

# COMP50006 Compilers

## Imperial College London

Boxuan Tang

Spring 2023

### Contents

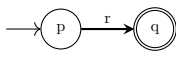
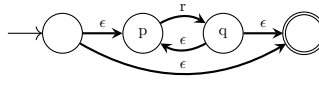
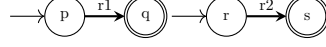

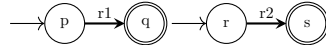
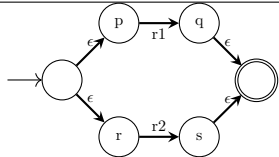
<b>1</b>	<b>Lexer</b>	<b>2</b>
1.1	Regex $\rightarrow$ NFA (Thompson's Construction)	2
1.2	NFA $\rightarrow$ DFA (Subset Construction)	2
<b>2</b>	<b>LR Bottom-Up/Shift-Reduce Parser</b>	<b>2</b>
2.1	LR Model	2
2.2	LR(0)	2
2.2.1	LR(0) $\rightarrow$ NFA	2
2.2.2	DFA $\rightarrow$ LR(0) Parsing Table	3
2.3	LR(1)	3
2.3.1	FIRST and FOLLOW	3
2.3.2	LR(1) $\rightarrow$ NFA	3
2.3.3	NFA $\rightarrow$ LR(1) Parsing Table	3
2.4	LALR(1)	3
2.5	Conflicts	3
2.5.1	Shift-Reduce Conflict	3
2.5.2	Reduce-Reduce Conflict	3
<b>3</b>	<b>LL Top-Down Parsing</b>	<b>3</b>
3.1	LL(1)	3
3.1.1	CFG $\rightarrow$ LL(1)	4
3.1.2	LL(1) to Parse Function	4
<b>4</b>	<b>Garbage Collection</b>	<b>4</b>
<b>5</b>	<b>Register Allocation</b>	<b>4</b>
5.1	Register Machine Strategy	4
5.2	Accumulator Strategy	4
5.3	Sethi-Ullman Weights	4
5.4	Graph Colouring	4
<b>6</b>	<b>Optimisation</b>	<b>5</b>
6.1	Live Variables	5
6.2	Reaching Definitions	5
6.3	Loops	5

## 1 Lexer

**Definition 1.1 (Context Free Grammar)** Contains a non-terminal start symbol, a set of productions, a set of terminals (tokens) and a set of non-terminals ( $S, P, t, nt$ ).

- A **sentence** is a derived string comprising only terminals.
- The **language** is the set of all sentences derived from the start symbol.
- It is **ambiguous** if the language contains strings that can be generated from two different ways

### 1.1 Regex $\rightarrow$ NFA (Thompson's Construction)

Regex	Given	Replace
(1) $r^*$		
(2) $r1 r2$		
(3) $r1 \mid r2$		

### 1.2 NFA $\rightarrow$ DFA (Subset Construction)

DFA start state =  $\epsilon$ -Closure(NFA start state)

foreach new subset state S of the DFA:

foreach unique symbol r leading out from any state of S:

add a transition r from S to S' where  $S' = \epsilon$ -Closure(states reached by r in 1 step)

Mark subset states accepting if any member state accepting in NFA

## 2 LR Bottom-Up/Shift-Reduce Parser

For LR, we add on an auxiliary rule with end-of-input symbol \$. For example,  $E' \rightarrow E \$$

### 2.1 LR Model

Parsing table contains all terminals under ACTIONS and non-terminals under GOTO

Push state 0 (start state) onto stack and then repeatedly perform

- **shift sN**: push state n onto stack, advance current token
- **goto gN**: not selected directly (reduce)
- **accept a**: accept input
- **reduce rN**: remove L elements from stack where  $L = \text{length rhs of rule N}$   
push Table[stack.top(), LHS of rule N]

### 2.2 LR(0)

**Definition 2.1 (LR(0) items)** are instances of the grammar rules with a  $\bullet$  on the rhs of the rule

#### 2.2.1 LR(0) $\rightarrow$ NFA

Given a state with item  $X \rightarrow A \bullet BC$ , add  $X \rightarrow A \bullet BC \xrightarrow{B} X \rightarrow AB \bullet C$

And if B is non-terminal, foreach initial item  $B \rightarrow \bullet D$ , add  $X \rightarrow A \bullet BC \xrightarrow{\epsilon} B \rightarrow \bullet D$

### 2.2.2 DFA $\rightarrow$ LR(0) Parsing Table

- state X with terminal transition  $X \xrightarrow{t} Y$ , add  $P[X,t] = sY$
- state X with non-terminal transition  $X \xrightarrow{N} Y$ , add  $P[X,N] = gY$
- state containing item  $R' \rightarrow \dots\bullet$ , add  $P[X,\$] = a$
- state containing item  $R \rightarrow \dots\bullet$ , add  $P[X,t] = rN$  for all terminals  $t$  where  $N$  is  $R$ 's rule number

