a CFG

Building a graph is relatively straight forward, with the input for every function declared parsed as a tuple type named `cfg` of the form:

```
type cfg = (lbl option * block) * (lbl * block) list
```

which consists of an optionally named entry block in the first part and a list of trailing blocks that are always named in the second. The reason for the first block to be optionally named is the fact that *some* LLVM programs transfer control back to the point of entry, whereas every subsequent block needs to be named because branching can only target it by its name.

From this, a CFG is constructed using the OCamlgraph library [13], which provides several approaches with various structures, but the one used in this case is `Imperative.Digraph.Abstract` where the index is parameterized over `int` which corresponds to the index of the instruction starting at 0 for the first instruction to be executed. At first, a graph is constructed with the input flattened such that labels, instructions and terminators are represented as a contiguous sequence as opposed to the basic blocks that is parsed from the input.

```
let flatten ((head, tail) : Ll.cfg) : insn list =
  let label l = Label l and insn i = Insn i in
  let block (b : Ll.block) = List.map insn b.insns @ [ Term b.terminator ] in
  let named_opt (n, b) = (Option.map label n |> Option.to_list) @ block b in
  let named (n, b) = Label n :: block b in
  named_opt head @ (List.map named tail |> List.flatten)
```

4 Liveness Analysis

Translating IR with an unbounded number of variables to a CPU with a bounded number of registers involves the process of assigning each variable a register such that no value that may be needed in the future

is overwritten. Variables that are in use at a given program point are considered *live*, and although variables can be assigned the same register, variables that are live at the same time (i.e. at intersecting program points) cannot, in which case they are also said to be in interference with one another. Variables that are not in interference can be assigned the same register, and finding the precise points at which any variable is live is trivial for linear sequences of instructions. However, when conditional branching is introduced, deriving the path of execution becomes undecidable.

For instance, suppose a function that calls another:

```

1  define i32 @countcall(i32 %x0) {
2      %x1 = add i32 %x0, 1
3      call void @printInt(i32 %x1)
4      call ptr @subproc()
5      ret i32 %x1
6  }
```

Figure 8: Increment, print and possibly return argument

Because of the halting problem, static analysis cannot necessarily determine if the call to `@subproc` will return. So when assigning `%x1` a register it is undecidable whether the variable must be live in the return instruction on line 5 in Figure 8. A register must be assigned to `%x1`, as it is used on line 3, although whether it must live across function calls can reduce the possible registers depending on the target platform and its calling conventions. Because of this, any sound approach to liveness analysis will be an approximation.

In some cases it is simply impossible to assign each variable its own register without conflict along interference edges, in which case one or more variables need to be assigned to memory instead. This is also referred to as *spilling* the variable/register to the stack, or the variable itself is referred to as *spilled*.

A sound albeit very naive approach is to consider every variable live at every program point, such that every variable is in interference with one another, producing a fully connected interference graph. This forces every variable to be assigned a different register, which is a viable approach for small programs with fewer variables than the set of assignable registers. Such a heuristic is greedy in the sense that it picks an assignment known to be safe for the least amount of preprocessing possible. However, this can be a very inefficient assignment at runtime. When the number of variables is greater than the number of assignable registers, the variable is spilled. A number of n variables greater than k working registers causes $n - k$ variables to be spilled to memory, which can greatly reduce performance, but allows for translation in linear time.

4.1 Dataflow Analysis

Another approach, which is a more precise approximation, is a specific variant of the dataflow analysis as described in *Modern Compiler Implementation in ML* ([3]) and *Compilers: Principles, Techniques, and Tools* ([4]). In general, dataflow analysis is the process of finding the possible paths in which data may propagate through various branches of execution. While several applications of this exist (like constant propagation, reaching definitions, available expressions etc.), one that is immediately beneficial in the case of liveness analysis is one that traverses a CFG in the reverse order of execution (i.e. *backwards* flow), and extracts any variable that *may* be used in execution (also referred to as *backwards may* analysis).

