

CS132: Compiler Construction

Professor Palsberg

Thilan Tran

Fall 2020

Contents

CS132: Compiler Construction	3
Introduction	3
Compiler Overview	3
Lexical Analysis	6
Parsing	7
Top-Down Parsing	8
Grammar Hacking	9
Achieving LL(1) Parsing	11
Predictive Parsing	11
Handling Epsilon	12
Formal Definition	13
JavaCC	13
Handling Syntax Trees	15
Visitor Pattern	15
Java Tree Builder	17
Type Checking	20
Simple Expressions	20
Statements	21
Declarations	22
Arrays	23
Methods	24
Classes	25
Entire Java Program	26

Sparrow	28
Grammar	28
Translation to IR	29
State and Transitions	29
Expressions	30
Statements	31
Examples	33
Arrays	34
Classes	35
Inheritance	36
Register Allocation	38
Liveness Analysis	38
Algorithm	40
Graph Coloring	41
NP-Completeness	44
Appendix	46
Practice Questions	46

CS132: Compiler Construction

Introduction

- a **compiler** is a program that *translates* an executable program in one language to an executable program in another language
- an **interpreter** is a program that *reads* an executable program and produces the results of running that program
 - usually involves executing the source program in some fashion, ie. *portions* at a time
- compiler construction is a *microcosm* of CS fields:
 - AI and algorithms
 - theory
 - systems
 - architecture
- in addition, the field is not a solved problem:
 - changes in architecture lead to changes in compilers
 - * new concerns, re-engineering, etc.
 - compiler changes then prompt new architecture changes, eg. new languages and features
- some compiler motivations:
 1. correct output
 2. fast output
 3. fast translation (proportional to program size)
 4. support separate compilation
 5. diagnostics for errors
 6. works well with debugger
 7. cross language calls
 8. optimization
- for new languages, how are compilers written for them?
 - eg. early compilers for Java were written in C
 - eg. for early, low level languages like C, **bootstrapping** is done:
 - * a little subset of C is written and compiled in machine code
 - * then a larger subset of C is compiled using that smaller subset, etc.

Compiler Overview

- abstract compiler system overview:
 - *input*: source code
 - *output*: machine code or errors
 - recognizes illegal programs, and outputs associated errors
- *two-pass* compiler overview:
 - source code eg. Java compiles through a frontend to an **intermediate representation (IR)** like Sparrow
 - * the **frontend** part of the compiler maps legal code into IR:
 - language *dependent*, but machine *independent*
 - allows for swappable front ends for different source languages
 - IR then compiles through a backend to machine code
 - * the **backend** part maps IR onto target machine:
 - language *independent*, but machine / architecture *dependent*
- frontend overview:
 - *input*: source code
 - *output*: IR
 - *responsibilities*:
 - * recognize legality syntactically
 - * produce meaningful error messages
 - * shape the code for the backend
 - 1. the scanner produces a stream of tokens from source code:
 - ie. *lexing* source file into tokens
 - 2. the parser produces the IR:
 - recognizes context free grammars, while guiding context sensitive analysis
 - both steps can be automated to some degree
- backend overview:
 - *input*: IR
 - *output*: target machine code
 - *responsibilities*:
 - * translate to machine code
 - * **instruction selection**:
 - choose specific instructions for each IR operation
 - produce compact, fast code
 - * **register allocation**:
 - decide what to keep in registers at each points
 - can move loads and stores
 - optimal allocation is difficult
 - more difficult to automate
- specific frontends or backends can be swapped
 - eg. use special backend that targets ARM instead of RISC, etc.
- middleend overview:
 - *responsibilities*:

- * optimize code and perform code improvement by analyzing and changing IR
 - * must preserve values while reducing runtime
- optimizations are usually designed as a set of iterative passes through the compiler
- eg. eliminating redundant stores or dead code, storing common subexpressions, etc.
- eg. GCC has 150 optimizations built in

Lexical Analysis

- the role of the **scanner** is to map characters into **tokens**, the basic unit of syntax:
 - while eliminating whitespace, comments, etc.
 - the character string value for a token is a **lexeme**
 - eg. `x = x + y;` becomes `<id,x> = <id,x> + <id,y> ;`
 - a scanner must recognize language syntax
 - how to define what the syntax for integers, decimals, etc.
1. use regular expressions to specify syntax patterns:
 - eg. the syntax pattern for an integer may be `<integer> ::= (+ | -) <digit>*`
 2. regular expressions can then be constructed into a **deterministic finite automaton (DFA)**:
 - a series of states and transitions for accepting or rejecting characters
 - this step also handles state minimization
 3. the DFA can be easily converted into code using a while loop and states:
 - by using a table that categorizes characters into their language specific identifier types or classes, this code can be language *independent*
 - as long as the underlying DFA is the same
 - a linear operation, considers each character once
- this process can be automated using **scanner generators**:
 - emit scanner code that may be direct code, or table driven
 - eg. `lex` is a UNIX scanner generator that emits C code

Parsing

- the role of the **parser** is to recognize whether a stream of tokens forms a program defined by some grammar:
 - performs context-free syntax analysis
 - usually constructs an IR
 - produces meaningful error messages
 - generally want to achieve *linear* time when parsing:
 - * need to impose some restrictions to achieve this, eg. the LL restriction
- context-free syntax is defined by a **context-free grammar (CFG)**:
 - formally, a 4-tuple $G = (V_t, V_n, S, P)$ where:
 - * V_t is the set of **terminal** symbols, ie. tokens returned by the scanner
 - * V_n is the set of **nonterminal** symbols, ie. syntactic variables that denote substrings in the language
 - * S is a distinguished nonterminal representing the **start symbol** or goal
 - * P is a finite set of **productions** specifying how terminals and non-terminals can be combined
 - each production has a single nonterminal on the LHS
 - * the **vocabulary** of a grammar is $V = V_t \cup V_n$
 - * the motivation for using CFGs instead of simple REs for grammars is that REs are not powerful enough:
 - REs are used to classify tokens such as identifiers, numbers, keywords
 - while grammars are useful for counting brackets, or imparting structure eg. expressions
 - factoring out lexical analysis simplifies the CFG dramatically
 - general CFG notation:
 - * $a, b, c, \dots \in V_t$
 - * $A, B, C, \dots \in V_n$
 - * $U, V, W, \dots \in V$
 - * $\alpha, \beta, \gamma, \dots \in V^*$, where V^* is a sequence of symbols
 - * $u, v, w, \dots \in V_t^*$, where V_t^* is a sequence of terminals
 - * $A \rightarrow \gamma$ is a production
 - * $\Rightarrow, \Rightarrow^*, \Rightarrow^+$ represent derivations of 1, ≥ 0 , and ≥ 1 steps
 - * if $S \Rightarrow^* \beta$ then β is a **sentential form** of G
 - * if $L(G) = \{\beta \in V^* | S \Rightarrow^* \beta\} \cap V_t^*$, then $L(G)$ is a **sentence** of G , ie. a derivation with all nonterminals
- grammars are often written in **Backus-Naur form (BNF)**:
 - non-terminals are represented with angle brackets