## 2.3 LR(1)

**Definition 2.2 (LR(1) items)** is a pair  $[LR(0) \text{ item}, \text{look-ahead token } t]$

### 2.3.1 FIRST and FOLLOW

**Definition 2.3 (FIRST)**  $FIRST(\alpha)$  is the set of all terminals that could start the derivation of  $\alpha$

$FIRST(\epsilon) = \{\epsilon\}$ ,  $FIRST(a) = \{a\}$

$FIRST(A) = \text{foreach } A \rightarrow \beta_1\beta_2\dots\beta_n$

include  $FIRST(\beta_1) - \{\epsilon\}$  in  $FIRST(A)$

if  $\epsilon$  in  $FIRST(\beta_1)$ : include  $FIRST(\beta_2) - \{\epsilon\}$  in  $FIRST(A)$  ...

... if  $\epsilon$  in  $FIRST(\beta_n)$ : include  $\epsilon$  in  $FIRST(A)$

**Definition 2.4 (FOLLOW)**  $FOLLOW(A)$  is the set of all terminals that follows non-terminal  $A$

foreach  $B \rightarrow CAD$

include  $FIRST(D) - \{\epsilon\}$  in  $FOLLOW(A)$

if  $\epsilon$  in  $FIRST(D)$ : include  $FOLLOW(B)$  in  $FOLLOW(A)$

if  $A$  ends the input include  $\$$  in  $FOLLOW(A)$

### 2.3.2 LR(1) $\rightarrow$ NFA

Add initial item  $[R' \rightarrow \bullet R, \$]$

Given a state with item  $[X \rightarrow A \bullet BC, t]$ , add  $[X \rightarrow A \bullet BC, t] \xrightarrow{B} [X \rightarrow AB \bullet C, t]$

And if  $B$  is non-terminal, foreach rule  $B \rightarrow \bullet D$ , foreach token  $u$  in  $FIRST(Ct)$ ,

add  $[X \rightarrow A \bullet BC, t] \xrightarrow{\epsilon} [B \rightarrow \bullet D, u]$

### 2.3.3 NFA $\rightarrow$ LR(1) Parsing Table

Like LR(0) but for states containing  $[X \rightarrow A \bullet, t]$ , only add reduction for column  $t$

## 2.4 LALR(1)

If any 2 LR(1) states have the same LR(0) items, combines the states

## 2.5 Conflicts

### 2.5.1 Shift-Reduce Conflict

$S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{other becomes}$

$S \rightarrow MS \mid UMS$

$MS \rightarrow \text{if } E \text{ then } MS \text{ else } MS \mid \text{other}$

$UMS \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } UMS \text{ else } UMS$

### 2.5.2 Reduce-Reduce Conflict

$\text{Expr} \rightarrow \text{Expr} + \text{Expr} \mid \text{Expr} * \text{Expr} \mid (\text{Expr}) \mid \text{int}$   
becomes

$\text{Expr} \rightarrow \text{Expr} + \text{Term} \mid \text{Term}$

$\text{Term} \rightarrow \text{Term} * \text{Factor} \mid \text{Factor}$

$\text{Factor} \rightarrow (\text{Expr}) \mid \text{int}$

## 3 LL Top-Down Parsing

### 3.1 LL(1)

**Definition 3.1 (LL(1) Grammar)** For all rules  $A \rightarrow \alpha \mid \beta$ ,  $FIRST(\alpha)$  and  $FIRST(\beta)$  are disjoint. If  $FIRST(\alpha)$  contains  $\epsilon$  then  $FIRST(\beta)$  and  $FOLLOW(\alpha)$  are disjoint. Vice-versa.

**Definition 3.2 (Extended BNF)** contains  $\{\alpha\}$  for 0 or more occurrences and  $[\alpha]$  for 0 or 1 occurrences

### 3.1.1 CFG $\rightarrow$ LL(1)

- **Left Factorisation**  $A \rightarrow BC \mid BD \rightarrow A \rightarrow B(C \mid D)$  and  $A \rightarrow BC \mid B \rightarrow A \rightarrow B[C]$
- **Substitution** replaces a rule with its alternatives to make conflicts direct
- **Left Recursion Removal**  $A \rightarrow X_1 \mid \dots \mid X_n \mid AY_1 \mid \dots \mid AY_n \rightarrow A \rightarrow (X_1 \mid \dots \mid X_n)\{Y_1 \mid \dots \mid Y_n\}$   
more specifically  $A \rightarrow X \mid AY \rightarrow A \rightarrow X\{Y\}$

### 3.1.2 LL(1) to Parse Function

- **A B:** A(); B();
- **A|B:** if token in FIRST(A) then A(); elif token in FIRST(B) then B(); else error();
- **{A}:** while token in FIRST(A): A();
- **[A]:** if token in FIRST(A): A();

