# Intro to Compilers

**What is a compiler?**

→ A program that translates software written in one language into another language.

→ Should improve the program in some way

→ Ex: C, C++ are both compiled

**What is an interpreter?**

→ A program that reads an executable program and produces the results of executing that program

→ Ex: Python, Scheme are both interpreted.

**Is Java interpreted or compiled?**

→ Both! Kinda

→ Java is compiled into bytecode (code for the JVM)
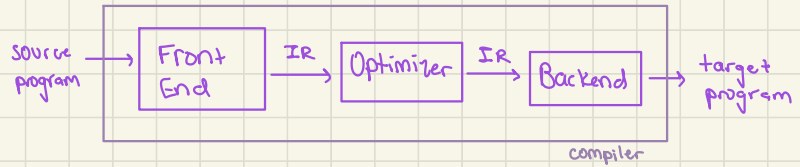
→ The bytecode is then interpreted

→ Some bytecode may be compiled (just-in-time compilation)

**What are the principles of compilation?**

→ Compiler must preserve the meaning of the input prog

→ Compiler must discernibly improve the input prog

**What is the structure of a compiler?**



source program → Front End → IR → Optimizer → IR → Backend → target program

compiler

**What does the Front End component do?**

→ Deals with the source code. Takes a stream of chars & converts into a stream of classified code.

→ INPUT: source-language program

→ Performs 3 types of tasks/operations:

1) Lexing (scanning): stream of characters → stream of words (e.g., removing whitespace)

2) Parsing: is stream of words a "sentence" in the source language?

3) Static semantic analysis: Is the sentence (statically) meaningful?

→ OUTPUT: Intermediate representation (IR)

**What does the Optimizer do?**

→ INPUT: IR

→ Focuses on improving program efficiency. multiple rewrite passes; reducing code size

→ OUTPUT: Intermediate representation (IR)

**What does the Backend Component do?**

→ Deals with the target machine
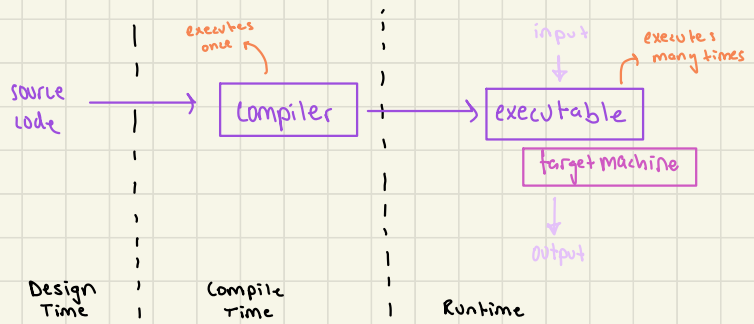
→ INPUT: Intermediate representation (IR)

→ Instruction selector: turns IR into target machine instructions

→ Register allocator: Fit the finite register set given to you by target machine, to the program that you have.

→ Instruction scheduler: reorder instructions for speed
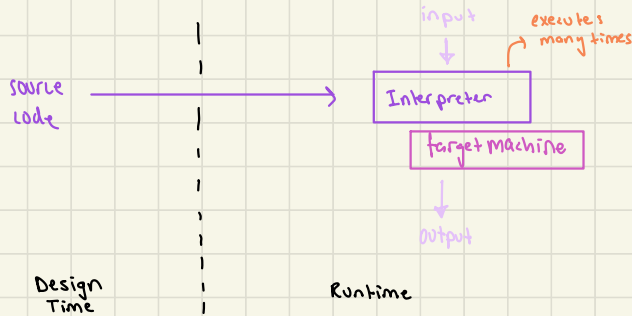
→ OUTPUT: target-language program

**What is the timing breakdown for <u>compiled</u> languages?**

executes once

input    executes many times

source code → Compiler → executable

target machine

output

Design Time | Compile Time | Runtime

---

**What happens at each stage?**

→ Design Time: design & implementation of the compiler

→ Build Time: using a compiler to compile the compiler

→ Compile Time: execute the compiler to translate source-language program to target-machine code

→ Runtime: target-machine code executes

→ Just-in-Time compiler (JIT): a runtime compile-time

---

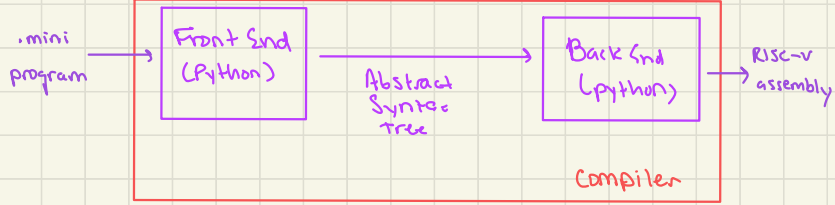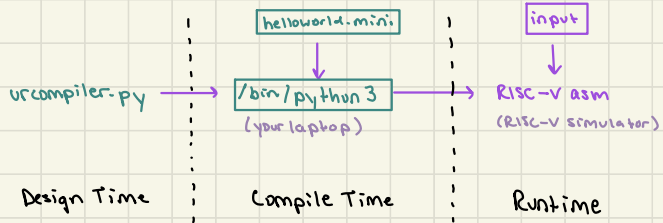**What is the timing breakdown for <u>interpreted</u> languages?**

input    executes many times

source code → Interpreter

target machine

output

Design Time | Runtime

---

**So, why study compilers?**

→ Fundamental in computation
  • key drivers of performance
  • enable use of programming langs

→ Broad applications
  ✶ classic compilation (e.g. gcc, rust)   ✶ embedded langs (e.g., macros in excel)
  ✶ scanning & parsing (e.g. HTML, SAT queries)   ✶ domain-specific langs

→ Fundamental in CS
  • computer organization/architecture, assembly programming, OS, algorithms, formal languages

→ Intrinsically interesting!

**What will my compiler look like?**

.mini program → Front End (Python) → Abstract Syntax Tree → Back End (Python) → RISC-V assembly

Compiler

**Timing for my compiler?**

urcompiler.py →

helloworld.mini → /bin/python 3 (your laptop) →

input → RISC-V asm (RISC-V simulator)

Design Time | Compile Time | Runtime

**RECALL:** What are the 3 main tasks of the Compiler Front End?

1. **Lexical Analysis (Scanning):**
   → concerned with microsyntax
   → maps: **stream of chars ~~~ stream of words**

2. **Parser:**
   → concerned with grammar
   → maps: stream of words → annotated sentences

3. **Semantic Elaboration**
   → checks that sentences are meaningful

**Example?**

→ Say that source lang is English.
   "Compilers arp fun" → error caught by scanner b/c "arp" isn't a valid word.
   "Compilers arc Fun" → All legal words; NOT caught by scanner
      ↳ Grammatically incorrect; caught by **parser**
   "Compilers walk Fun" → All legal words; NOT caught by scanner
      ↳ Grammatically correct; NOT caught by **parser**
      ↳ Semantically non-meaningful; caught by **semantic analysis**

**What does the Scanner do?**

→ The only part of the compiler that **reads every character**
→ Maps streams of chars to streams of tokens.
   • **Token : <syntactic category, lexeme>**
   • e.g.: INPUT: "Compilers are Fun."
     OUTPUT:                              Stream of tokens
     <NOUN, "Compilers">, <VERB, "are">, <ADJ, "Fun"> <ENDMARK, ".">
→ Finds microsyntax errors: unrecognized lexemes, spelling mistakes
   • **microsyntax:** specifies correct spelling of all **words** in the lang.

**What is the difference between Scanning and Parsing?**

→ Scanning: identifies & classifies words, recognizing words in the lang.
→ Parsing: identifies sentences
   • **grammar:** specifies all **legal sentences** in the language

**Why separate the two?**

→ greater efficiency; parsing is harder than scanning; better to give the parser only well-formed tokens
→ Scannerless parsing — e.g. combining the 2 — is possible.
→ **Scanner Construction is based on Finite Automata**

**What is the plan for lexical analysis (e.g. the compiler's plan)?**

| Design Time | Build Time | Build Time | Runtime |
|---|---|---|---|
| • specify microsyntax using REs | • Convert RE to FA | • Generate code to implement FA | • Code analyzes input to produce <category, lexeme> pairs |

**RECALL: What is a DFA?**

→ A deterministic finite automaton (DFA) is a 5-tuple
$(S, \Sigma, \delta, s_0, S_A)$, where

1. $S$ is a finite set, called the states
2. $\Sigma$ is a finite set, called the alphabet
   - represents the "input alphabet"
   - indicates the allowed input symbols
   - for ex, with binary string inputs, $\Sigma = \{0, 1\}$
3. $\delta: S \times \Sigma \longrightarrow S$ is the transition function
4. $s_0 \in S$ is the start state
5. $S_A \subseteq S$ is the set of accept states.

→ A DFA accepts a string iff :
- Starting in $s_0$ & consuming each char in the string according to the transition function, leaves it in an accepting state.

→ We regard **Recognizers** as DFAs

**Example of a Recognizer?**

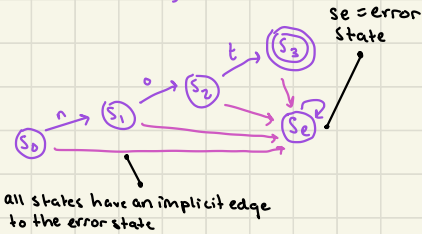→ Recognizing the word "not" :

```
c ← next char
If c = 'n' then {
    c ← next char
    if c = 'o' then {
        c ← next char
        if c = 't'
            then return <NOT, "not">
            else report error
    } else report error
} else report error
```
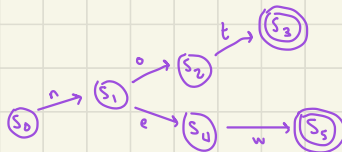
Code

Recognizer

se = error state



all states have an implicit edge to the error state
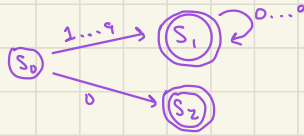
→ EX2 : Recognizing "not" or "new"

How would you recognize the infinite set of integers?



What are some key terms & facts regarding Finite Automata (FA)s?

→ **language**: set of words accepted by an FA. Countable set of strings over a fixed alphabet

→ **L( A )** : lang. recognized by an F.A. $\underline{A}$

→ An FA recognizes exactly one language

→ The transition diagram of FA specifies how to spell every word in L(FA)

→ 2 languages cannot be recognized by the same FA _or_ the same RE.

## — Specifying Languages : Regular Expressions —

Why use REs?

→ FAs are precise but not concise. REs are precise & concise

→ Set of langs that can be recognized by an FA = " by an RE

**RECALL**: what are the rules for REs?

RE Syntax?

→ ε is an RE, and $L(ε) = \{ε\}$ (empty string)

→ if $a \in \Sigma$, then $a$ is an RE and $L(a) = \{a\}$

→ Let $X = \{a, b\}$ and $Y = \{c, d\}$

1. UNION: $X | Y = X \cup Y = \{a, b, c, d\}$

2. CONCATENATION: $X \cdot Y = \{xy | x \in X \text{ and } y \in Y\} = \{ac, ad, bc, bd\}$

3. KLEENE CLOSURE: $X^* = \{x_1 x_2 \dots x_K | K \geq 0 \text{ and each } x_i \in X\}$
$= \{ε, a, b, ab, ba, aab, abba, ba, \dots\}$

What is some syntactic sugar for REs?

→ Positive closure: $r^+$ (equivalent to $rr^*$). Specifies "1 or more rs"

→ $[0 \dots 9] = 0|1|2|3|4|5|6|7|8|9$

→ $[a \dots z] = a|b|c|\dots|z$

→ Complement: $^c = \Sigma \setminus \{c\}$

What are **the steps for building a Scanner**?

1. Build REs for each category in the programming language.
   → e.g., Cat1 = integers (like $[0 \dots 9]^*$ or smthn)
   cat2 = keywords ( if | while | then | for | ...)
   cat3 = identifiers $(a-z | A-Z)(a-z | A-Z | 0-9)^*$

2. Combine into single RE (union) (cat1 | cat2 | cat3)

3. Convert RE ⟶ NFA using Thompson's Construction (1)

4. Convert NFA ⟶ DFA using subset construction (2)

5. DFA minimization: DFA ⟶ DFA (3)

6. Scanner Generator: DFA ⟶ code (4)
   • skeleton scanner + generated table
   • generated scanner

RECALL: What is an NFA? → in nondeterministic finite automatons (NFAs), several choices may exist for the next state at any point.

> A nondeterministic finite automaton is a 5-tuple $(S, \Sigma, \delta, s_0, s_A)$, where
> 1. $S$ is a finite set of states
> 2. $\Sigma$ is a finite alphabet
> 3. $\delta: S \times \Sigma_\varepsilon \longrightarrow P(S)$ is the transition function
> 4. $s_0 \in S$ is the start state
> 5. $s_A \subseteq S$ is the set of accept states.

→ $\delta: S \times \Sigma_\varepsilon \longrightarrow P(S)$, where $\delta$ takes an input of the cartesian set of all possible combos of states ($S$) and input symbols plus the empty string ($\Sigma_\varepsilon$), and produces the power set of $S$, aka the set of all possible next states.

→ $P(Q)$ = power set of $Q$ = collection of all possible subsets of $Q$.
  - ex: $Q = \{1, 2\}$ ... $P(Q) = \{\{\}, \{1\}, \{2\}, \{1, 2\}\}$

so
How do NFAs compute? → When the NFA arrives at a state with multiple ways to proceed (like if we were at $q_2$ and the next input symbol is a $1$, we can either stay in $q_2$ or move to $q_2$), the machine splits into multiple copies of itself and then follows all of the possibilities in parallel.
  - each copy of the machine takes one of the possible paths & then continues on reading the input
  - Every time there are choices, the machine "splits" again

→ NFAs are like a parallel computation where multiple independent "threads" can be running concurrently.

What happens when the NFA arrives on a state that has $\varepsilon$ on an exit arrow? → Similar: without/before reading any further input, the machine splits into at least two but possibly more copies — one that stays at the current state, and one following each of the exiting $\varepsilon$-labeled arrows.
  - For ex, when it arrives at $q_2$ and the next symbol is a $0$, $N_1$ splits into 2 copies befor reading the next symbol — one that stays at $q_2$ and one that advances to $q_3$

What are the 2 types of NFA Computation models? → Omniscient NFA: at each nondeterministic choice, follow the transition that leads to an accepting state for the input string

→ Cloning NFA: the one described above (from COMP 455)
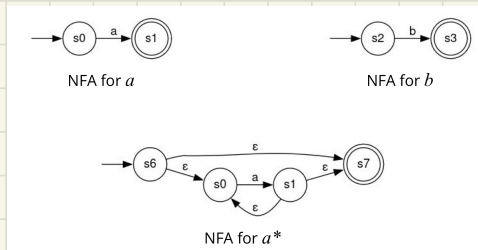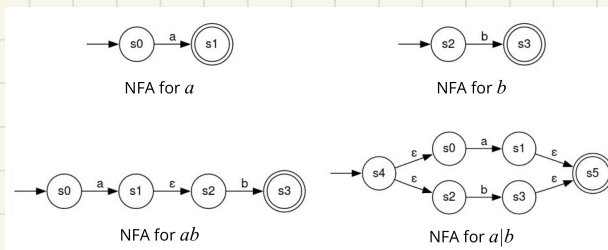
**What is Thompson's Construction?**

→ Method to generate an NFA from an RE

→ For each symbol in $\Sigma$ and each R.E. operator, there is an **NFA template**

1) Join templates together to form larger REs
- Join with $\varepsilon$ transitions
- Join in precedence order

2) Adjust set of accepting states

**What is the precedence order?**

$()$, $*$, concatenation, union    — left-associative

**What are the templates?**



NFA for $a$          NFA for $b$

NFA for $ab$          NFA for $a|b$

NFA for $a$          NFA for $b$

NFA for $a*$

**Other properties?**

→ Single start & single accepting state

→ No transitions return to start state

→ No transitions leave accepting state

→ Each state has <u>max</u> 2 transitions per symbol; one exiting and one entering

→ Each state has <u>max</u> 4 $\varepsilon$-transitions; 2 exiting & 2 entering

**Example of using Thompson's Construction?**

→ build NFA for $a(b|c)*$

1) 

2) 

(wait this might be wrong? version from textbook on next page)
* have to check these with prof

| | |
|---|---|
| Why use Thompson's construction if we can just make a minimal DFA? | → e.g., for $a(b\|c)^*$ :  |
| | → ANS: b/c it can't be auto-generated; we had to come up with it in our brains. |
| | → Thompsons Construction process can be coded/automated |

<span style="background:yellow">— Subset Construction: NFA → DFA — (2)</span>

| | |
|---|---|
| What do we do in this step? | → Convert the NFA into a DFA |
| | → Each state in DFA = set of states in NFA |
| | → Transitions in DFA represent the possible transitions from the set of NFA states. |
| | → <span style="background:violet">Subset construction</span>: the alg to convert NFA → DFA |
| How do we execute subset construction? | 1) start state $q_0$ in DFA = set of NFA states including: |
| | • n0, NFA start state |
| | • all states reachable from n0 by following $\varepsilon$ |
| | • $q_0 = s_0$ |
| | 2) Given state q in DFA: |
| | • For each NFA state $\underline{n}$ in q, for each char c: follow all possible NFA transitions, including $\varepsilon$ |
| | • the set of reachable states from doing this = new $q'$ |
| | • add $q'$ and transition to DFA |
| | 3) iterate step 2 until no more states are added |
| | 4) accepting states in DFA: contains an NFA accepting state |
| | → So basically, make a table & start listing out states. |
| | → EX: NFA for $a(b\|c)^*$ : |



Set of reachable states for char...

*Why not including s2 in the set for "b" in row for q1?? (red annotation)

| DFA State Name | NFA States | a | b | c |
|---|---|---|---|---|
| $q_0$ | {s0} | {s1, s2, s9, s3, s4, s6, } | –none– | –none– |
| | | *this becomes our next state q!* | | |
| $q_1$ | {s1, s2, s9, s3, s4, s6, } | –none– | {s3, s9, s4, s6, s5, s8} | {s3, s9, s4, s6, s7, s8} |
| $q_2$ | {s3, s4, s4, s6, s5, s8} | –none– | $q_2$ | $q_3$ |
| $q_3$ | {s3, s9, s4, s6, s7, s8} | –none– | $q_2$ | $q_3$ |

if you write out the states, it ends up being = to set of NFA states listed for q2!

**What is the pseudocode for subset construction?**

→ When alg halts, each $q_i \in Q$ corresponds to a state $d_i \in D$ in the DFA

→ Alg builds elements of $Q$ by following the transitions that the NFA can make on a given input.

→ each element of $Q$ is a subset of $N$ (the set of states in the NFA)

```
q0 = follow_epsilon ({s0})
Q = {q}
workList = {q0}
while (workList != {}):
    workList.remove (q)
    for char in Σ:
        tmp = follow_epsilon (delta (q,c))
        if tmp != {}:
            if tmp not in Q:
                Q.append ({tmp})
                workList.append ({tmp})
            T[q, char] = tmp
```

each $q$ represents a valid config of the original NFA

helper function that applies the NFA's transition function to each element of $q$, and returns $\bigcup\limits_{s \in q_i} \delta_N (s,c)$

**What is follow_epsilon?**

→ helper function that takes a set $S$ of NFA states. Then, for each state $s_i \in S$, it adds to $S$ any NFA states reachable by following one or more $\varepsilon$-transitions from $s_i$. It returns the "updated" version of $S$, with all of those states added to it.

**What type of algorithm is this?**

→ Fixed-point computation : guaranteed to terminate
 • Monotonic construction of finite set ($Q$)
 • exists a finite # of possible states to add.
 • Halts when it stops adding to the set.