This algorithm calculates which program points each variable may be accessed from with some conservative constraints known to maintain correctness. Specifically, these are the transfer and control-flow constraints. The transfer constraint is based on a *transfer function* that describes how liveness is affected across instructions. For each instruction, there is a transfer function that describes how liveness changes from one point to the one immediately after. For example, as an arithmetic operation needs to be assigned a new temporary variable, the liveness of a new variable is propagated to all instructions executed subsequently. This is done by applying the transfer function to the current *live-out* set to stop further propagation of variables defined by the currently visited instruction.

$$in[n] = use[n] \cup (out[n] - def[n]) \quad (1)$$

with the $use[n]$ set being defined as any variables that may be used and the $def[n]$ as the set of variables defined in node n . Because the chosen IR is in SSA form the $def[n]$ either consists of one variable or equal to \emptyset . The $use[n]$ set is effectively unbounded as some instructions take any m variables as parameters.

Control-flow constraints on the other hand propagate the use of variables to previously executed instructions, expecting these to be defined somewhere further up the CFG. This is done by propagating the union of the *live-in* set associated with all immediate successor nodes. This is also referred to as the *meet operator*, whose operator depends on the type of dataflow analysis as well, but for liveness analysis a union is performed on the *live-in* sets of any successive nodes:

$$out[n] = \bigcup_{s \in succ[n]} in[s] \quad (2)$$

Initially, two empty sets are associated with each instruction: the *live-in* and *live-out* sets, which are the sets of variables that are live respectively before and after execution. Then the two equations above are applied iteratively until a fixed point is reached, i.e. an invariant point where neither $in[n]$ or $out[n]$ is changed for all instructions n .

The simplest equations to implement are the `def` and `use` functions, as all of the values of interest are located immediately within the instruction itself and not hidden behind some layer of indirection:

```

1 let def (s : S.SS.t) (insn : Cfg.insn) =
2   match insn with Insn (Some dop, _) -> S.SS.add dop s | _ -> s
3
4 let use (s : S.SS.t) (insn : Cfg.insn) =
5   let op o s = match o with Ll.Id i -> S.SS.add i s | _ -> s in
6   let po = Fun.flip op in
7   match insn with
8   | Insn (_, Allocan (_, (_, o))) (* | Bitcast _ | ... | Zext _ *) ->
9     op o s
10  | Insn (_, Binop (_, _, l, r)) (* | Icmp _ | Store _ *) ->
11    op l s |> op r
12  | Insn (_, Call (_, _, args)) -> List.map snd args |> List.fold_left po s
13  | Insn (_, Gep (_, bop, ops)) -> List.fold_left po (op bop s) ops
14  | Insn (_, Select (c, (_, l), (_, r))) -> op c s |> op l |> op r
15  | Insn (_, PhiNode (_, ops)) -> List.map fst ops |> List.fold_left po s
16  | Term (Ret (_, Some o) | Cbr (o, _, _)) -> op o s
17  | _ -> s

```

Where `S.SS` is a `Set.S` module built over the `symbol` type found in `lib/symbol.ml`:

```

1 type symbol = string * int
2 (* ... *)
3 module SS = Set.Make (struct
4   type t = symbol
5   let compare (_, n1) (_, n2) = compare n1 n2
6 end)
7 type set = SS.t

```

With the actual dataflow analysis performed recursively as follows:

```

1 let dataflow (insns : Cfg.insn list) (ids : Cfg.G.V.t array) (g : Cfg.G.t) =
2   let insns = List.mapi (fun i v -> (i, v)) insns |> List.rev in
3   let in_ = Array.init (List.length insns) (fun _ -> S.SS.empty) in
4   let out = Array.init (List.length insns) (fun _ -> S.SS.empty) in
5   let rec dataflow () =
6     let flowout = (* ... *)

```

```

7   let flowin = (* ... *)
8   let flow changed insn = changed || flowout insn || flowin insn in
9   if List.fold_left flow false insns then dataflow () else (in_, out)
10  in
11  dataflow ()

```

The flowin function corresponds to the *live-in* equation (1) and is implemented as follows:

```

1  let flowin (i, insn) =
2    let newin = S.SS.union (use insn) (S.SS.diff out.(i) (def insn)) in
3    let changed = not (S.SS.equal newin in_.(i)) in
4    if changed then in_.(i) <- newin;
5    changed

```

And the flowout function which corresponds to the *live-out* equation (2) implemented as follows:

```

1  let flowout (i, _) =
2    let newout =
3      let succ = Cfg.G.succ g ids.(i) in
4      List.fold_left
5        (fun s v -> S.SS.union s in_.(Cfg.G.V.label v))
6        S.SS.empty succ
7    in
8    let changed = not (S.SS.equal newout out.(i)) in
9    if changed then out.(i) <- newout;
10   changed

```

To

4.2 Interference Graph

The purpose for conducting dataflow analysis as above is finding variables that may be assigned the same register. This is done by building an interference graph, which is an undirected graph, whose nodes represent variables and edges that signify interference between them, i.e. variables a and b live at overlapping program points is represented with an edge (a, b) .