## 4 Garbage Collection

**Heap Compaction:** Mark live blocks  $\rightarrow$  Relocate live blocks  $\rightarrow$  Update pointers to relocated blocks

- **Reference-Counting:** Requires special techniques for cyclic data structures
- **Mark-Sweep:** **Mark:** Mark all blocks reachable from non-heap references as live  
**Sweep:** Scan all blocks to reclaim dead blocks, unmark live blocks for next sweep  
Use pointer reversal technique to visit all nodes of a Directed Graph without additional stack space  
Provides the largest possible block available when combined with compaction
- **Two-Space:** Allocate blocks to **From-Space**. When it's exhausted, copy live blocks to **To-Space**  
Automatic compaction, very fast to allocate objects, wastes half of memory, relocate long-lived objects
- **Generational:** Heap divided areas based on block age. Perform GC on younger generations more

## 5 Register Allocation

transExp (Binop Minus e1 e2) r

### 5.1 Register Machine Strategy

```
transExp e1 r ++
transExp e2 (r+1) ++
[Sub r (r+1)]
```

### 5.2 Accumulator Strategy

```
transExp e2 r ++
[Push r] ++
transExp e1 r ++
[SubStack r]
```

### 5.3 Sethi-Ullman Weights

```
if weight e2 > weight e1 then
  transExp (nxtreg:dstreg:regs) ++
  transExp (dstreg:regs) ++
  [Sub dstreg nxtreg]
```

```
weight (Binop Minus e1 e2) = min [
  max [weight e1, (weight e2) + 1],
  max [weight e2, (weight e1) + 1]
]
```

### 5.4 Graph Colouring

1. Generate intermediate three-address code where values are always saved in named locations
2. Construct inference graph where nodes are locations and nodes are linked if their live ranges overlap
3. Try to colour the nodes so no connected nodes have the same colour

## 6 Optimisation

### 6.1 Live Variables

$$\begin{aligned} \text{LiveOut}(n) &= \bigcup_{s \in \text{succ}(n)} \text{LiveIn}(s) \\ \text{LiveIn}(n) &= \text{uses}(n) \cup (\text{LiveOut}(n) - \text{defs}(n)) \end{aligned}$$

```
foreach n in CFG:
  LiveIn(n) := {}; LiveOut(n) := {}
while true: foreach n in CFG (backwards):
  LiveIn(n) = uses(n) + (LiveOut(n) - defs(n))
  LiveOut(n) = union [LiveIn(s) | s <- succ n]
```

### 6.2 Reaching Definitions

$$\begin{aligned} \text{ReachIn}(n) &= \bigcup_{s \in \text{pred}(n)} \text{ReachIn}(s) \\ \text{ReachOut}(n) &= \text{Gen}(n) \cup (\text{ReachIn}(n) - \text{Kill}(n)) \end{aligned}$$

For  $n : t = u1 \oplus u2$  ( $\text{defs}(n) = \{t\}$ ,  $\text{uses}(n) = \{u1, u2\}$ )

- $\text{Gen}(n)$  is definitions generated by  $n$ ,  $\{n\}$
- $\text{Kill}(n)$  is all definitions of  $t$  except  $n$

### 6.3 Loops

**Definition 6.1 (Loop)** A set of nodes  $S$  including a header node  $h$  such that

- Any node in  $S$  has a path leading to  $h$
- There is a path from  $h$  to any node in  $S$
- There is no edge from any node outside  $S$  to any node in  $S$  other than  $h$

**Definition 6.2 (Dominator)** A node  $d$  dominates node  $n$  if every path from the start node to  $n$  must go through  $d$ . Every node dominates itself.

$$\begin{aligned} \text{Doms}(s) &= \{s\}, \text{ for start node } s \\ \text{Doms}(n) &= \{n\} \cup \left( \bigcap_{p \in \text{preds}(n)} \text{Doms}(p) \right), \text{ otherwise} \end{aligned}$$

**Definition 6.3 (Back Edge)** An edge from node  $n$  to node  $h$  that dominates  $n$  is called a back edge. For every back edge, there is a loop. Two loops can share the same header.

**Definition 6.4 (Loop Invariants)** For hoisting a node  $d$ :  $t = a \oplus b$ ,

- **Reaching Definitions:** All reaching defs used by  $d$  occur outside loop
- **Dominators:**  $d$  dominates all loop exits
- **Count:** There must only be one def of  $t$  in the loop
- **Live Variables:**  $t$  must not be  $\text{LiveOut}$  from the loop's preheader

**Definition 6.5 (Single State Assignment)** Introduces a new name each time a variable is assigned. At control-flow joins, insert a dummy operator  $t = \phi(t_1, t_2)$  which magically picks either value depending on what path is taken.

In the generated code, we push the assignments of  $t$  backwards into the two predecessor paths.

To hoist a node after SSA conversion, only need to check **Reaching definitions**.