**How do we construct the DFA after generating subsets & DFA states?**

→ We now have our set of DFA states $Q$, and their transitions:

| DFA State | NFA States | a | b | c |
|---|---|---|---|---|
| $q_0$ | {s0} | $q_1$ | -none- | -none- |
| $q_1$ | {s1, s2, s9, s3, s4, s6,} | -none- | $q_2$ | $q_3$ |
| $q_2$ | {s3, s9, s4, s6, s5, s8} | -none- | $q_2$ | $q_3$ |
| $q_3$ | {s3, s9, s4, s6, s7, s8} | -none- | $q_2$ | $q_3$ |

→ $q_0$ becomes start state; any $q_i \in Q$ who contains a state $\{n_i \in N \,(NFA) \mid n_i = \text{accepting state}\}$, becomes an accepting state.

**What is a partition of a set?**

→ For a set $Q$: a collection of subsets of $Q$ which are disjoint & whose union gives you back the original set. Basically dividing elements of $Q$ up into groups.

**What is Hopcroft's Algorithm?**

→ Idea: detect when 2 states produce the same behavior on any input string

1) Build a partition of DFA states, $P = \{P_1, P_2, \dots P_n\}$ s.t.:
   - for all $c$ in $\Sigma$: if $x_a, x_b$ in $P_s$, $\delta(x_a, c) = x_i$, and $\delta(x_b, c) = x_j$, then $x_i, x_j$ must be on the same partition.
   
   (basically, if P1 has $a$ and $b$, and char $c$ transitions a & b to c & d, respectively — THEN, c and d must be in the same partition.)

2) Start with the initial partition $P_0$: $\{S_A, \{S - S_A\}\}$   *(accepting state)*

3) Split partitions to satisfy the requirements

4) Iterate until reqs satisfied. Iterate for each $c \in \Sigma$

5) Create minimized DFA with 1 state for each partition group.

**Example?**

→ Final req:
   - 2 DFA states $a, b \in P_s$ have the same behavior in response to all input chars.
   - Property holds for every $P_s \in P$, every pair of states $a, b \in P_s$, and every input char $c$.
   - When examining a partition $\overset{P_x}{\{a, b, c\}}$, if on a given char, all states transition to a state inside $P_x$ ... the states can remain in this group

→ Non-minimized DFA recognizing $fee | fie$ :



| Step | Current Partition | char | set | action |
|---|---|---|---|---|
| 0 | $\{\{s_0, s_1, s_2, s_4\}, \{s_3, s_5\}\}$ | — | — | — |
|  | $\{\{s_0, s_1, s_2, s_4\}, \{s_3, s_5\}\}$ | e | $\{s_0, s_1, s_2, s_4\}$ | Split off $\{2, 4\}$ |

*why? because on char "e", s2 & s4 both transition to states in the partition $\{s_3, s_5\}$. So they must be in the same partition.*

| | | | | |
|---|---|---|---|---|
|  | $\{\{s_2, s_4\}\{s_3, s_5\}, \{s_0, s_1\}\}$ | f | $\{s_0, s_1\}$ | split off $\{s_0\}$ |
|  | $\{\{s_2, s_4\}, \{s_3, s_5\}, \{s_0\}, \{s_1\}\}$ | | | |

→ Minimized DFA:

| | |
|---|---|
| How do we generate a scanner from a DFA? | → Skeleton scanner + DFA tables |
| | → Direct-coded scanner |
| | → Handwritten scanner |
| What is a Skeleton Scanner? | → Scanner generator: RE → DFA, DFA → table |
| | → Skeleton Scanner interprets the tables to simulate the DFA |
| | → Generated Scanner uses the same skeleton for every RE |
| | → Tables encode specifics of given RE |
| Example? | → Skeleton Scanner: → Table: |

Skeleton Scanner:

```
char = next_char
State = S_0
while (char ≠ EOF):
    state = δ(state, char)
    char = next_char
if (state = final state):
    //report success
else:
    report failure
```

| δ | a | b | other |
|---|---|---|---|
| s0 | s1 | s0 | se |
| s1 | s0 | s1 | se |
| se | se | se | se |

| | |
|---|---|
| What is a Direct-Coded Scanner? | → Scanner Generator: RE → DFA, DFA → code |
| | → Transitions are compiled into conditional logic |
| | → Generated Scanner is different for each RE |
| | → Generated Code encodes specifics of RE |
| | → Low overhead per char |
| High level: How is the scanner generated? | → Build REs for each category in programming language |
| | → Combine into a single RE (union) |
| | cat1 | cat2 | cat3 |
| | • Ex: cat1 = (if|while|for|then)   cat2 = (a-z |A-Z|0-9)* |
| | → Use associated DFA to build scanner |
| Recap: How does the scanner work? | → Reads enough input to recognize a word — DFAs read all the input |
| | → Reads until it reaches state s, next char c |
| | • if δ(s, c) = se and s in S_A ... the scanner found a word |
| | • else, scanner backs up until it finds an accepting state or it exhausts the current lexeme. |
| | → Returns a token <lexeme, category> |

**Summary / Recap?**

→ FAs are good for recognizing words

→ REs are good for specifying the set of words to recognize

**What is an Abstract Syntax Tree (AST)?**

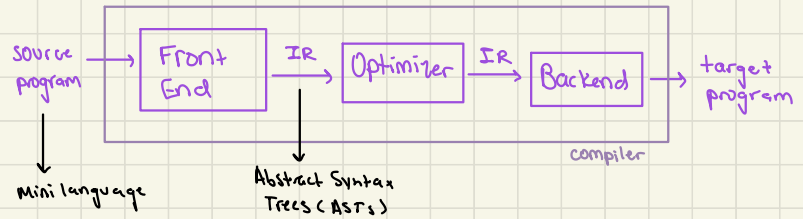→ Represents the structure of a language, as a tree



→ the structure that shows up here, is given by your parser.

→ parser generates a *parse tree*, which has tokens & stuff & isn't the same as an AST

→ Lab 0: Visit AST & pretty print an expression

Front End → IR → Optimizer → IR → Backend

source program → Front End

compiler

↓ Mini language

Abstract Syntax Trees (ASTs)

→ Scanning: Mini.g4 grammar → generated parse tree (ANTLR4)
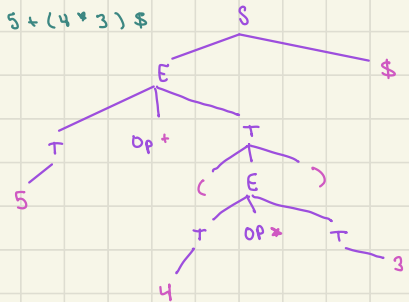
→ Parsing: generated parse tree → AST

→ A graph that describes the structure of your input source code.

→ AST design is closely linked with that of a compiler

→ Construct the syntax tree for  5 + (4 * 3) $

**What is an Abstract Syntax Tree?**

**Example?**

```
S ::= E $
E ::= T (Op T)*
T ::= (E) | num
Op ::= + | *
```

("::=" means "is defined as")

S
├── E
│   ├── T
│   │   └── 5
│   ├── Op +
│   └── T
│       ├── (
│       ├── E
│       │   ├── T
│       │   │   └── 4
│       │   ├── Op *
│       │   └── T
│       │       └── 3
│       └── )
└── $

→ Precedence rules (eg. order of ops) must be enforced for the correct AST to be generated.

→ We need to be able to traverse our ASTs

**What is the Visitor Model/ Design Pattern?**

→ MiniToAST Visitor handles construction of the AST. It builds AST from the parse tree.

→ Subsequent passes implement the ASTVisitor to traverse the AST constructed.

**What does ASTVisitor provide?**

→ abstract base classes with visit methods for all AST node types.

→ Establishes a concrete visitor pattern interface (All visitors must implement)

RECAP: Where are we?

1. Lexical Analysis (Scanning): ✓
   → concerned with microsyntax
   → maps: stream of chars ⤳ stream of words
2. Parser:
   → concerned with grammar
   → maps: stream of words → annotated sentences
   → looks at every word
3. Semantic Elaboration
   → checks that sentences are meaningful

} Front End

**What does the parser do?**

→ Given an input of **tokens** $\langle$ category, lexeme $\rangle$, maps it into a sentence.

→ Ex: int x; ⟶ $\langle$type,"int"$\rangle$ $\langle$ID,"x"$\rangle$

**What is parsing?**

→ Given a stream of words, s, and a grammar G, find a **derivation** in G that produces S.

→ Goal: recognize syntactically valid sentences in the language.

   1. Specify syntax of the lang. using grammar G.
   2. Test membership in L(G)

**What are the 2 methods for parsing?**

→ Top-Down parsing: Using G, try to recreate the sentence S. Start with G.
   • begin with root node & build down
   • root node ≃ grammar's start symbol

→ Bottom-Up parsing: Starting with S, identify rules to verify that it fits in G.
   • begin with leaf nodes & build up
   • leaf nodes ≃ string of terminals

**How do we specify the syntax of a lang?**

RECALL: What is a CFG?

→ Context Free Grammars! COMP 455!!

A context-free grammar is a 4-tuple $(V, T, P, S)$, where
   1. V is a finite set of syntactic **variables**, aka nonterminals
   2. T is a finite set, disjoint from V, called the **terminals**,
      • disjoint = no common elements between V and Σ. T = Σ
   3. P is a finite set of **rules**, with each rule being a variable and a string of variables & terminals, and
   4. S∈V is the start variable.

→ The language of a given CFG = set of sentences that can be derived from S.

**Example?**

→ Grammar G1:    S ⟶ aSb
                  Rules →  | SS
                          | ε

   V→ Variables = {S}
   T→ Terminals = {a, b, ε}
   S→ Start variable = S

→ Ex strings in $G_1$: ε, aSb, aaSbb, aaaSbbb, aSSb, ...

**RECALL: What is a derivation?**

→ A sequence of applications of rules to transform the start symbol in the grammar, into a sentence in the language.

**What is the difference between Context-Free & Regular languages?**

→ Regular Languages: (e.g., $a^*b^*$)
- Cannot count to an unbounded number
- recognized by FAs

→ CF Languages: (e.g., $a^n b^n$)
- Can count, but cannot recognize context
- recognized by nondeterministic PDAs

→ (set of all Reg. langs) $\subseteq$ (set of all CF langs)

**What are the key components of a derivation?**

→ STEPS to derive:

$$Start \rightarrow \gamma_0 \rightarrow \gamma_1 \rightarrow \gamma_2 \rightarrow \cdots \rightarrow \gamma_{n-1} \rightarrow \gamma_n \rightarrow Sentence$$

1.) select a non-terminal $A \in \gamma_i$

2.) Use rule $A \rightarrow \beta$ to expand A to get $\gamma_{i+1}$ (replace A with $\beta$)

→ IF $\gamma_i$ contains only terminal symbols, it is a sentence in L(G)

→ IF $\gamma_i$ contains 1 or more non-terminals, it is a **sentential form**.

→ Parsing = finding a derivation for a sentence

**What is a Sentential Form?**

→ A sequence of symbols that occurs at one step in a valid derivation.

**What is a parse tree?**

→ A diagram that illustrates how a string gets generated. Alternative to derivations, e.g. another way to represent a sentence

**Example?**

→ From $G_2$ on prev page ... tree for string $aaSSbb$



**What is an Ambiguous grammar?**

→ A grammar that contains at least one sentence which can be derived in multiple ways via **multiple parse trees.**
- If the same string can be made from 2 diff parse trees, the grammar is ambiguous
- If same string can be made from 2 diff derivations, doesn't mean anything

**Formal definition?**

→ A string w is derived **ambiguously** in CFG G if it has 2 or more leftmost derivations.

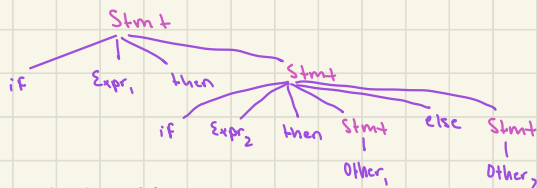→ A grammar G is ambiguous if it generates some string ambiguously.
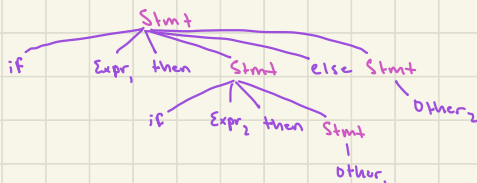
**Example of an Ambiguous Grammar?**

→ Grammar $G_2$ :

| Stmt ⟶ if Expr then Stmt |
| if Expr then Stmt else Stmt |
| Other |

→ Ex: deriving sentence S : if $Expr_1$ then if $Expr_2$ then $Other_1$ else $Other_2$

→ PARSE TREE 1 :



→ PARSE TREE 2 (Start with Rule 2 ):



**What is the difference between Parse Trees & Derivations?**

→ Derivations : capture the order in which rules are expanded ; 2 diff derivations may produce the same parse tree

→ Parse Trees : capture the structure & meaning of the sentence
  • For ambiguity, we are concerned with structure & meaning

→ To focus on parse trees, we <u>fix</u> the derivation order

**What are the derivation orders?**

→ **Leftmost Derivation** : replace the leftmost variable at each step
  • A left-setential form occurs in a leftmost derivation

→ **Rightmost Derivation** : replace the rightmost variable at each step
  • A right-setential form occurs in a leftmost derivation

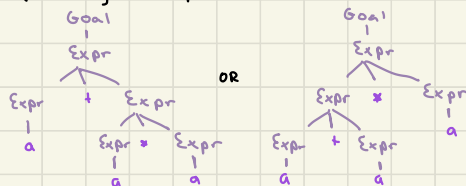**What is an unambiguous grammar?**

→ Only one possible parse tree (regardless of derivation order)

→ If derivation order is fixed, there should only be one possible derivation.

→ Ambiguous grammar $G_5$, and generating the expression $a + a * a$ :
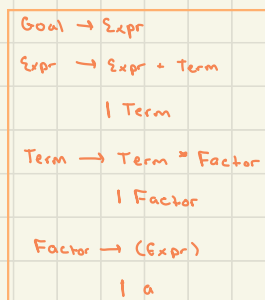
**Example of an ambiguous vs unambiguous expression?**

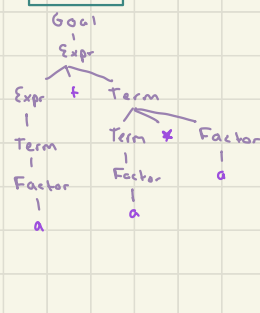| Goal ⟶ Expr |
| Expr ⟶ Expr + Expr |
| Expr * Expr |
| (Expr) |
| a |



OR



CTD next page

**Example of an ambiguous vs unambiguous expression?**

→ Unambiguous grammar $G_6$ :

| $a + a * a$ | $(a + a) * a$ |

```
Goal  → Expr
Expr  → Expr + Term
      | Term
Term  → Term * Factor
      | Factor
Factor → (Expr)
      | a
```



→ $G_6$ preserves PEMDAS / operation order unambiguously

**What is semantic ambiguity?**

→ When ambiguity in the grammar reflects true ambiguity in the language (e.g., overloading)

→ Ex: Fortran Grammar:
```
Term  → name
      | Call
      | ArrayRef
Call  → name (ExprList)
ArrayRef → name (ExprList)
```

- Given $a = F(17)$ ... we don't know whether this is a call to function $F$, or a reference to element 17 of array $F$.

**How do we handle semantic ambiguity?**

→ Solve outside the grammar :
- Use declaration info to return diff categories from the scanner
- Parse an imprecise grammar & fix it during type checking.

→ Change language to be unambiguous
- e.g., using different syntax for array refs, like []

→ In Summary :
- handle language ambiguity through type information
- handle grammar ambiguity by rewriting the grammar.

**Final notes on ambiguity?**

→ ==The grammar for a programming lang. should be unambiguous.==

→ Some ambiguous grammars can be rewritten to be ambiguous & still capture the same language.

→ But some langs are inherently ambiguous. For ex, $\{a^i b^j c^k \mid i = j \text{ OR } j = k\}$

→ A CFG $G$ is ambiguous if, for some sentence $l$ in $L(G)$:
- it can produce more than one leftmost derivation for $l$
- it can produce more than one rightmost derivation for $l$
- it can produce more than one parse tree for $l$.

→ In an unambiguous grammar, the leftmost & rightmost derivations may differ. BUT, they must have the same parse tree.

How do you add precedence to
a grammar ?

1. Decide how many levels of precedence are needed (e.g. PEMDAS)
2. Create a nonterminal (NT) for each level
3. Make higher precedence operators go through the productions for
lower-precedence operators

→ Ex : ( ), +, -, *, and / :

- Level 1 = ( )
- Level 2 = *, /
- Level 3 = +, -

Example grammar implementing
precedence

Goal → Expr

Expr → Expr + Term
| Expr - Term
| Term

Term → Term * Factor
| Term ÷ Factor
| Factor

Factor → (Expr)
| num
| name

Final note on CFGs ?

→ the language = the set of strings
→ The grammar :
  - defines the set of strings
  - encodes info about the Structure
→ 2 grammars may yield the same language
  - with or w/o ambiguity
  - with different heirarchical Structure

**How can we classify CFGs?**

→ based on the difficulty of parsing their grammars.

→



Regular
LL(1)
LR(1)
Unambiguous
Arbitrary context-Free
(All CFGs)

**What are the classifications?**

→ **Regular (RG):** CFGs that generate regular langs

- equivalent to REs (encode langs recognizable by DFAs)
- Primary use in compiler construction = for specifying scanners

→ **Arbitrary CFGs:** require more time to parse than the more restricted LR(1) or LL(1) grammars (for ex, $O(n^3)$).

- A lang is context free if it can be recognized by an NPDA

→ **LR(1) grammars:**

- able to be parsed **bottom-up**, in a linear L-to-R scan, looking at most one word ahead of the current input symbol

- A language is **LR(k) if it can be recognized by a DPDA**

- LR(k) and LR(1) grammars can express the same set of langs.

- L R (1)
  - Scan input left-to-right
  - construct a rightmost derivation
  - → lookahead of 1 symbol

→ **LL(1) grammars:**

- able to be parsed **top-down**, in a linear L-to-R scan, looking at most one word ahead of the current input symbol

- can be parsed w/ either a hand-coded recursive-descent parser, or a generated one

- Many programming langs can be expressed in an LL(1) grammar.

- L L (1)
  - scan input left-to-right
  - construct a leftmost derivation
  - → lookahead of 1 symbol

→ Almost all programming lang. constructs can be expressed in LR(1), and often LL(1) form.

# Top-Down Parsing

What does the pseudocode for a top-down parser look like?

→ RECALL: Top-down parsing = try to <u>recreate</u> sentence & using G.

(let NT = Nonterminal, aka symbol (not in Σ))

```
root = node for start symbol S
focus = root
stack.push(null)  push null onto the stack
word = NextWord()
While (true):
    if (focus is NT):
        pick rule to expand focus, A → β₁β₂β₃...βₙ
        build nodes for β₁ β₂β₃...βₙ as children of focus
        push βₙβₙ₋₁βₙ₋₂...β₂ onto the stack
        focus = β₁
    else if (word matches focus):
        word = NextWord()
        focus = stack.pop()  pop value off of stack
    else if (word ==eof && focus == null):
        accept input
        return root
    else:
        backtrack (set focus = parent node & try
        a diff production. If no untried ones remain, return error
```

What is the "lower fringe"?  → In a tree: it ≃ a sentential form in the table.

What is the focus?  → The node in a tree that is working to expand

What is the lookahead symbol?  → The next word in the string to match.

What are the requirements for a grammar to be Top-Down parsable?

→ Not left recursive: left recursion leads to infinite expansions
  • transform grammar to be right recursive

→ Backtrack-free (aka predictive OR LL(1))
  • backtracking is expensive
  • transform grammar to be predictive
  • not all languages have an LL(1) grammar

What is the focus symbol?

→ In a top-down parser, you keep a <u>stack</u> of symbols you still need to match/expand.