Constructing an interference graph only depends on the *live-out* set and type of instruction. If the instruction defines a variable, said variable is in interference with all variables in the *live-out* set. There is one exception however: according to the Appel text, move instructions (i.e. phi nodes in the case of SSA form) are given special consideration. The purpose of phi nodes is to copy/move a certain value from a certain predecessor, so they are not necessarily in conflict for being live at the same time. Rather it would often benefit if they were assigned the same register to spare unnecessary moves.

Because of this, for any phi node of the form

$$a = \Phi(b_1, \dots, b_n) \tag{3}$$

add edges to all *live-out* variables not in $\{b_1, \dots, b_n\}$:

$$\forall b_j \in out[i] \setminus \{b_1, \dots, b_n\}, add_edge(a, b_j)$$

For any other instruction that defines a variable a

$$\forall b_j \in \{b_1, \dots, b_n\}, add_edge(a, b_j).$$

Although the Appel text ([3]) describes interference of variables with concrete registers as well as overlapping variables, this isn't considered in this implementation.

5 Graph Coloring

Once an interference graph is constructed, the actual assignments can be found using graph coloring. Although this has long been known to be NP-complete, the heuristic as introduced in both the Appel and Aho et al. texts is a linear time approximation to this problem. It is based on an iterative approach wherein nodes known to be colorable are removed until either an empty graph remains, in which case the original graph G is k -colorable or nodes with more than k neighbours remain. In this case a node is chosen to be spilled to the stack and removed. This is then repeated.

Let k be the number of working registers available on the target architecture. After building the interference graph G , a node n with fewer than k neighbors is chosen. As n has at most $k - 1$ neighbors, it can be removed safely, effectively simplifying G . It is pushed to a stack to preserve the order in which they are removed so that G can be rebuilt and the corresponding registers can be assigned correctly once a k -colorable assignment is found.

5.1 Coloring by simplification

5.2 Coalescing

Coalescing is the process of eliminating moves/copies of data from one GPR to another by combining their interference graph nodes. This is similar to but not the same as the lack of interference edges between variables that are subject to move operations. This is because variables a and b may still be assigned different registers or even spilled if, for instance, either of them are of significant degree. Coalescing joins nodes a and b to node ab preserving the edges of both to maintain soundness.

Since all edges are preserved, the resulting node ab may be of a much higher degree. Because of this, only strategies that produce a k -colorable graph are worth considering, as additional spills negate the purpose entirely.

6 Linear Scan

Another approach that does not depend on liveness analysis is the linear scan algorithm. It is a simple and efficient approach that performs a linear pass over IR to produce a register assignment significantly faster than all of the steps involved with graph coloring.

The primary text is *Linear Scan Register Allocation* ([14]), the execution of which also relies on prior liveness analysis.

7 Instruction Selection

7.1 LLVM-- instruction set

The instruction set used in this paper will be a union of the sets used in the 2022 and 2023 compilers courses in order to work as a drop-in replacement of LLVM for either of the two respective source languages: Tiger and Dolphin. This instruction set is a subset of the one used in practice, as, for instance, neither of the languages implemented support exception handling, floating point operations and so on, and instead only strive to cover the basics of compilers.

The instructions included are `trunc` has only been added in order to help cover more generated LLVM. The branching for most of these is trivial: after successful execution the flow of all but `br`, `ret` and `unreachable` will unconditionally attempt to execute the next instruction in memory (i.e. increment the instruction pointer by instruction length).

1. binary operations `add`, `and`, `ashr`, `lshr`, `mul`, `or`, `sdiv`, `srem`, `shl`, `sub` and `xor`
2. integer comparison `icmp` (conditions being `eq`, `ne`, `sge`, `sgt`, `sle` and `slt`)

3. memory/address operations `alloca`, `gep`, `load` and `store`
4. `mov` operations `gep`, `phi`, `ptrtoint`, `trunc`, and `zext`
5. control flow operations `br`, `call` and `ret`
6. a nullary block terminator `unreachable`

7.2 Translating LLVM-- to x86

Translating each IR instruction to x86 correctly is a matter of eliminating unintended side-effects. Each LLVM- instruction is defined to have only one purpose, as concepts such as calling conventions, stack frames, or a `FLAGS` register are completely abstracted over in order to remain platform independent.