- terminals are represented in monospace font or underlined
- productions follow the form `<nont> ::= ...expr...`
- the productions of a CFG can be viewed as rewriting rules:
 - by repeatedly rewriting rules by replacing nonterminals (starting from goal symbol), we can **derive** a sentence of a programming language
 - * **leftmost derivation** occurs when the *leftmost* nonterminal is replaced at each step
 - * **rightmost derivation** occurs when the *rightmost* nonterminal is replaced at each step
 - this sequence of rewrites is a **derivation** or **parse**
 - *discovering* a derivation (ie. going backwards) is called **parsing**
- can also visualize the derivation process as construction a tree:
 - the goal symbol is the root of tree
 - the children of a node represents replacing a nonterminal with the RHS of its production
 - note that the ordering of the tree dictates how the program would be *executed*
 - * can multiple syntaxes lead to different parse trees depending on the CFG used?
 - parsing can be done **top-down**, from the root of the derivation tree:
 - * picks a production to try and match input using backtracking
 - * some grammars are backtrack-free, ie. *predictive*
 - parsing can also be done **bottom-up**:
 - * start in a state valid for legal first tokens, ie. start at the leaves and fill in
 - * as input is consumed, change state to encode possibilities, ie. recognize valid prefixes
 - * use a stack to store state and sentential forms

Top-Down Parsing

- try and find a linear parsing algorithm using top-down parsing
- general top-down parsing approach:
 1. select a production corresponding to the current node, and construct the appropriate children
 - want to select the right production, somehow guided by input string
 2. when a terminal is added to the *fringe* that doesn't match the input string, backtrack
 3. find the next nonterminal to expand
- problems that will make the algorithm run worse than linear:

- too much backtracking
- if the parser makes the wrong choices, expansion doesn't even terminate
 - * ie. top-down parsers *cannot* handle left-recursion
- top-down parsers may backtrack when they select the wrong production:
 - do we need arbitrary **lookahead** to parse CFGs? Generally, yes.
 - however, large subclasses of CFGs *can* be parsed with *limited* lookahead:
 - * LL(1): left to right scan, left-most derivation, 1-token lookahead
 - * LR(1): left to right scan, right-most derivation, 1-token lookahead
- to achieve LL(1) we roughly want to have the following initial properties:
 - no left recursion
 - some sort of *predictive* parsing in order to minimize backtracking with a lookahead of only one symbol

Grammar Hacking

Consider the following simple grammar for mathematical operations:

```
<goal> ::= <expr>
<expr> ::= <expr> <op> <expr> | num | id
<op> ::= + | - | * | /
```

- there are multiple ways to rewrite the same grammar:
 - but each of these ways may build *different* trees, which lead to *different* executions
 - want to avoid possible grammar issues such as precedence, infinite recursion, etc. by rewriting the grammar
 - eg. classic precedence issue of parsing $x + y * z$ as $(x+y) * z$ vs. $x + (y*z)$
- to address **precedence**:
 - additional machinery is required in the grammar
 - introduce extra **levels**
 - eg. introduce new nonterminals that group higher precedence ops like multiplication, and ones that group lower precedence ops like addition
 - * the higher precedence nonterminal cannot reduce down to the lower precedence nonterminal
 - * forces the *correct* tree

Example of fixing precedence in our grammar:

```
<expr> ::= <expr> + <term> | <expr> - <term> | <term>
<term> ::= <term> * <factor> | <term> / <factor> | <factor>
<factor> ::= num | id
```

- **ambiguity** occurs when a grammar has more than one derivation for a single sequential form:
 - eg. the classic **dangling-else** ambiguity `if A then if B then C else D`
 - to address ambiguity:
 - * rearrange the grammar to select one of the derivations, eg. matching each `else` with the closest unmatched `then`
 - another possible ambiguity arises from the context-free specification:
 - * eg. **overloading** such as `f(17)`, could be a function or a variable subscript
 - * requires context to disambiguate, really an issue of type
 - * rather than complicate parsing, this should be handled separately

Example of fixing the dangling-else ambiguity:

```
<stmt> ::= <matched> | <unmatched>
<matched> ::= if <expr> then <matched> else <matched> | ...
<unmatched> ::= if <expr> then <stmt> | if <expr> then <matched> else <unmatched>
```

- a grammar is **left-recursive** if $\exists A \in V_n$ s.t. $A \Rightarrow^* A\alpha$ for some string α :
 - top-down parsers fail with left-recursive grammars
 - to address left-recursion:
 - * transform the grammar to become right-recursive by introducing new nonterminals
 - eg. in grammar notation, replace the productions $A \rightarrow A\alpha|\beta|\gamma$ with:
 - * $A \rightarrow NA'$
 - * $N \rightarrow \beta|\gamma$
 - * $A' \rightarrow \alpha A'|\epsilon$

Example of fixing left-recursion (for `<expr>`, `<term>`) in our grammar:

```
<expr> ::= <term> <expr'>
<expr'> ::= + <term> <expr'> | - <term> <expr'> | E // epsilon

<term> ::= <factor> <term'>
<term'> ::= * <factor> <term'> | / <factor> <term'> | E
```

- to perform **left-factoring** on a grammar, we want to do repeated prefix factoring until no two alternatives for a single non-terminal have a common prefix:
 - an important property for LL(1) grammars
 - eg. in grammar notation, replace the productions $A \rightarrow \alpha\beta|\alpha\gamma$ with:
 - * $A \rightarrow \alpha A'$
 - * $A' \rightarrow \beta|\gamma$
 - note that our example grammar after removing left-recursion is now properly left-factored

Achieving LL(1) Parsing

Predictive Parsing

- for multiple productions, we would like a *distinct* way of choosing the *correct* production to expand:
 - for some RHS $\alpha \in G$, define $FIRST(\alpha)$ as the set of tokens that can appear *first* in some string derived from α
 - key property: whenever two productions $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar, we would like:
 - * $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$, ie. the two token sets are disjoint
 - this property of **left-factoring** would allow the parser to make a correct choice with a lookahead of only *one* symbol
 - if the grammar does not have this property, we can hack the grammar
- by left factoring and eliminating left-recursion can we transform an *arbitrary* CFG to a form where it can be predictively parsed with a single token lookahead?
 - no, it is undecidable whether an arbitrary equivalent grammar exists that satisfies the conditions
 - eg. the grammar $\{a^n 0 b^n\} \cup \{a^n 1 b^{2n}\}$ does not have a satisfying form, since would have to look past an arbitrary number of a to discover the terminal
- idea to translate parsing logic to code:
 1. for all terminal symbols, call an `eat` function that *consumes* the next char in the input stream
 2. for all nonterminal symbols, call the corresponding function corresponding to the production of that nonterminal
 - perform predictive parsing by looking *ahead* to the next character and handling it accordingly
 - * there is only one valid way to handle the character in this step due to the left-factoring property
 - how do we handle epsilon?
 - * just do nothing ie. consume nothing, and let recursion handle the rest
 - creates a mutually recursive set of functions for each production
 - * the name is the LHS of production, and body corresponds to RHS of production

Example simple recursive descent parser:

```
Token token;
void eat(char a) {
```

```

    if (token == a) token = next_token();
    else error();
}

void goal() { token = next_token(); expr(); eat(EOF); }
void expr() { term(); expr_prime(); }
void expr_prime() {
    if (token == PLUS) { eat(PLUS); expr(); }
    else if (token == MINUS) { eat(MINUS); expr(); }
    else { /* noop for epsilon */ }
}
void term() { factor(); term_prime(); }
void term_prime() {
    if (token == MULT) { eat(MULT); term(); }
    else if (token == DIV) { eat(DIV); term(); }
    else { }
}
void factor() {
    if (token == NUM) eat(NUM);
    else if (token == ID) eat(ID);
    else error(); // not epsilon here
}

