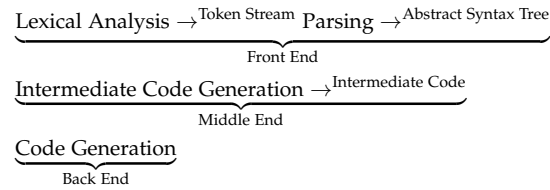


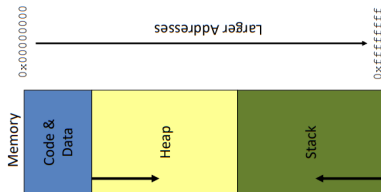
Compiler Design - Cheatsheet

Simplified Compiler Structure

Source Code →



Stack layout



Suffixes

1. q = quadword (4 words)
2. l = long (2 words)
3. w = word (16-bit)
4. b = byte (8-bit)

movq-instruction

movq SRC, DST = DST ← SRC

Flags

e (equality)	ZF is set
ne (inequality)	(not ZF)
g (greater than)	(not ZF) and (SF = OF)
l (less than)	SF <> OF
ge (greater or equal)	(SF = OF)
e (less than or equal)	SF <> OF or Z

Addressing

-8(%rsp)	rsp - 8
(%rax, %rcx)	rax + 8·rcx
8(%rax, %rcx)	rax + 8·rcx + 8

leaq vs. movq

In leaq we just compute the address, in movq we dereference the address.

Callee vs. Caller saved

Caller saved register (Rdi, Rsi, Rdx, Rcx, R09, R08, Rax, R10, R11) are saved by the caller before calling the function. Callee saved register (Rbx, R12, R13, R14, R15) are saved by the called function.

Parameter

1. 1...6: rdi, rsi, rdx, rcx, r8, r9
2. 7+: on the stack (in right-to-left order), nth arg.
 $((n - 7) + 2) \cdot 8 + rbp$

Why Intermediate Representations?

1. resulting code quality is poor (direct translation)
2. Richer source language features are hard to encode (Structured data types, Objects, ...)
3. hard to optimize
4. Control-flow is not structured

Basic Blocks

1. Starts with a label that names the entry point of the basic block
2. Ends with a control-flow instruction
3. Contains no other control-flow instructions
4. Contains no interior label used as a jump target

CFGs

1. Nodes are basic blocks
2. There is a directed edge from node A to node B if the control flow instruction at the end of block A might jump to the label of block B
3. No two blocks have the same label

getelementptr instruction

The first argument is always a type used as the basis for the calculations. The second argument is always a pointer or a vector of pointers, and is the base address to start from. The remaining arguments are indices that indicate which of the elements of the aggregate object are indexed. **GEP never dereferences the address it's calculating.**

```
struct RT {
    char A;
    int B[10][20];
    char C;
};
struct ST {
    int X;
    double Y;
    struct RT Z;
};
```

```
int *foo(struct ST *s) {
    return &s[1].Z.B[5][13];
}
```

is translated in:

```
%arrayidx = getelementptr %struct.ST, ptr %s, i64 1,
i32 2, i32 1, i64 5, i64 13
ret ptr %arrayidx
```

Regular Expressions

Regular expressions precisely describe sets of strings. A regular expression R has one of the following forms:

1. ϵ : Epsilon stands for the empty string
2. 'a': An ordinary character stands for itself
3. $R_1|R_2$: Alternatives, stands for choice of R_1 or R_2
4. R_1R_2 : Concatenation, stands for R_1 followed by R_2
5. R^* : Kleene star, stands for zero or more repetitions of R

Useful extensions:

1. "foo": Strings, equivalent to 'f' 'o' 'o'
2. R^+ : One or more repetitions of R, equivalent to RR^*
3. $R?$: Zero or one occurrences of R, equivalent to $(\epsilon|R)$
4. $[a' - 'z']$: One of a or b or c or ... z, equivalent to $(a|b| \dots |z)$
5. $[^'0' - '9']$: Any character except 0 through 9
6. R as x: Name the string matched by R as x