→ Focus symbol: the symbol at the point you're about to act on.
  • If its terminal: try to match it with the next input token (the lookahead)
  • If its nonterminal: expand it using one of its productions

**What is a left-recursive grammar?**
→ a rule in a CFG is left-recursive if the 1$^{st}$ symbol on its right-hand side is the symbol on its left-hand side, or can derive that symbol.

→ Ex:

| Fee → Fee a | · Fee is on |
|---|---|
| \| β | left-hand side |

Direct Left Recursion

| a → β | · β can be used to derive a |
|---|---|
| β → γ | |
| γ → a S | |

Indirect Left Recursion

→ A CFG is left-recursive if,
for some $A \in NT$ and $B \in (T \cup NT)^*$, $A \rightarrow^+ AB$ ⎤ Formal
defn

**How do we eliminate left-recursion?**
**What are the steps?**

→ Transform grammar from left to right-recursive :
→ Naive rewriting changes the associativity
1) Forward substitution: change indirect LR → direct LR
2) Transformation: change direct LR → RR ...assumes that the grammar has no cycles & no ε- productions

**Example?**

→ Classic expression grammar $G_1$ :

| Goal → Expr |
|---|
| Expr → Expr + Term |
| \| Expr - Term |
| \| Term |
| Term → Term * Factor |
| \| Term ÷ Factor |
| \| Factor |
| Factor → (Expr) |
| \| num |
| \| name |

left-recursive rule: 1$^{st}$ symbol on right (Expr) = symbol on left

→ Rewritten R-recursive rules for $G_1$:

| Goal → Expr |
|---|
| Expr → Term Expr' |
| Expr' → + Term Expr' |
| \| - Term Expr' |
| \| ε |
| Term → Factor Term' |
| Term' → x Factor Term' |
| \| ÷ Factor Term' |
| \| ε |

→ by moving recursion to right-side, infinite loop is avoided :

| | Rule | Transformed |
|---|---|---|
| Expr | 2 | Term Expr' |
| | 3 | Term + Term Expr' |
| | 4 | Term * Term - Term Expr' |
| | 5 | Term + Term - Term |

→ Why? Because left-recursive productions don't generate a leading terminal symbol.

**What do we have to watch out for when making a grammar right-recursive?**

→ ==Preserving Associativity==: eliminating left recursion must maintain associativity.
→ Ex : Converting left-recursive $G_1$: Expr → Expr - name | name into
$G_2$: Expr → name - Expr | name would be a naive attempt.

**What is backtracking?**

→ When the parser expands the lower fringe with the wrong production, eventually encounters a mismatch between that fringe & the parse tree's leaves (aka the terminals/string of the sentence itself), and then has to undo all of the actions that built the wrong fringe. And try other productions.

→ Backtracking – expanding, retracting, & re-expanding the fringe – wastes time & effort.

**How do we make a parser backtrack-free?**

→ A top-down parser may choose the wrong production & need to backtrack

→ Better: look ahead in the input to pick the right production

→ In general, an arbitrarily large amount of look-ahead may be needed

  • Look-ahead handled by Cocke-Younger, Kasami, or Earley's algorithms

**What is the effect of using look ahead?**

→ Large subclasses of CFGs can be parsed with limited lookahead

  • e.g., LL(1) & LR(1) grammars. Most prog. lang constructs fall in these subclasses

→ When parser goes to select next rule (in trying to recreate sentence), it can consider both the focus symbol & the next input (lookahead) symbol.

→ w/ one-symbol lookahead, parser can disambiguate all of the choices that arise in parsing the right-recursive expression grammar.

**What is predictive parsing?**

→ Given $A \rightarrow a \mid \beta$ and the next token in the stream, the parser can correctly choose between $a$ and $\beta$

→ Additional definition needed: FIRST sets, FOLLOW sets

→ Backtrack-free grammar = **Predictive grammar.**

**What is a FIRST set?**

→ Let $a$ = some RHS symbol in grammar G.

RECALL:
  T = terminals = $\Sigma$ (alphabet)

→ FIRST($a$) = the set of tokens that appear as the first symbol in some string that derives from $a$:

  • FIRST: $T \cup NT \cup \{\epsilon, eof\} \rightarrow$ PowerSet $(T \cup \{\epsilon, eof\})$
  • $x \in$ FIRST($a$) iff $a \Rightarrow^* x\gamma$, for some $\gamma$

→ For every terminal $\beta$, FIRST($\beta$) = $\{\beta\}$. e.g., FIRST(num) = $\{num\}$

→ Consider grammar $G_1$:

**Example of FIRST sets?**

1. Compute FIRST sets for terminal symbols

  FIRST(Expr) = $\{$ (, name, num $\}$   ← 2↑↑

  'Expr goes to Term Expr'. Term goes to Factor Term'.
  Factor goes to (Expr), num, name
  FIRST(Expr') = $\{$ +, −, $\epsilon$ $\}$
  FIRST(Term) = $\{$ (, num, name$\}$

2. For each production $A \rightarrow \beta$, compute FIRST($\beta$)

3. Iterate until a fixed point is reached.

| Goal | → | Expr |
|------|---|------|
| Expr | → | Term Expr' |
| Expr' | → | + Term Expr' |
| | | − Term Expr' |
| | | $\epsilon$ |
| Term | → | Factor Term' |
| Term' | → | × Factor Term' |
| | | ÷ Factor Term' |
| | | $\epsilon$ |
| Factor | → | (Expr) |
| | | num |
| | | name |

Example to understand intuition behind predictive (LL(1)) parsing with FIRST sets?

→ Take this grammar: $S \to c \mid aSb$

→ Say the parser has string **aacbb** that it is trying to "recreate"

**0. START**

Stack (topmost = lastin)

$\$$
$S$

Input = a a c b b $\$$

lookahead = a

→ Build the LL(1) FIRST table: For each production $A \to \alpha$, compute FIRST($\alpha$):

| $\alpha$ | Rule No. | FIRST($\alpha$) |
|---|---|---|
| $c$ | 1 | $\{c\}$ |
| $aSb$ | 2 | $\{a\}$ |

**1.** → Focus symbol = top of stack = Stack.pop() = $S$ ... nonterminal.

→ Find the production $S \to x$ s.t. $a \in$ FIRST($x$) : Parser checks lookup table, and now knows that it should pursue Rule 2, NOT Rule 1.

→ Pop $S$ ; push RHS $aSb$ in reverse so that leftmost $a$ ends up on top:

Stack.push(b) ; stack.push(S) ; stack.push(a);

Stack

$\$$
$a$
$S$
$b$

Input: a a c b b $\$$

lookahead = a

**2.** → Focus symbol = Stack.pop() = $a$ ... terminal.

→ Since focus symbol is terminal, match it with lookahead $\underline{a}$

→ Pop $a$ ; advance input

Stack

$\$$
$S$
$b$

Input: ~~a~~ a c b b $\$$

lookahead = a

**3.** → Focus symbol = $S$ ... nonterminal. Lookahead = a So choose Rule 2.

→ Pop $S$, push $aSb$ in reverse (same as step 1)

Stack

$\$$
$a$
$S$
$b$
$b$

Input: a c b b $\$$

lookahead = a

**4.** → Focus symbol = $a$ ... Same as step 2. Pop $a$ b/c matches lookahead.

Stack

$\$$
$S$
$b$
$b$

Input = ~~a~~ c b b $\$$

lookahead = c

**5.** → Focus symbol = $S$. Parser checks lookup table & chooses Rule 1! $S \to c$

→ Pop $S$ and push RHS $c$ to stack

Stack

$\$$
$c$
$b$
$b$

Input = c b b $\$$

lookahead = c

**6.** Focus symbol = c (terminal). Match w/ lookahead. Pop c & advance input.

| Stack |
|-------|
| $ |
| b |
| b |

Input: b b $

lookahead = b

**7.** Focus symbol = b (terminal). Match w/ lookahead. Pop b & advance input.

| Stack |
|-------|
| $ |
| b |

Input: b̶ b $

lookahead = b

**8.** Focus symbol = b (terminal). Match w/ lookahead. Pop b & advance input.

| Stack |
|-------|
| $ |

Input: b̶ $

lookahead = ∅

→ All done! With no backtracking

→ If $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar,

then $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$

- Allows parser to make correct choice with single look-ahead token

- Does not account for ε-productions

What is the **predictive/LL(1) property**?

What are **ε-productions**?

→ Rules that expand to ε, e.g. $a \rightarrow \varepsilon$.

if $A \rightarrow a | B$ and $\varepsilon \in FIRST(a)$, how should the parser choose which rule to apply?

→ Ex:

| $A \rightarrow a | \beta$ |
|---|
| $a \rightarrow cd | \varepsilon$ |

$FIRST(a) = \{c, \varepsilon\}$

$FIRST(A) = \{c, \varepsilon\}$

→ Apply ε when the look-ahead symbol appears in no other production's FIRST sets

→ FOLLOW sets are needed.

What is a FOLLOW set?

→ $FOLLOW : NT \rightarrow PowerSet(T \cup \{eof\})$

→ For a nonterminal $a$, $FOLLOW(a)$ is the set of **terminal** symbols that can appear immediately after $a$ in some sentential form.

→ $FOLLOW(a)$ = set of all words that can occur to the immediate right of a string derived from $a$

→ Purpose: parser needs to know which words can appear as the leading symbol after a valid application of a rule.

How does the parser compute FOLLOW sets?

**1.** For the topmost rule (ex: $Goal \rightarrow Expr$): Set $FOLLOW(Goal) = \{eof\}$

**2.** For each NT, if $\beta_i \beta_{i+1}$ appears in the RHS, then $FOLLOW(B_i)$ contains $FIRST(\beta_{i+1})$

**3.** Iterate until a fixed point is reached.

→ Parser utilizes pre-calculated lookup table of all FIRST sets.

**What is the pseudocode for computing FOLLOW sets?**

```
for each A ∈ NT:
    FOLLOW(A) = ∅          — S = start symbol
FOLLOW(S) = { eoF }
while (FOLLOW sets are still changing):
    for each p ∈ P of the form A → β₁β₂...βₖ:
        TRAILER = FOLLOW(A)
        For i = k, k-1, ... 1: (aka from right-to-left of the RHS of a production)
            if βᵢ ∈ NT:
                FOLLOW(βᵢ) = FOLLOW(βᵢ) ∪ TRAILER
                if ε ∈ FIRST(βᵢ):
                    TRAILER = TRAILER ∪ (FIRST(βᵢ) - ε)
                else: TRAILER = FIRST(βᵢ)
            else:
                TRAILER = FIRST(βᵢ) (aka { βᵢ } since FIRST(x)={x}
                                     for some terminal x )
```

**Example?**

→ Classic expression grammar:

$$Goal \to Expr$$
$$Expr \to Term\ Expr'$$
$$Expr' \to +\ Term\ Expr'$$
$$\qquad |\ -\ Term\ Expr'$$
$$\qquad |\ \varepsilon$$
$$Term \to Factor\ Term'$$
$$Term' \to \times\ Factor\ Term'$$
$$\qquad |\ \div\ Factor\ Term'$$
$$\qquad |\ \varepsilon$$
$$Factor \to (Expr)$$
$$\qquad |\ num$$
$$\qquad |\ name$$

| NT | FIRST(NT) | FOLLOW(NT) |
|---|---|---|
| Goal | (, num, name | eof |
| Expr | (, num, name | eof, ) |
| Expr' | +, −, ε | eof, ) |
| Term | (, num, name | eof, ), +, − |
| Term' | ×, ÷, ε | eof, ), +, − |
| Factor | (, num, name | eof, ), +, −, ×, ÷ |

**What is the (Full) Predictive, or LL(1) Property?**

→ Let START(A → α) be defined as:
- FIRST(α) ∪ FOLLOW(A), if ε ∈ FIRST(α), or
- FIRST(α) otherwise.

→ Given a grammar G and its FIRST and FOLLOW sets, we say that G has the LL(1) property i.f.f. A → α and A → β implies

$$START(A \to α) \cap START(A \to β) = ∅$$

→ A grammar with the LL(1) property is a **predictively parsable grammar.**

# Example of a Non-LLC(1) Grammar?

| | |
|---|---|
| Goal → S | |
| S → A \| B | |
| A → ( A ) \| a | |
| B → ( B ) \| b | |

| NT | FIRST(NT) |
|---|---|
| S | ( , a , b |
| A | ( , a |
| B | ( , b |

→ Let A → α ≅ S → A and A → β ≅ S → B

→ START(S → A) = FIRST(A) = { ( , a }

→ START(S → B) = FIRST(B) = { ( , b }

→ START(S → A) ∩ START(S → B) = ( . Since it isn't ∅, this grammar doesn't satisfy LLC(1) property!

## What is Left Factoring?

→ Transforming a non-LLC(1) grammar into an LLC(1) grammar

→ Not all grammars can be transformed

### How does it work?

→ Assume a grammar G with productions $A → αβ_1$, and $A → αβ_2$

→ If α derives anything other than ε, G is not LLC(1) b/c doesn't satisfy property.

1. Left-Factor to pull the common prefix, α, into a separate production:

$A → αA'$ , $A' → β_1$, and $A' → β_2$

2. IF START(A' → β_1) ∩ START(A' → β_2) = ∅, G MAY be LLC(1)

## Example of Left Factoring?

→ Non-LLC(1) Grammar G:

Factor → name \| name [ ArgList ] \| name ( ArgList )

ArgList → Expr MoreArgs

MoreArgs → , Expr MoreArgs \| ε

→ Left-Factored:

Factor → name Args

Args → [ ArgList ] \| ( ArgList ) \| ε

ArgList → Expr MoreArgs

MoreArgs → , Expr MoreArgs \| ε

## How does automatic parser generation work?

→ Write a CFG and generate an excellent parser

→ The parser generator analyzes the CFG & encodes the information into:

1) Tables to drive a skeleton parser ; OR

2) A direct-code recursive descent parser

→ Generated parsers are efficient (asymptotic) & fast (constants)

→ Parser generators handle most prog. lang. features

• LR(1) is more general

• LLC(1) accepts a strictly smaller class of langs & grammars

→ Standard work-arounds cover the features that do not immediately work

## RECAP: What is Backtrack-Free Top-Down Parsing?

→ Requires backtrack free (aka predictive) grammar

→ Build FIRST and FOLLOW sets

→ Use the LLC(1) condition/property to know whether a grammar is backtrack-free

→ Use left-factoring to rewrite a grammar to make it backtrack free

# Top-Down Parsing

→ Recursive Descent

→ Table Driven

→

(SKIPPED — see end of
bottom-up lecture. AND top-down lecture)

# Intermediate Representations

Todays topic!

RECAP: Where are we now, in the compiler?

source program → **Front End** → *IR → **Optimizer** → IR → **Backend** → target program

compiler

→ **Scanning** → **Parsing** → **Semantic Elaboration** →

→ Stream of chars → stream of words
→ outputs <category, lexeme> pairs
→ COVERED ✓

→ is stream of words a sentence in the source language?
→ maps stream of words → sentences in source lang
→ COVERED ✓

→ Static semantic analysis
→ is the sentence (statically) meaningful?
→ Output: I.R.
→ NOT YET COVERED

**What is the IR?**

→ The authoritative version of the program

| Front End | Optimizer | Backend |
|---|---|---|
| ↓ | ↓ | ↓ |
| Produces IR | traverses & transforms IR | transforms IR into native code |

→ Encompasses program representation, and derived knowledge:
  * Symbol tables          * Storage maps

→ The compiler can manipulate only those details that are represented in the IR

**What is the taxonomy (classification) of an IR?**

→ **Structural Organization**: graphical, linear, hybrid

→ **Level of Abstraction**: near-source form, near-machine code form

→ **Mode of Use**: definitive, derivative

**What are the ways to structurally organize an IR?**

→ **Graphical IR** (e.g. parse trees, ASTs)
  → aka MSO
  * encodes program structure
  * useful in optimization & in source-to-source translators
  * Graphs can be large if not careful



load

@A

×

⑩  ⓙ  ①

i  ①

good for instruction selection

→ **Linear IRs** (e.g. 3-address code, stack machine code)
  * simple, compact data structures
  * easy to rearrange
  * no structural information

load2DArray A, i, j
good for code relocation

| loadI 1 | => | $r_1$ |
| sub $r_j, r_1$ | => | $r_2$ |
| loadI 10 | => | $r_3$ |
| mult $r_2, r_3$ | => | $r_4$ |
| sub $r_i, r_1$ | => | $r_5$ |
| add $r_4, r_5$ | => | $r_6$ |
| loadI @A | => | $r_7$ |
| add $r_7, r_6$ | => | $r_8$ |
| load $r_8$ | => | $r_{Aij}$ |

good for address optimization

→ **Hybrid IRs** (e.g. control-flow graphs)
  * combine graphs & linear code

**What are the modes of use?**

→ **Definitive IR**: primary representation of the code, transmitted from one pass to the next

→ **Derivative IR**: built for a specific temporary purpose. Used within a single pass of the compiler.

| What are the key properties of an IR? | → Ease of generation          → Cost of manipulation |
| | → Ease of manipulation        → procedure size |
| | → IR design affects the speed & efficiency of a computer |

| How do Parse Trees weigh up, as IRs? | → RECALL: Parse Tree ∴ concrete rep. of the derivation |
| | → Captures the structure of the grammar |
| **Pros?** | → Includes syntactic details of the input program |
| | → Useful in source-to-source compilers |
| **Cons?** | → Tend to be inefficient |
| | → Explicitly states what could be represented implicitly |
| **Example?** | → Ex: |

Grammar box:
```
Goal → Expr
Expr → Term Expr'
Expr' → + Term Expr'
      | - Term Expr'
      | ε

Term → Factor Term'
Term' → x Factor Term'
      | ÷ Factor Term'
      | ε

Factor → (Expr)
       | num
       | name
```

$x - 2 * y$

Parse tree:
```
          Goal
           |
          Expr
        / | \
    Expr  -  Term
     |      /  |  \
   Term  Term  *  Factor
     |    |        |
  Factor Factor    y
     |    |
     x    2
```

| How are ASTs, as IRs? | → AST = a parse tree minus most nodes for non-terminal symbols (many fewer nodes than a syntax tree) |
| **Pros?** | → Could regenerate source code in a treewalk. |
| | → Often used as initial IR |
| **Cons?** | → AST node sizes tend to grow large |
| **Example?** | → For the expression   $c = x - 2 * y$; |

AST tree:
```
          Program
             |
          Functions
             |
          Fun main
         /        \
  Return Type     Body
     |             |
    void         Assign
              /        \
          Target       Source
            |            |
        (valueID) c      OP -
                        /    \
                     Id x    OP *
                            /    \
                           2    Id y
```

| How do we use Directed Acyclic Graphs (DAGs) as IRs? | → DAG IR = AST with a unique node for each value.   • e.g., AST with sharing |
| **Pros?** | → Given 2 instances of the same expr, the compiler may be able to arrange the code to evaluate it only once. (For ex, only performing evaluation of $a*2$ once but using the result twice, in an expression $a*2 + a*2*b$ ) |
| | → Encodes redundancy. e.g., with this example, the AST would contain 2 distinct copies of the expression $a*2$.   • DAG avoids this duplication by reusing identical subtrees & allowing nodes to have multiple parents. |
| | → Makes sharing explicit |
| **Cons?** | → Compiler must **prove** an expressions value doesn't change (in order to reuse)   • easy if there are assignments or function calls |
| **Example?** | → $a*2 + a*2*b$ → |

DAG:
```
        +
       / \
      x   x
     /|   |\
    a 2    b
```

| How do the trees get implemented? | → Nodes of various types & arities, connected by edges |
|---|---|

→ Allocating nodes of diff sizes complicates allocation & fragmentation in the heap (malloc())

→ A compiler executes just often enough to justify a careful implementation.