```

Handling Epsilon

- handling epsilon is not as simple as just ignoring it in the descent parser
- for a string of grammar symbols α , $NULLABLE(\alpha)$ means α can go to ε :
 - ie. $NULLABLE(\alpha) \iff \alpha \Rightarrow^* \varepsilon$
- to compute $NULLABLE$:
 1. if a symbol a is terminal, it cannot be nullable
 2. otherwise if $a \rightarrow Y_1 \dots Y_n$ is a production:
 - $NULLABLE(Y_1) \wedge \dots \wedge NULLABLE(Y_n) \Rightarrow NULLABLE(a)$
 3. solve the constraints
- again, for a string of grammar symbols α , $FIRST(\alpha)$ is the set of terminal symbols that begin strings derived from α :
 - ie. $FIRST(\alpha) = \{a \in V_t \mid \alpha \Rightarrow^* aB\}$
- to compute $FIRST$:
 1. if a symbol a is a nonterminal, $FIRST(a) = \{a\}$

2. otherwise if $a \rightarrow Y_1 \dots Y_n$ is a production:
 - $FIRST(Y_1) \subseteq FIRST(A)$
 - $\forall i \in 2 \dots n$, if $NULLABLE(Y_1 \dots Y_{i-1})$:
 - * $FIRST(Y_i) \subseteq FIRST(A)$
 3. solve the constraints, going for the \subseteq -least solution
- for a nonterminal B , $FOLLOW(B)$ is the set of terminals that can appear immediately to the *right* of B in some sentential form:
 - ie. $FOLLOW(B) = \{a \in V_t \mid G \Rightarrow^* \alpha B \beta \wedge a \in FIRST(\beta \$)\}$
 - to compute $FOLLOW$:
 1. $\{\$ \} \subseteq FOLLOW(G)$ where G is the goal
 2. if $A \rightarrow \alpha B \beta$ is a production:
 - $FIRST(\beta) \subseteq FOLLOW(B)$
 - if $NULLABLE(\beta)$, then $FOLLOW(A) \subseteq FOLLOW(B)$
 3. solve the constraints, going for the \subseteq -least solution

Formal Definition

- a grammar G is **LL(1)** iff. for each production $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$:
 1. $FIRST(\alpha_1), \dots, FIRST(\alpha_n)$ are pairwise *disjoint*
 2. if $NULLABLE(\alpha_i)$, then for all $j \in 1 \dots n \wedge j \neq i$:
 - $FIRST(\alpha_j) \cap FOLLOW(A) = \emptyset$
 - if G is ε -free, the first condition is sufficient
 - eg. $S \rightarrow aS | a$ is not LL(1)
 - * while $S \rightarrow aS', S' \rightarrow aS' | \varepsilon$ accepts the same language and is LL(1)
- provable facts about LL(1) grammars:
 1. no left-recursive grammar is LL(1)
 2. no ambiguous grammar is LL(1)
 3. some languages have no LL(1) grammar
 4. an ε -free grammar where each alternative expansion for A begins with a distinct terminal is a simple LL(1) grammar
- an LL(1) **parse table** M can be constructed from a grammar G as follows:
 1. \forall productions $A \rightarrow \alpha$:
 - $\forall a \in FIRST(\alpha)$, add $A \rightarrow \alpha$ to $M[A, a]$
 - if $\varepsilon \in FIRST(\alpha)$:
 - * $\forall b \in FOLLOW(A)$, add $A \rightarrow \alpha$ to $M[A, b]$ (including EOF)
 2. set each undefined entry of M to an error state
 - if $\exists M[A, a]$ with multiple entries, then the grammar is *not* LL(1)

JavaCC

- the **Java Compiler Compiler (JCC)** generates a parser automatically for a given grammar:
 - based on LL(k) vs. LL(1)
 - transforms an EBNF grammar into a parser
 - can have embedded (additional) action code written in Java
 - `javacc fortran.jj` → `javac Main.java` → `java Main < prog.f`

JavaCC input format:

```
TOKEN :
{
  < INTEGER_LITERAL: ( ["1"-"9"] (["0"-"9"])* | "0" ) >
}

void StatementListReturn() :
{}
{
  ( Statement() )* "return" Expression() ";"
}
```

Handling Syntax Trees

Visitor Pattern

- parsers generate a syntax tree from an input file:
 - this is an aside on design patterns in order to facilitate using the generated tree
 - see Gamma's *Design Patterns* from 1995
- for OOP, the **visitor pattern** enables the definition of a *new* operation of an object structure *without* changing the classes of the objects:
 - ie. new operation *without* recompiling
 - set of classes must be fixed in advance, and each class must have a hook called the `accept` method

Consider the problem of summing up lists using the following list implementation:

```
interface List {}

class Nil implements List {}

class Cons implements List {
  int head;
  List tail;
}
```

First approach using type casts:

```
List l;
int sum = 0;
while (true) {
  if (l instanceof Nil)
    break;
  else if (l instanceof Cons) {
    sum += ((Cons) l).head;
    l = ((Cons) l).tail;
  }
}
```

- *pros*:
 - code is written without touching the classes
- *cons*:

- code constantly uses type casts and `instanceof` to determine classes

Second approach using dedicated methods (OO version):

```
interface List { int sum(); }

class Nil implements List {
    public int sum() { return 0; }
}

class Cons implements List {
    int head;
    List tail;
    public int sum() { return head + tail.sum(); }
}
```

- *pros*:
 - code can be written more systematically, without casts
- *cons*:
 - for each new operation, need to write new dedicated methods and re-compile
- **visitor pattern** approach:
 - divide the code into an object structure and a **visitor** (akin to functional programming)
 - insert an `accept` method in each class, which takes a Visitor as an argument
 - a visitor contains a `visit` method for each class (using *overloading*)
 - * defines both actions and access of *subobjects*
 - *pros*:
 - * new methods without recompilation
 - * no frequent type casts
 - *cons*:
 - * all classes need a hook in the `accept` method
 - used by tools such as JJTree, Java Tree Builder, JCC
 - summary, visitors:
 - * make adding new operations easily
 - * gather *related* operations
 - * can accumulate state
 - * can break encapsulation, since it needs access to internal operations

Third approach with visitor pattern:

```
interface List {
    // door open to let in a visitor into class internals
    void accept(Visitor v);
}
```



```

}
interface Visitor {
    void visit(Nil x); // code is packaged into a visitor
    void visit(Cons x);
}

class Nil implements List {
    // 'this' is statically defined by the 'enclosing' class
    public void accept(Visitor v) { v.visit(this); }
}

class Cons implements List {
    int head;
    List tail;
    public void accept(Visitor v) { v.visit(this); }
}

class SumVisitor implements Visitor {
    int sum = 0;
    public void visit(Nil x) {}
    public void visit(Cons x) {
        // take an action:
        sum += x.head;
        // handle subobjects:
        x.tail.accept(this); // process tail 'indirectly' recursively

        // The accept call will in turn call visit...
        // This pattern is called 'double dispatching'.
        // Why not just visit(x.tail) ?
        // This 'fails', since x.tail is type List.
    }
}

```

Using `SumVisitor` :

```

SumVisitor sv = new SumVisitor();
l.accept(sv);
System.out.println(sv.sum);

```

Java Tree Builder

- the produced JavaCC grammar can be processed by the JCC to give a parser

that produces syntax trees:

- the produced syntax trees can be traversed by a Java program by writing subclasses of the default visitor
- JavaCC grammar feeds into the **Java Tree Builder (JTB)**
- JTB creates JavaCC grammar with embedded Java code, syntax-tree-node classes, and a default visitor
- the new JavaCC grammar feeds into the JCC, which creates a parser
- `jtb fortran.jj` → `javacc jtb.out.jj` → `javac Main.java` → `java Main < prog.f`