Chomsky Hierarchy

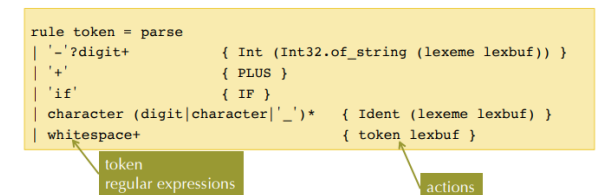
Regular \subset Context-Free \subset Context Sensitive \subset Recursively Enumerable

Matching Rule for Lexer

Most languages choose "longest match".

Lexer Generator

1. Reads a list of regular expressions: R_1, \dots, R_n , one per token
2. Each token has an attached "action" A_i (just a piece of code to run when the regular expression is matched)

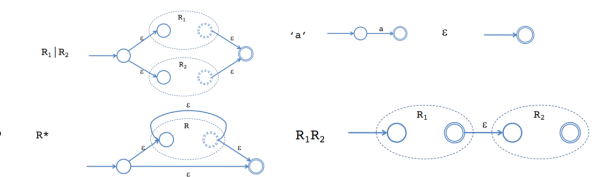


Generates scanning code that:

1. Decides whether the input is of the form $(R_1| \dots |R_n)^*$
2. After matching a (longest) token, runs the associated action

Implementation Strategies

We can use a finite automaton since one must exist for a regular expression.



DFA vs. NFA

DFA:

1. Action of the automaton for each input is fully determined
2. Accepts if the input is consumed upon reaching an accepting state
3. Obvious table-based implementation

NFA:

1. Automaton potentially has a choice at every step
2. Accepts an input if there exists a way to reach an accepting state
3. Less obvious how to implement efficiently

Parsing

Def.: Finding Syntactic Structure

Limits of regular expressions:

1. DFA's have only finite # of states (i.e., finite memory)
2. So, DFA's can't count

Therefore we need something more powerful than DFA's.

CONTEXT FREE GRAMMARS

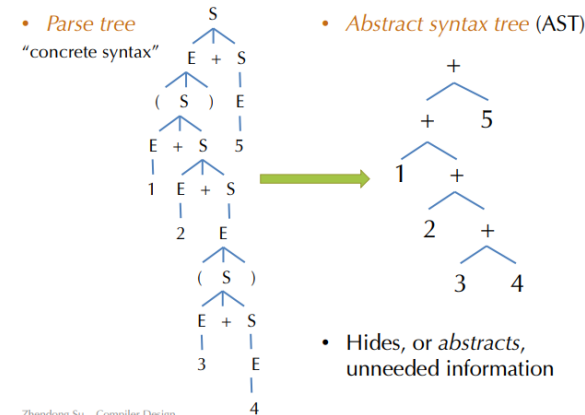
Def. recursive: S mentions itself, e.g. $S \mapsto (S)S$ A Context-free Grammar (CFG) consists of

1. A set of terminals (e.g., a lexical token, ϵ is not a terminal)
2. A set of nonterminals (e.g., S and other syntactic variables)
3. A designated nonterminal called the start symbol
4. A set of productions: $LHS \mapsto RHS$

single step: For arbitrary strings α, β, γ and production rule $A \mapsto \beta$ a single step of the derivation is $\alpha A \gamma \mapsto \alpha \beta \gamma$

Parse Tree: Leaves are terminals and Internal nodes are nonterminals.

Parse Trees to Abstract Syntax



Derivation Orders

Productions of the grammar can be applied in any order. Both strategies (and any other) yield the **same** parse tree!

1. Leftmost derivation: Find the left-most nonterminal and apply a production to it e.g. $S \mapsto \underline{E} + S$
2. Rightmost derivation: Find the right-most nonterminal and apply a production there e.g. $S \mapsto E + \underline{S}$

productive and nonproductive: This grammar has nonterminal definitions that don't mention any terminal symbols.

finite: There is a finite derivation starting from start symbol.