**What is a dependence graph?**

→ Encodes the flow of values in a block

→ Nodes = operations

→ Edges represent a flow of values

→ ==Often built as a secondary IR for scheduling or optimization==

→ Used as a primary IR in some JITs (just-in-time compilation)

**What are call graphs?**

→ Represent the flow of control between procedures

→ Nodes = procedures

→ Edges represent individual calls (each edge represents a different calling environment)

→ Used to optimize multiple procedures together

→ Complex call patterns create complex graphs (e.g., recursion can create cycles)

→ Some interprocedural optimizations traverse & change the call graph (e.g., inline substitution)

## — Linear IRs —

**What is 3 address code?**

→ Statements of form    $x \leftarrow y\ op\ z$

  • 1 operator (op)    • at most 3 names (x, y, & z)

→ Used in many modern compilers

→ Resembles many real machines

→ Introduces a new set of names:

  e.g.    $z \leftarrow x - 2^* y$    becomes   $t \leftarrow 2^* y$,  $z \leftarrow x - t$

**How do we implement 3-address codes?**

→ As a set of quadruples: A table of k*4 small integers

→ Each quadruple is represented w/ 4 fields: operator, 2 operands (or sources), and a destination

→ Opcodes & operands stored as small integers

$t_1 \leftarrow 2$
$t_2 \leftarrow b$
$t_3 \leftarrow t_1 \times t_2$
$t_4 \leftarrow a$
$t_5 \leftarrow t_4 - t_3$

3-address code for $a - 2 \times b$

**Pros of quadruples?**

→ Simple record structure    → Easy to reorder

→ Explicit names

**Example?**

→ For $a - 2 \times b$:

```
load  r1  b
loadi r2  2
mult  r3  r2  r1
load  r4  a
sub   r5  r4  r3
```

equivalent RISC assembly

←

| Opcode | Op₁ | Op₂ | Op₃ |
|---|---|---|---|
| load | 1 | b | |
| loadi | 2 | 2 | |
| mult | 3 | 2 | 1 |
| load | 4 | a | |
| sub | 5 | 4 | 3 |

Quadruples

| | |
|---|---|
| What are control flow graphs? | → Models the transfer of control in the procedure |
| | → Nodes = basic blocks |
| | → Edges = transfers of control |
| | • branches, jumps, predicated instructions |
| How do we use them? | → Use CF Graph (CFG) to represent control flow, & a linear IR (like 3-addr code) to represent code in blocks |
| | → Graph benefit: easy navigation between blocks |
| | → Linear IR benefit: explicit, low-level detail & operation sequence |
| | → Useful for program analysis |

Example?

```
if (x = y)
    a ← 2
    b ← 5
else
    a ← 3
    b ← 4
c ← a + b
```

$\longrightarrow$

$n_5$  if (x = y)

$a \leftarrow 2$
$b \leftarrow 5$

$a \leftarrow 3$
$b \leftarrow 4$

$n_6$    $n_7$

$n_8$   $c \leftarrow a + b$

| | |
|---|---|
| How do CFGs compare to ASTs? | → ASTs show grammatical structure, while CFGs show **control flow** |
| | → Eg: |

```
while i < 100
    stmt₁
stmt₂
```

| | |
|---|---|
| What are namespaces & scopes? | → A fundamental abstraction in most programming languages, that allow us to build real systems |
| | • E.g., use libraries w/o concern that a name was already used |
| | → New scopes create a clean slate for naming |
| How does the IR track them? | → IR includes a collection of annotations; it records names, values, constants, & locations for all program entities |
| How are the IRs records used? | → Mapping of an instance of a name ⟶ a specific declaration |
| | → Searching hierarchies to model scope rules & inheritance |
| What are some key capabilities? | → Compile-time map from a name to a specific entity (declaration, definition) |
| | • Scopes, visibility          • binding & general organization |
| | → Runtime mechanism to find & access a given entity |
| | • storage assignment |
| | • data-area addressability |

**How does the compiler map a name to an entity (at compile-time)?**

→ Compiler needs info to generate an access (load, store, or call)

→ To translate a reference to an entity, the compiler needs to know about the entity

→ SOLN: Derive info in the compiler's front-end, & store that info in some handy form for later use.

   • This is what the IR is for! IR is where we "store that info"

**What are symbol tables?**

→ Compile-time data structure for storing info for mapping names to entities

→ Maps names (symbols) to their attributes

   • This is how the compiler resolves names

1. Initial info is given where the entity is defined. Aka:

   • A variable's declaration or first use

   • A procedure's definition or its prototype

   • An object's class declaration (plus its superclass, super-superclass...)

2. Compiler constructs a type signature to represent that entity

   • It may also assign storage to the entity

**How are symbol tables (STs) structured?**

→ entry exists for every symbol, & contains associated attributes:

   • lexeme     • type     • address     • length

   • Storage class (local, static, global, constant, etc)

**conceptual example of an ST?**

| Name | Type | Addr | Len | Storage Class |
|------|------|------|-----|---------------|
| y | int | 8 | 4 | local |
| w | float | @_w | 4 | static |
| x | char | 12 | 0 | static |
| z | double | 0 | 8 | local |

What is **scope**?
→ DEFN: The region of a program where a given name can be accessed
→ A contiguous set of statements in which a name is declared, visible, & accessible

What are the 2 main types?
→ Lexical Scope                              → Inheritance heirarchies

What is **name resolution**?
→ Mapping each var reference to its specific declaration
→ Compiler does this using either lexically-scoped STs, or Inheritance hierarchy STs

## —Lexical Scopes —

What are **lexical scopes**?
→ DEFN: Scopes that nest in the order in which they appear in the source code
→ Main rule:

An occurrence of name n in the prog. at point p, refers to the entity named n that was created — implicitly or explicitly — in the scope that is lexically closest to p.

How do scopes present in programs?
→ Nested regions of code — e.g. procedures, blocks, structs, modules
→ Each scope creates a new name space
→ Scoping is tied to program semantics — languages may differ in scoping rules
   • e.g. rules for nesting, inheriting, or obscuring names from outer scopes

What are **Algol-like languages**?
→ ALGOL = Algorithmic language (name comes from the family of languages developed in 1958)
→ Examples: Algol, Pascal, Fortran, C, Scheme, Java, C++

How do Algol-like languages handle scoping?
→ Each scope creates its own namespace (e.g. procedure, block, structure)
   • Procedure scopes typically contain formal parameters
→ Almost any name can be declared locally
→ Local names obscure identical non-local names
   • e.g. in this code:

```
int i = 10;
for (int i = 0; i < 3; i++) {
   ... }
```
everything inside the scope of the for-loop only recognizes the variable named "i" that was declared inside of it (in the guard line)

→ Local names cannot be seen outside their scope
→ Scopes are searched in nesting order, inner scope to outer scope

How do common Algol-like langs behave?
→ Algol, Pascal : nested procedures, often with deep nesting
→ Fortran : • local, static, & named global scopes
            • no nested procedures
→ C : added a little to Fortran
→ Scheme : global, procedure-wide, & nested scopes

**How does Python handle lexical scopes?**

→ Python has 3 nested lexical scopes : • builtin • global (module)
  • local (function); functions can be nested
→ Declaration before use not required — first use is defining use
→ If defining use is :
  • Assignment : creates new local symbol with inferred type
  • Reference : binds to or creates a global symbol
  • Nonlocal declaration : ensures name is in the global scope

**How does the compiler model lexical scopes?**

→ Compiler builds model of the scope structure of the code
→ The ==parser builds a set of tables & search paths==
  • ==one table for each scope==
  • names declared in a scope go into the scope's table
  • ==search paths link the tables together==
→ The parser builds a ==distinct path as it enters or leaves a scope==
  • Path : instantiates the namespace according to scope rules
  • Preserve the search paths with the code
→ Compiler preserves the tables for later use (debugging)

**How do Scoped symbol tables work?**

→ New table for each scope ; Individual tables are hash tables with $O(1)$ lookups
→ ==Chain tables together into a "search path"==
  • Search starts with $1^{st}$ table, & "fails" its way up the path until it finds the name, or falls off the end of the path.



Search Path

| x | int |
| y | char* |
| z | float |

| v | int |
| b | char* |
| x | double |
| w | float |

| b | int |
| a | char* |
| c | short |

→ "Scope Model" of the code :

Scope J :
  int a, b, c
  Scope k :
    char c, d
  Scope L :
    float b, d

**Example of scoped STs?**

→ Tables to model the name space :

Scope : K
Parent : J

| c | char |
| d | char |

Scope : L
Parent : J

| b | float |
| d | float |

Scope : J
Parent : None

| a | int |
| b | int |
| c | int |

| How are the scoped tables built? | → Source language defines the syntactic constructs that demarcate (set the limits/boundaries of) scope |
|---|---|
| | • e.g. new procedure, block, class decl, struct decl |
| | → **Demarcated scope entry & exit:** |
| | • Table is created & added to front of search path UPON block entry |
| | • Table is finalized on block exit |
| | → Variable Declaration: |
| | • An entry for the declared variable, along with its attributes, is created in the current table |
| | • langs that don't require type declarations infer type info (e.g. Python) |
| | → **Variable Reference:** a lookup upon the current search path (e.g. from current scope, moving outwards) is triggered |

## — Inheritance Hierarchies —

| What are inheritance hierarchies? (I4) | → Another type of naming environment |
|---|---|
| | → A data-centric naming scheme: data & code is accessed relative to the **object** they belong to, rather than relative to the **current procedure** |
| When are they used? | → In **object-oriented languages** (eg C++, Java, Python) |
| | → subclasses & superclasses define the inheritance hierarchy, which is orthogonal to the lexical hierarchy |
| RECALL: what are OO-langs? | → Object: an abstraction with members (data, code, other objects); i.e. **a class instance.** maintains internal state |
| | → Class: A collection of objects w/ the same structure. Defines data & code members of each instance |
| | → **Inheritance:** A relationship defining a partial order on the namespaces of classes |
| | → Receiver: methods are invoked relative to an object -- the method's receiver |
| | • e.g. "self." for Python, or "this." for Java |
| How does inheritance work in OO-langs? | → For ex: `Class Alpha {...}` `Class Beta extends Alpha {...}` |
| | → Beta objects inherit code, data members from Alpha definition |
| | → Beta may redefine names from Alpha |
| | → Methods in Alpha must work correctly on Beta objects (provided the method is visible in Beta) |
| | → Beta may itself be extended — this creates a **tree of inheritance** |
| | → Beta may inherit from multiple immediate superclasses — this creates an **acyclic graph of inheritance** |

→ Given a named member, compiler searches the inheritance hierarchy to resolve the name.

**How does hierarchical scoping work?**

→ Each class definition creates a new scope

→ Data members, code members are in the scope
  - Visibility parameters can affect scope

→ Scope tables & search paths work as in lexical scoping

**Example?**

```
Class Point {
    public int x;
    private int y;
        public void draw () {...}
        public void move() {...}
Class ColorPoint extends Point {
    private Color c;
    public void draw () {...}
    public void setc (Color x) {
        this.c = x };
    }
```

Scope Tables

Class: Point
Superclass: Class

| x | int | public |
|---|---|---|
| y | int | private |
| draw | void() | public |
| Move | void() | public |

Class: Class
Superclass: None

|  |  |  |
|---|---|---|
|  |  |  |

Class: ColorPoint
Superclass: Point

| c | Color | private |
|---|---|---|
| draw | void() | public |
| setc | void() | public |

**What is compile-time name resolution?**

→ At compile time, the compiler looks at STs and etc. & maps each name to its address in memory.

→ Runtime name resolution: The lookups & etc. are done during execution.

**What is a closed class structure?**

→ An OO language has a closed-class structure if it requires that class definitions/structures must be present & fixed at compile-time.

→ OOLs with closed class structure are able to resolve names at compile-time

**Characteristics?**

→ Hierarchy known at compile time

→ Compiler can build models & emit code for all references     `EX: C++`

→ Virtual functions force runtime binding

**What is an open class structure?**

→ Characteristic of OOLs that allow the running program to change its class structure at runtime (for ex, importing classes in Java)

**Characteristics?**

→ OOLs with open class structure must resolve names at runtime.     `EX: Java`

→ Hierarchy can change at runtime

→ Compiler can build models & emit code for some references.

→ Compiler must emit code to resolve some references at runtime

**How does dynamic (runtime) name resolution work?**

→ Runtime creates & maintains the namespace tables
  - higher cost for references     • New tables & search paths created as hierarchy changes

→ Frequency of change is critical

<span style="color:red">(unfinished notes from slide)</span>

| | |
|---|---|
| RECAP: Lexical vs Hierarchical scoping? | → Lexical: Given an unqualified name $\underline{n}$ in procedure $\underline{p}$ :<br>• compiler first searches lexical scope for n, then searches hierarchical scope<br>→ Hierarchical: Given a qualified name my_obj.n<br>• compiler first resolves the name my_obj (aka an object) lexically; eg looking for where it was instantiated<br>• Compiler then resolves the name n using my_obj's hierarchical scope |

## — Visibility —

| | |
|---|---|
| What is visibility? | → A name is **visible** at point p in the src code if it can be referenced at point p.<br>→ Visibility is orthogonal to lexical & heirarchical scoping<br>→ Lang. defines programmer's level of control over visibility<br>→ Visibility info included in ST entries |
| EX: How is visibility defined in C? | → C: Static variables: Lifetime = the entire execution<br>• Inside a procedure, value is preserved across invocations<br>→ Visibility restricted to current & nested scopes<br>• Outside a proc: visibility restricted to current file |
| Ex: How is visibility defined in Java? | → Java: public, private, protected, & default methods or data members<br>• public: visible anywhere in the prog<br>• private: visible only in the enclosing class<br>• protected: visible in enclosing class, other classes in same package, & any subclasses in any package<br>• default: visible in enclosing class, other classes in same package<br>→ Ex: Given MyClass, member MyMethod(), and instance Obj, Obj.myMethod() can access:<br>• local vars declared in myMethod()    • data members of both obj & MyClass<br>• public & protected vars of any of MyClass' superclasses<br>• classes defined in same pkg as MyClass, or in any explicitly imported package<br>• public class & instance vars of imported classes    • package class & instance vars in the pkg containing MyClass |

**What are types?**
→ An abstract category defined by a set of properties. All members of the category have the same set of properties.

**What is a type system?**
→ The set of types, their properties, & the lang's semantic definitions that use types
→ Provides additional context & rules of behavior over the CFG of the lang, which are useful to the compiler, e.g.
  - Checking that operator & operands are conformable (e.g M81!)
  - checking function sigs
  - garbage collection

**What does a type system consist of?**
→ Base Types
→ Rules for constructing types
→ Tests for type equivalence
→ Rules for type inference

**What are the Base Types?**
→ Common: numbers, characters, booleans
  - processor support exists        • various varieties
→ Additional (e.g., not universal; language dependent):
  - pointers, recursive lists, rationals, strings, maps

**What are constructed types?**
→ Aggregate types that allow the programmer to organize & use base data as higher-level constructs.
→ E.g., arrays, strings, enums, structs, unions, records

**What is an array (high level)?**
→ Groups multiple entities of the same type. Indices give each entity a unique name (e.g A[2][3] is unique).

**What are strings (high level)?**
→ More than just an array of characters
→ Concatenation (e.g. "hello" + "world"); Comparison (e.g. "Fee" < "Fie")
→ length of string vs length of allocated memory

**What are enums (high level)?**
→ A set of self-documenting names for a small set of constants. e.g.
  enum Day { Mon, Tues, Wed, ... }
→ Compiler maps distinct elements to distinct values.

**What are structures?**
→ Group multiple entities of arbitrary type
→ Elements are given explicit names, e.g. "int pid;"
→ Type is the ordered product of the element types.

Ex:
```
struct Student {
    int pid;
    char [] name; }
```

**What are unions?**
→ The disjunction of multiple entities
→ Type is disjunction of the type of component entities

Ex:
```
union Node {
    struct N1 one;
    struct N2 two; };
```

| | |
|---|---|
| How is type equivalence tested? | → Straightforward for base types |
| | → Language-dependent for constructed types |
| | • name equivalence • structure equivalence |
| How is type inference conducted? | → Simple components (num, name): declare before use; infer by context (may require second pass of the IR) |
| | → Simple expressions: operator & operand type determine resulting type |
| | → Expressions involving function calls: |
| | • requires type sig. of function |
| | • function type may be parameterized |
| Example of parametric polymorphism? | → filter : (a → bool) × list of a → list of a |
| | filter (lambda x: x < 5, (1,4,6,7,8)) → [6,8] |

RECAP: What does the IR do? → Building the IR encodes information about named entities that is needed by subsequent phases of the compiler :

- Name space
- Storage
- Type

→ RECALL : Definitive IRs:

- AST
- 3-address code

What is syntax-driven translation?
( S-DT )

→ A collection of computations tied to the grammatical structure of the source code

→ GOAL: Translate the source code into an IR

→ Computations may include :

- building the AST
- building the symbol tables & type checking
- emmitting 3-address code
- evaluating expressions

What are the 2 strategies for
S-DT ?

→ Build computation into initial parse — e.g., action rules in the grammar.

- the grammar's structure determines where the computation happens . e.g., taking the code where you are parsing, and adding in code that also performs relevant computations for the variables being parsed
- tying semantic actions to grammar productions

→ Walk the parse tree to compute

## — Using Action Rules in the Grammar —

What is our "toy example" grammar? → Grammar for positive integers, $G_2$ ( terminals = {digit}, NTs = {Number, DigitList}):

Number → digit DigitList

DigitList → digit DigitList | ε

• Ex: "123" →



What might its parser look like?

```
boolean Number() {
    if (word == digit) :
        word = NextWord()
        return DigitList()
    else: return False }
boolean DigitList() {
    if (word == digit) :
        word = NextWord()
        return DigitList()
    else: return true
```

| How would we use SDT here? | → SDT : The grammar determines the computation |
|---|---|

→ Use points in the grammar to convert stream of lexemes into corresponding integer value . E.g., turning $\{<1,num>, <2,num>, <5,num>\}$ into an actual integer with value 125.

**How would we modify our parser for $G_2$ to perform translations?**

→ Same control-flow as before, but now we are passing/carrying an attribute (value) through the calls in order to simultaneously accumulate/compute the integer value

```
(boolean,int) Number(value){
    if (word == digit):
        value = CtoI(word)
        word = NextWord()
        return DigitList(value)
    else: return (False, invalid_value)}
(boolean,int) DigitList(value){
    if (word == digit):
        value = (value * 10) + CtoI(word)
        word = NextWord()
        return DigitList(value)
    else: return (true, value)
```

CtoI = function that returns int value of a single char

**What is another way to use action rules to build computation into initial parse?**

→ Embed computations in the grammar!

→ Consider a left-recursive grammar $G_3$ for 1 or more digits:

```
Number → DigitList
DigitList → DigitList digit | digit
```

↳ eg "blocks"/code snippets

**Example?**

→ For each production, define semantic actions. When the parser reduces a rule – for EX recognizing DigitList digit & reducing it to DigitList – it then runs that block, using the already-computed value of the left DigitList, and the current digit, to compute the new DigitList value.

→ Implementation: Let $$ = the NT on the LHS , and $i = the $i^{th}$ symbol on the RHS

```
Number → DigitList           {return $1 ; }
DigitList → DigitList digit   { $$ = $1 * 10 + CtoI($2); }
         | digit              { $$ = CtoI($1); }
```

## – Walking the Parse Tree –

**What is the other strategy for SDT?**

→ Creating PT first, then separately traversing it & performing computations.

→ e.g., for grammar $G_3$:   Value (root).

```
if (root is Number): return Value(root.child)
elif (root is DigitList):
    return 10 * Value(root.left) + Value(root.right)
elif (root is digit): return CtoI(root.lexeme)
```

**RECAP: What is physical &**
**Virtual memory?**

→ Physical addr. space: $2^{64}$ linearly addressable bytes in memory

→ Virtual memory: each processor/process gets its own linearly addressable bytes in mem.