Translating a grammar production with JTB:

```
// .jj grammar
void assignment() :
{
{ PrimaryExpression() AssignmentOperator() Expression() }

// jtb.out.jj with embedded java code that builds syntax tree
Assignment Assignment () :
{
    PrimaryExpression n0;
    AssignmentOperator n1;
    Expression n2; {}
}
{
    n0 = PrimaryExpression()
    n1 = AssignmentOperator()
    n2 = Expression()
    { return new Assignment(n0, n1, n2); }
}
```

JTB creates this syntax-tree-node class representing `Assignment` :

```
public class Assignment implements Node {
    PrimaryExpression f0;
    AssignmentOperator f1;
    Expression f2;

    public Assignment(PrimaryExpression n0,
        AssignmentOperator n1, Expression n2) {
        f0 = n0; f1 = n1; f2 = n2;
    }

    public void accept(visitor.Visitor v) {
```

```
    v.visit(this)
  }
}
```

Default DFS visitor:

```
public class DepthFirstVisitor implements Visitor {
    ...
    // f0 → PrimaryExpression()
    // f1 → AssignmentExpression()
    // f2 ⇒ Expression()
    public void visit(Assignment n) {
        // no action taken on current node,
        // then recurse on subobjects
        n.f0.accept(this);
        n.f1.accept(this);
        n.f2.accept(this);
    }
}
```

Example visitor to print LHS of assignments:

```
public class PrinterVisitor extends DepthFirstVisitor {
    public void visit(Assignment n) {
        // printing identifier on LHS
        System.out.println(n.f0.f0.toString());
        // no need to recurse into subobjects since assignments cannot be nested
    }
}
```

Type Checking

- a program may follow a grammar and parse properly, but other problems may remain
 - eg. type errors, use of undeclared variables, cyclical inheritance, etc.

Simple Expressions

- eg. in Java, `5 + true` gives a type error:
 - this type rule could be expressed as:
 - * if two expressions have type `int`, then their addition will also be of type `int`
 - in typical notation:

$$\frac{a : \text{int} \quad b : \text{int}}{a+b : \text{int}}$$

- * in this notation, the **conclusion** appears under the bar, with multiple **hypotheses** above the bar
 - ie. if hypotheses are true, then conclusion is true
- * to *check* this rule, recursively check if e_1 is type `int`, and then if e_2 is also of type `int`
- when given `5: int` and `true: boolean`, type checker will see that the types don't obey the rule, and should throw an error

Implementing a simple type checker:

```
class TypeChecker extends Visitor {
    // f1: Expression (e_1)
    // f2: Expression (e_2)
    Type visit(Plus n) { // should either return Plus or throw an error
        // recursive calls for subexpressions that *could* individually fail
        Type t1 = n.f1.accept(this);
        Type t2 = n.f2.accept(this);

        // making a decision, ie. actual type checking
        if (t1 == "int" && t2 == "int") { // need some mechanic to check type equality
            return "int";
        } else {
            throw new RuntimeException();
        }
    }
}
```

- what about for more complex compound type rules, eg. for parsing `3 + (5 + 7)` :
 - simply by following the recursive calls, the previous type checker would still successfully check the type
 - * ie. type checking in a DFS manner

$$\frac{5 : \text{int} \quad 7 : \text{int}}{5+7 : \text{int}} \rightarrow \frac{3 : \text{int} \quad 5+7 : \text{int}}{3+(5+7) : \text{int}}$$

- handling the simple nonterminal `true` :

$$\frac{}{\text{true} : \text{boolean}}$$

- handling boolean negation:

$$\frac{e : \text{boolean}}{!e : \text{boolean}}$$

- handling ternary expressions:

$$\frac{a : \text{boolean} \quad b : t \quad c : t}{(a ? b : c) : t}$$

- t is a type variable since b and c should have the same type

Statements

- different types of statements in MiniJava:
 - `System.out.println(e)` , assignments, `if` and `while` statements
- unlike expressions, statements don't *return* anything:
 - they may have side effects, but do not have their *own* types
 - type checkers would only either return silently or throw an error
 - no overall type value to return
 - in typical notation \vdash means that a sentence type-checks
 - * not necessarily that it *is* a particular type
- handling `System.out.println` :

$$\frac{\vdash e : \text{int}}{\vdash \text{System.out.println}(e)}$$

- handling `if` statements:

$$\frac{\vdash e : \text{boolean} \quad \vdash a \quad \vdash b}{\vdash \text{if } (e) \ a \ \text{else } b}$$

- handling `while` statements:

$$\frac{\vdash e : \text{boolean} \quad \vdash s}{\vdash \text{while } (e) \ s}$$

Declarations

- for declared variables, how can we track what specific identifiers represent?
 - create and add to a **symbol table** that caches the declaration of variables with types
 - maps identifiers to types
- type checking rule for an assignment statement, given the symbol table A :

$$\frac{A(x) = t \quad A \vdash e : t}{A \vdash x = e}$$

- by convention, A should go before each \vdash in all type checking rules
 - * this is because any subexpressions may contain variables
- ie. $A \vdash e : t$ can be read as expression e *given* program context from table A can *lead* to conclusion t when type checking

Symbol table example:

```
class C {
  boolean f;
  t m(int a) {
    int x,
    ...
    x = a + 5;
    ...
  }
}
```

```
// symbol table contains:
// id | type
// -----
// f  | boolean
// a  | int
// x  | int
```

- to type check a variable, just look it up in the symbol table:

$$\frac{A(x) = t}{A \vdash x : t}$$

- rewriting our symbol table lookup using the previous notation
- to type check `x = a + 5` :

$$\frac{A \vdash a : \text{int} \quad A \vdash 5 : \text{int}}{A \vdash a+5 : \text{int}} \rightarrow \frac{A \vdash x : \text{int} \quad A \vdash a+5 : \text{int}}{A \vdash x=a+5}$$

Example of variable shadowing:

```
class a {
  int a;
  int a(int a) {
    // int a;    // compiler error, can't redeclare parameter
    ... a + 5 ... // checks closest `a`, ie. read from bottom of symbol table
    return 0;
  }
}
```

Arrays

- type checking arrays:
 - expressions like `arr[idx]`, `new int[len]`, `arr.length`
 - statements like `arr[idx] = val`
 - note that it is *unreasonable* for the type checker to check out of bound indices
 - * would require the type checker to do some arithmetic in this stage
- handling array indexing:

$$\frac{A \vdash a : \text{int}[] \quad A \vdash b : \text{int}}{A \vdash a[b] : \text{int}}$$

- in MiniJava, the only array type is `int` , but arrays can generally hold any type *t*
- note that the `int` in the hypotheses refers to the array type while the other refers to the index type
- note that this rule is an **elimination rule** since the `int[]` type is *consumed* into an `int` in the conclusion

- handling array assignments:

$$\frac{A \vdash x : \text{int[]} \quad A \vdash i : \text{int} \quad A \vdash v : \text{int}}{A \vdash x[i] = v : \text{int}}$$

- handling array constructions:

$$\frac{A \vdash e : \text{int}}{A \vdash \text{new int}[e] : \text{int[]}}$$

- handling the array length property:

$$\frac{A \vdash e : \text{int}[]}{A \vdash e.\text{length} : \text{int}}$$

Methods

- what needs to be checked in a method call in Java?
 - the method being called needs to *exist*
 - the type of the actual parameter should match the formal parameter

