

Compiler Design

by dcamenisch

A compiler translates one programming language to another. The simplified compiler has the following structure:

- Lexical Analysis: Source Code \rightarrow Token Stream
- Parsing: Token Stream \rightarrow AST
- Intermediate Code Generation: AST \rightarrow Intermediate Code
- Code Generation: Intermediate Code \rightarrow Target Code

The first two steps are the frontend and machine independent, the last step is the backend and machine dependent.

x86lite

x86lite memory consists of 2^{64} bytes numbered 0x00000000 through 0xffffffff, split into 8-byte quadwords (has to be quadword-aligned).

The stack grows from high addresses to low addresses, **rsp** points to the top of the stack, **rbp** points to the bottom of the current stack frame.

The stack sits at the top of memory space, at the bottom we have code and data followed by the heap.

Registers: **rax, rbx, rcx, rdx, rsi, rdi, rsp, rbp, rip, r08 - r15**

Flags: **OF** overflow/underflow, **SF** sign (1 = negative), **ZF** zero

Condition Codes:

Code	Condition
e (equality)	ZF
ne (not equals)	not ZF
g (strictly greater)	not ZF and SF = OF
l (strictly less)	SF \neq OF
ge (greater or equal)	SF = OF
le (less or equal)	SF \neq OF or ZF

Instructions: **INSTR SRC DEST** (AT&T syntax), prefix register with % and immediate values with \$. Note that **subq** is **DEST - SRC**.

Operands:

- **Imm**: 64-bit literal signed integer
- **Lbl**: label representing a machine address
- **Reg**: one of the registers, the value is its content
- **Ind**: machine address

Ind is **offset(base, index)** is calculated **base + index * 8 + offset**.

Thus, **%rax** refers to the contents of the register, while **(%rax)** refers to either the memory address or the contents of the memory address, depending on whether its used as location or value.

x86 assembly is organized into labeled blocks, indicating code locations used by jumps, etc. Program begins execution at designated label (**main**).

Calling Conventions

- Setup Stack Frame: **pushq %rbp movq %rsp, %rbp**
- Teardown: **popq %rbp**

- Caller Save - freely usable by the called code.
- Callee Save - must be restored by the called code (**rbp, rsp, rbx, r12-15**).
- Arguments: In **rdi, rsi, rdx, rcx, r08, r09** and starting with $n = 7$ in $(n - 7) * 8 + rbp$
- Return value in **rax**.
- 128 byte "red zone" - scratch pad for the callee (beyond **rsp**), this means a function can use up to 128 byte without allocating a stack frame..

Intermediate Representations

Direct translation is bad as it is hard to optimize the resulting assembly code. The representation is too concrete, as it already committed to using certain registers etc. Further retargeting the compiler to a new architecture is hard. Finally control-flow is not structured, arbitrary jumps from one code block to another. Implicit fall-through makes sequences non-modular.

Using a universal IR means that for p programming languages and q ISA's, we only need $p + q$ compilers instead of $p * q$.

IR's allow machine independent code generation and optimization.

Multiple IR's: get program closer to machine code without losing the information needed to do analysis and optimizations (high / mid / low level IR).

Good IR: Easy translation target, easy to translate, narrow interface (fewer constructs means simpler phases / optimizations).

Basic Blocks are a sequence of instructions that are always executed from the first to last instruction. They start with a label and end with a control-flow instruction (no other control-low instruction or label).

Basic blocks can be arranged into a control-flow graph (CFG): Nodes are basic blocks - directed edges represent potential jumps.

LLVM (Low Level Virtual Machine)

Storage Types: local variable **%uid**, global variable **@gid**, abstract locations (stack-allocated with **alloca**), heap-allocated structures (**malloc**).

Each **%uid** appears on the left-hand side of an assignment only once in the entire control flow graph (SSA).

The entry block of the CFG does not have to be labeled, the last instruction of a block is called the terminator.