→ Mapping from physical to VM:

• Policy: OS determines mapping via page tables

• mechanism: HW walks the page tables on every memory reference

→ OS + HW enforce process isolation

**RECALL: How is the Address**
**Space defined?**

| OS | 0xFFFFFFFC |
|---|---|
| Stack | → holds procedure activation records |
| Heap | → memory allocated explicitly by program |
| Static (.data) | • statically defined entities • size determined at link time |
| Code (.text) | → executable code; fixed size |
| Reserved | 0x0000 |

**How does the compiler handle**
**Storage layout?**

→ Compiler decides where each entity will live at runtime (requires namespace & type info).

→ Assign each named entity a runtime home:

a) the logical data area (e.g. stack, static, etc.) ; depends on lifetime & visibility

b) the offset within the data area

**What are the Storage classes?**

→ Automatic, Static, irregular

→ Depends on entities lifetime

**What are automatic variables?**
**Where are they stored?**

→ Lifetime matches that of declaring scope

→ Either the scope's local data area, or a register

→ e.g., if x declared in procedure fee(), then compiler stores x in fee()'s local data
area. x becomes part of fee()'s activation record (AR).

→ If x is local, scalar, unambiguous: compiler can store x in a register.

**What is an unambiguous value?**

→ A value that can be accessed by only one name.

**What are some examples of**
**ambiguous values?.**

→ Pointer-based vars               → Call by reference parameters

→ Array-element references

→ Ambiguous values cannot be kept in a register across assignment statements.

**What are static variables?**

→ Lifetime spans first definition to last use (or, depending on the implementation,
lifetime may span entire execution of program)

→ EX: global vars ; static vars inside procedures

**where are they stored?**

→ Single .static data area per source file

| | |
|---|---|
| What are irregular variables? | → Lifetime is under program's control – program code explicitly allocates space ( e.g. malloc() in C ) |
| | ↳ Compiler's runtime support library |
| Where are they stored? | → The variable itself is stored in the runtime heap |
| | → The handle (pointer?) to the variable is stored in local or static space. |
| What are temporary variables? | → Lifetime is short |
| | → e.g., $x - 2*y$ |
| Where are they stored? | → Placement depends on size. If variable size is... |

|  |  |  |
|---|---|---|
| known, and small | ────→ | Compiler places in register |
| large, but bounded | ────→ | compiler can place in local data area |
| unknown | ────→ | compiler emits code for runtime heap allocation |

## – Storage Assignment –

| | |
|---|---|
| How is storage assigned by the compiler ? | 1. Assign all small, local, scalar, unambiguous variables to a unique virtual register (reg) |
| | 2. Assign each ambiguous value an offset within its data area. |
| How are 1D arrays laid out in storage ? | → Let $w$ = width of each element (in bytes) ; high = max valid index ; low = min valid index |
| | → 1D Arrays laid out as $\boxed{(\text{high} - \text{low} + 1) * w}$ contiguous bytes of memory |
| | → Then, the address of V[i] is $@V + (i - \text{low}) * w$ |
| What is one method for laying out 2*D arrays in storage ? | → Let col = num. columns and row = num. rows |
| | → Method 1 : laid out as $\text{col} * \text{row} * w$ bytes of contiguous storage |
| | • layout can be row-major or column-major order ; language dependent |
| | → Entire array stored in one big block of memory. |
| What is the other method? | → Array of pointers : laid out as one contiguous block of $\boxed{\text{row}}$ pointers, each of which points to $\boxed{\text{col} * w}$ bytes of contiguous storage |
| | → E.g, A [i][j][k]. <u>A</u> points to an array of pointers, aka a contiguous block of i elements. Each of these elements is a pointer to an array (aka contiguous block) of j 'row-pointers.' Each of these elements is a pointer that points to an array (contiguous block) of <u>K</u> data elements. |
| | • The array has a total of $i * j * k$ data elements. |
| | → Total storage needed for the array : $\boxed{(\text{row} * \text{ptr}) + (\text{row} * \text{col} * w)}$ , where ptr = width of each pointer (in bytes) |

Example of storing an array using pointers / indirection vectors?

→ $A[0:1][0:2][0:3]$ :



How are strings laid out in storage?

→ Size of layout is not fully determined by length of string
- variable may hold diff strings over its lifetime

→ 2 main storage options:
- Null-terminated: length computation is $O(n)$, where $n$ = length of the string
- Explicit length field: length computation is $O(1)$

How are structures laid out in storage?

→ Space allocated for each field in the struct

→ 2 options for assigning a field's offset within the struct: declaration layout, or compiler controlled

What is an **object record**? (OR)

→ Data structure / block; each instance of a class has an OR.

→ Contains:
- data members specified by the object's class
- pointer to its class
- pointer to array of the class' methods

→ Additionally, inheritance requires containing data members specified by superclasses & access to superclasses' methods

Where are they stored?

→ On the heap

Example?

```
Class Point {
    public int x;
    private int y;
    public void draw() { ...};
    public void move() {...};
Class ColorPoint extends Point {
    private Color c;
    public void draw() { ...};
    public void setc( Color x ) {
        this.c = x
    };
}
```

| | |
|---|---|
| How are object records laid out in storage? | → Layout in subclass instances must meet superclass method expectations<br>→ Single inheritance: prefix layout<br>→ Multiple inheritance:<br>   • compiler imposes an order on superclasses<br>   • subclass instance ORs contain members for all superclasses<br>   • compiler adjusts OR pointer as needed, when calling superclass' methods |
| How does the compiler assign storage layout in the IR? | → e.g. how is storage layout info embedded in the IR<br>→ Depends on the type of language |
| How is this done for langs with 'all declarations before executables'? | → Compiler builds up symbol tables while processing declarations<br>→ Compiler performs storage layout before processing executables<br>→ Code gen. uses concrete references |
| How is this done for languages with 'declare before use'? | → Compiler builds up symbol tables<br>→ Code generation uses abstract references<br>→ compiler performs type inference, storage layout<br>→ References are refined |
| What about langs that don't require declarations? | → Same as for langs that do, but with possible additional passes required |

<u>— Alignment & Padding —</u>

| | |
|---|---|
| What is alignment? | → Address of a value is divisible by the size of the value (in bytes)<br>→ ISAs require aligned values, e.g., a 32-bit int begins on a 32-bit addr boundary. |
| How does the compiler enforce alignment? | → Ensure data area begins at aligned address<br>→ Ensure layout of data area is internally aligned |
| What is the effect of caching? | → For 2 vars with spatial proximity (close together in the code): <small>declared</small><br>   • brought into cache at the same time, on the same cache line<br>→ For 2 vars with temporal proximity (accessed close-ly in time):<br>   • want them both in cache at the same time<br>   • non-conflicting cache lines<br>→ Compilers control depends |

**What is an activation record (AR)?**

→ Block of mem. for the control & data values of a procedure instance
  - often stack based, e.g. a "stack frame"

→ Key data structure for the implementation of procedures, e.g.: calls & returns, lexically scoped name spaces

**What is a procedure?**

→ Control-flow based containment of a unit of execution; an abstraction to create a controlled execution environment

→ A callee may be called by many diff caller

→ A Procedure (proc) defines 3 abstractions:
  1. call: transfer of control to a proc
  2. namespace: a new, protected namespace per call
  3. interface: a standardized interaction b/w caller and callee.

**How would a proc's SF look?**

→
| ... |
|-----|
| arg 10 |
| arg 9 |
| Saved registers |
| local vars |
| ... |
| Heap |

**What state is saved across procedure calls?**

→ AR pointer (Frame pointer, base pointer)

→ Stack pointer                    → RA of caller

→ Callee-saves registers

**What state is NOT saved?**

→ temporary registers

## — Procedure Linkage —

**What are the phases of a procedure call?**

1. Precall ———————————— caller
2. Prologue ─┐
            ├ callee
3. Epilogue ─┘
4. Postreturn ————————————

**What does the caller do during precall?**

→ Save any necessary temporary registers in AR

→ Allocate space for callee's AR

→ Save actual parameter values or references (in registers and/or designated AR slots)

→ Save current AR pointer in AR slot

→ Save address to return to (in dedicated register or AR slot)

→ Jump to callee

**What does the callee do during the prologue?**

→ Save any needed callee-saves registers :
- return address
- s0 - s11

→ Allocate space for local data area

→ Initialize local variables

**RECALL: What are caller- & callee- saves registers?**

→ callee-saved : registers that caller expects callee to NOT overwrite. If callee wants to use them, they have to preserve the values in those registers (e.g. on the stack), & then restore them before returning

→ caller-saved : registers that caller expects callee to overwrite. If caller wants to preserve those values, it must save them on stack before making function call.

**What does the caller do during epilogue?**

→ Save return values (in dedicated AR slot and/or registers)

→ Restore RA

→ Free space for local data area ; e.g., move the sp/frame pointer back up

**What does the callee do during postreturn?**

→ Free callee's AR

→ Restore any needed caller-saves registers

→ Continue execution

**Example?**

```
fun plusplus(int n) int {
    return n+1; }

fun add2(int a, int b) int {
    int c;
    b = plusplus(b);
    c = a + b;
    return c; }

fun main() int {
    int w, y, val;
    w = 4;
    y = 8;
    val = add2(w, y);
    print val;
    return 0; }
```

→ main :
```
lw a0 w
lw a1 y
jal add2
jal print
jal zero exit
```

add2 :
```
// prologue :
addi sp sp -12   // make room on stack for RA and params a,b
sw   ra  0(sp)   // push old RA onto stack
sw   a0  4(sp)   // push the params  a and b on stack
sw   a1  8(sp)   // b
lw   a0  8(sp)   // now that param saved, load the param for plusplus in reg
jal plusplus     // result/ret val will be stored in a0

lw t0 4(sp)      // t0 = a
add t0 t0 a0     // t0 = a + b
addi a0 t0 0     // a0 = t0 = a + b = return val for add2
```

(add2 ctd...)

//epilogue:

lw ra 0(sp) //restore old RA

addi sp sp 12 //restore SP back up

ret

```
fun plusplus(int n) int {
    return n+1; }

fun add2(int a, int b) int {
    int c;
    b = plusplus(b);
    c = a + b;
    return c; }

fun main() int {
    int w, y, val;
    w = 4;
    y = 8;
    val = add2(w, y);
    print val;
    return 0; }
```

plusplus: //plusplus makes no calls, so no prologue or epilogue needed

addi a0 a0 1

ret

**Summary of the phases of a proc call?**

→ Precall & prologue:
- Save state of caller
- set up AR for callee
- transfer params & control to callee

→ Epilogue & postreturn:
- restore State of caller
- tear down AR for callee
- transfer ret vals & control to caller

— Parameter Binding —

**What is parameter binding?**

→ Mapping caller's actual parameters to callee's formal ones

→ Key to enforcing the procedure abstraction
- proc can be called from multiple call sites
- caller does not need details of callee's implementation

→ 2 conventions:
- call-by-value
- call-by-reference

**What is call-by-value param binding?**

→ Caller copies value into appropriate location — e.g. AR slot, or register
- copy of param val is passed to the callee

→ In the callee, the value has one name: the name of the formal parameter

→ The parameter's value is evaluated & initialized by caller
- Modifications to the param value inside the callee are not visible to the caller

**Example?**

→
```
int fee (int x) {
   x = x * 2;
   return x; }

int main {
   a = 3;
   b = fee(a); }
```
→
```
main:
   lw a0 a
   jal fee
fee:
   mul a0 a0 2
   ret
```

↳ basically, the caller copies the value of the param into another reg, & then 'passes' on the val via the reg. Upon return, it is caller's responsibility to 'save' the ret value — e.g., by storing it back in memory.

| | |
|---|---|
| What is call-by-reference param binding? | → Caller passes the addresses of parameters to the callee (either via a reg or a slot in the callee's AR)<br>    • e.g., rather than `lw a0 param1`, it would be `la a0 param1`<br>→ In callee, value may have more than one name<br>    • the name of the formal param<br>    • the value itself, e.g. if callee has direct access, or if 2 formal params point to the same actual param<br>→ Changes to the value inside the callee, ARE visible to caller<br>→ Callee has extra level of indirection for each use of the param |

Example?

```
int fee (int x, int y) {            main:
    x = 2 * x                          la a0 a
    y = x + y;                         la a1 b
    return y;                          jal fee
int main() {                          la a1 a
    int a = 2;                         jal fee
    int b = 2;                      fee:
    int c = fee(a,b);                  lw t0 0(a0)  → callee directly loads in the val
    int d = fee(a,a);                  lw t1 0(a1)
                                       mul t0 t0 2
                                       add t1 t0 t1
                                       sw t0 0(a0)  → callee directly modifies the val
                                       sw t1 0(a1)
                                       ret
```

→ In this ex, fee(a,b) will return 8. but fee(a,a) will return 16 - even though both times, we were basically doing fee(2,2).

→ With fee(a, a), both params are pointing to the same param.

| | |
|---|---|
| How are values returned? | → Caller sets return value<br>    • in caller's AR        • in designated register<br>→ Ret val may be a value, or a pointer to a value on the heap. |

<u>— Procedures in OO-languages —</u>

| | |
|---|---|
| How do procedures work for OO langs? | → AR: manages control & local storage info for methods<br>→ Object Record:<br>    • Objects' lifetimes are orthogonal to methods' lifetimes<br>    • OR manages persistent state of object<br>    • OR instantiates the inheritance hierarchy |

**RECALL: What is an object record?**

→ A runtime data structure, typically on the heap

→ Pointed to by the Object Record Pointer (ORP)

→ Contains:
- pointer to class
- pointer to class' method vector
- class-specific members

→ OR allocated when an object is created

→ OR deallocated when an object is no longer reachable

**Example?**

What would the symbol tables look like?

→ Single inheritance, open class structure example:

→ RECALL: heirarchical scoping

```
Class Point {
    public int x, y;
    private int z;
    public void draw() { ...};
    public void move() {...};
Class ColorPoint extends
    Point {
    private Color c;
    public void draw() { ...};
    public void setc( Color x )
        { this.c = x };
}
```

Class: Class
Superclass: None

Class: Point
Superclass: Class

| x | int | public |
|---|-----|--------|
| y | int | public |
| z | int | private |
| draw | void() | public |
| Move | void() | public |

Class: ColorPoint
Superclass: Point

| c | Color | private |
|------|--------|--------|
| draw | void() | public |
| setc | void() | public |

**What would the OR for a Point obj. instance look like?**



aPoint

Class Point

draw:
Move:

**What if we add a ColorPoint obj. instance?**



aPoint

a ColorPoint

Class ColorPoint

class Point

class Class

Class Point

draw:
setc:
draw:
Move:

| | |
|---|---|
| What is the process of emitting code for method invocation? | → Invocation happens relative to an obj. instance<br>→ Instance must be visible at point of invocation<br>→ Compiler does the following:<br>  • Symbol-table lookup to find instance's ORP<br>    1. lexical hierarchy<br>    2. class hierarchy<br>    3. global scope<br>  • emit code to obtain ORP<br>  • emit code to obtain code-vector pointer<br>  • emit code to obtain method pointer<br>  • make procedure call (first param = instance's ORP) |
| What is a closed class structure? | → Method impl. is known at compile time<br>→ ST lookup produces the label of the method<br>→ Compiler emits code to generate the proc. call |

— Addressing Variables within the Procedure —

| | |
|---|---|
| How are variables addressed? | → 2 parts: base addr, and offset<br>→ local vars & static + global vars have static base addresses<br>→ Variables with dynamic base addresses:<br>  • those belonging to enclosing lexical scope<br>  • heap-allocated objects |
| Example of a static base address? | ```
.data
  w: .word 4  → Static vars
  y: .word 8

main:
  lw a0 w  → uses the compiler-generated label
  lw a1 y
``` |
| How does the compiler handle dynamic base addresses? | → For dynamic vars, base addr. is NOT known at compile time<br>→ Compiler does the following:<br>  1) ST lookup to find variable's coordinates (lexical hierarchy, global scope)<br>  2) emit code to obtain correct ARP<br>  3) emit code to add offset to pointed-to. AR<br>  4) emit code to load the variable @ the calculated address |

| | |
|---|---|
| What is code generation ? | → Mapping source-language constructs to the target machine's instruction set |
| |      * Expression evaluation |
| |      * Accessing vars & aggregate data structures |
| |      * Control-Flow constructs |
| |      * procedure calls |
| What is the register-to-register model? | → emit code using virtual registers for unambiguous, scalar values |
| | → later : register allocation assigns virtual regs to physical ones |
| What is the memory-to-memory model ? | → emit code using memory as primary home for all values |
| | → later : optimizations to increase register use |
| What is the emitted code ? | → 3-addr code IR ; array of tuples |
| What is code shape ? | → The explicit & derived knowledge about the source prog., encoded in the IR |
| |    • live values               • registers in use |
| |    • memory references       • type information |
| |    • instruction mapping |
| | → For ex, consider x+y+z ... code shape determines later opportunities for optimization |
| | → Deciding code shape often requires 2 passes : |
| |      1. learn the surrounding context |
| |      2. decide code shape & emit code |

## — Expression Evaluation —

| | |
|---|---|
| RECALL : What is an expression? | → e.g. x+y+z |
| | → Represent a significant portion of code |
| | → Efficient evaluation improves runtime of overall program |
| | → Target machine with 3-addr form operators |
| |    • allows compiler to name the result of operations for later reuse |
| |    • most RISC ISAs, including RISC-V |
| How is code generated for expressions ? | → Walk AST to : emit code to put operands in registers , emit code for expressions |

What would the program
to generate code for expressions
look like?