In contrast, the x86 *instruction set architecture* (ISA) is targeting a *complex instruction set computer* (CISC) family of processors, the instructions of which perform a much broader set of operations [3, p. 190]. This is in part due to pipelining, i.e. an abstraction over the concrete implementation of the actual processor, which in turn is made for the sake of performance.

An example of this would be the division/remainder operation: since integer division is a non-trivial iterative process wherein both the quotient and remainder is needed throughout, the result of both of these is stored in the `%rax` and `%rdx` registers respectively. This means two operations are performed simultaneously regardless of which value is used, hence they need to be restored before executing the next instruction, as any variables assigned to `%rax` or `%rdx` will be overwritten.

<code>add</code>	<code>addx</code>	<code>{ }</code>
<code>mul</code>	<code>imul</code>	<code>{ %rax }</code>
<code>sdiv</code>	<code>idivx</code>	<code>{%rax, %rdx}</code>
<code>srem</code>	<code>idivx</code>	<code>{%rax, %rdx}</code>
<code>call</code>	<code>callx</code>	<code>{%rax, caller saved registers}</code>

7.2.1 Implementing `call` instructions

The `call` operation is by far the most complicated one. Mostly because of calling conventions, The `%rax` register is zeroed before

7.2.2 Implementing `phi` instructions

Another irregular instruction is the `phi` instruction representing the Φ -functions. Their intended purpose is to copy data from whichever predecessor node was executed. As opposed to the other instructions that are only concerned with transferring the current state to another, this operation needs to concern itself with what was executed immediately prior. Fortunately there is no need to further complicate the scenario with additional state for the past executed block, as simply moving from `assign[bi]` to `assign[a]` immediately prior to branching into the block whose header contains a phi node of the form (3) will suffice.

```

1  define i32 @main (i32 %0) {
2  1:
3    %2 = icmp sgt i32 %0, 0
4    br i1 %2, label %3, label %9
5  3:
6    %4 = phi i32 [%0, %1], [%7, %3]
7    %5 = phi i32 [0, %1], [%6, %3]
8    %6 = add i32 %5, %0
9    %7 = sub i32 %4, 1
10   %8 = icmp ugt i32 %7, 0
11   br i1 %8, label %3, label %9
12  9:
13   %10 = phi i32 [0, %1], [%6, %3]
14   ret i32 %10
15 }

```

Figure 9: `cat tests/square0.ll`

```

1  _main$3:
2  # ...
3  # br i1 %8, label %3, label %9
4  movq %rdx, %rax
5  cmpq $0, %rax
6  movq %rbx, %rax
7  movq %rax, %rdx
8  je _main$9
9  movq %rsi, %rax
10  movq %rax, %rsi
11  movq %rbx, %rax
12  movq %rax, %rdx
13  jmp _main$3
14  _main$9:
15  # ...

```

Figure 10: `dune exec build -- -t x86`

Notice for instance Figure 9 which has three phi nodes: two in block 3 and one in block 9, i.e. `%0` needs to be copied to `%4` when control is transferred from block 1 and `%7` when transferred from block 3. Likewise the immediate value `$0` must be copied to `%5` when transferred from block 1 and `%6` when transferred from block 3. All of the relevant translation can be seen in the fragment on Figure 10. After the comparison operation on lines 4-6 that check if the boolean condition of `%2` is nonzero and setting `FLAGS` register accordingly, `$0` is moved into `assign[%1]` in preparation of transferring to `_main$9` if the zero flag indeed is set. If not, execution is continued in which case both `assign[%6]` and `assign[%7]` are moved into `assign[%4]` and `assign[%5]` respectively, in preparation for continuing execution at `_main$3`.

As `mov` operations do not change the `FLAGS` register [15], having `mov` instructions between the branching instructions of `je` and `jmp` do not affect the branching. Likewise, because of the `use[n]` semantics of Φ -functions, all b_i are in interference in one another, so may not be overwritten by the `movs` inserted.

7.3 Assessing Correctness

It is difficult to assert the correctness of a compiler. While writing simple unit tests is easy, guaranteeing the correctness of code generated for any given input is on a different order of magnitude in terms of complexity.

