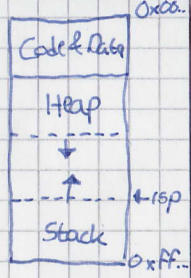


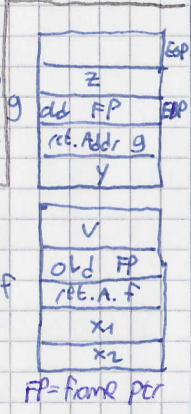
Compiler Design

X86-Lite: Rsp = Top of stack
 Heap: Stores dynam. alloc. Objects
 Stack: Stores local vars & return addr.



Calling Convention SystemV AMD64 ABI

- Callee save: rbp, rbx, r12-r15
- Params 1..6: rdi, rsi, rdx, rcx, r8, r9
- 7+: on stack, right-to-left

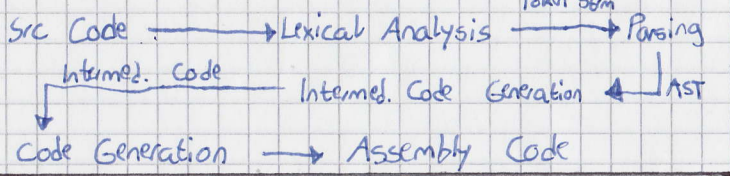


Function Call Frames:

$f(x_1, x_2)$ with local variable v
 v calls $g(y)$ with local variable z

Basic Block: Seq. of instr. that execute together, starting at first instr. & end at last instr.

Compiler Structure



Lexing "Character stream" -> Tokens

Token is a datatype that represents "chunks" of text; e.g. Identifiers, Keywords, Integers etc.

- Regex R:**
- R^* -> zero or more
 - R^+ -> one or more
 - $R?$ -> zero or one
 - $['a' - 'z']$ -> Any in Range
 - $[^ '0' - '9']$ -> Any except in Range

Lexer Generator

- Reads list of Regex R_1, \dots, R_n , one per Token
- Each token has an "Action" A_i (piece of code to run when R_i is matched)

DFA for Lexing: Regex can be represented by a DFA (Graph or Transition Table Represent.)

NFA for Lexing: Has states & transitions like DFA, can also have trans. lbb ϵ which does not consume input. Two arrows leaving the same state may have same lbb (nondet.)

Lexer Behaviour

- 1) Take each Regex R_i with Action A_i
- 2) Compute NFA formed by $(R_1 | \dots | R_n)$
- 3) Compute DFA for this NFA
- 4) Compute minimal equivalent DFA
- 5) Produce transition table
- 6) Implement longest match

Parsing: Transforms a stream of tokens into an abstract syntax tree (AST)

- Strategy: Parse token stream to traverse "concrete" syntax; during traversal build a tree representing the "abstract" syntax

Context-free grammars (CFG)

- Terminals (e.g. Lexical token)
- Set of Nonterminals
- Designated Nonterminal called start symbol

• Set of productions LHS -> RHS

Parse-Tree is a tree representation of the derivation. Leaves are terminals (in-order traversal yields the input sequence).

The internal nodes are the nonterminals

Leftmost derivation: Find the leftmost nonterminal & apply a production to it

Rightmost derivation: Find the rightmost nonterminal & apply a production to it

LL(1) Grammar -> Top-down in tree

- Left-to-right scanning
- Leftmost derivation • 1 lookahead symbol

Remove Left-Recursion

Rewrite $S \rightarrow S\alpha_1 | \dots | S\alpha_n | \beta_1 | \dots | \beta_m$ as $S \rightarrow \beta_1 S' | \dots | \beta_m S'$ & $S' \rightarrow \alpha_1 S' | \dots | \alpha_n S'$

First-Set

First(X) for a grammar symbol X is the set of terminals that begin the strings derivable from X (Also need to recursively derive non-terminals if they are first in a production.)

Follow-Set

- $Follow(A) = \{ t \in T \mid S \rightarrow \dots \rightarrow \beta A t \gamma \}$
 ↳ Contains all terminals that follow on A .
- 1) $Follow(S) += \{ \$ \}$ (S = start, $\$$ = EOL)
 - 2) $A \rightarrow \alpha B \beta \Rightarrow Follow(B) += First(\beta) \setminus \{ \epsilon \}$
 - 3) $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ with $\beta^* \rightarrow \epsilon$: $Follow(B) += Follow(A)$