Example of right associative: $S \mapsto E + S | E, E \mapsto \text{number} | (S)$

Example of Ambiguity: $S \mapsto S + S | (S) | \text{number}$, accepts the same set of strings as the previous one but there are two leftmost derivations. Moreover, if there are multiple operations, ambiguity in the grammar leads to ambiguity in their precedence.

Eliminating Ambiguity:

1. by adding nonterminals and allowing recursion only on the left (or right)
2. Higher-precedence operators go farther from the start symbol

Example of Eliminating Ambiguity: $S \mapsto S + S | (S) | \text{number}$ becomes $S_0 \mapsto S_0 + S_1 | S_1, S_1 \mapsto S_2 * S_1 | S_2$ and $S_2 \mapsto \text{number} | (S_0)$

LL(1) GRAMMARS

Top-down: Start from the start symbol (root of the parse tree), and go down. Not all grammars can be parsed top-down with a single lookahead.

LL(1) means:

1. Left-to-right scanning
2. Left-most derivation
3. 1 lookahead symbol

Making a grammar LL(1)

Problem: We can't decide which S production to apply until we see the symbol after the first expression **Solution:** "Left-factor" the grammar. There is a common S prefix for each choice, so add a new non-terminal S' at the decision point.



Also need to eliminate left-recursion.

$S \mapsto S\alpha_1 | \dots | S\alpha_n | \beta_1 | \dots | \beta_m$ becomes $S \mapsto \beta_1 S' | \dots | \beta_m S'$
 $S' \mapsto \alpha_1 S' | \dots | \alpha_n S' | \epsilon$



Predictive Parsing

For a given nonterminal, the lookahead symbol uniquely determines the production to apply. Therefore we get a table with nonterminal \times input token \rightarrow production

Construction of the parse table:

Consider a given production: $A \mapsto \gamma$.

1. (Case 1) Construct the set of all input tokens that may appear first in strings that can be derived from γ . Add the production $\mapsto \gamma$ to the entry (A, token) for each such token.
2. (Case 2) If γ can derive ϵ (the empty string), then we construct the set of all input tokens that may follow the nonterminal A in the grammar. Add the production $\mapsto \gamma$ to the entry (A, token) for each such token.

Example: We have

$T \mapsto S \$, S \mapsto ES', S' \mapsto \epsilon, S' \mapsto +S, E \mapsto \text{number} | (S)$. We get

1. $\text{First}(S\$) = \text{First}(S) = \text{First}(E') = \text{First}(E) = \{ \text{number}, '(' \}$
2. $\text{First}(\epsilon) = \{ \epsilon \}$
3. $\text{First}(+S) = \{ + \}$
4. $\text{First}(\text{number}) = \{ \text{number} \}$
5. $\text{First}((S)) = \{ '(' \}$
6. $\text{Follow}(S') = \text{Follow}(S)$
7. $\text{Follow}(S) = \{ \$, ')' \} \cup \text{Follow}(S')$

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

LR GRAMMARS

LR(k) parser

1. Bottom-up Parsing
2. Left-to-right scanning
3. Rightmost derivation
4. k lookahead symbols

LR grammars are more expressive than LL. Can handle left-recursive (and right recursive) grammars.

Shift/Reduce Parsing

Shift: Move look-ahead token to the stack

Reduce: Replace symbols γ at top of stack with nonterminal X s.t. $X \mapsto \gamma$ is a production

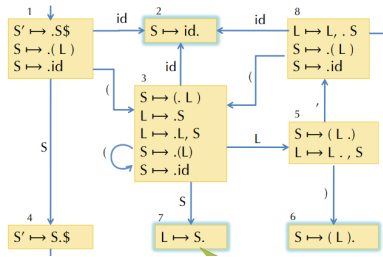
Stack	Input	Action
	$(1 + 2 + (3 + 4)) + 5$	shift (
($1 + 2 + (3 + 4)) + 5$	shift 1
(1	$+ 2 + (3 + 4)) + 5$	reduce: $E \mapsto \text{number}$
(E	$+ 2 + (3 + 4)) + 5$	reduce: $S \mapsto E$

LR(0) state: items to track progress on possible upcoming reductions

LR(0) item: a production with an extra separator "." in the RHS, e.g. $S \mapsto .(L)$ or $S \mapsto (.L)$. Idea is stuff before the "." is already on the stack and stuff after the "." is what might be seen next.