Declaring and calling methods in Java:

```
u2 m(u a) { // declaration
  s
  return e2;
}

m(e); // call
```

- when type checking the method body, need to ensure:

$$A \vdash s, A \vdash e_2 : u_2$$

- to type check the method call:

$$\frac{A \vdash e : t \quad A \vdash m : u \rightarrow u_2 \quad t = u}{A \vdash m(e) : u_2}$$

- note that we need a mechanism to find the types of the method, eg. `m` takes parameter of type `u` and returns type `u2`

Classes

- expressions such as `new C()`, `this`, `(C) e` refer to classes:
 - need a new type C , representing some class:
 - C is another contextual record similar to the symbol table A
 - $A, C \vdash e : t$
 - ie. expression e given program context from table A and class C can lead to conclusion t when type checking
 - `this` refers to the lexically enclosing class type
- handling `new C()` :
 - need to check that C exists
- handling casts:

$$\frac{A, C \vdash e : t}{A, C \vdash (D)e : D}$$

- again, may need a runtime check (like array indexing) to actually *perform* the type cast
- in Java, *upcasts* will always succeed, while *downcasts* need a runtime check
- handling subtyping:
 - the notation $C \leq D$ indicates the class `C` is a subclass of the class `D`, perhaps transitively or reflexively
 - eg. for primitives as well, `char ≤ short ≤ int ≤ long ≤ float ≤ double`
 - note that:

$$\frac{}{C \leq C'}, \frac{C \leq D \quad D \leq E}{C \leq E}$$

- when $u \leq t$:
 - for polymorphic assignments:

$$\frac{A, C \vdash x : t \quad A, C \vdash e : u}{A, C \vdash x = e}$$

- for method parameters:

$$\frac{A, C \vdash a : D \quad A, C \vdash e : u \quad D.m : t \rightarrow t_2}{A, C \vdash a.m(e) : t_2}$$

- note that to check subclassing relationships, we can either use visitors to traverse the `extends` relationship on the fly, or cache a previous result

Polymorphism in Java:

```
// class B and C both extend class A
A x = (e ? new B() : new C());
// compiler checks *both* B and C have correct subclassing relationship to A
// at runtime, the type of x will be related to either B or C
```

Extended method example:

```
class D {
  t1 f1
  t2 f2
  u3 m(u1 a1, u2 a2) {
    t1_1 x1
    t2_2 x2
    s
    return e;
  }
}
```

- hypotheses for type checking the above method:
 1. $a1, a2, x1, x2$ are all different
 2. $A, C \vdash s$, where $C = D$:
 - but what should the symbol table A hold?
 - * $f1:t1, f2:t2, a1:u1, a2:u2, x1:t1_1, x2: t2_2$
 - * when using A , we should search from *bottom* of table to handle variable shadowing
 3. $A, C \vdash e : u'_3$, where $u'_3 \leq u_3$
- conclusion from the above hypotheses:
 - $A, C \vdash D.m$ ie. the method $D.m$ type-checks

Entire Java Program

Sample Java program:

```
class Main {...}

class C1 extends D1 {...}
...
class Cn extends Dn {...}
```

- type checking responsibilities:
 1. `main` must exist
 2. `C1...Cn` need to all be different

3. `D1...Dn` need to all exist
4. no `extends` cycle in classes

Sparrow

- Sparrow is the intermediate language used in CS132
- characteristics:
 - no classes
 - * methods in classes have concatenated names, eg. `Fac.ComputeFac` becomes `FacComputeFac`
 - uses `goto` to implement `if else`
 - uses brackets to indicate heap loads or stores
 - * no global variables, only heap and functions have global visibility
 - functions may have extra parameters added

Grammar

- a program `p` is a series of functions, `p ::= F1...Fm`
- a function declaration `F` has syntax `F ::= func f(id1...idf) b`
- a block `b` is a series of instructions, `b ::= i1...in return id`
- an instruction can be:
 - a label `l:`
 - an assignment `id = c`, where `c` is an integer literal:
 - * or `id = @f`, where `f` is a function
 - an operation `id = id + id`, `id = id - id`, `id = id * id`
 - a less-than test `id = id < id`
 - a heap load or store:
 - * `id = [id + c]`, `[id + c] = id`, where `c` is the offset and heap addresses are valid
 - an allocation `id = alloc(id)`:
 - * eg. `v0 = alloc(12)` allocates 12 bytes
 - * creates a 3-tuple of addresses to values accessible by `[v0 + 0]`, `[v0 + 4]`, `[v0 + 8]`
 - a print `id = print(id)`
 - an error print `id = error(s)`
 - an unconditional goto `goto l`, where `l` is a label
 - a conditional goto `if0 id goto l`, jumps if `id` contains 0
 - a function call `id = call id(id1...idf)`, where `id` contains a function name
- identifiers can be any reasonable identifiers, except:
 - `a2...a7`, `s1...s11`, `t0...t5` which are all RISC register names

Translation to IR

- overall translation pipeline is:
 1. MiniJava
 - mostly *unbounded*, eg. number of variables, parameters, methods, classes, etc.
 2. Sparrow
 - still unbounded
 - may even generate *new* variables for simplicity / ease of translation
 3. Sparrow-V:
 - number of registers becomes bounded
 4. RISC-V:
 - still bounded register count

State and Transitions

- program *state* consists of the tuple (p, H, b^*, E, b) :
 - p is the program
 - H is the *heap* that maps from heap addresses to *tuples* of values
 - * the tuple can be *indexed* into
 - b^* is the body of the function that is executing right now
 - * can only perform a `goto` within this function block
 - E is the environment that maps from identifiers to values
 - b is the remaining part of the block that is executing right now
 - * ie. b^* contains the entire function block, while b only contains the current and remaining statements in the block
- in a state transition, we want to *step* from one state to the next:

$$(p, H, b^*, E, b) \rightarrow (p, H', b^{*'}, E', b')$$

- assignment state transition:

$$(p, H, b^*, E, \text{id=c} \ b) \rightarrow (p, H, b^*, E \cdot [\text{id} \mapsto c], b)$$

- arithmetic state transition:

$$(p, H, b^*, E, \text{id=id1+id2} \ b) \rightarrow (p, H, b^*, E \cdot [\text{id} \mapsto (c_1 + c_2)], b)$$

- where $E(\text{id}_1) = c_1$ and $E(\text{id}_2) = c_2$

- ie. this transition requires a runtime check in the environment
- heap load state transition:

$$(p, H, b^*, E, \text{id}=[\text{id1}+\text{c}] \quad b) \rightarrow (p, H, b^*, E \cdot [\text{id} \mapsto (H(a_1))(c_1 + c)], b)$$
 - where $E(\text{id}_1) = (a_1, c_1)$ such that a_1 is a heap address and c_1 its offset, and $(c_1 + c) \in \text{domain}(H(a_1))$
 - ie. the new computed offset is a valid index into the tuple on the heap
- heap allocation state transition:

$$(p, H, b^*, E, \text{id}=\text{alloc}(\text{id1}) \quad b) \rightarrow (p, H \cdot [a \mapsto t], b^*, E, b)$$

- where $E(\text{id}_1) = c$ and c is divisible by 4, a is a fresh address, and $t = [0 \mapsto 0, 4 \mapsto 0, \dots (c - 4) \mapsto 0]$
- unconditional `goto` state transition:

$$(p, H, b^*, E, \text{goto } l \quad b) \rightarrow (p, H, b^*, E, b')$$

- where $\text{find}(b^*, l) = b'$
- $\text{find}(b, l)$ is used to find the label l inside the block b
- conditional `goto` state transition:

$$(p, H, b^*, E, \text{if } 0 \text{ id goto } l \quad b) \rightarrow (p, H, b^*, E, b')$$