```
expr(node) {
    int result, t1, t2
    switch (type(node)) {
        case x, /, +, - :
            t1 = expr(node.left_child)
            t2 = expr(node.right_child)
            result = next_register()
            emit(node.op, t1, t2, result)
            break
        case NAME:
            result = access_value(node.st_lookup())
            break
        case NUM:
            result = access_number(node.st_lookup())
            break
    } return result; }
```

→ Ex: a-b×c
→ AST:

Goal
Expr
Expr — Term
Term    Term × Factor
Factor  Factor  (name)c
(name)a  (name)b

Example walkthrough?

→ For a-b×c, with a,b,c on the stack:

| | | | | |
|---|---|---|---|---|
| lw | t0 | 0(sp) | → | t0 = a |
| lw | t1 | 4(sp) | → | t1 = b |
| lw | t2 | 8(sp) | → | t2 = c |
| mul | t3 | t1 t2 | → | t3 = b×c |
| sub | t4 | t0 t3 | → | t4 = a-b×c |

What is the convention for
register usage in RISC-V?

| Reg Name | Reg Num | Description |
|---|---|---|
| zero | $0 | constant 0 |
| ra | $1 | return address |
| sp | $2 | stack pointer |
| t0-t6 | $5-7, $28-31 | temporaries (caller-saves) |
| s0-s11 | $8-9, $18-27 | callee-saves (safe for caller) |
| a0-a1 | $10-11 | function args & return values |
| a2-a7 | $12-17 | function args |

→ RECALL COMP 311

| | |
|---|---|
| How does the compiler reduce register pressure? | → In a binary operation, the operand requiring more registers should be computed first<br>    • results in the same number of virtual registers (VRegs)<br>    • later, reg. allocation will produce code using fewer registers<br>→ Requires 2 passes by compiler:<br>    1. compute (virtual) register demand for each operand<br>    2. decide order of subexpression evaluation & emit code |
| How are expression operands handled? | → Variables:<br>    • In memory (AR; labeled area in .data)<br>    • In (virtual) register<br>→ Constants:<br>    • labeled area in .data<br>    • Immediate values<br>→ Temporary values (aka the result of a subexpression): in a register<br>→ Procedure calls:<br>    • Postreturn sequence puts result in appropriate location<br>    • Subexpression evaluation order cannot be reordered around procedure calls |
| What about for mixed-type operands? | → Source languages allow some mixed-type expressions (eg int + float)<br>→ Compiler converts operand(s) to the appropriate type<br>→ In compiler-time conversion, compiler looks up conversion rules & applies them at the point of emitting code<br>→ In Run-time conversion:<br>    • compiler adds type tag to value<br>    • runtime case analysis jumps to appropriate subroutine for mixed types |
| How does the compiler handle assignments? | → e.g. a = b + c<br>→ For RHS (eg b + c): expression evaluation; result is a value; compiler emits code<br>→ For LHS (eg a): evaluation to access a value<br>    • result is a location (eg addr where a is stored)<br>    • compiler emits code to move RHS value to LHS location |

## — Accessing Values —

| | |
|---|---|
| How are values kept in registers accessed? | → REGALLOC: values kept primarily in registers: scalar, unambiguous, & temporary values<br>→ IR may assign register-based values to virtual registers<br>→ Registers can be read & written to in the same cycle<br>→ Ideal: values kept in registers wherever they are used frequently<br>    • optimization, register allocation |

| | |
|---|---|
| How are values kept in memory accessed? | → RECALL: values kept primarily in memory:<br> • ambiguous names<br> • names that have global visibility<br> • names that have static lifetimes<br>→ IR represents the name with its coordinates<br> • lexical level (e.g. global = 0), & offset<br>→ Type info in the IR is needed for emitting correct operations |
| Which values are kept on the stack? | → Local values<br>→ Spill values from registers<br>→ Storage allocated in the activation record (AR)<br>→ These values are referenced by offset (to ARP or SP) |
| How are variables of surrounding scopes accessed? | → E.g., values from enclosing lexical scopes (excluding the global scope)<br>→ IR represents the name with its coordinates (lexical level, offset)<br>→ Compiler emits code to traverse the chain of scopes |
| How are static & global vars accessed? | → Compiler emits code to access using the label |
| How are parameters accessed? | → Compiler emits code according to the ISAs linkage convention |
| How are values on the heap accessed? | → Virtual address must be maintained<br> • in a register      • on the stack<br> • as a field in another heap-based entity |

## – Accessing Aggregate Data Structures –

| | |
|---|---|
| What are the memory access instructions in RV32I? | → load word: lw rd imm (rs1) ... reg[rd] = mem[reg[rs1] + imm]<br> • loads 4 bytes (16 bits)        • Immediate is 12 bits, sign-extended<br>→ load word pseudo-inst: lw rd label<br>→ load byte: lb rd imm (rs1) ... reg[rd] = mem[reg[rs1] + imm]<br> • loads 1 byte, sign-extended<br>→ load byte unsigned: lbu rd imm (rs1)... same as lb but loaded in byte is zero-extended<br>→ Store word: sw rs2 imm(rs1)<br>→ Store byte: sb rs2 imm (rs1)... stores LSB of rs2 |
| RECALL: What are aggregate objects?<br><br>How are they accessed? | → Structs, class objects, vectors, multi-dimensional arrays<br>→ Accessing indiv. elements of the aggregate: accessible via start addr of the object + offset within the aggregate<br>→ Elements have individual type information |

| How are struct elements accessed? | → EX : Struct Ex { int x; int y; }<br><br>1. ST lookup to find start address<br>   • struct's storage class determines location: local AR, enclosing scope, static/global label, or heap<br>2. Emit code to put start address in register, eg `la t0 struct1`<br>3. ST lookup to obtain element offset & type<br>4. Emit code to add offset to start address, eg `addi t0 t0 4`, for int y<br>5. Emit load operation (s) ... use type info to select appropriate instructions |
|---|---|
| How are object members accessed? | 1. ST lookup to find start address<br>   • Object record pointer (ORP) is associated with the object's source-language name<br>2. Emit code to put ORP in register<br>3. ST lookup to obtain object member's offset & type<br>4. Emit code to add offset to OR's start address<br>5. Emit load operation (s) ... use type info to select appropriate instructions |
| How are vector elements accessed? | 1. ST lookup to find start address & element type<br>   • Vector's storage class determines location: local AR, enclosing scope, static/global label, or heap<br>2. Emit code to put start address in register<br>3. Emit code to compute offset in vector<br>4. Emit code to add offset to start address<br>5. Emit load operation (s) ... use type info to select appropriate instructions |
| Example for vector elements? | → ```
# Vector [0:9]... a0 holds the address
# of V and a1 holds the index i
addi t1 zero 4      ── move length of each el. into t1
mul t0 a1 t1        ── t0 = index × size of el.
add a0 a0 t0        ── a0 = base addr + offset
lw t0 0(a0)
``` |
| How are multidimensional arrays accessed? | → RECALL: 2 ways to store : row/col-major order, OR indirection vectors ⟍ for 2D arrays<br>→ Num. of mem. accesses for each type is different; indirection vectors can be slow<br>   • Mem. access is expensive compared to calculations<br>   • May be harder for compiler to take advantage of locality |
| <span style="color:magenta">From → 'Storage Layout' lec</span><br><br>RECALL: What are indirection vectors? | → Innermost level is a set of vectors containing data elements<br>→ Outer level (s) are a (set of) vectors containing pointers to the next inner level |

| | |
|---|---|
| How are arrays as parameter arguments handled ?. | → Different call sites might pass arrays with different dimensions<br>→ Callee needs dimension information<br>→ Passed via a Dope vector, which is a fixed size & layout and contains a descriptor for the array |
| What are range checks ? | → Test that the given array index falls within the allocated range<br>→ Range checking is determined by the language design...either :<br>    • compiler does check statically<br>    • compiler builds code to check dynamically<br>    • compiler does no checking (e.g. C ...there are buffer overflow & other security vulnerabilities) |

## — Boolean & Relational Operators —

| | |
|---|---|
| How are boolean & relational operators used ?. | → For specifying control flow, e.g. conditionals & loops<br>→ HW support : assembly instructions for comparing, performing boolean ops (e.g. AND, OR), and acting upon the result (e.g. beq, bne, blt, bne) |
| What is short-circuit evaluation? | → Evaluating a boolean expression only until the result is determined<br>    • e.g. for if (a<b) && (c<d)...if (a<b) evaluates to false, then (c<d) won't even be evaluated<br>    • If operand is false/true for an AND/OR operation, $2^{nd}$ operand doesn't need to be evaluated<br>→ Evaluate the operand that requires fewer operations first – aka, the more shallow subtree in an AST |

## — Control-Flow Constructs —

| | |
|---|---|
| What is a basic block ? | → DEFN : Maximal-length sequence of straight-line, unlabeled, unpredicated instructions<br>→ Control-flow transfer ends the block<br>→ Labeled statement ends the block (assuming label is a possible target branch)<br>→ Block starts with labeled statements & the fall-thru instruction of a branch |
| How is an if-then-else section constructed in assembly? | → Build the basic blocks (if block, else block), then insert instructions to implement the control flow .<br>→ E.g. : |

```
if (condition) {
    block₁
} else { block₂ }
} block₃
```

⟶

```
La:
    b {eq, ne, ge, lt } rs1 rs2 Lb
    # code for block 1
    jal zero Lc
Lb:
    # code for block 2
Lc:
    # code for block 3
```

statement that checks for OPPOSITE of condition; if condition false, jump to block₂. Else, execute block₁.

| | |
|---|---|
| What logical & branching instructions does RV32I have? | → slt, slti, beq, bne, bge, blt |
| | → jump and link : jal rd label ... sets rd = PC+4, then sets PC = PC + offset |
| | (eg offset given by label) ... label is 20 bits |
| What instructions can be used to avoid branching? | → Sample Code : |

```
if (e < b) {
    x = c + d 3
else { x = e + f }
```

→ NOT in RISC-V, but :

- conditional move operations : move the value in one of 2 registers into dest reg, based on condition. e.g., for sample code above :

```
add  t0  rc  rd
add  t1  re  rf
slt  t2  ra  rb
mvcon rx  t2  t0 t1
```

- predicated operations : instr. executes or not, based on the condition, e.g. :

```
slt   t0  ra  rb
xori  t1  t0  0x01
addeqz t1  rx  rc  rd
addeqz t0  rx  re  rf
```

| | |
|---|---|
| What are condition codes ? | → Provide detailed results from comparison instructions |
| | → E.g., in x86 : — zero flag : result is 0 |
| | — sign flag : result is negative |
| | — carry flag : a borrow was needed (unsigned) |
| | — overflow flag : an overflow occurred (signed) |
| | → Enables a variety of branch instructions |
| What are loops (high level)? | → Perform iteration, w/ common features being : |
| | • Initialization   • Conditional   • update state used in guard |
| | → In the basic case (eg no nested loops) : one initial branch + one branch per iteration |
| How are for loops implemented in assembly? | → e.g. |

```
for ( expr1 ; expr2 ; expr3 ) {
    // loop body 3
    // post loop
```

1. evaluate expr1
2. if NOT (expr2) go to no.5
3. loop body
4. evaluate expr3. if expr2 go to no.3
5. post-loop

| | |
|---|---|
| How are while loops implemented ? | → Similar structure to for loops |
| | → Condition expression, but no initialization |
| | → double test loops are better |

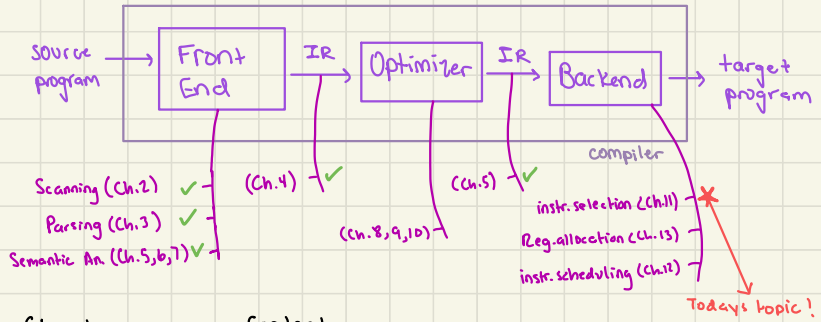| | |
|---|---|
| What are until loops? | → Guard is at the end of the loop body (execute until guard is true) |
| | → Execute at least once |
| | → More compact than while loop |
| What is a tail call? | → A call that is the last statement in a procedure |
| What is tail recursion? | → A recursive call, that is a tail call |
| | → Optimizations allow eliminating much of the procedure set-up & tear-down w/ tail recursion for iteration ... allowing resulting code to be as efficient as traditional loop |
| How are break statements handled? | → Compiler emits immediate jump to labeled breakpoint, or labeled stmt following the loop or case |
| How are skip/continue statements handled? | → Compiler emits immediate jump to code evaluating the loop condition |
| How are case statements handled? | → Basic case: evaluate the controlling expr, branch to the selected case, execute the code for that case |
| | → Break stmts exit the case. Otherwise, end of exec of a case falls through to the next case |
| | → Challenge: efficiently branching to the selected case |
| | • Compiler uses value of controlling expr to locate the corresponding case |
| What are the methods for efficiently branching to the selected case? | → Linear Search |
| | → Computed address |
| | → Binary Search |
| What is linear search? | → Implement the case as a series of nested if-then-else statements |
| | → Cost of branching to correct case depends on order of cases |
| | • Compiler should order cases according to estimated execution frequency |
| | → Works well for small number of cases |
| What is computed address? | → Compiler builds a jump table, indexed by case: vector of block-address labels |
| | → Efficient when case labels form a compact set |
| What is binary search? | → Compiler builds ordered table of case labels + block-address labels |
| | → Compiler emits binary search: |
| | • Find the right label      • branch to associated block addr |
| | → Useful when: num. of cases is small and/or label set is sparse |
| How are string operations handled? | → String length: null-terminated storage, or explicit length storage |
| | → String assignment: HW support for byte-oriented mem. operations is useful |
| | → String concatenation: |
| | • appending = string length + assignment |
| | • new string = 2 assignments |

**RECAP: Where are we now?**

Source program → Front End →[IR]→ Optimizer →[IR]→ Backend → target program

compiler

Scanning (Ch.2) ✓   (Ch.4) ✓      instr. selection (Ch.11)  ✗ ← Today's topic!
Parsing (Ch.3) ✓                  Reg. allocation (Ch.13)
Semantic An (Ch.5,6,7) ✓   (Ch.8,9,10)   (Ch.5) ✓   instr. scheduling (Ch.12)

| Chapter | Content |
|---------|---------|
| 5 | Namespaces, Naming Env., Type Info, SDT, Storage Layout |
| 6 | Procedures |
| 7 | Code Shape |

**What is the role of the backend?**

→ Optimizer has primary responsibility for efficiency

→ Backend also plays a role:

    \* **Instruction Selection**: local mapping from optimized IR to ASM

    \* **Reg. Allocation**: mapping virtual registers to ISA ones    (TM)

    \* **Inst. Scheduling**: scheduling of ISA Inst.s on target machine functional units

**What is instruction selection?**

→ Rewrite IR code into TM's ASM (assembly)

→ Challenge: many possible ASM implementations for each IR statement

⌐ Goal: efficient, local rewriting

→ Ex: $x = 1$ / $y = 1$   could turn into   addi t0 $0 1 / add t1 t0 $0   or   add t0 $0 1 / slt t1 t0 $0 , . . . . etc.

    I.S. ↑

**What is the instruction selector generator?**

→ Aka backend generator, code-generator generator

→ Compiler writer produces specifications: IR, ISA, & mapping b/w them

→ Then, backend generator (BG) builds inst. selection engine for the given specs (similar to how parser generators work)

→ Basically, instruction selection must take IR representation (s) of a piece of code — like an AST and/or ST — & produce assembly language program.

**What is our running example?**

→ $a = b - 2 \times c$, where:   a = local var in AR, w/ base addr = ARP  & offset = 4
    b = call-by-ref var in AR, w/ base addr = ARP  & offset = -16
    c = global var, w/ base addr = label "GL" & offset = 12

**What would its IRs look like?**

**What is our running example?** → $a = b - 2 \times c$, where: $a$ = local var in AR, w/ base addr = ARP & offset = 4

$b$ = call-by-ref var in AR, w/ base addr = ARP & offset = -16

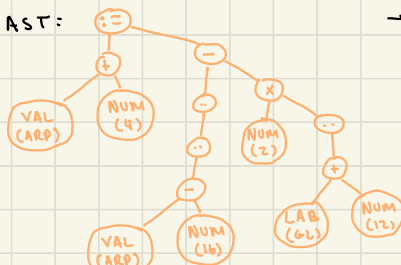$c$ = global var, w/ base addr = label "GL" & offset = 12

**What would its IRs look like?** → AST:



→ Table of quadruples:

| Result | Arg$_1$ | Arg$_2$ | Op |
|--------|---------|---------|-----|
| t0 | 2 | c | $\times$ |
| a | b | t0 | $-$ |

**How would the I.S. use these IRs?** → Linear IR (e.g. quadruples):

- text-based pattern matching
- based on peephole optimization

→ Tree-based IR (eg AST): pattern matching on trees

## — Peephole Optimization —

**What is peephole optimization?**
→ One approach to inst. selection

→ Using pattern matching over short inst. sequences to find local improvement

→ Traditionally compiler's last pass — aka, reading & writing ASM

→ Idea: 1) Translate IR into lower-level IR (LLIR)

2) Systematically improve LLIR via peephole optimization

3) Map optimized LLIR to TM instructions (ASM)

→ Basically, the compiler examines the code in small adjacent/overlapping sequences/
windows. In a given window (small sequence of code), it examines the operations in
this window, looking for specific patterns that it can improve. When it recognizes
a pattern, it rewrites it w/ a better instruction sequence.

**Examples of peephole optimizations?** → EX: a store & then load of the same value:
**(PpO)**

```
sw t0 8(sp)
lw t0 8(sp)
```

→ P.O. can replace the load w/ a copy inst., which is less expensive:

```
sw t0 8(sp)
mv t1 t0
```

→ EX: unnecessary branching:

```
jal zero La
jal zero Lb
```
→
```
jal zero Lb
```

**What are the phases of PpO?**

IR → **Expansion** → LLIR → **Simplification** → LLIR → **Matching** → ASM

- Rewrite IR into LLIR
- Capture all the IR effects
- Template driven
- Bottom-up rewrite (calculate liveness information)

- 'Sliding window' of 2-4 instructions
- Top-down rewrite:
  - Forward substitution
  - algebraic simplification
  - evaluation of constant-valued exprs
  - elimination of dead effects

- Library of patterns
- Top-down rewrite
- Final ASM (that may use virtual registers):
  - preserves all vals in LLIR (mem, regs, condition codes)
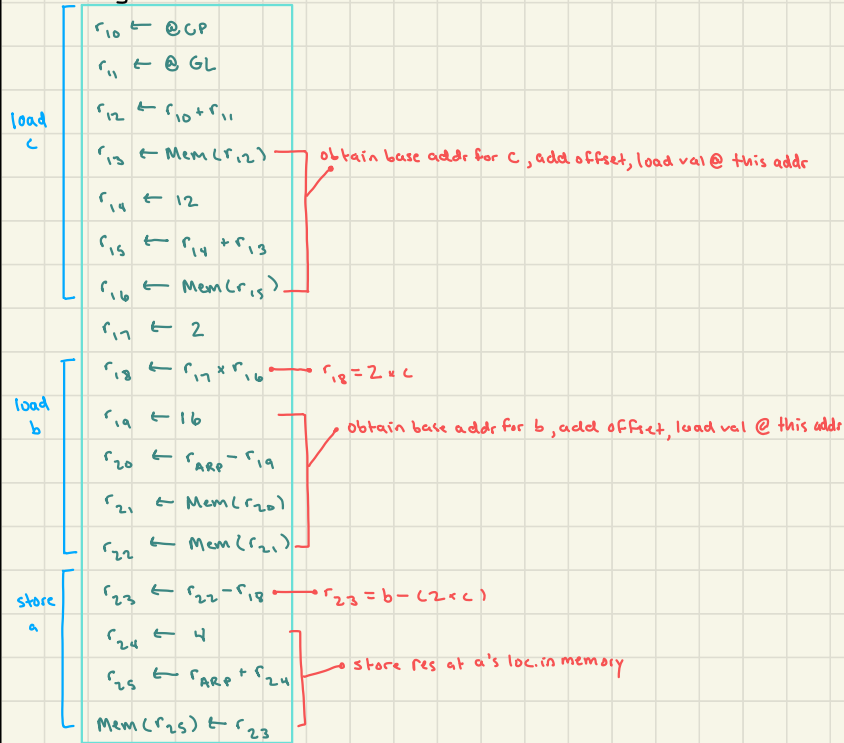
What does the IS do during expansion?

→ Rewrite IR into an LLIR, capturing all memory details ...Bottom—up rewrite

→ Template - driven : Expander uses a template for each IR op, & substitutes appropriate reg names, constants, & labels in the template.

→ Expander rewrites IR, op by op, into a sequence of LLIR ops that represent all the direct effects of an op.

Example of the LLIR that expansion might produce?

→ GIVEN : The linear IR for $a = b - 2 \times c$ :

| Result | Arg$_1$ | Arg$_2$ | Op |
|--------|---------|---------|-----|
| t0 | 2 | c | × |
| a | b | t0 | — |

→ Resulting LLIR :

load c
$r_{10} \leftarrow @GP$
$r_{11} \leftarrow @GL$
$r_{12} \leftarrow r_{10} + r_{11}$
$r_{13} \leftarrow Mem(r_{12})$ — obtain base addr for c, add offset, load val @ this addr
$r_{14} \leftarrow 12$
$r_{15} \leftarrow r_{14} + r_{13}$
$r_{16} \leftarrow Mem(r_{15})$

$r_{17} \leftarrow 2$

load b
$r_{18} \leftarrow r_{17} \times r_{16}$ — $r_{18} = 2 \times c$
$r_{19} \leftarrow 16$
$r_{20} \leftarrow r_{ARP} - r_{19}$ — obtain base addr for b, add offset, load val @ this addr
$r_{21} \leftarrow Mem(r_{20})$
$r_{22} \leftarrow Mem(r_{21})$

store a
$r_{23} \leftarrow r_{22} - r_{18}$ — $r_{23} = b - (2 \times c)$
$r_{24} \leftarrow 4$
$r_{25} \leftarrow r_{ARP} + r_{24}$ — store res at a's loc. in memory
$Mem(r_{25}) \leftarrow r_{23}$

What is 'last use' ?

