

Introduction

- Compiler: Translator between programming lang.
- Operators:
 - numeric: 64-bit signed integer w. exponent
 - Lbl: 16-bit representing some machine address
 - Reg: 16-bit of 1 of 16 registers w. presented %
 - Loc: Machine address: base + Reg [index: 16bit displacement in 32bit] base + index & disp.
 - Instructions: • only quadword • loc before dest
 - move Src, Dest • reg • mem • [mem]
 - neg, dest • dest • reg • add
 - addf Src, dest • dest • dest = dest - Src
 - subf Src, dest • dest = dest + Src
 - movq Src, Reg • Reg = Reg + Src • swapped to 128-bit
 - movl Src, Reg • Reg = Reg + Src • swapped to 128-bit
 - Logic / Bit Manipulation Instructions
 - not of dest • logical negation • because not
 - and Src, dest • dest = dest & Src • bitwise and
 - orf Src, dest • dest = dest | Src • bitwise or
 - xor Src, dest • dest = dest xor Src • bitwise xor
 - setf cond, dest • dest = dest >= 0 and • arithmetic shift right
 - Shifl dest, dest • dest = dest < 0 and • logic shift right
 - shr and, dest • dest = dest > 0 and • logic shift right
 - Conditional Instructions
 - cmpq Src1, Src2 • dest = Src1 - Src2 and set condition flags
 - setb(C) dest • dest's last byte = 1 if CC else 0
 - jcc Src • rip = Src if CC else fallthrough
 - Condition Flags
 - OF • overflow
 - SF • sign • ZF • zero
 - Condition Codes
 - Z (equality) → ZF
 - NE (inequality) → (not ZF)
 - G (greater than) → (not ZF) and (SF = OF)
 - GE (greater equal) → SF = OF
 - LE (less equal) → SF <= OF or ZF
- Structure:
 - Same code as C/C Stream
 - Char Stream
 - Lexical Analysis
 - Token Stream
 - Parsing
 - Abstract Syntax Tree
 - Compile-dependent Code Generation
 - Intermediate Code
 - Machine-dependent Code Generation
 - Code Generation
 - Assembly Code
- Typical Stages:
 - Lexing
 - Parsing
 - Disambiguation
 - Semantic analysis
 - Translation
 - Control-flow Analysis → Control-Flow Graph
 - Data-flow Analysis → Interference Graph
 - Register Allocation → Assign bkg
 - Code Generation
- Optimisations may be done at any stage

<p>• Label Block: Jam → label · begin of main</p> <p>• Registers: • RIP • RAX • RBX • RCX • RSI • RDI • RBP • RSP • R0-R15</p> <p>• Arguments: • RDI • RSI • RDX • RCX • R0-R9</p> <p>• Call Save: • RIP • RDX • R12 • R13 • R14 • R15</p> <p>• Return: • RAX • Instruction Pts: • RIP</p>	<p>$\S = \text{global }$ i32, i64</p> <p>declare void @main()</p> <p>define void @main()</p> <p>i64 @foo(i64 a, i64 b) {</p> <p> i64 = add i64 i64, a2</p> <p> i64 = i64 set i64, a2</p> <p> br i1 %cond, label %then, label %else</p> <p> call void @use(i64 %sum)</p> <p> ret i64 %sum</p> <p> } else:</p> <p> store i64 %sum, i64 *%b</p> <p> }</p>
<p>• Expressive: matches human idea • Redundant to catch errors • Abstract: exact compilation</p> <p>• Low-level Code: Optimized for Machine:</p> <p>• Reducers of: • Readability • Redundancy • Ambiguity • Abstraction • Information</p> <p>• Translation goals: • Correctness • Maintainable</p> <p>• Code • Performance Code • Efficient Translation</p> <p>• Source level expressiveness of task</p>	<p>• High-level (HL): <ul style="list-style-type: none"> • Arithmetic instructions • Negf dest • dest = dest - dest • Addf Src, dest • dest = dest + dest • Subf Src, dest • dest = dest - dest • Movq dest, Src • dest = dest • Src • Imovq dest, Src • dest = dest • Src • Logic / Bit Manipulation Instructions <ul style="list-style-type: none"> • not of dest • logical negation • because not • and Src, dest • dest = dest & Src • bitwise and • orf Src, dest • dest = dest Src • bitwise or • xor Src, dest • dest = dest xor Src • bitwise xor • setf cond, dest • dest = dest >= 0 and • arithmetic shift right • Shifl dest, dest • dest = dest < 0 and • logic shift right • Conditional Instructions <ul style="list-style-type: none"> • cmpq Src1, Src2 • dest = Src1 - Src2 and set condition flags • setb(C) dest • dest's last byte = 1 if CC else 0 • jcc Src • rip = Src if CC else fallthrough • Condition Flags <ul style="list-style-type: none"> • OF • overflow • SF • sign • ZF • zero • Condition Codes <ul style="list-style-type: none"> • Z (equality) → ZF • NE (inequality) → (not ZF) • G (greater than) → SF < DF • GE (greater equal) → SF = DF • LE (less equal) → SF > DF or ZF </p>
<p>• Optimisations may be done at any stage</p>	<p>• High-level Type: <ul style="list-style-type: none"> • Abstract separate & new node types not gen. by parser • Struct Point \Rightarrow Struct Point { i64, i64 } • Represents higher-level language constructs • High-level optimisations based on program structure <ul style="list-style-type: none"> • Use fall for semantic analysis like type checking • Low-level (LL): <ul style="list-style-type: none"> • Machine dep assembly code + extra pseudo-instructions • Source structure of program is lost • Source structure of program is lost • Low-level optimisations based on target architecture </p>