Validating the correctness of each instruction translation is a complicated task given how low-level it is. One approach is attaching a debugger and examining the state before and after the translated sequence of x86 instructions is executed. The result of applying a binary operation, for instance, should only affect the assigned register.

7.3.1 Debugger Harness

A more elaborate way of testing correctness of translation is to assert that there are no effects except for those intended. One way of achieving this is to check if a group of x86 instructions representing a single LLVM- operation change any values except those expected. An approach to this is to attach a debugger, insert a breakpoint before every LLVM- instruction, and ensure that only the intended registers are altered.

To achieve this, only a few changes need to be made to the codebase. A debug flag is introduced, which inserts a label before every instruction by adding a `Breakpoint` variant of the `ins` type before the assembly emitted of that specific operation. A unique identifier as well as a list of changes expected to be made is encoded. This includes a bitmask of GPRs that can possibly be changed (aside from the scratch registers of course) as well as which parts of the stack may be modified. In addition to this, each of the labels are output with a `.globl` directive in the file header. A file is built with these extra symbols using the `-d` flag.

To validate that this works as intended,

8 Evaluation

While there are many approaches to assessing the quality of translation, one of the primary means is to simply measure the time taken to execute the code generated. Assuming the translation is sound given the previous steps taken to validate correctness in translating `LLVM--` to x86, it is a solid foundation for future optimizations. While there are certainly other factors worth measuring, only runtime performance is considered for this project.

8.1 Benchmarking

Since the target architecture is x86-64, all benchmarks are executed natively on the fastest processor available to be used consistently at the time of writing. This is to maximize the sample size and in turn reduce uncertainty, but the clock speed at which the benchmarks are performed is irrelevant to how fast the generated code is. Because of this, the metric recorded is the amount of CPU cycles spent executing from start to finish instead of actual time in seconds, as the time measure only works as a more imprecise estimate of the amount of underlying work/execution steps performed on the processor.

CPU cycles isn't a perfect measurement either, but it works to eliminate the clock speed as well as negate much of the time the scheduler allocates for other processes. Scheduling still has a negative impact on execution because of the cache misses caused by context switching, but outside of executing each program in immediate mode (which isn't possible on Linux or macOS) this is a sensible estimate. To further reduce context switching, benchmarks are made on a fresh restart with minimal processes running (Xorg etc.).

Additionally, all benchmarks are deterministic, meaning relevant values at all points across all runs are consistent. So while not completely equivalent due to address randomization, such noise should never leak to value space, and assuming values allocated are initialized before use, no remnants from other processes should leak to value space either. This means that the cycles measured over infinitely many runs will converge towards a 'perfect' run with no cache misses as caused by context switching. Although realistically this will never happen in a lifetime, across n runs the run least affected by context switching will be the one with the least CPU cycles measured, so is the one most representative of the actual performance without noise.

Measurements are made with the `perf` utility, which relies on the *hardware counters* (HC) registers of modern microprocessors, which are special purpose registers meant to record performance data during execution. Because the analysis is integrated into the chip on which it is executing, very little overhead is incurred, and is therefore the go-to for estimating native performance on Linux. The `perf` subcommand is used, as it starts a process and attaches from the very beginning, in addition to padding a `-e` flag specifying that only cycles are meant to be recorded and the `-x` flag to help parse the output. Both min, avg and max are tracked and reported but only min are represented in the matrices below.

Each `.11` file present in the `benches` directory is listed here in a table denoting the input parameter on the left column and type of allocator used.

8.1.1 `benches/fib.11`

The 'dumb' fib is a very naive approach to finding the n th number of the Fibonacci sequence by calling itself recursively twice (decrementing n before each), which causes an exponential growth in function calls.