Constructing the DFA: Start state & Closure

1. Start state of the DFA = empty stack, so it contains the item $S' \mapsto .S\$$
2. Closure of a state: Adds items for all productions whose LHS nonterminal occurs in an item in the state just after the $'.'$ (e.g., S in $S' \mapsto .S\$$). The added items have the $'.'$ located at the beginning (no symbols for those items have been added to the stack yet), e.g.
 $\text{CLOSURE}(\{S' \mapsto .S\}) = \{S' \mapsto .S, S \mapsto .(L), S \mapsto .id\}$
3. Next we add the transitions, e.g. after reading id we get $S \mapsto id$.
4. Finally, for each new state, we take the closure
5. If a reduce state is reached, reduce otherwise, if the next token matches an outgoing edge, shift



Implementing the Parsing Table

Entries for the “action table” specify two kinds of actions: Shift and go to state n and Reduce using reduction $X \mapsto \gamma$. We only have reduction when we don’t have any outgoing edges otherwise we shift.

	()	id	,	\$	S	L
1	s3		s2			g4	
2	S→id	S→id	S→id	S→id	S→id		
3	s3		s2			g7	g5
4					DONE		

LR(0) Limitations

shift/reduce: $S \mapsto (L)$, and $L \mapsto .L, S$ yield a problem in a state since with (L) , we are allowed to reduce but with $.L$ we are also allowed to shift. **reduce/reduce:** $S \mapsto L, S$, and $S \mapsto .S$, yield a problem in a state since there exists two different reduction rules for S .

LR(1) Parsing

1. LR(1) state = set of LR(1) items
2. An LR(1) item is an LR(0) item + a set of look-ahead symbols $A \mapsto \alpha.\beta, \mathcal{L}$

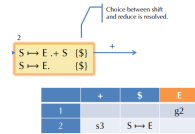
LR(1) closure:

1. Form the set of items just as for LR(0) algorithm
2. Whenever a new item $C \mapsto .\gamma$ is added because $A \mapsto \beta.C\delta, \mathcal{L}$ is already in the set, we need to compute its look-ahead set \mathcal{M}

- (a) The look-ahead set \mathcal{M} includes $\text{FIRST}(\delta)$
- (b) If δ is or can derive ϵ , then the look-ahead \mathcal{M} also contains \mathcal{L}

Example Closure: Assume we have
 $S' \mapsto S\$, S \mapsto E + S | E, E \mapsto \text{number} | S$

- Start item: $S' \mapsto S\$, \{\}$
- Since S is to the right of a $'.'$, add
 $S \mapsto .E + S, \{\$ \}; S \mapsto .E, \{\$ \}$
- Need to keep closing, since E appears to the right of a $'.'$ in $.E + S, E \mapsto .\text{number}, \{+, \}; E \mapsto .(S), \{+, \}$
- Because E also appears to the right of $'.'$ in $.E$ we get:
 $E \mapsto .\text{number}, \{\$ \}; E \mapsto .(S), \{\$ \}$
- All items are distinct, so we’re done



Ambiguity with if else:

if (E1) if (E2) S1 else S2

This is known as the **dangling else problem**. What should the right answer be? We know two solutions: the simple one would just require $\{\}$ and another one could use the grammar