→ A reference to a name (eg a reg holding a certain var/value) after which the value represented by that name is no longer 'live' (used)

→ Simplifier needs expansion to "mark last use". For ex, with

```
li t0 2
add t1 t0 t0
```
, we could easily simplify this to `addi t1 zero 4`. BUT, expander cannot eliminate the first op. unless it knows that t0 is not live after its original use in the second op.

How does the expander mark last use ?

→ Compute LiveOut sets for each block; then, in a backwards pass over the block, track which vals are live at each operation

→ Then, walk each block from bottom—up & build the LiveNow set: set of values that are live at each operation.

| How is the LiveNow set computed? | → Initialize LiveNow with either all of the names used in more than one block, or the LiveOut set for the given starting block. |
| | → Then, while processing each operation $r_i \leftarrow r_j$ op $r_k$, update LiveNow by: |
| | 1) removing $r_i$ from it |
| | 2) adding $r_j$ and $r_k$ to it |
| | → RESULT : alg produces, at each step, a LiveNow set that is as precise as the initial info used at the bottom of the block. |

| What does the IS do during Simplification? | → Sliding window : considers 3-4 instructions at a time |
| | → Pattern-match to identify simplifications : |
| | • forward substitution •  algebraic simplification |
| | • evaluation of constant-valued expressions • elimination of dead effect |
| | → Emit simplified instructions ; discard dead instructions |
| | → Basically, Simplifier makes a pass over the LLIR, examining the ops in a small window on the LLIR & systematically trying to improve them. |
| What is forward substitution? | → Substituting defined values into uses later in the window |
| | → Example : |

```
la  r_5  L1          la  r_5  L1
li  r_6  12    ⟶     lw  r_8  12(r_5)
add r_7 r_5 r_6
lw  r_8  0(r_7)
```

→ Rather than defining (12 + L1) as a new value, we can directly do a load word.

| What is algebraic simplification? | → Using algebraic identities to simplify, e.g. : |

```
subi t0 t1 0     ⟶   mul t3 t1 t2
mul t3 t0 t2
```

→ Algebraic simplifications to look for :

| | | | |
|---|---|---|---|
| • a + 0 = a | a × 1 = a | a && a = a | a << 1 = a + a |
| a - 0 = a | a × 2 = a + a | a ll a = a | |
| a - a = 0 | a/1 = a | a >> 0 = a | |
| a ∧ 0 = 0 | a/a = 1 | a << 0 = a | |

**What is constant folding?** → Evaluate constant-valued expressions & substitute the value

$$\boxed{\begin{array}{l} \text{li } r_2 \ 2 \\ \text{add } r_2 \ r_2 \ r_2 \end{array}} \longrightarrow \boxed{\text{addi } r_2 \ \text{zero } 4}$$

**How does the simplifier eliminate dead operations?** → Simplifier can only eliminate defining operations if the defined value is dead

→ Simplifier uses the LiveNow computation from the Expander

**Example?**



$$\boxed{\text{li } r_2 \ 2} \quad \bullet \ \text{LN} = \emptyset$$
$$\quad \bullet \ \text{LN} : \{r_2\}$$
$$\boxed{\text{add } r_2 \ r_2 \ r_2} \quad \bullet \ \text{LN} : \{r_2\}$$

**What would oor example LLIR look like after simplification?**

$$
\begin{array}{l}
r_{10} \leftarrow @CP \\
r_{11} \leftarrow @GL \\
r_{12} \leftarrow r_{10} + r_{11} \\
r_{13} \leftarrow \text{Mem}(r_{12}) \\
r_{14} \leftarrow 12 \\
r_{15} \leftarrow r_{14} + r_{13} \\
r_{16} \leftarrow \text{Mem}(r_{15}) \\
r_{17} \leftarrow 2 \\
r_{18} \leftarrow r_{17} \times r_{16} \\
r_{19} \leftarrow 16 \\
r_{20} \leftarrow r_{ARP} - r_{19} \\
r_{21} \leftarrow \text{Mem}(r_{20}) \\
r_{22} \leftarrow \text{Mem}(r_{21}) \\
r_{23} \leftarrow r_{22} - r_{18} \\
r_{24} \leftarrow 4 \\
r_{25} \leftarrow r_{ARP} + r_{24} \\
\text{Mem}(r_{25}) \leftarrow r_{23}
\end{array}
\longrightarrow
\begin{array}{l}
r_{13} \leftarrow \text{Mem}(@CP + @G) \\
r_{16} \leftarrow \text{Mem}(r_{13} + 12) \\
r_{18} \leftarrow r_{16} + r_{16} \\
r_{21} \leftarrow \text{Mem}(r_{ARP} - 16) \\
r_{22} \leftarrow \text{Mem}(r_{21}) \\
r_{23} \leftarrow r_{22} - r_{18} \\
\text{Mem}(r_{ARP} + 4) \leftarrow r_{23}
\end{array}
$$

## — IS $P_eO$: Matching —

**How does matching work?** → Matcher compares simplified LLIR against the pattern library, looking for the pattern that best captures all the effects in the LLIR

→ Some operations are a one-to-one match; others may require multiple ISA instructions for one LLIR inst. (this may introduce temporary names)

→ Matcher must know:
- size of constants
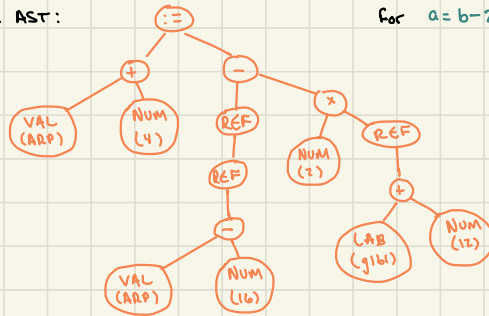- address modes
- available memory operations

**RECAP ?**

→ Peephole optimizations use linear IRs to do instruction selection

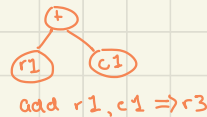→ Another way to do IS is using tree-based IRs, like ASTs (which can be an LLIR)

→ Example AST:                                        for  a = b - 2 × c



**What is an operation tree ?**

→ Tree-representations of operations in target ASM . e.g. :



add r1, c1 => r3

**What is the library of tree patterns?**

→ Rewrite rules; each pattern maps constructs in the IR to operations in target ASM

→ Each pattern contains:
- low-level IR pattern tree            • cost
- ASM code template

**What is tiling?**

→ a Tile is <x, y>, where :
- x is a subtree in the AST
- y is a tree pattern from the Library of Tree Patterns
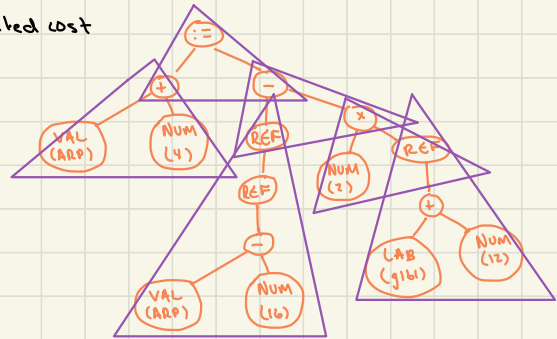- pattern y can implement the subtree at node x

→ Tiling : a set of tiles that implement the AST
- covers every node in the AST
- overlap between any 2 trees occurs at a single node
- where 2 pattern trees overlap, they agree in storage class & value type
- root of each pattern tree overlaps a leaf of another pattern tree, unless it covers the root of the AST

**What is the goal?**
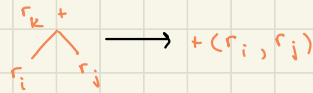
→ Find a locally optimal tiling

→ Triangles represent tree patterns: each has an associated ASM code template & associated cost



**How do we describe an LLIR AST using prefix notation?**

→ Write the op. tree for <u>add</u> for ex, as:

$$r_k\ + \qquad \longrightarrow\ +(r_i, r_j)$$
$$r_i \quad r_j$$

→ Combine them all. Ex for the AST above:

$:= ( + (VAL(ARP), 4), - (- (VAL(ARP), 16), \times (2, + (LAB(glbl), 12))))$

# Register Allocation

| | |
|---|---|
| How does register allocation work? | 1. Determine, at each point in the code. <br>     • which values reside in registers <br>     • which values reside in memory <br> 2. Rewrite the code to implement the allocation <br>     • spill & restore values when there are too many live values <br>     • keep ambiguous values in memory |
| RECALL: What are registers? | → Defined by the ISA <br> → Provide fast access b/c physically on the CPU, and no address translation |
| What is register allocation? | → Input: IR program that is <u>nearly</u> executable ASM, but uses an arbitrary number of virtual registers <br> → Output: equivalent program that fits into the ISA's register set <br>     • aka executable ASM <br> → Can be followed by: <br>     • instruction scheduling |
| RECALL: What is the register-to-register model? | → IR uses as many virtual registers (VRs) as needed <br> → Register allocator maps VRs to physical registers (PRs) |
| What is the memory-to-memory model? | → IR establishes a home memory location for all named values <br> → Register allocator can promote a value to a register for a longer range of its lifetime |
| What are 2 methods for register allocation? | → Local register allocation: <br>     • each basic block analyzed independently (scope of analysis = basic block) <br>     • a single class of physical registers (PR) is given <br>     • assumes analysis begins with no values in PRs <br>     • namespace for values: Live Range <br>     • key algorithm: compute most distant next use to determine spills <br> → Global register allocation: <br>     • each procedure is analyzed as a whole <br>     • may be multiple classes of physical registers <br>     • key algorithm: uses graph coloring to compute allocation <br>     • namespace for values: Live Range |
| How does local register allocation work? | 1. Build a new namespace of Live Ranges (LLR) <br> 2. Perform register allocation (calculate interference; generate spill code) <br> 3. Rewrite the code using the namespace of physical registers |

| What are Live Ranges? | → DEFN: A closed set of related definitions and uses that serves as the base name |
| (LLR) | space for register allocation |

**What are Live Ranges? (LLR)**

→ DEFN: A closed set of related definitions and uses that serves as the base name space for register allocation

→ Each LR corresponds to a single value
- from the value's definition
- to the value's last use

→ In a basic block, each LR
- starts with one definition
- extends until that definition's last use

→ An LR starts: with the first operation after it is defined
- differentiates b/w a use and a definition in the same operation

→ An LR ends:
- at last use of the name, OR
- when the name gets a new value

**What is liveness?**

→ A variable $v$ is live at point $p$ if it has been defined along a path from the procedure's entry to $p$, and there exists a path from $p$ to a use of $v$ along which $v$ is not redefined.

→ Anywhere that $v$ is live, its value must be preserved because subsequent execution might use $v$.

**Example?**

→ For $a = a \times 2 \times b \times c \times d$, we have this IR block, and each distinct live range in the block:

| | IR | Register | Interval |
|---|---|---|---|
| 1 | $r_{arp} \leftarrow Mem(ARP)$ | $r_{arp}$ | $[1, 11]$ |
| 2 | $r_a \leftarrow Mem(r_{arp} + offset)$ | $r_a$ | $[2, 7]$ |
| 3 | $r_2 \leftarrow 2$ | $r_a$ | $[7, 8]$ |
| 4 | $r_b \leftarrow Mem(B)$ | $r_a$ | $[8, 9]$ |
| 5 | $r_c \leftarrow Mem(C)$ | $r_a$ | $[9, 10]$ |
| 6 | $r_d \leftarrow Mem(D)$ | $r_a$ | $[10, 11]$ |
| 7 | $r_a \leftarrow r_a \times r_2$ | $r_2$ | $[3, 7]$ |
| 8 | $r_a \leftarrow r_a \times r_b$ | $r_b$ | $[4, 8]$ |
| 9 | $r_a \leftarrow r_a \times r_c$ | $r_c$ | $[5, 9]$ |
| 10 | $r_a \leftarrow r_a \times r_d$ | $r_d$ | $[6, 10]$ |
| 11 | $Mem(A) \leftarrow r_a$ | | |

→ basically, a reg's LR is the interval over which it represents (& is used as) a given value. E.g., $r_a$ has an LR from $(2, 7)$ b/c line 7 uses the val from operation 2. Op. 8 uses the value from Op. 7 (rather than op.2), SO, $r_a$ has an ^diff LR from $[7, 8]$

**More about LRs?**

→ A reg. allocator doesn't have to keep each distinct LR^in the same reg... *(for one register. e.g. $r_a$)*

it can treat each LR in the block as an independent val for allocation

& assignment. For ex, the ILOC ex from prev. page could become:

```
lw   r_a  @a(sp)
addi r_2  zero 2
lw   r_b  @b(sp)
lw   r_c  @c(sp)
lw   r_d  @d(sp)
mul  r_1  r_a r_2
mul  r_3  r_b r_1
mul  r_4  r_3 r_c
mul  r_5  r_4 r_2
sw   r_5  @a(sp)
```

*in local Reg Alloc, 2 LRs interfere if their spans overlap* ↑

**What is interference?**

→ 2 LRs interfere if there exists an operation for which both are live and they (probably) don't have the same value.

→ Two LRs can share a physical reg. i.f.f. they do not interfere, and use the same PR set... 2 LRs using distinct PRs cannot interfere

  • in local Reg Alloc, all LRs use the same PR set

## — Local Register Allocation —

**What is Best's Algorithm?**

1. Compute live information (bottom-up pass)

2. Add interference info (top-down pass)

→ Idea: Most distant next use

  → *aka a store in mem*

  • Spill the PR with the val whose next use is furthest in the future

  • reduce demand for PRs over the longest interval

→ Optimal choice if spill costs are uniform (in general tho this isn't the case)

**What are non-uniform spill costs?**

→ Dirty values, clean values, rematerializable values

**What is a dirty value?**

→ A value that does not exist in memory

  • STRAT: store at spill points (out of registers (?.)), load at restore points

**What is a clean value?**

→ A value that is already in memory

  • STRATEGY: load at restore points

**What is a rematerializable value?**

→ A value that can be recomputed at each use

  • STRATEGY: recompute at restore points (e.g. load immediate operations)

**How are the live ranges determined?.**

→ Maintain 2 maps :

- • VRtoLR : initialized to INV for all virtual registers (VRs)
  → VRtoLR[x]=y means that the curr. val in register x belongs to LR y.
- • PrevUse : initialized to inf for all VRs
  → maps virtual register to index of closest future use
  → PrevUse(v) = the next inst. that uses VR v (while scanning bottom-up).

→ Bottom-up Pass

- • visit definitions first, <u>then</u> uses
- • update maps
- • fill in the LR and Next Use (NU) fields

**What would this pseudocode look like?.**

```
For each operand of the inst.:
    For each op defined:
        if VRtoLR [op.VR] == INV, then VRtoLR[op.VR] := next LR
        op.LR = VRtoLR (op.VR)
        op.NU = PrevUse [op.VR]
        PrevUse [op.VR] = inf
        VRtoLR (op.VR) = INV

    For each op used:
        if VRtoLR [op.VR] == INV, then VRtoLR[op.VR] := next LR
        op.LR = VRtoLR (op.VR)
        op.NU = PrevUse [op.VR]

    For each op used:
        PrevUse (op.VR) = index of current instruction
```

**So how are registers assigned?.**

→ Process uses before definitions

→ For each operation:

1. For each use, look for an existing PR assigned to the LR
   a. maintain LRtoPR map
   b. call getPR if new PR is needed
2. Determine if a use is the last use of the LR
   a. if so, free the PR, making it avail. for reassignment
3. For each definition, assign a PR to the LR
   a. call getPR
   b. update the LRtoPR map

| | |
|---|---|
| How would the getPR function work? | 1. Find a PR for a reference v |
| | 2. If there exist free PRs, use the next available (maintain a stack of free PRs) |
| | 3. Else, choose an LR to evict from its PR: |
| | • spill the LR val to mem/stack |
| | • reassign the PR to v |
| | • if the ref to v is a use rather than a definition, restore its value from its mem. location, to the newly assigned PR |
| RECAP: What are the steps for RegAlloc? | 1) Translate to namespace of live ranges |
| | 2) Identify interferences b/w live ranges |
| | 3) Find an allocation |
| | 4) Decide what to spill |
| What is the "Max Live" value? | → Maximum demand for registers in a basic block |
| | → Given by the max. num. of concurrently live ranges |
| | → If Max Live is greater than the num. of PRs, a spill is needed |

**— Global RA —**

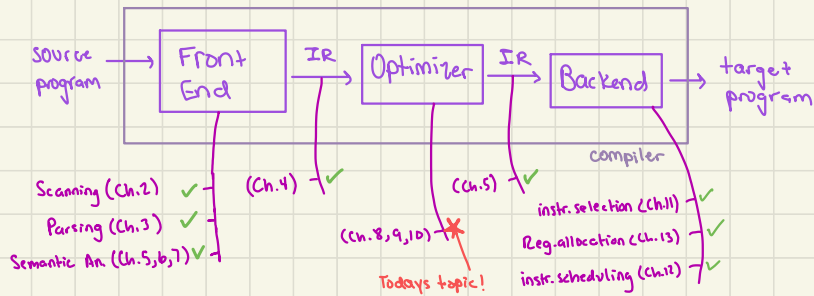| | |
|---|---|
| What are the challenges with local RA? | → LRs have multiple definitions & uses |
| | → RA must coordinate across basic blocks |
| | → Spill cost estimation must account for blocks that execute multiple times |
| How does global RA work? | 1. Find LRs |
| | 2. Build interference graph |
| | 3. Coalesce copies — loop to step 2 |
| | 4. Estimate spill costs |
| | 5. Find a coloring |
| | 6. Insert spills — loop back to step 1 |

# Optimization

RECAP: Where are we now?



Scanning (Ch.2) ✓
Parsing (Ch.3) ✓
Semantic An (Ch.5,6,7) ✓

(Ch.4) ✓

(Ch.8,9,10) ✗
Today's topic!