arg(s)	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
8	139964	1.00041x	0.98897x	1.00518x	1.00397x	0.99405x	1.00683x	1.02014x
10	141471	1.00025x	1.00662x	1.00701x	1.01013x	1.01107x	0.99899x	1.03276x
12	143729	1.01286x	1.00840x	1.02403x	0.97714x	1.00222x	1.02376x	1.03188x
14	152717	1.00032x	1.03967x	1.03472x	1.03637x	1.02969x	1.03224x	1.17572x
16	173395	1.08069x	1.05265x	1.08045x	1.08271x	1.06793x	1.06925x	1.40948x
18	227988	1.15964x	1.15565x	1.16758x	1.16818x	1.16349x	1.16611x	1.80801x
20	375309	1.24200x	1.23989x	1.24662x	1.24412x	1.24000x	1.24497x	2.26177x
22	760328	1.30708x	1.31289x	1.31045x	1.31230x	1.30795x	1.30885x	2.63144x
24	1763965	1.35080x	1.35029x	1.35039x	1.35129x	1.35274x	1.35185x	2.83868x
26	4321515	1.38984x	1.39041x	1.38939x	1.39050x	1.39175x	1.39096x	2.98115x
28	11291147	1.37185x	1.37282x	1.37280x	1.37369x	1.37371x	1.37320x	2.96675x
30	28225827	1.42877x	1.42993x	1.42896x	1.43085x	1.43172x	1.43069x	3.09771x
32	73749658	1.43033x	1.42973x	1.42936x	1.42954x	1.43088x	1.42974x	3.10158x

Table 1: Benchmark of benches/fib.ll output by dune `exec bench -- -f fib -n 1000`

The measures in 1 are significantly equivalent to one another, with all allocators (except greedy 0) seeming to converge towards being 37% slower than Clang. This isn't particularly surprising, as the x86 translation is more concerned with preserving callee saved registers than necessary in pushing each and every one despite not being in use. Also, instead of pushing with the `push` instruction it would be faster to simply subtract from the `%rsp` register directly and using `mov` directly.

What's more interesting and related to the allocation is the fact that greedy 0 is more than twice as slow. This is presumably because the `fib` function that has 7 variables can assign each of them to a register for the first six allocators listed. However, the last one (greedy 0), has no assignable registers, so it must necessarily start spilling from the very beginning. To elucidate this hypothesis, here is the same benchmark but only for the greedy of decreasing k assignable registers:

arg(s)	clang	greedy 12	greedy 8	greedy 6	greedy 4	greedy 2	greedy 1	greedy 0	tiger
8	139813	0.99868x	1.00816x	1.00434x	1.00410x	1.01222x	1.00127x	1.01345x	1.01526x
10	141482	0.99391x	1.01001x	1.01231x	1.01469x	1.00529x	1.02271x	1.03397x	1.02248x
12	145091	1.00774x	0.99627x	1.01584x	1.01047x	1.01694x	1.07056x	1.06622x	1.04396x
14	149598	1.06186x	1.06462x	1.06184x	1.06665x	1.08504x	1.20101x	1.20679x	1.14215x
16	170028	1.09372x	1.09738x	1.12226x	1.13603x	1.15900x	1.42535x	1.43162x	1.29614x
18	230073	1.15757x	1.15498x	1.17517x	1.21433x	1.25770x	1.73753x	1.79427x	1.51243x
20	376308	1.24275x	1.24539x	1.28023x	1.35349x	1.41074x	2.22527x	2.26082x	1.81290x
22	762044	1.30899x	1.30974x	1.35873x	1.45410x	1.52442x	2.57640x	2.62467x	2.05193x
24	1765499	1.34932x	1.34929x	1.40559x	1.51369x	1.59160x	2.78143x	2.83690x	2.18519x
26	4399288	1.36595x	1.36628x	1.42456x	1.53853x	1.62035x	2.87046x	2.92749x	2.24179x
28	11290936	1.37293x	1.37379x	1.43301x	1.54853x	1.63198x	2.90889x	2.96589x	2.27035x
30	28716847	1.40466x	1.40503x	1.46656x	1.58655x	1.67357x	2.98560x	3.04538x	2.32981x
32	76582019	1.37696x	1.37647x	1.43704x	1.55462x	1.64074x	2.92847x	2.98745x	2.28358x

Table 2: Benchmark of benches/fib.ll output by dune `exec bench -- -f fib -n 1000`

8.1.2 benches/ackermann.ll

Another benchmark that derives highly from the recursive structure of `fib.ll` above is the Ackermann function.