Example Program:

```
@s = global i32 42
```

```
declare void @use (i64)
```

```
define i64 @foo(i64 %a, i64* %b) {
    %sum = add nsw i64 %a, 42
    %cond = icmp sgt i64 %sum, 100
    br i1 %cond , label %then , label %else
```

```
then:
    call void @use(i64 %sum)
    ret i64 %sum

else:
    store i64 %sum, i64* %b
    ret i64 %sum
}
```

GEP

LLVM supports structured data with the use of **types**, e.g.:

```
%struct.Node = type {i64, %struct.Node*}
```

To computer pointer values of structs or index into arrays, LLVM provides the **getelementptr** instruction. Given a pointer and a path through the structured data pointed to by that pointer, GEP computes an address - analog of LEA.

```
getelementptr <ty>* <ptrval> {, <ty> <idx>}*
```

GEP never dereferences the address it is calculating.

Lexing

Lexing is the process of taking the source code as an input and producing a token stream as output. The problem is to precisely define tokens and matching tokens simultaneously.

One way of implementing a lexer is, using regular expressions. Regex rules precisely describe a sets of strings. But regex alone can be ambiguous if we have multiple matching rules. Most languages therefore choose the longest match or have another specified order.

Regex can be implemented by forming an NFA and then transforming it to a DFA.

Parsing

In this part we take the token stream and generate an abstract syntax tree (AST). Parsing itself does not check things such as variable scoping, type agreement etc.

Parsing uses a more powerful tool than regex - context free grammars (CFG).

Chomsky Hierarchy:

- Regular - Productions have at most one nonterminal and it is at the start or end of the word
- Context-Free (CFG) - LHS of productions only have a single nonterminal
- Context-Sensitive
- Recursively Enumerable

An example for a non CFG would be $a^n b^m c^n d^m$. This corresponds to methods having matching parameters.

A CFG consists of a set of terminals, a set of nonterminals, a start symbol and a set of productions. A production consists of a single nonterminal LHS and an arbitrary RHS.

- Leftmost derivation: Find the left-most nonterminal and apply a production to it
- Rightmost derivation: Find the right-most nonterminal and apply a production there

In CFGs ambiguity can (often) be removed by adding nonterminals and allowing recursion only on one side. For example:

```
S_0 -> S_0 + S_1 | S_1
S_1 -> S_2 * S_1
S_2 -> n | (S_0)
```

$$\begin{aligned} S &\rightarrow b_1 S' \mid \dots \mid b_m S' \\ S' &\rightarrow a_1 S' \mid \dots \mid a_n S' \mid \text{epsilon}\$ \end{aligned}$$

For a given production $A \rightarrow \gamma$:

- Construct the **first set** of A , this set contains all terminals that begin strings derivable from the nonterminal. For each nonterminal of the first set, add the corresponding production to the table.
- Construct the **follow set** of A , this set contains all terminals that can appear immediately to the right of the given nonterminal. If ϵ is derivable by the production, add the corresponding production to the table.

Note: we want the *least* solution to this system of set equations... a *fixpoint* computation. More on these later in the course.

	number	+	()	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{number}$		$\mapsto (S)$		

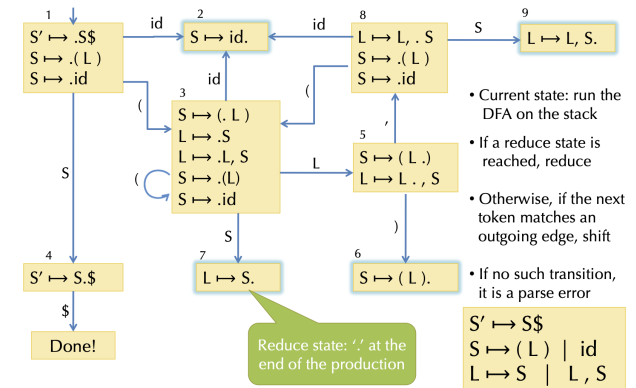
Bottom-up parsing is a sequence of **shift** and **reduce** operations:

- Action Selection Problem:

- Given a stack σ and a lookahead symbol b , should the parser **shift** b onto the stack (new stack is σb), or **reduce** a production $X \mapsto \gamma$, assuming that $\sigma = \alpha\gamma$?
- Sometimes the parser can reduce, but should not, sometimes the stack can be reduced in different ways.