(Ch.5) ✓

instr. selection (Ch.11) ✓
Reg.allocation (Ch.13) ✓
instr. scheduling (Ch.12) ✓

What are the goals & process in optimization?
→ Goal: improve performance of compiled code
→ Process: Analysis, Transformation

What are the considerations?
→ **Safety**: transformed code produces same results as original code
→ **Profitability**: (trnsfm) transformation improves performance
→ **Risk**: Trnsfms intended to improve performance can make it harder for the compiler to generate good code

What are the different scopes where we can perform optimization?
→ **Local**: one basic block
→ **Regional**: a subset of the CFG
→ **Global**: one procedure
→ **Interprocedural**: 2 or more procs

## — Local Value Numbering —

What is the "local scope" exactly?
→ A single basic block = scope for local OPT methods; so each basic block

RECALL: What is a basic block?
→ RECALL: basic block = maximal-length sequence of branch-free code
  - Starts w/ a label or fell-through instruction
  - Ends with a branch or jump
  - Instructions execute sequentially
  - If one inst. in a BB executes, they all do

What is Local Value Numbering (LVN)?
→ A type of local optimization for eliminating redundancy
→ Finds & eliminates redundantly computed expressions:
  - exprs previously computed in the basic block, AND
  - exprs where no operands have been redefined since 1st computation

Examples of redundant expressions? →

```
a ← b + c
b ← a - d
c ← b + c    → repeats line 1
d ← a - d    → repeats line 2
```

```
a ← b + c
d ← b
e ← d + c
```
→ L2 is redundant b/c we could just do "e ← b + c"

**What is the algorithm for LVN?**

→ Traverse block **top-to-bottom**, and map each value to a distinct number
　(each "value" ≈ each operand or expression)
→ Two vals map to same number i.f.f. they are equal in all possible executions
→ Use a hash of each val to index into the mapping table (or use a dictionary)
→ Start with an empty hash table mapping value names to their distinct
　numbers / IDs :　　HT =

| Val number | name |
|---|---|
| ... | ... |

**What is the pseudocode?**

```
num = 0
For each instruction T ← L Op R :
      VN_L = getNumber(L)
      VN_R = getNumber(R)
      if < VN_L, Op, VN_R >  is an entry in HT :
          x = getName(< VN_L, Op, VN_R >)
          VN_x = getNumber(< VN_L, Op, VN_R >)
          if getNumber(x) == VN_x :
              rewrite instruction to be "T ← x"
          else :
              set the "name" associated w/ the same
              value, < VN_L, Op, VN_R > (aka name whose
              value = value for < VN_L, Op, VN_R > entry )
              to be INVALID
      else:
          HT[< VN_L, Op, VN_R >] = num
          num += 1
      VN_T = HT[< VN_L, Op, VN_R >]
      HT[T] = VN_T
def getNumber(operand) :
   if HT[operand] :   // if entry in HT alr exists for this value
      return HT[operand]
   else :
      HT[operand] = num
      num += 1
      return HT[operand]
```

find the number
that < VN_L, Op, VN_R >
is assoc. with, aka
HT[< VN_L, Op, VN_R >]
Then, find the other
entry in HT that has
the _same_ number.
the key for this
entry — aka the other
name sharing that num—
is returned

Example walkthrough?

| | HT: | num | name |
|---|---|---|---|
| a ← b+c | | ~~0~~ 4 | b |
| b ← a-d | | ~~1~~ 5 | c |
| c ← b+c | | 2 | 0+1 |
| ~~d ← a-d~~ → d ← b | | 2 | a |
| | | ~~3~~ 4 | d |
| | | 4 | 2-3 |
| | | 5 | 4+1 |

1. a ← b+c
   - create new entries for b & c → HT[b]=0 and HT[c]=1
   - create new entry for "0+1" → HT[0+1]=2
   - create new entry for a, same num as "0+1" → HT[a]=2

2. b ← a-d
   - create new entry for d → HT[d]=3
   - create new entry for "2-3" → HT[2-3]=4
   - set HT[b]= HT[2-3]=4

3. c ← b+c
   - create new entry for "4+1" → HT[4+1]=5
   - set HT[c]=HT[4+1]=5

4. d ← a-d
   - entry already exists for "2-3"!
   - x = getName [2-3] = entry where num is 4 = b
   - $VN_x = 4$        • getNum (x) = 4
   - rewrite Instruction to: d ← b
   - set HT[d]= HT[2-3]=4

How can the LVM algorithm be extended for more optimizations?

→ Commutative operations                    → Algebraic identities

→ Constant folding

How would the LVN use constant folding?

→ Constant Folding: Store values for constants, rather than names.
   - For instructions (that are operations) where all operands are constants, LVN can evaluate/perform the operation, & then replace the instruction with a load imm. inst. that loads in the evaluated result.
   - Ex:   a ← 3          → Rewrite L2 using →    a ← 3
           b ← a+2            constant folding       b ← 5

How would the LVM use commutative operations?

→ Commutative operations: ensure that operations that are commutative, eg a×b and b×a, receive the same value numbers/ are treated as the same value

**How would the LVN use algebraic identities?**

→ <u>Algebraic identities</u>: Keep a dictionary **AID** of algebraic identities that the LVN can apply to simplify the code, e.g. :

| Operator | Expr | Simplified identity |
|---|---|---|
| add | $a + 0$ | $= a$ |
| sub | $a - 0$ | $= a$ |
| sub | $a - a$ | $= 0$ |
| mul | $2 \times a$ | $= a + a$ |
| mul | $1 \times a$ | $= a$ |
| mul | $0 \times a$ | $= 0$ |
| div | $a \div 1$ | $= a$ |
| div | $a \div a$ | $= 1$ |

| Operator | Expr | Simplified identity |
|---|---|---|
| exp | $a^1$ | $= a$ |
| exp | $a^2$ | $= a \times a$ |
| srl | $a >> 0$ | $= a$ |
| sll | $a << 0$ | $= a$ |
| AND | $a \& a$ | $= a$ |
| OR | $a \parallel a$ | $= a$ |
| MAX | $max(a, a)$ | $= a$ |
| min | $min(a, a)$ | $= a$ |

• here, each "simplified identity" represents an operation which is equivalent to its corresponding $Expr$, BUT is simpler/faster/more efficient. For ex, $sub\ \$b\ \$a\ \$a$ is equivalent to but less efficient than doing $li\ \$b\ 0\times0$

→ For each instruction $T_i \leftarrow L_i\ Op\ R_i$, iterate thru all the entries in **AID** where "operator" is $Op$. If any of the $Exprs$ match the instruction, replace the instruction with the corresponding simplified version.

→ ==Performing LVN over an IR in SSA prevents lost values==

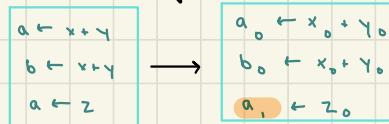## — Static Single Assignment (SSA) —

**What is SSA Form?**

→ Linear IR + Namespace ; encodes control flow & value flow into names

→ Each name corresponds to one definition point in the code
  • each defn. has a distinct name
  • each use refers to a single defn.

**How do we translate a Linear IR into SSA Form?**

→ At each definition, generate a **new name** − e.g., by adding a subscript:

$$a \leftarrow x + y$$
$$b \leftarrow x + y$$
$$a \leftarrow z$$

$\longrightarrow$

$$a_0 \leftarrow x_0 + y_0$$
$$b_0 \leftarrow x_0 + y_0$$
$$a_1 \leftarrow z_0$$

→ At each control-flow join, add a $\emptyset$-function − a "fake" instruction that lives at the top of a basic block that has multiple predecessors
  • the $\emptyset$-function basically allows compiler to choose which previous definition's value should get assigned as the value of a new defn., when that value is dependent on the control flow that led up to it.

→ Compiler inserts $\emptyset$-functions at points where different control-flow paths merge, and it then renames variables to make the single-assignment property hold.

**How do the ∅-functions work?**

→ At control-flow merge points, a var may come from differing previous definitions depending on which path was taken

→ The ∅ operation selects which prev. defn will be used.

**small example?**

→ Ex:  $a_2 \leftarrow \emptyset(x_1, x_2)$

- When execution enters this block, if we came from predecessor 1, then $a_1$ shld take the value of $x_1$. If we came from predecessor 2, then $a_2$ should take the value of $x_2$.

→ ==SSA attaches each argument of the ∅ to a particular incoming edge in the control-flow graph (CFG)==; the ∅-function selects the arg. corresponding to the CFG edge just taken
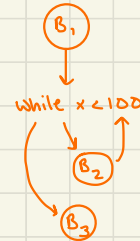
**How would we turn a while block into SSA form?**

→ Example linear IR:

```
x ← ...
y ← ...          ] ─ B₁
while (x<100)
   x ← x+1
   y ← y+z       ] ─ B₂
... (B₃)
```

→ CFG:



→ There are 2 merge points:

1) Top of the loop body, where control could either enter from above (e.g. the first iteration), OR from the bottom of the loop, when starting next iter.

2) The block AFTER the loop, where the val of $x$ could either come from the while loop, or the prev. block if we never entered the loop

→ SSA encoding:

```
x₀ ← ...
y₀ ← ...
bge x₀ 100 continue
Loop:
   x₁ ← ∅(x₀, x₂)
   y₁ ← ∅(y₀, y₂)
   x₂ ← x₁ +1
   y₂ ← y₁ + x₂
Continue:
   x₃ ← ∅(x₀, x₂)
   yₑ ← ∅(y₀, y₂)
```

**What is Superlocal Valve Numbering (SLVN)?**

→ An extension of LVN, applied to a regional scope

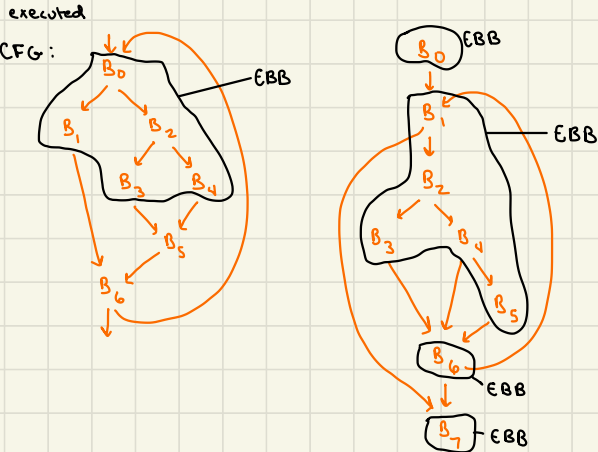→ The regional scope = Extended Basic blocks (EBBs)

**What is an EBB?**

→ A set of basic blocks $\{B_0, B_1, \ldots B_n\}$ where :

- $B_0$ = the entry node, OR $B_0$ has multiple CFG predecessors, AND
- all other $B_i$ have only one predecessor, which is in the set

→ If one instruction executes, every prior inst. on the CFG path in the EBB also executed

**Examples of EBBs?**

→ CFG :



**How does SLVN work?**

→ Apply LVN to paths through EBBs : treat each path thru an EBB as straight-line code & find redundancies that cross basic block boundaries

→ Challenge : eliminating redundant LVN work

**Example?**

$B_0:$ $m_0 = a_0 + b_0$
$\quad n_0 = a_0 + b_0$
$\quad (a_0 > b_0) \longrightarrow B_1, B_2$

$B_1:$ $p_0 = c_0 + d_0$
$\quad r_0 = c_0 + d_0$
$\quad \longrightarrow B_6$

$B_2:$ $q_0 = a_0 + b_0$
$\quad r_1 = c_0 + d_0$
$\quad (a_0 > b_0) \longrightarrow B_3, B_4$

$B_3:$ $e_0 = b_0 + 18$
$\quad s_0 = a_0 + b_0$
$\quad u_0 = e_0 + f_0$
$\quad \longrightarrow B_5$

$B_4:$ $e_1 = a_0 + 17$
$\quad t_0 = c_0 + d_0$
$\quad u_1 = e_1 + f_0$
$\quad \longrightarrow B_5$

$B_5:$ $e_2 = \emptyset(e_0, e_1)$
$\quad u_2 = \emptyset(u_0, u_1)$
$\quad v_0 = a_0 + b_0$
$\quad w_0 = c_0 + d_0$
$\quad x_0 = e_2 + f_0$
$\quad \longrightarrow B_6$

$B_6:$ $r_2 = \emptyset(r_0, r_1)$
$\quad y_0 = a_0 + b_0$
$\quad z_0 = c_0 + d_0$

**What are the EBBs in this example?**

→ The CFG for this code:



→ EBBs:

$\{B_0, B_1, B_2, B_3, B_4\}$ ✗

$\{B_5\}$

$\{B_6\}$

---

**How would basic SLVN work on the starred EBB?**

→ By treating each path as if it were a single block, & performing LVN on each possible path.

→ Possible paths: $(B_0, B_1)$, $(B_0, B_2, B_3)$, $(B_0, B_2, B_4)$

→ For each path, it would perform LVN on each block in the path, & then simply carry the hash table forward to the next block in the path

---

**How would the basic SVN alg work on this entire block?**

1. Create scope (e.g., HT) for $B_0$
2. Apply LVN to $B_0$ — eliminates redundant $n_0$ assignment
3. Create scope for $B_1$: e.g., carrying forward $B_0$'s HT; apply LVN to $B_1$
4. Add $B_6$ to Worklist
5. Delete $B_1$'s scope: e.g., revert changes made to o.g. HT    from $B_0$ when we analyzed $B_1$
6. Create scope for $B_2$; apply LVN to $B_2$
7. Create scope for & apply LVN to $B_3$
8. Add $B_5$ to Worklist
9. Delete $B_3$'s scope
10. Create scope for $B_4$; apply LVN to $B_4$
11. Delete $B_4$'s scope
12. Delete $B_2$'s scope
13. Delete $B_0$'s scope
14. Create scope for $B_5$; apply LVN to $B_5$
15. Delete $B_5$'s scope
16. Create scope for $B_6$; apply LVN to $B_6$
17. Delete $B_6$'s scope

---

**What are the drawbacks of this naive approach?**

→ Analyzing basic blocks once per path = lots of repeated work; inefficient

---

**What is a better approach?**

→ Basic blocks are analyzed once, and ValueNumber table is reused

  • checkpoint state at BB boundaries to restore table's state
  • unwind effects on VN table by walking BB backward, OR
  • use linked list of VN tables & delete leaf tables to start a new path

→ Requires SSA form

# Data Flow Analysis

| | |
|---|---|
| What is Data-flow analysis? (DFA) | → A way of reasoning (at compile time) about the flow of values at runtime |
| | → Goal: analyze prog. in support of optimizations |
| | → Process: Iterative DFA |
| What is the scope? | → global analysis, or Full CFG |
| What are the considerations? | → Precision, correctness, & termination |
| What is iterative Data Flow Analysis? <br> 9'.01 | → Iterative fixed-point algorithms where the fixed point exists & is unique <br> → For each node of the CFG: <br>    1. evaluate an equation <br>    2. update the set defined by the eqn. <br>    3. go back to step 1 until fixed point is reached <br> → BENEFIT: simple & robust <br>    • works directly on any valid CFG <br>    • guaranteed to terminate |
| Why is precision important to consider? | → DFA is a static analysis; it may be imprecise, & transformations must be conservative <br> → Sources of imprecision: <br>    • control flow <br>    • ambiguous references (analysis must assume any feasible value is possible) <br>    • reading from an external source (adds non-determinism to program) <br>    • external procedures |
| Why is control flow a source of imprecision? | → Analysis assumes all paths through the CFG are feasible <br> → Analysis is precise up to symbolic execution |
| Why are procedure calls a source of imprecision? | → Analysis must assume a callee will use & modify any accessible variable |

## — Dominance —

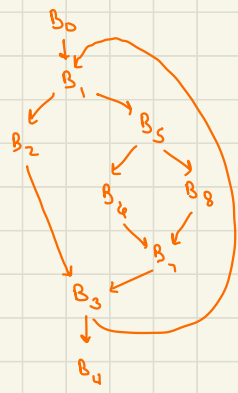| | |
|---|---|
| What is dominance? | → Given a CFG with nodes $\{b_0, \dots b_n\}$ where $b_0$ = entry node, node $b_i$ dominates node $b_j$ (written $b_i \gg b_j$) if & only if $b_i$ lies on every path from $b_0$ to $b_j$ <br> → By defn, $b_i \gg b_i$ for all $b_i$ <br> → Dominance = a tool used by compiler for reasoning about the shape & structure of a CFG <br>    • Important in construction of SSA form |
| What is a DOM set? | → Each node $b_i$ in a CFG has a set $Dom(b_i)$, which contains the names of all nodes that dominate $b_i$ |

**Example of DOM sets given a CFG?**



- $DOM(B_0) = \{B_0\}$
- $DOM(B_1) = \{B_0, B_1\}$
- $DOM(B_2) = \{B_0, B_1, B_2\}$
- $DOM(B_3) = \{B_0, B_1, B_3\}$
- $DOM(B_4) = \{B_0, B_1, B_3, B_4\}$
- $DOM(B_5) = \{B_0, B_1, B_5\}$
- $DOM(B_6) = \{B_5, B_0, B_1, B_6\}$
- $DOM(B_7) = \{B_5, B_0, B_1, B_7\}$
- $DOM(B_8) = \{B_5, B_0, B_1, B_8\}$

**What is immediate dominance?**

→ Given a CFG with nodes $\{b_0, \ldots b_n\}$ where $b_0$ = entry node :

- node $b_i$ is the immediate dominator of $b_j$ iff :
  → $b_i$ dominates $b_j$
  → $b_i$ is the dominating node closest to $b_j$
  → $b_i \neq b_j$
- $b_0$ has no IDOM

**Example finding IDOM sets given a CFG?**

→ From same CFG as above :

| node | IDOM | node | IDOM |
|------|------|------|------|
| 0 | $\emptyset$ | 5 | $B_1$ |
| 1 | $B_0$ | 6 | $B_5$ |
| 2 | $B_1$ | 7 | $B_5$ |
| 3 | $B_1$ | 8 | $B_5$ |
| 4 | $B_3$ | | |

**How does the compiler compute dominance, conceptually?**

→

$$DOM(n) = \{n\} \cup \left( \bigcap_{m \in preds(n)} DOM(m) \right)$$

(and $DOM(b_0) = \{b_0\}$)

→ MEANING: for node $n$, compute $DOM(n)$ :

1. Add $\{n\}$ to $DOM(n)$
2. For each node that precedes $n$ (e.g. every node by which $n$ is reachable), get the intersection of the DOM sets of each of those nodes. Add this to $DOM(n)$

**How is this computation /algorithm implemented?**

→ Iterative Fixed-point ALG for a CFG with $|N|$ nodes:

1. initialize $DOM(b_0) = \{b_0\}$

2. For each node $b_1, \ldots b_n$, initialize $DOM(n) = \{b_0, b_1, \ldots b_n\}$

   (aka, initialize each node's DOM set as the set of all nodes)

3. For $i=1$ to $|N|-1$:

   $$DOM(b_i) = \{b_i\} \cup \left( \bigcap_{b_m \in preds(b_i)} DOM(b_m) \right)$$

4. Repeat step 3 until, for all $i$, $DOM(b_i)$ does not change at all (aka fixed point reached)

What is the proof that this alg considers termination?

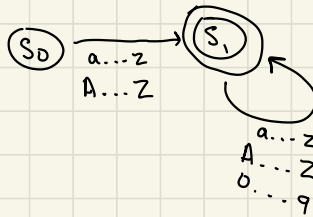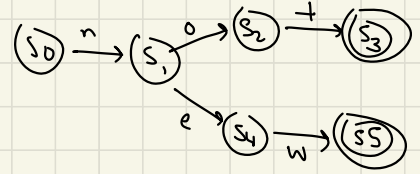DFA to accept "not" and "new"    let    X = error state    <span style="color:green">Yay COMPLSS review!</span>
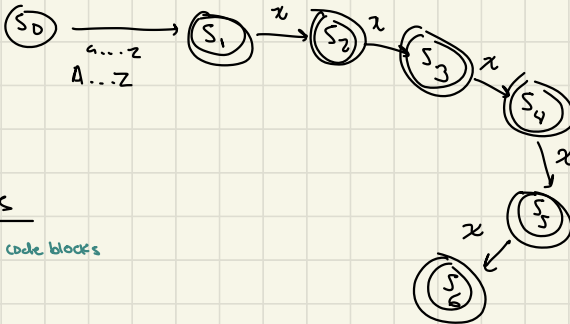
$\delta : \{S_0, S_1, S_2, S_3, S_4, S_5\}$

$\Sigma : \{n, o, t, e, w\}$

$\delta :$

|    | n  | o  | t  | e  | w  |
|----|----|----|----|----|----|
| S0 | S1 | x  | x  | x  | x  |
| S1 | x  | S2 | x  | S4 | x  |
| S2 | x  | x  | S3 | x  | x  |
| S3 | x  | x  | x  | x  | x  |
| S4 | x  | x  | x  | x  | S5 |
| S5 | x  | x  | x  | x  | x  |





Identifier with one alphabetic char followed 0 or more alphanumeric chars

let set $X = a \ldots z$
              $A \ldots z$
              $0 - 9$



Identifier with one alphabetic char followed by up to 5 alphanumeric chars

COLORS

<span style="color:teal">• Main secondary</span>   <span style="color:teal">+ code blocks</span>

<span style="color:red">• Main secondary</span>

<span style="color:purple">• Main purple</span>

<span style="color:violet">• Secondary purple</span>

<span style="color:purple">• neutral purple</span>