1

2

Lexer

<ul style="list-style-type: none"> Discrimination: Two ways: <ul style="list-style-type: none"> Top-Down: Start from start symbol and go down Need to guess which production to use. gives good parse Errors 	<ul style="list-style-type: none"> Input: Char Stream Object: Token Stream Token: Datatype to represent invisible channel Lexical Grammar: Describes sets of strings Lexical Analyzer: Rules consisting of regular expressions Implementation: Table based
<ul style="list-style-type: none"> Bottom-up: Start from final string Left-to-right parse, Left-most derive, A look-ahead (LL(1)) Parse Grammar LL(1): Left-Facture Grammar: 	<ul style="list-style-type: none"> If 3 common prefix for each choice, then add a new non-terminal \$' at decision pt.: Ex: $\begin{array}{l} S \xrightarrow{*} E + S \mid E \\ \downarrow \\ E \xrightarrow{*} num \mid \\$' \end{array}$ Elision rule Left - derivation: Neutralizes $E \xrightarrow{*} num \mid \\$'$ so it contains $S \xrightarrow{*} S\\$'$ Close of a State: Adds items for all productions whose LHS non-terminal occurs in the state just after the * (e.g. $S \xrightarrow{*} S\\$). The added items have the * located at the beginning. The newly added items may cause yet more items to be added to the state \Rightarrow elide till it's point already in the set. (complete lookahead set \hat{M}):

Parser

<ul style="list-style-type: none"> Input: Token Stream Output: Abstract Syntax Tree 	<ul style="list-style-type: none"> Context-free grammars (C=GL): Define String by applying some rules to starting token. Consider GL: Set of Terminals: element w/o derivation role. Set of Nonterminals: elem. w/ derivation rules of start symbol: designated non-terminal Set of Productions: LHS \rightarrow RHS Accept String: If 3 derivation rule to get to this string Parse Tree: Represents derivation. (carries terminals) Inland Nodes: non-terminals. E.g., $y \rightarrow A$ to convert to AST Discrimination: Two ways: <ul style="list-style-type: none"> Leftmost: apply production to left-most non-terminal Rightmost: apply production to right-most non-terminal Both produce the same parse tree. Ambiguity: Ex. of a right-associative grammar
<ul style="list-style-type: none"> Left-Right Scanning: right most recursive right-hand grammar Fix: Add non-terminals and allow recursion only on left or right. High-priority operators go farther from the start symbol Ex: $S \xrightarrow{*} S + S \mid S \xrightarrow{*} S^2 \mid S \xrightarrow{*} num \mid \\$ 	<ul style="list-style-type: none"> Left-Right Scanning: right most recursive grammar is left and right associative derivation. Shift-Reduce: Work bottom-up. Sequence of shift and reduce operations: Shift: move lookahead token to stack Replace: replace symbols at top of stack w/ non-terminal X. St. $X \rightarrow Y$ is a production Parse Stack: Stack: terminals and non-terminals Object: String of terminals

2

- Reaching Definition Analysis:** What var definition reaches a particular use of the variable.
- Constant propagation & copy propagation
- Def:** $\{n\}$: Set of nodes that define var n
- $\text{fan}(n) = \{j\}$ if j quadratic in n is an assignment
- $\text{Lis} = \{\}$
- Constraint:** $i[n] \cdot \text{out}[n] \geq 2$ and $i[n] \in \{n\}$
- $i[n]$ is predecessor of n • $\text{out}[n] \cup \text{in}[n] \geq 2$ in $\text{in}[n]$
- $\text{in}[n] = \text{out}[n] = \{\}$ $\forall n$
- Update:** $\text{in}[n] := \cup_{i \in \text{pred}[n]} \text{out}[i]$
- $\text{out}[n] := \text{gen}[n] \cup \text{in}[n] \setminus \text{out}[n]$

- $\text{in}[n]$: set of nodes defining some var such that the definition may reach the beginning of node n
- $\text{out}[n]$: set of nodes defining some var such that the definition may reach the end of node n
- Available Expressions:** For common subexpression elimination
- Define:** infl : set of nodes whose values are available on entry to n
- $\text{out}[n]$: set of nodes whose values are available on exit of n
- Available Expressions:** For common subexpression elimination
- Define:** infl : set of nodes whose values are available on entry to n
- $\text{out}[n]$: set of nodes whose values are available on exit of n
- Available Expressions:** For common subexpression elimination