Constructing the DFA:

- Add new production: $S' \rightarrow S\$$, this is the start of the DFA.



To form the LR(1) closure, we first do the same as for LR(0). Additionally for each item $C \rightarrow \cdot \gamma$ we add due to a rule $A \rightarrow \beta.C\gamma$, \mathcal{L} , we compute its look-ahead set \mathcal{M} including $\text{FIRST}(\gamma)$ and if γ can derive ϵ also \mathcal{L} .

$$\begin{aligned} S' &\mapsto S\$ \\ S &\mapsto E + S \mid E \\ E &\mapsto \text{number} \mid (S) \end{aligned}$$

Note: {\$} is FIRST(\$)

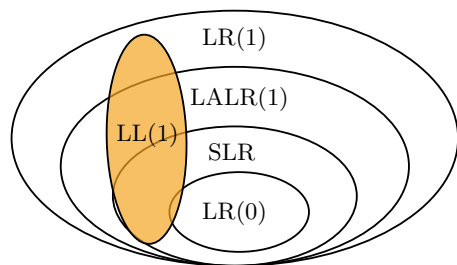
Note: + added for reason 1

$\text{FIRST}(+ S) = \{+\}$

Note: \$ added for reason 2
 δ is ε

All items are distinct, so we're done

For LR(1) we have a shift-reduce conflict if the shifted token is contained in the follow set of the reduction.



FirstClass Functions

Consider the Lambda Calculus. It has variables, functions and function application. The only values are (closed) functions. Instead of `(fun x -> e)` we write: $\lambda x.e$

$x\{v/x\}$	$= v$	<i>(replace the free x by v)</i>
$y\{v/x\}$	$= y$	<i>(assuming $y \neq x$)</i>
$(\text{fun } x \rightarrow \text{exp})\{v/x\}$	$= (\text{fun } x \rightarrow \text{exp})$	<i>(x is bound in exp)</i>
$(\text{fun } y \rightarrow \text{exp})\{v/x\}$	$= (\text{fun } y \rightarrow \text{exp}\{v/x\})$	<i>(assuming $y \neq x$)</i>
$(e_1 e_2)\{v/x\}$	$= (e_1\{v/x\} e_2\{v/x\})$	<i>(substitute everywhere)</i>

Function application is interpreted by substitution. In `fun y -> x + y`, `x` is said to be free and `y` is bound by `fun y`.

A term without free variables is **closed**, else it is **open**. Two terms that differ only by consistent renaming of bound variables are alpha equivalent.

To avoid accidentally capturing a free variable by a substitution $e_1\{e_2/x\}$, we first pick an alpha equivalent version of e_1 such that the bound variables do not mention the free variables of e_2 .

Some special Lambda Calculus terms:

- Omega Term (infinite loop):
 $(\lambda x. x x)(\lambda x. x x)$
- Y-Combinator (computes fixed point so $Yg = g(Yg)$):
 $\lambda f. (\lambda x. f(x x))(\lambda x. f(x x))$

Operational Semantics is a way to give meaning to a program (interpreter) using inference rules. $exp \Downarrow v$ means exp evaluates to v .

Inference rules are of the form $G; L \vdash e : t$. This means in the global environment G and local environment L the expression e is of type t . Sometimes we include a third symbol on the LHS referring to the return type.

With this we can build up derivation or proof trees. Leaves of the tree are axioms.

Typing

Applying a set of inference rules allows for type checking of a program. For simply typed lambda calculus this implies termination, for well-typed expressions we are a bit more general.

A well-typed program either terminates in a well-defined way, or it continues computing forever.

If we view types as sets of values, there is a natural inclusion relation $\text{Pos} \subseteq \text{Int}$. This gives rise to a subtype relation $P <: \text{Int}$ and to a subtyping hierarchy.

The LUB (least upper bound) is defined for two types $T_1 \vee T_2$.

A subtyping rule is sound if it approximates the underlying subset relation, i.e. if $T_1 <: T_2$ implies $[[T_1]] \subseteq [[T_2]]$. It follows that $[[T_1]] \cup [[T_2]] \subseteq [[LUB(T_1, T_2)]]$.