1. $S \mapsto M | U // M = \text{matched}, U = \text{unmatched}$
2. $U \mapsto \text{if (E)} S // \text{Unmatched if}$
3. $U \mapsto \text{if (E)} M \text{ else } U // \text{Nested if is matched}$
4. $M \mapsto \text{if (E)} M \text{ else } M // \text{Matched if}$
5. $M \mapsto X = E // \text{Other statements}$

Lambda Calculus

The lambda calculus is a minimal programming language. It has variables, functions, and function application. It’s Turing Complete. Concrete syntax:

```
exp ::=
  | x                //variables
  | fun x -> exp      //functions
  | exp1 exp2        //function application
  | ( exp )          //parentheses
```

The only values of the lambda calculus are (closed) functions that is $\text{val} ::= \text{fun } x \rightarrow \text{exp}$.

substitute: Replace all free occurrences of x in e by v also written as $e\{v/x\}$. If we try to substitute a variable which is not free, the expression remains the same.

$x\{v/x\} = v$ (replace the free x by v)
 $y\{v/x\} = y$ (assuming $y \neq x$)
 $(\text{fun } x \rightarrow \text{exp})\{v/x\} = (\text{fun } x \rightarrow \text{exp})$ (x is bound in exp)
 $(\text{fun } y \rightarrow \text{exp})\{v/x\} = (\text{fun } y \rightarrow \text{exp}\{v/x\})$ (assuming $y \neq x$)

$(e_1 e_2)\{v/x\} = (e_1\{v/x\} e_2\{v/x\})$ (substitute everywhere)

free variable: We say variable x is free in $\text{fun } y \rightarrow x + y$. Free variables are defined in an outer scope

bound variable: We say variable y is bound by $\text{fun } y$. Its scope is the body $x + y$ in $\text{fun } y \rightarrow x + y$

closed: A term with no free variables is called closed.

open: A term with one or more free variables is called open. **Free Variable Calculation:**

$fv(x) = \{x\}$
 $fv(\text{fun } x \rightarrow \text{exp}) = fv(\text{exp}) \setminus \{x\}$
 $fv(\text{exp}_1 \text{ exp}_2) = fv(\text{exp}_1) \cup fv(\text{exp}_2)$

Variable Capture: In

$(\text{fun } x \rightarrow (xy))\{(\text{fun } z \rightarrow x)/y\} = \text{fun } x \rightarrow (x(\text{fun } z \rightarrow x))$ x is captured. Usually not the desired behavior. This property is sometimes called dynamic scoping. The meaning of x is determined by where it is bound dynamically.

Alpha Equivalence: Two terms that differ only by consistent renaming of bound variables are called alpha equivalent, e.g. $(\text{fun } x \rightarrow yx)$ the same as $(\text{fun } z \rightarrow yz)$

Operational Semantics

\Downarrow
 $\frac{\text{exp}_1 \Downarrow (\text{fun } x \rightarrow \text{exp}_3) \quad \text{exp}_2 \Downarrow v \quad \text{exp}_3\{v/x\} \Downarrow w}{\text{exp}_1 \text{ exp}_2 \Downarrow w}$

Y Combinator & Factorial

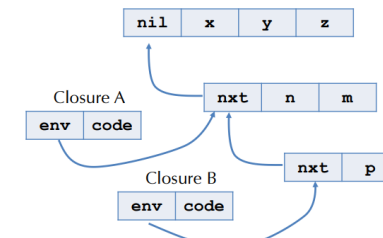
Y combinator: $Y == \lambda f. (\lambda x. f(x x)) (\lambda x. f(x x))$. $Y F = F (Y F)$ for any term F . **Example factorial function:**

- Typical recursive definition:
 $\text{fac} = \lambda n. \text{if}(n = 0) 1 (n * \text{fac}(n - 1))$
- Abstract it: $F = \lambda f. n. \text{if}(n = 0) 1 (n * f(n - 1))$
- Thus, $\text{fac} = F \text{ fac}$ (i.e., fac is a fixpoint of F)
- $Y F$ being fixpoint of F can thus be viewed as the factorial function!