8.1.3 benches/factori32.ll

Prime factorization of `i32`.

n	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
268435399	1.00000x	1.11499x	1.11025x	1.11595x	1.11943x	1.66718x	2.47709x	2.44535x
536870909	1.00000x	0.64555x	0.64645x	0.63974x	0.64095x	1.04505x	1.60979x	1.61997x
1073741789	1.00000x	1.18438x	1.19584x	1.19047x	1.18472x	2.09863x	3.36121x	3.37087x
2147483647	1.00000x	1.20872x	1.20665x	1.20405x	1.20826x	2.32725x	3.71253x	3.71693x
4294967291	1.00000x	1.17597x	1.18070x	1.18067x	1.18159x	2.37710x	4.04637x	4.06492x
8589934583	1.00000x	1.28534x	1.29317x	1.29145x	1.28941x	2.73679x	4.74271x	4.76719x
17179869143	1.00000x	1.31790x	1.32297x	1.32585x	1.32452x	2.96961x	5.15949x	5.07055x
34359738337	1.00000x	1.34905x	1.34969x	1.34780x	1.38507x	3.08190x	5.37886x	5.39434x

Table 3: Benchmark of benches/fib.ll output by dune `exec bench -- -f fib -n 1000`

8.1.4 benches/factori64.ll

n	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
268435399	1.00000x	1.11499x	1.11025x	1.11595x	1.11943x	1.66718x	2.47709x	2.44535x
536870909	1.00000x	0.64555x	0.64645x	0.63974x	0.64095x	1.04505x	1.60979x	1.61997x
1073741789	1.00000x	1.18438x	1.19584x	1.19047x	1.18472x	2.09863x	3.36121x	3.37087x
2147483647	1.00000x	1.20872x	1.20665x	1.20405x	1.20826x	2.32725x	3.71253x	3.71693x
4294967291	1.00000x	1.17597x	1.18070x	1.18067x	1.18159x	2.37710x	4.04637x	4.06492x
8589934583	1.00000x	1.28534x	1.29317x	1.29145x	1.28941x	2.73679x	4.74271x	4.76719x
17179869143	1.00000x	1.31790x	1.32297x	1.32585x	1.32452x	2.96961x	5.15949x	5.07055x
34359738337	1.00000x	1.34905x	1.34969x	1.34780x	1.38507x	3.08190x	5.37886x	5.39434x

Table 4: Benchmark of benches/fib.ll output by dune `exec bench -- -f fib -n 1000`

8.1.5 benches/sieven.ll

arg(s)	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
128	159130	1.00777x	1.00293x	1.00644x	1.01562x	1.27104x	1.29856x	1.29894x
256	163800	1.00726x	1.01200x	1.00852x	1.01440x	1.53225x	1.58048x	1.59431x
512	172096	1.01902x	1.02488x	1.01789x	1.01876x	2.03762x	2.13731x	2.15792x
1024	188591	1.03687x	1.03706x	1.02964x	1.03689x	2.93413x	3.09867x	3.15326x
2048	221875	1.04971x	1.06586x	1.04806x	1.06306x	4.34620x	4.64145x	4.72466x
4096	288132	1.08649x	1.09864x	1.08082x	1.10033x	6.25886x	6.69614x	6.84597x
8192	424179	1.10601x	1.12721x	1.10683x	1.12700x	8.26470x	8.85856x	9.07318x
16384	697378	1.12603x	1.15191x	1.12581x	1.15269x	10.00521x	10.73828x	10.98249x

Table 5: Benchmark of benches/sieven.ll output by dune `exec bench -- -f fib -n 1000`

8.1.6 benches/subset.ll

n	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
14	1.00000x	1.08770x	1.15188x	1.08832x	1.15813x	1.92564x	2.59504x	2.95798x
15	1.00000x	1.08304x	1.14485x	1.08536x	1.14973x	1.93885x	2.60300x	2.96001x
16	1.00000x	1.08622x	1.14309x	1.08345x	1.14388x	1.94770x	2.62809x	2.99463x
17	1.00000x	1.08533x	1.14335x	1.08652x	1.14605x	1.94922x	2.64844x	3.01363x

Table 6: Benchmark of benches/sieven.ll output by dune `exec bench -- -f fib -n 1000`

8.1.7 benches/sha256.ll

n	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
100	1.00000x	3.17253x	3.35873x	3.17637x	3.37014x	19.50758x	19.44391x	19.61921x
1000	1.00000x	4.54197x	4.83710x	4.54593x	4.83359x	31.35583x	31.26454x	31.56475x
10000	1.00000x	4.80821x	5.11478x	4.80648x	5.12019x	33.50659x	33.41428x	33.72138x

Table 7: Benchmark of benches/sha256.ll output by dune `exec bench -- -f fib -n 1000`