Argument type is contravariant (it is okay if a function takes more arguments), output type is covariant (it is okay if a function returns less arguments).

$$\frac{S_1 <: T_1 \quad T_2 <: S_2}{(T_1 -> T_2) <: (S_1 -> S_2)}$$

For records, we have to decide between **width** and **depth subtyping**. In width subtyping, a record is a subtype of another record if it has more (or the same number of) fields (order matters!). In depth subtyping a record is a subtype of another if every elements type is a subtype of the others.

Mutable structure need to be invariant - else one can break type-safety. Thus, $T \text{ ref } <: S \text{ ref} \implies T = S$.

OAT Type System

Primitive (non-reference) types: `int`, `bool`

Definitely non-null reference types: `R` (named) mutable structs with width subtyping, `strings`, `arrays`

Possibly-null reference types: `R?`

Subtyping: $R <: R?$

Compiling Objects

The dispatch problem occurs when the same interface is implemented by multiple classes. In the client program, it may be necessary to dynamically choose which implementation to use. In order to do this, objects contain a pointer to a **dispatch vector** (vtable) with pointer to method code.

For extension / inheritance, the dispatch vector gets extended at the end.

For multiple inheritance there are different approaches:

- Allow multiple DV tables (C++), choose which DV to use based on static types, casting requires runtime operations.
- Use a level of indirection: Map method identifiers to code pointers using a hash table, search up through the class hierarchy.
- Give up separate compilation: Use sparse dispatch vectors or binary decision trees.

Multiple Dispatch Vectors: Objects may have multiple entry points with individual DVs, casts change entry point of a variable

Optimizations

There are different kinds of optimization: Power, Space, Time.

- **Constant Folding:** If operands are statically known, compute value at compile-time. More general algebraic simplification: Use mathematical identities.
- **Constant Propagation:** If x is a constant replace its uses by the constant.
- **Copy Propagation:** For $x = y$ replace uses of x with y
- **Dead Code Elimination:** If side-effect free code can never be observed, safe to eliminate it.
- **Inlining:** Replace a function call with the body of the function (arguments are rewritten to local variables).
- **Code Specialization:** Create specialized versions of a function that is called from different places with different arguments.
- **Common Subexpression Elimination:** It is the opposite of inlining, fold redundant computations together.
- **Loop Optimizations**
 - Hot spots often occur in loops (esp. inner loops)
 - Loop Invariant Code Motion (hoist outside)
 - Strength Reduction (replace expensive ops by cheap ones by creating a dependent induction variable)
 - Loop Unrolling

Dataflow Analysis

Almost every dataflow analysis is a variation of the following algorithm.

Forward Must Dataflow Analysis

```
for all n, in[n] = T, out[n] = T
repeat until no change in 'in' or 'out'
  for all n
    in[n] = intersect out[n'] for all n' in pred[n]
    out[n] = gen[n] union (in[n] \ kill[n])
```

Backward: swap `in` and `out` and `pred` with `succ`.

May: swap \top with \perp or \emptyset and replace \cup with \cap .

For each dataflow analysis we only need to define the set `gen`, `kill` as well as the domain of dataflow values \mathcal{L} and a combining operator \cup or \cap .

Liveness (Backward, May)

We can use the same registers for multiple `%uids` if they are not alive at the same time. We define `gen[s]` as all the variables used and `kill[s]` as all the variables defined by statement s . \mathcal{L} corresponds to the variables and the combination operator to the set union.

It holds: $\text{in}[n] \supseteq \text{gen}[n]$, $\text{in}[n] \supseteq \text{out}[n] \setminus \text{kill}[n]$ and $\text{out}[n] \supseteq \text{in}[n']$ if $n' \in \text{succ}[n]$.

Reaching Definition (Forward, May)

What variable definitions reach a particular use of a variable? Used for constant and copy propagation. `in` / `out` is the set of nodes defining some variable such that the definition may reach the beginning resp. end of the current node. `gen` is the current statement, if it defines a variable, and `kill` is all the other definitions of the defined variable.

Available Expressions (Forward, Must)