- where $E(\text{id}) = 0$ and $\text{find}(B^*, l) = b'$
- function call state transition:

$$(p, H, b^*, E, \text{id} = \text{call id0}(\text{id1} \dots \text{idf}) \quad b) \rightarrow (p, H', b^*, E \cdot [\text{id} \mapsto E'[\text{id}']], b)$$

- where $E(\text{id}_0) = f$, p contains `func f (id1'...idf') b'`, and $E' = [\text{id}'_1 \mapsto E(\text{id}_1), \text{id}'_2 \mapsto E(\text{id}_2), \dots]$
- the function f is then *called* through the following state transition:
 - * $(p, H, b', E', b') \rightarrow (p, H', b', E', \text{return } \text{id}')$
- note the intermediate transfer of control to the callee

Expressions

-
- want to translate $e, k \rightarrow c, k'$ where e is some expression in MiniJava and c is the output code in Sparrow, while managing:
 1. additional variables (that could have been added during translation)

2. labels used in jumps in IR

- k is a “fresh” or new integer number that has not yet been used
 - * can be utilized to generate variable and label names
- after translation, the number k' should be the next *new* number
- by convention, the *result* of the expression e is stored in the variable t_k

• simple expression:

$$5, k \rightarrow tk = 5, k + 1$$

• expression with a local variable:

$$x, k \rightarrow tk = x, k + 1$$

• recursive translations:

$$e_1 + e_2, k \rightarrow c_1 c_2, k_2$$

- given that $e_1, k + 1 \rightarrow c_1, k_1$ and $e_2, k_1 \rightarrow c_2, k_2$
- note that by convention, c_2 will be stored in t_{k_1} and c_1 will be stored in t_{k+1}
- then, $t_k = t_{k+1} + t_{k_1}$

• example addition translation:

$$7+9, 3 \rightarrow t4=7 \ t5=9 \ t3=t4+t5, 7$$

- after $7, 4 \rightarrow t4=7, 5$ and $9, 5 \rightarrow t5=9, 6$

• example nested addition translation:

$$(7+9)+11, 3 \rightarrow t5=7 \ t6=9 \ t4=t5+t6 \ t7=11 \ t3=t4+t7, 8$$

- $7, 5 \rightarrow t5=7, 6$ and $9, 6 \rightarrow t6=9, 7$
- $7+9, 4 \rightarrow t5=7 \ t6=9 \ t4=t5+t6, 7$
- $11, 7 \rightarrow t7=11, 8$
- variables $t4...t6$ handle $7+9$, $t7$ handles 11 , and $t3$ holds the entire expression

Statements

-
- want to translate $s, k \rightarrow c, k'$ where s is some statement in MiniJava and c is the output code in Sparrow
 - note that k' may differ from k when a subexpression is contained within the statement

- simple recursive statements:

$$s_1 s_2, k \rightarrow c_1 c_2, k_2$$

- where $s_1, k \rightarrow c_1, k_1$ and $s_2, k_1 \rightarrow c_2, k_2$
- no return value since statements only have a side effect, so only ordering of output code matters

- simple assignment statement:

$$x=e, k \rightarrow c \quad x=tk, k'$$

- where $e, k \rightarrow c, k'$ and the result of c is held in variable t_k by convention

- if else statement:

$$\text{if}(e) \ s1 \ \text{else} \ s2, k \rightarrow c_k, k_2$$

- $e, k + 1 \rightarrow c_e, k_e$
- $s_1, k_e \rightarrow c_1, k_1$
- $s_2, k_1 \rightarrow c_2, k_2$
- to generate control code, need to *generate* labels and use jumps:
 - * can use k in order to avoid label duplication
 - * also need to *linearize* the code while allowing s_1 to execute without s_2 , and vice versa
 - * uses an unconditional jump and a conditional jump

c_k , the generated code from translating the if else statement:

```
c_e // result stored in temporary t_{k+1}
if0 t_{k+1} goto else_k
c1
goto end_k
else_k:
  c2
end_k:
```

- while statement:

$$\text{while}(e) \ s, k \rightarrow c_k, k_s$$

- $e, k + 1 \rightarrow c_e, k_e$
- $s, k_e \rightarrow c_s, k_s$
- again, need to generate control code
- note that the compiler does not care about how many times the loop runs
 - * thus the label subscript k simply depends on static code structure
- conventionally, for loops are usually reduced into while loops

c_k , the generated code from translating the `while` statement:

```
loop_k:
c_e // result stored in temporary t_{k+1}
if0 t_{k+1} goto end_k
c_s
goto loop_k
end_k:
```

Examples

To translate $s, k \rightarrow c, k1$ for a simple sequence s of two substatements using visitors:

```
class Result {
    String code;
    int k1;
}

Result visit(Seq n, int k) {
    Result res1 = n.f1.accept(this, k);
    Result res2 = n.f2.accept(this, res1.k);
    return new Result(res1.code + res2.code, res2.k);
}
```

Translating the following code starting from $k = 0$:

```
while (true) {
    if (false)
        x = 5;
    else
        y = 7;
}
```

Translation process:

```
// handling `while (e) s`:
true, 1 → t1 = 1, 2 // e

// handling s, which has the form `if (e1) s1 else s2`:
false, 3 → t3 = 0, 4 // e1

5, 4 → t4 = 5, 5 // s1
x=5, 4 → x = t4, 5
```

```

7, 5  → t5 = 7, 6 // s2
y=7, 5 → x = t5, 6

if (false) x=5 else y=7, 2 → { // s
  t3 = 0
  if0 t3 goto else2
  t4 = 5
  x = t4
  goto end2
else2:
  t5 = 7
  y = t5
end2:
}, 6

while (true) {if (false) x=5 else y=7}, 0 → {
  loop0:
    t1 = 1
    if0 t1 goto end0
    // code generated from s, the if-else statement
    goto loop0
  end0:
}, 6

```

Arrays

- need to generate Sparrow IR for the following MiniJava code involving arrays:
 - expressions like `arr[idx]`, `new int[len]`, `arr.length`
 - statements like `arr[idx] = val`
- in Sparrow, arrays will be represented on the heap:
 - can create an array by allocating space on the heap
 - need to *track* and store the array length somewhere:
 - * by *convention*, simply store the length of the array at position 0 on the heap, and shift array positions in the heap accordingly
 - * thus in Sparrow, first array element is stored at offset 4, and last element is at offset $4n$ where n is array length
 - * total amount to allocate is $4 \times (n + 1)$ to store length
- array length:

`e.length, k` → `c tk=[t1+0], k'`

- after evaluating `e, k + 1` → `c, k'` and where $t_l = t_{k+1}$

- array allocation:

`new int[e], k → c tk'=4*(t1) tk=alloc(tk') [tk+0]=t1, k' + 1`

- after evaluating $e, k + 1 \rightarrow c, k'$ and where $t_l = t_{k+1}$
- t'_k is a new temporary to store $4 * (e + 1)$ while saving the original array length to store back in array
- in the Sparrow spec, `alloc` will initialize all heap entries to zero
- array indexing:

`e1[e2], k → c1c2 tk2=4*(tk1+1) t1=t1+tk2 tk=[t1+0], k2 + 1`

- after evaluating $e_1, k + 1 \rightarrow c_1, k_1$ and $e_2, k_1 \rightarrow c_2, k_2$, and where $t_l = t_{k+1}$
- note that the second part of a heap lookup must be a constant so we build up the entire dynamic offset in the first parameter
- to check bounds, need to ensure that $0 \leq e_2 < n$ where n is the length of the array in e_1
- similar process for array assignments in the form `a[e1] = e2`