CLOSURE CONVERSION

The closure is a pair of the environment and a code pointer: $\text{Code}(\text{env}, y, \text{body})$. We get $\text{fun}(\text{env}, y) \rightarrow \text{let } x = \text{nth env } 0 \text{ in body}$, where y is possibly in body .

Array-based Closures with N-ary Functions: We have $(\text{fun } (x y z) \rightarrow (\text{fun } (n m) \rightarrow (\text{fun } p \rightarrow (\text{fun } q \rightarrow n + z) x)))$



Theorem: (simply typed lambda calculus with integers): If

$\vdash e : t$, then there exists a value v such that $e \Downarrow v$

Theorem: (Type Safety) If $\vdash P : t$ is a well-typed program, then either

1. the program terminates in a well-defined way, or

- the program continues computing forever

type: A type is just a predicate on the set of values in a system. E.g., the type `int` can be thought of as a boolean function that returns true on integers and false otherwise. Equivalently, we can think of a type as just a subset of all values.

subtype: This subset relation gives rise to a subtype relation: if $Pos <: Int$, then `Pos` is a subtype of `Int`

LUB: least upper bound, for statically unknown conditionals, we want the return value to be the LUB of the types of the branches. LUB is also called the join operation.

subtyping relation is a **partial order**, that is:

- Reflexive: $T \vdash T$ for any type T
- Transitive: $T_1 \vdash T_2$ and $T_2 \vdash T_3$ then $T_1 \vdash T_3$
- Antisymmetric: If $T_1 \vdash T_2$ and $T_2 \vdash T_1$ then $T_1 = T_2$

Soundness of Subtyping Relations: A subtyping relation $T_1 \vdash T_2$ is sound if it approximates the underlying semantic subset relation, that is let be $\llbracket T \rrbracket = \{v \mid v : T\}$. If $T_1 \vdash T_2$ implies $\llbracket T_1 \rrbracket \subseteq \llbracket T_2 \rrbracket$, then $T_1 \vdash T_2$ is sound. Whenever we have a sound subtyping relation, it follows that $\llbracket LUB(T_1, T_2) \rrbracket \supseteq \llbracket T_1 \rrbracket \cup \llbracket T_2 \rrbracket$. Using LUBs in the typing rules yields sound approximations of the program behavior (as if the IF-B rule)

Subsumption Rule: $\frac{E \vdash e : T \quad T <: S}{E \vdash e : S}$

Subtyping for Function Types: Need to convert an S_1 to a T_1 and T_2 to S_2 , so the argument type is contravariant and the

output type is covariant: $\frac{S_1 <: T_1 \quad T_2 <: S_2}{(T_1 \rightarrow T_2) <: (S_1 \rightarrow S_2)}$

reference types: Are not covariant because of

```
Int bad(NonZero ref r) {
  Int ref a = r; (* OK because NonZero ref <: Int ref *)
  a := 0; (* OK because 0 : Zero <: Int *)
  return (42 / !r) (* OK because !r has type NonZero *)
}
```

and not contravariant because of

```
Assume: NonZero <: Int => ref Int <: ref NonZero
Int ref a;
a := 0;
NonZero ref b;
b = a;
return (1 / !b);
```

Immutable Record Subtyping: it holds

$\{x : int, y : int\} \neq \{y : int, x : int\}$ and
 $\{x : int, y : int, z : int\} <: \{x : int, y : int\}$

Mutable Structures: Mutable structures are invariant that is

$T_1 ref <: T_2 ref \implies T_1 = T_2$

Structural vs. Nominal Typing: Checking against the name is nominal typing and checking against the structure is structural typing.

Compiling Objects

Objects contain a pointer to a dispatch vector (also called a virtual table or vtable) with pointers to method code

Single Inheritance: every method has its own small integer index, Index is used to look up the method in the dispatch vector.

Each interface and class gives rise to a dispatch vector layout. DV layout of new method is appended to the class which is being extended