LL(1) Parser Table

- Create first & follow set for all productions. Then:
 - for every production $A \rightarrow \alpha$:
 - 1) For every a in $\text{First}(\alpha)$: α in $T[A, a]$
 - 2) If ϵ in $\text{First}(\alpha)$:
 - for every terminal b in $\text{Follow}(A)$:
 - α in $T[A, b]$

LL(1) Criteria

- A grammar is LL(1) iff for 2 different productions $A \rightarrow \alpha$ and $A \rightarrow \beta$:
- 1) $\text{First}(\alpha) \cap \text{First}(\beta) = \emptyset$
 - 2) If $\epsilon \in \text{First}(\alpha)$: $\text{First}(\beta) \cap \text{Follow}(A) = \emptyset$
 - If $\epsilon \in \text{First}(\beta)$: $\text{First}(\alpha) \cap \text{Follow}(A) = \emptyset$

Shift/Reduce Parsers

- Parser State:
 - Stack of terminals & non-terminals
 - Unconsumed input is a string of terminals
 - Current derivation step is stack + input
- Parsing = Shift & Reduce operations
- Shift = Move lookahead token to the stack
- Reduce = Replace symbols x at top of stack with nonterminal X , such that $X \rightarrow x$ is a production (pop x , push X)

LR(0) has a state, which is a set of items, keeping track of possible upcoming reductions
 → LR(0) item is a production with separator \bullet

somewhere in the R.H.S. → Intuition:

- Stuff before \bullet is on stack (poss. x to be reduced)
- Stuff after \bullet is what might be seen next
- LR(0) DFA** start state: New prod. $S' \rightarrow S\ \$$ with item $S' \rightarrow \bullet S\ \$$.
- Closure of state: Add items for all productions whose L.H.S nonterminal occurs in an item of the state just after the \bullet (e.g. $S' \rightarrow \bullet S\ \$$: Add all with $S \rightarrow \dots$)
- Transitions: Outgoing edges = Terminals & nonterminals that appear after \bullet in src state
 - Target state includes all items that have edge symbol after \bullet in src → advance \bullet to simulate shifting on stack
 - Reduce state: \bullet at end of rule → reduce when reached

- LR(1) Parsing** state = set of LR(1) items
 → LR(1) item = LR(0) item + set of lookahead symbols
 When new item $C \rightarrow \bullet x$ is added b/c $A \rightarrow \beta \bullet C S, L$ (closure), we need to compute new look-ahead set M :
- 1) M includes $\text{First}(S)$
 - 2) If S can derive ϵ : $M += L$

First class functions & Lambda calculus

- Write fun $x \rightarrow x$ as $\lambda x. x$
- $e\{v/x\}$: Subst. all free x in e by v
- Free variable: Defined in an outer scope
- Bound variable: Defined in function
- Alpha Equivalence**: 2 terms that only differ by consistent renaming of bound vars. are called alpha equivalent.

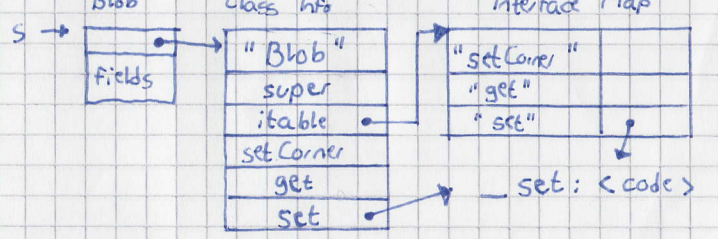
Inference Rules

- $G; L \vdash e : t = e$ is well typed and has type t
- Context: Has bindings like $x : \text{int}, y : \text{bool}, \dots$
- $\text{exp} \Downarrow \text{val} = \text{exp}$ evaluates to val
- Subtyping** $t_1 <: t_2 = t_1$ is subtype of t_2

- Depth Subtyping**: Corresponding fields may be subtypes (Number of fields needs to be equal)
- Width subtyping**: Subtype record may have more fields

Multiple Inheritance

- Dispatch vector: table in class with ptrs to methods
- Option 1: **Multiple Dispatch Vectors**
 - choose DV based on static type
 - Casting from/to class may require run-time operations
- Option 2: **Use a level of indirection**
 - Map method identifiers to code pointers
 - Use a hash table
 - Map need to search up the class hierarchy



- Option 2.1: **Use a Hashtable**: Fill Dispatch vector sparsingly, using a hash function. Don't use an interface map.

Option 3: Give up separate compilation

- Get "single class performance"
- Use "sparse" DV or binary decision trees
- Must know the entire class hierarchy

Register Allocation

Accessing Spilled Registers:

- Option 1: Reserve registers specifically for moving from/to memory
- Option 2: Rewrite the program to use a new temporary variable, with explicit moves from/to mem.