Constraints:

- $\text{in}[n]$ is predecessor of n • $\text{out}[n] \cup \text{in}[n] \geq 2$ in $\text{in}[n]$
- initial : $\text{in}[n] = \text{out}[n] = \text{set of all nodes}$, if nodes
- Update:** $\text{in}[n] := \cap_{i \in \text{pred}[n]} \text{out}[i]$
- $\text{out}[n] := \text{gen}[n] \cup \text{in}[n] \setminus \text{out}[n]$
- Update:** $\text{infl} = \text{out}[n]$
- Forward Reaching defns Available Expressions**
- Backward Live vars Very basic definitions**

- (Forward) Data Flow Analysis Framework:**
 - Domain of data flow values \mathcal{D}
 - For each node n , a flow function $F_n: \mathcal{D} \rightarrow \mathcal{D}$
 - $\text{I.E. } F_n(n) = \text{gen}(n) \cup \text{in}[n] \setminus \text{out}[n]$
 - Combine operator $\sqcap: \mathcal{D} \times \mathcal{D} \rightarrow \mathcal{D}$
 - Initial:** $\text{in}[n] = \text{out}[n] = \text{Top}$ in Top
- Maximal Information** \sqsubseteq : \mathcal{D}_1 provides of least as much information as \mathcal{D}_2 if \mathcal{D}_2 is a poset under \sqsubseteq
- General:** $\text{Wf}(n) = \text{gen}(n) \cup \text{in}[n] \setminus \text{out}[n]$
- Forward:** $\text{Wf}_{in} = \text{out}[n] = \text{pred} \cdot \text{Max}$, $\text{Wf}_{out} = \text{in}[n] = \text{succ}$
- Backward:** $\text{Wf}_{in} = \text{out}[n] = \text{pred} \cdot \text{Min}$, $\text{Wf}_{out} = \text{in}[n] = \text{succ}$

- Best Possible Solution:** Suppose we have CFA
- Best solution of the outfit is $\text{Sol} \subseteq \mathcal{P}$
- V paths \mathcal{P}
- Meet-one-path (MOP) Solution:** [] best path to n [] low functions F_i : dominants i over []
- Distributivity:** $[i:j] F_i(n) = F_i([i:j])$
- all 4 solutions compute MOP solution
- Calculating MOP costs $O(2^n)$ for $O(n)$ nodes and edges
- Relative Costs $O(n^2)$ or so
- Code Analysis / Control-Flow Ana.**
- loop with same header one merged
- loop: set of nodes in CFA • Header: Only entry point • Body nodes: either first points directly connected component (SCC) or loop
- Domination:** n dominates m iff only way to reach m from start node is via n
- Properties:** • Transitive • Anti-Symmetric
- Back Edge:** Edge b target dominates the source
- Dominator Tree:** Exists δ CFAs • Node: Node in CFA • Edge: If several dominants destination directly (i.e. last dominator)
- Dominator Domination Analysis:** $\text{dom}[n]$: set of nodes that dominate n • $\text{infl}[n] = \cap_{i \in \text{dom}[n]} \text{out}[i]$
- $\text{out}[n] = \text{in}[n] \cup \text{out}[n]$ • $T = \text{set of all nodes}$
- $F_n(e) = X \cup \{n\}$ • \sqcap is MOP
- Optimization:** $\text{dom}[b] = \text{nodes along short node to } n \text{ path, in domination order}$ • $\text{dom}[b] = \text{immediate dominator of } n \cdot \text{calcable dom}[n]$
- Domination analysis identifies back edges by walking $\text{dom}[n]$
- SSA:** Problem: How to name variable x to which we assign a value in both branches of if/else? x pointer to x
- Phi Function:** It takes operator used for analysis
- Choose version of x variable by the path below control takes phi node
- $\text{Yield} = \text{phi } \sqcup_{i=1}^k v_i, \text{label } i, \dots, v_k, \text{label } n$
- $\text{Yield} = \text{phi } \sqcup_{i=1}^k v_i, 2, \dots, k, \text{label } n$
- Mark and Sweep:** When out of Mem, exec. 2 phase
- Mark Phase:** Trace reachable objects • Add all root pointers to todo list • Set mark bits of all objects to 0 • Go through todo list • Set mark bit to 1 when visiting obj. • Recursively add pointers to todo
- Sweep Phase:** Collect garbage • Add all objects with mark bit 0 to free list • Reset all mark bits
- Efficiency:** Need to store todo and free list in the objects itself, since we are out of Mem • Phi
- Reverse:** when phi followed, reverse it
- \oplus no need to copy obj Θ fragments Mem