Implementing a bounds check:

```
t_{k2+1} = -1 < tk1
t_{k2+2} = [t_{k+1} + 0]
t_{k2+3} = tk1 < t_{k2+2}
t_{k2+4} = t_{k2+1} * t_{k2+3} // multiplication acts as logical AND
// including a bounds check uses up to k2+5 for temporaries

if0 t_{k2+4} goto error
... // loading or storing code here
goto end
error:
    error("out of bounds")
end:
```

Classes

MiniJava method example:

```
class C {
    Y m(T a) {
        S;
        return e;
    }
}
```

```

}
}

```

In Sparrow, there are no more classes:

```

func C.m(this a) { // name mangling
  c
  return x
}

```

- how to implement objects in Sparrow?
 - like arrays, represent objects on the heap
 - object fields will be given by offsets into the heap
 - * eg. to access the field $f, k \rightarrow tk=[this+L]$ where L is the static offset remembered by the compiler associated with the field f
 - but how can we handle `this` ie. the current object?
 - * need to translate the method call `e1.m(e2)` into `C.m(t1 t2)`
 - * where `t1, t2` hold the expression results of `e1, e2` respectively
 - * note that we can *identity* the full mangled name for `e1.m` by finding the *type* `C` of `e1`

Inheritance

- when we allow for inheritance and `extends`, we no longer know exactly which methods we are calling
 - similarly, object fields can also be inherited and accessed in subclasses

MiniJava inheritance example:

```

class B {
  t p() { this.m(); }
  u m() { ... }
}

class C extends B {
  u m() { ... } // overrides B.m
}

```

In Sparrow:

```

B.p(this) {
  B.m(...) // this is now *wrong*, could be called with C instead of B
}
B.m(this) { ... }
C.m(this) { ... }

```

- to handle inheritance, need to add an additional level of *indirection*:
 - add a **method table** for all objects:
 - * by convention, method table stored at position 0, similarly to where length is stored in arrays
 - * shift locations of object fields by 4
 - method table contains entries pertaining to all inherited visible methods
 - * each entry holds the address of a function
 - for the above example:
 - * a **B** object's method table should contains references to **B.p**, **B.m**
 - * while a **C** object's method table should contains references to **B.p**, **C.m**
 - * note that the offsets across related objects should line up
 - method tables are *shared* across objects of the same class, so each class can have a single method table allocated and all objects will point to that single copy
 - to call a class function in OOP:
 1. *load* method table
 2. *load* function name
 - * additional indirection to handle inheritance
 3. *call*
- consider building method tables for the following inheritance relationship:
 - class **A** has methods **m**, **n**
 - * method table holds refs. to **A.m**, **A.n**
 - class **B** extends **A** has no methods
 - * method table holds refs. to **A.m**, **A.n**
 - class **C** extends **B** with methods **p**, **m**
 - * method table holds refs. to **C.m**, **A.n**, **C.p** , regardless of method definition order
 - class **F** extends **C** with methods **q**
 - * method table holds refs. to **C.m**, **A.n**, **C.p**, **F.q**
 - thus compiler needs to associate the method **n** with offset 4 into method tables

Register Allocation

- in Sparrow, all variables are on the stack, which is allows for relatively *slow* access:
 - in the Sparrow-V IR, *some* variables are held in registers for *faster* access
 - want to fit as many variables as possible into the registers in a process called **register allocation**
- register allocation can be broken down into two subproblems:
 1. liveness analysis
 2. graph coloring (such that no adjacent nodes are the same color)
 - interface between the two steps is an **interference graph**
 - an approach using heuristics
- incremental steps for saving stack variables:
 1. no registers used
 2. use registers for storing all subcomputations
 3. use registers to pass function parameters
 - have to save all registers before function calls in case they are *clobbered*
 4. avoid saving registers that are not utilized after function calls
 5. using liveness analysis to reuse registers

Liveness Analysis

Sparrow to Sparrow-V example:

```
// Sparrow
a = 1      // prog. pt. 1 (point lies just before instruction)
b = 2      // prog. pt. 2
c = a + 3  // prog. pt. 3
print b + c // prog. pt. 4

// Sparrow-V with two registers (s0, s1)
s0 = 1
s1 = 2     // cannot assign to s0, would clobber `a`
c = s0 + 3 // leave `c` on the stack
// but `a` is no longer used, can reuse a register
print s1 + c

// Sparrow-V attempt 2
s0 = 1
```

```
s1 = 2
s0 = s0 + 3 // load before store
print s1 + s0
```

- **program points** lie between lines of instructions:
 - ie. where labels can be introduced
 - first program point lies before the first instruction
- in liveness analysis, want to examine:
 - *when* is a variable *live*, or what is its **live range**?
 - * **a** is live in the range $[1, 3]$, ie. from 1 *up to* 3
 - * **b** is live in $[2, 4]$
 - * **c** is live in $[3, 4]$
 - note that the live ranges for **a, b** and **b, c** overlap, while the live ranges for **a, c** do *not*
 - * ie. overlapping live ranges *conflict*
- in an **interference graph**:
 - nodes represents variables
 - edges represent conflicts
 - allocating registers is similar to coloring this graph

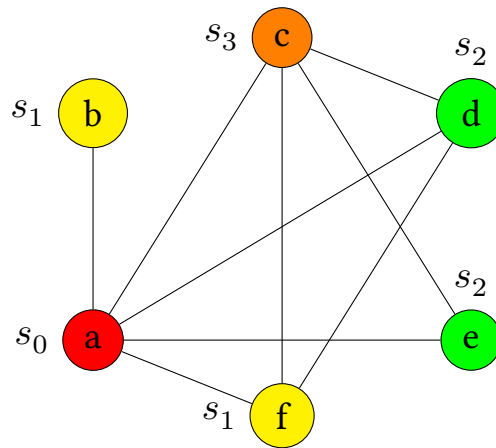
What is the fewest number of registers to entirely allocate the following code?

```
a = 1      // p.1
b = 10     // p.2
c = 8 + b  // p.3
d = a + c  // p.4
e = a + d  // p.5
f = c + e  // p.6
d = a + f  // p.7
c = d + f  // p.8
f = a + c  // p.9
return c + f // p.10
```

Live ranges, from inspection:

```
a : [1, 9]
b : [2, 3]
c : [3, 6] *and* [8, 10] (c can be reused before it is declared again)
d : [4, 5] and [7, 8]
e : [5, 6]
f : [6, 8] and [9, 10]
```

Corresponding interference graph:



Sparrow-V output code:

```

s0 = 1
s1 = 10
s3 = 8 + s1
s2 = s0 + s3
s1 = s0 + s2
s1 = s3 + s1
s2 = s0 + s1
s3 = s2 + s1
s1 = s0 + s3
return s3 + s1

```

Algorithm

- given some statement n :
 - $def[n]$ includes the variables defined or assigned in n
 - $use[n]$ includes the variables used in n
 - $in[n]$ includes the variables live coming *in* to n
 - $out[n]$ includes the variables live coming *out* of n
- in a **control-flow graph**:
 - nodes are statements
 - directed edges are possible control flow directions
 - a node's successor(s) are the nodes reachable from it
- defining **liveness equations**:

$$out[n] = \bigcup_{s \in succ(n)} in[s]$$

$$in[n] = use[n] \cup (out[n] - def[n])$$

Example code:


```

a = 0          // 1
L1: b = a + 1   // 2
c = c + b      // 3
a = b * 2      // 4
if a < 10 goto L1 // 5
return c       // 6

// in control-flow graph:
// 1→2, 2→3, 3→4, 4→5, 5→6, 5→2