Kempe's Algorithm: k Color this Graph

- 1) Find a node with degree $< k$ and cut it out of the graph \rightarrow simplifying the graph
- 2) Recursively k-color the remaining subgraph
- 3) When remaining graph is colored, there must be at least one free color available for the deleted node. Pick such a color

Coalescing: "Merge" nodes of the interference graph if they are connected by non-related edges

Briggs' Strategy: It's safe to coalesce x & y if the resulting node will have fewer than k neighbors that have degree $\geq k$

Georges' Strategy: We can safely coalesce x & y if for every neighbor t of x, either t already interferes with y or t has degree $< k$

Reaching Definition Analysis: What variable definitions reach a particular use of the var.?

- Constraints:
- $out[n] \supseteq gen[n]$
 - $in[n] \supseteq out[n']$ if n' in $pred[n]$
 - $at[n] \cup kill[n] \supseteq in[n]$

- Available Definitions**
- $out[n] \supseteq gen[n]$
 - $in[n] \subseteq out[n']$ if n' in $pred[n]$
 - $out[n] \cup kill[n] \supseteq in[n]$

Dataflow Analyses

Liveness (backward, may) *

$out[n] := \bigcup_{n' \in succ[n]} in[n']$ $in[n] := gen[n] \cup (at[n] - kill[n])$

Reaching Definitions (forward, may) *

$in[n] := \bigcup_{n' \in pred[n]} out[n']$ $at[n] := gen[n] \cup (in[n] - kill[n])$

Available Expressions (forward, must) *

$in[n] := \bigcap_{n' \in pred[n]} out[n']$ $out[n] := gen[n] \cup (in[n] - kill[n])$

Very busy Expressions (backward, may) *

Generic Iterative (Forward) Analysis

for all n: $in[n] = T$ $out[n] = T$

repeat for all n until no change:

$in[n] := \bigcap_{n' \in pred[n]} out[n']$

$out[n] := F_n(in[n])$

Meet operator \sqcap : greatest lower bound

Join operator \sqcup : Least upper bound

* These are all distributive

Loops: A loop is a set of nodes in the CFG, with one distinguished entry point

- called the header.
- Every node is reachable from header & header is reachable from every node
 - Nodes with outgoing edges are called loop exit nodes

Domination

- Node A dominates B if the only way to reach B from start is through A.
- **Back Edge** if target node (of back node) dominates the source node

• Domination is transitive & anti-symmetric ($A \text{ dom } B \ \& \ B \text{ dom } A \rightarrow A = B$)

Dominator Dataflow Analysis

$in[n] := \bigcap_{n' \in pred[n]} out[n']$

$out[n] := in[n] \cup \{n\}$

Strictly dominates: A sd B if A dominates B but also $A \neq B$

Dominance Frontier of node n:

- 1) Calc set of nodes it dominates
 - 2) Calculate succ of those nodes (K)
 - 3) Remove strict dominations from K
- \rightarrow This is $DF[n]$

Phi Placement (Eazy Version)

Place Nodes "maximally" i.e. at every node with \geq predecessors

Garbage Collection

Garbage: X is reachable iff

- A register contains a pointer to X or
- Another reachable object y contains a pointer to X

↳ Unreachable objects are garbage

Mark and Sweep

• Mark Phase: traces reachable objects

↳ Has mark bit to use in mark phase

• Sweep phase: collects garbage objects

⇒ Can fragment the memory

⇒ Advantage: objects are not moved during GC

Stop and Copy

Memory is organized into two areas

- Old space: used for allocation

- New space: used as reserve for GC

- Heap pointer points to the next free word in the old space

Garbage Collection:

- Copy all reachable objects from old into new

- After copy, the roles of old and new are reversed and the program resumes

⇒ When copying the elements, need to store forwarding pointer to the new copy (when we reach such an object, we know it's already copied)

⇒ Partition new space into 3 regions:

copied & scanned		copied		empty
↑ start	1)	↑ scan	2)	↑ alloc

1) Objects whose pointer fields were scanned & fixed

2) Objects whose pointer fields weren't scanned

⇒ Must copy objects around → Expensive

Conservative Garbage Collection

• If a memory word looks like a pointer it is considered a pointer

• All such "pointers" are followed and we overestimate the reachable objects

Reference Counting

• Store number of pointers to each object

• Each assignment operation has to manipulate the reference count

+ Easy to implement

+ Collects garbage incrementally without large pauses in the execution

- Manip. ref. counts at each assignment is slow

- Cannot collect circular structures