8.1.8 benches/fannkuch-redux.ll

arg	clang	simple 12	simple 2	briggs 12	briggs 2	linear	greedy 12	greedy 0
4	197975	0.93202x	1.08772x	0.94359x	1.37648x	1.16456x	1.09012x	1.14835x
5	193800	1.05122x	1.10079x	1.07640x	1.11231x	1.72450x	1.78685x	1.73329x
6	307184	1.28426x	1.32022x	1.25973x	1.33286x	3.81882x	3.85798x	3.90989x
7	1312573	1.50533x	1.56355x	1.51619x	1.53672x	6.40226x	6.44050x	6.32867x
8	11252155	1.60713x	1.63817x	1.52961x	1.61191x	7.25974x	7.25458x	7.28693x
9	115513112	1.57771x	1.62155x	1.55882x	2.57223x	7.60550x	7.87522x	7.69289x
10	1352747472	1.54700x	1.74875x	1.58381x	1.67440x	7.77269x	7.74449x	7.82170x

Table 8: Benchmark of benches/sieven.ll output by dune `exec bench -- -f fib -n 1000`

8.2 Comparison to other work and ideas for future work

9 Conclusion

References

- [1] Y. Shi, K. Casey, M. Ertl, and D. Gregg, “Virtual machine showdown: Stack versus registers,” eng, *ACM transactions on architecture and code optimization*, vol. 4, no. 4, pp. 1–36, 2008, ISSN: 1544-3566.
- [2] J. Dean and P. Norvig, “Latency comparison numbers.” (), [Online]. Available: <https://web.archive.org/web/20240101182413/https://gist.github.com/jboner/2841832> (visited on 01/01/2024).
- [3] A. W. Appel, *Modern Compiler Implementation in ML*, eng. Cambridge University Press, 1997, ISBN: 0521582741.
- [4] A. V. Aho, M. S. Ian, R. Sethi, and J. D. Ullmann, *Compilers : principles, techniques, and tools*, eng, 2nd ed., internat. ed. Harlow: Pearson Education Limited, 2014, ISBN: 1292024348.
- [5] L. Foundation. “Clang: A c language family frontend for llvm.” (), [Online]. Available: <https://web.archive.org/web/20231230181618/https://clang.llvm.org/> (visited on 12/30/2023).
- [6] D. L. Foundation. “Ldc – the llvm-based d compiler.” (), [Online]. Available: <https://web.archive.org/web/20231230181935/https://github.com/ldc-developers/ldc> (visited on 12/30/2023).
- [7] A. Inc. “Ldc – the llvm-based d compiler.” (), [Online]. Available: <https://web.archive.org/web/20231230182842/https://developer.apple.com/swift/> (visited on 12/30/2023).
- [8] A. Inc. “Ldc – the llvm-based d compiler.” (), [Online]. Available: <https://web.archive.org/web/20231230185321/https://rustc-dev-guide.rust-lang.org/overview.html> (visited on 12/30/2023).

- [9] L. Foundation. “Writing an llvm backend.” (), [Online]. Available: <https://web.archive.org/web/20231230200453/https://llvm.org/docs/WritingAnLLVMBackend.html> (visited on 12/30/2023).
- [10] . F. Martin Woodward Executive Director. “Announcing llilc - a new llvm-based compiler for .net.” (2015), [Online]. Available: <https://web.archive.org/web/20211212184833/https://dotnetfoundation.org/blog/2015/04/14/announcing-llilc-llvm-for-dotnet> (visited on 09/12/2020).
- [11] K. Foundation. “Kotlin native.” (), [Online]. Available: <https://web.archive.org/web/20231230192142/https://kotlinlang.org/docs/native-overview.html> (visited on 12/30/2023).
- [12] *System v application binary interface amd64 architecture processor supplement (with lp64 and ilp32 programming models) version 1.0.*
- [13] “Ocamlgraph.” (), [Online]. Available: <https://github.com/backtracking/ocamlgraph>.
- [14] M. Poletto and V. Sarkar, “Linear scan register allocation,” eng, *ACM transactions on programming languages and systems*, vol. 21, no. 5, pp. 895–913, 1999, ISSN: 0164-0925.
- [15] “X86 opcode and instruction reference 1.12.” (), [Online]. Available: <https://web.archive.org/web/20240104142515/http://ref.x86asm.net/geek64-abc.html#M>.