```

Building liveness equations for example code:

n	def	use	in_1	out_1	in_2	out_2	in_3	out_3
1	a			a		a, c	c	a, c
2	b	a	a	b, c	a, c	b, c	a, c	b, c
3	c	b, c	b, c	b	b, c	b	b, c	b, c
4	a	b	b	a	b	a, c	b, c	a, c
5		a	a	a, c	a, c	a, c	a, c	a, c
6		c	c		c		c	

- algorithm steps:
 1. initially, in and out are both \emptyset
 2. use liveness equations as update steps, iteratively
 - in the table, in_k represents in at iteration k
 3. keep iterating until *no* change
 - give the final in and out values, we can determine interfering variables:
 - * every variable in the same box interferes with each other, pairwise
 - given n statements and $O(n)$ variables, there are $O(n^2)$ iterations:
 - * ie. at each stage, at least one element will be added to out
 - * at each iteration, n set-unions are performed each with a runtime of $O(n)$
 - thus, the algorithm has $O(n^4)$ runtime

Graph Coloring

- **graph coloring** ie. **liveness allocation** is the problem of allocating colors / registers given live ranges
- liveness analysis in *linear* time:
 - instead of the full $O(n^4)$ liveness analysis algorithm, is there a linear time approximation?

- instead of dividing up live ranges, simply take the entire interval from a variable's first declaration to its last use
 - * intervals can be retrieved in linear time
 - * ignores gaps in live ranges, and interpolates them
 - * seems to be a reasonable approximation that still indicates interfering variables
- **linear scan liveness allocation:**
 - now that we have a linear liveness analysis algorithm, want achieve liveness *allocation* in linear time as well:
 - * will not be as thorough of an allocation as running the full $O(n^4)$ liveness analysis and a complete graph coloring in exponential time
 - * but should be a good estimate, using heuristics
- **linear scan algorithm:**
 1. sort live range intervals by starting point
 2. perform a *greedy* coloring in a left-to-right scan:
 - tentatively assign different color ranges to colors ie. registers while the ranges overlap
 - once a previous interval *expires* before the start time of the current interval, *finalize* its coloring
 - when we have run out of registers, we have to **spill** a variable onto the stack:
 - * ie. store it on the stack instead of in registers
 - * as a heuristic, spill the one that extends the *furthest* into the future (based on end time)
 - * note that we can *only* take over ie. spill tentatively assigned registers
 - thus, this algorithm has $O(n)$ time
 - * there are only a constant number of registers to track of tentative assignments for

Liveness analysis and allocation example using linear algorithms:

```

a = 1      // 1
b = 10     // 2
c = 9 + a  // 3
d = a + c  // 4
e = c + d  // 5
f = b + 8  // 6
c = f + e  // 7
f = e + c  // 8
b = c + 5  // 9
  
```

```
return b + f // 10
```

Live range intervals (using linear approximation):

```
a : [1, 4]
b : [2, 9]
c : [3, 8]
d : [4, 5]
e : [5, 8]
f : [6, 10]
// note here that the intervals in this problem are already in sorted order
```

- given three registers `r1`, `r2`, `r3` to color with:
 1. assign `r1` to `a`
 2. assign `r2` to `b`
 3. assign `r3` to `c`
 4. `a` has expired, finalize `r1` for `a`
 - assign `r1` to `d`
 5. `d` immediately expires, finalize `r1` for `d` as well
 - assign `r1` to `e`
 6. no more registers, have to spill:
 - spill `f` onto the stack
 - coloring is complete

Allocation results using two registers:

```
a : r1
b : r2
c : r3
d : r1
e : r1
f : <mem>
// r1 can be reused multiple times without inteferences
```

- given two registers `r1`, `r2` to color with:
 1. assign `r1` to `a`
 2. assign `r2` to `b`
 3. no more registers, have to spill:
 - spill `b` onto the stack
 - assign `r2` to `c`
 4. `a` has expired, finalize `r1` for `a`
 - assign `r1` to `d`
 5. `d` immediately expires, finalize `r1` for `d` as well
 - assign `r1` to `e`
 6. no more registers, have to spill:

- spill `f` onto the stack
- coloring is complete

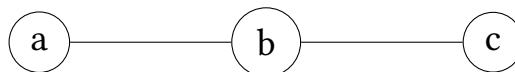
Allocation results using two registers:

```
a : r1
b : <mem>
c : r2
d : r1
e : r1
f : <mem>
// r1 can be reused multiple times without interferences
```

NP-Completeness

- in a full interference graph, generated from running the full liveness analysis algorithm:
 - $k \geq s$, where k is the minimum number of colors and s is the size of the max clique
 - *however*, if all the live ranges are presented as intervals, $k = s$
 - in general, graph coloring is NP-complete
 - * by using intervals, graph coloring becomes a polynomial time operation
- liveness analysis transforms the register allocation problem into a graph coloring problem:
 - do we *lose* anything in this transformation?
 - * ie. is the register allocation problem itself also an NP-complete problem?
 - * can transform in the other direction, from a coloring to an allocation problem
 - * eg. transform a certain graph into a problem, and run it through liveness allocation in order to color it

Graph coloring example:



From graph coloring to a representative program:

```
a = 1
b = 2
c = 3
```

```
// consider the graph as an interference graph
print(a + b) // `a` and `b` overlap, ie. both are live
print(b + c) // `b` and `c` overlap
// meanwhile, `a` and `c` are never live at the same time
```

- because this problem is transformable, register allocation is indeed an NP-complete problem:
 - thus, if the compiler should perform anything *smart* in its execution, it should do register allocation
 - register allocation by hand is unfeasible considering the very good linear approximation the compiler can quickly perform

Appendix

Practice Questions

1. given the following grammar:

- $A ::= \varepsilon | zCw$
- $B ::= Ayx$
- $C ::= ywz | \varepsilon | BAx$
- then:
 - $FIRST(A) = \{z\}$
 - $FIRST(B) = \{y, z\}$
 - $FIRST(C) = \{y, z\}$
 - $NULLABLE(A) = true$
 - $NULLABLE(B) = false$
 - $NULLABLE(C) = true$
- we can make the following observations for each nonterminal on the RHS:
 - $w \in FOLLOW(C)$
 - $y \in FOLLOW(A)$
 - $FIRST(A) \subseteq FOLLOW(B)$
 - $x \in FOLLOW(B)$
 - $x \in FOLLOW(A)$
- thus:
 - $FOLLOW(A) = \{x, y\}$
 - $FOLLOW(B) = \{x, z\}$
 - $FOLLOW(C) = \{w\}$
- therefore the grammar is *not* LL(1), since for C :
 - $FIRST(ywz) \cap FIRST(BAx) \neq \emptyset$

2. What is the fewest number of registers to entirely allocate the following code?

```

a = 1      // 1
b = 10     // 2
c = a + a  // 3
d = a + c  // 4
e = c + d  // 5
f = b + 8  // 6
c = f + e  // 7
f = e + c  // 8
b = c + 5  // 9
return b + f // 10

```

Live ranges, from inspection:

a : [1, 4]
b : [2, 6], [9, 10]
c : [2, 6], [7, 9]
d : [4, 5]
e : [5, 8]
f : [6, 7], [8, 10]

Corresponding interference graph:

