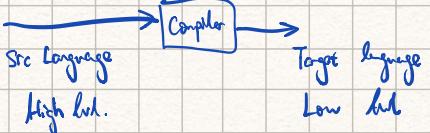
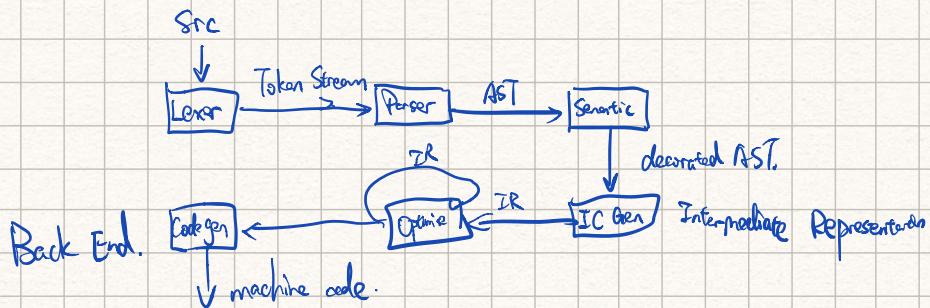


Compiler



Compilers: bridge semantic gap

Front End.

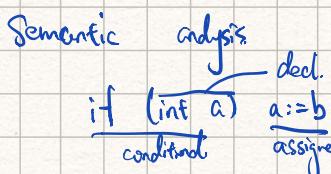
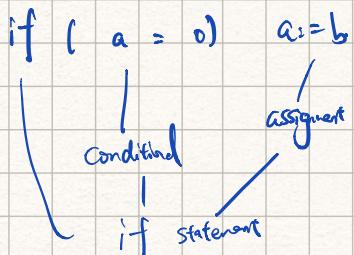


Lexical Analysis.

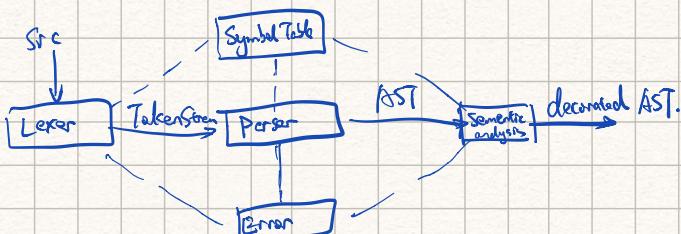
Delimiters define words ↳ In. It /.../

Lexemes Token Types

Syntax Analysis



Error Handling & Symbol Table.



Regular Expressions

Lexeme.

Def.: a contiguous seq. of chars that form a lexical unit in the grammar of language.

Token Types:

Lexemes	Token Type	descrip
if, while	IF, WHILE, FOR	Keywords
is, true	ID	identifiers
10, 0.75	INT	integers
2.0, 0.3	REAL	floats
(,), ,	LBRACK, RBRACK	Symbols

We use Regex to specify lexemes.

Defining a Language

↪ finite set of symbols

Ex. $\Sigma = \{a, b\}$

Regex R.	LCR.
a	$\{a\}$
ab	$\{a, b\}$
ab	$\{a, b\}$
(ab)*	$\{\epsilon, ab, abab, \dots\}$
(a ε)b	$\{b, ab\}$
(a b*)a*	$\{aa, a, ba, bba, \dots\}$
(ΣΣΣ)*	$\{www \text{ mod } 3 = 0\}$

Associativity

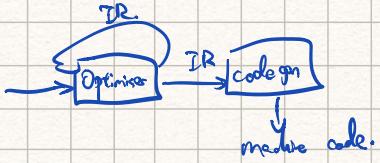
Operator priority $* > . > |$

Syntactic Sugar

a^+	aa^*
$a?$	$a \epsilon$
$[abc]$	$a b c$
$[a-z]$	$a b \dots z = \Sigma$
*	$a b \dots z = \Sigma$
$[a-c]$	$\Sigma - \{a, b, c\}$

Ambiguity

All phase.



To compile or to interpret.



REPL.

Read - Evaluate - Print Loop.
scripting shell languages

$\text{if } a := 0$

$\text{if } a := 0 \text{ or if } a := 0$
If IDENTIFIER

Disambiguation:

Longest match / rule priority

What's Next?

PSA FSM Finite State Automata

Exercise

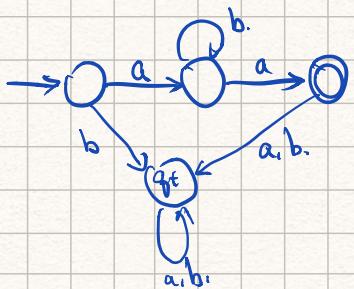
1. $(aa)^*$

2. $/(["^"]^* | "d" | /["^"]^*) //$

'
Lexical Analysis.

input \rightarrow Acceptor for L . \rightarrow Yes if $t \in LCR$
No if $t \notin LCR$

Example DFA for ab^*a over $\Sigma = \{a, b\}$.



NFA. $(\Sigma, Q, \delta, q_0, F)$

\downarrow states
 \downarrow alphabets
 \downarrow transitions
 $Q \times \Sigma \xrightarrow{\delta} P(Q)$
power set.

$f \subseteq Q$

Set of acc. states

$\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$ ϵ -transition.

$2^Q \Leftrightarrow$ power set of Q .

If $Q = \{a, b, c\}$,

$2^Q / P(Q) = \{a, b, c\}$.

$$2^Q = \{ \emptyset, \{a\}, \{b\}, \{c\}, \{ab\}, \{ac\}, \{bc\}, \{abc\} \} \geq 2^{101}$$

NFA always get the accept route if accepted
parallel interpretation / perfectly lucky

DFA vs. NFA.

for each symbol fully determined.	choose i) Σ -transition. ii) which transition.
Accept ends in q_F	contains a path from $q_0 \rightarrow q_F$.

Implementation	Table driven Space \uparrow Time \uparrow .	possible states Space \downarrow Time \uparrow .
Accept ends in q_F	Table driven Space \uparrow Time \uparrow .	possible states Space \downarrow Time \uparrow .

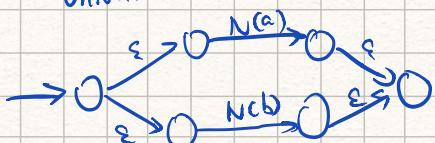
Regex \rightarrow NFA. (Thompson's Construction)



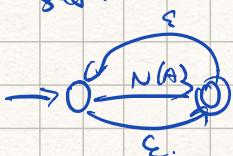
Concat.



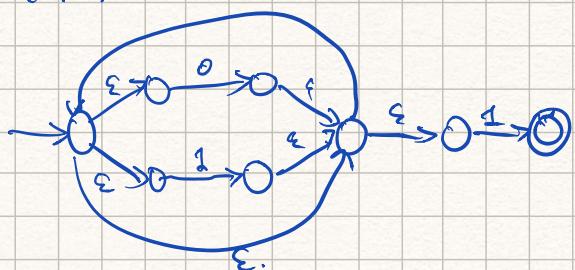
Union.



Star.



$(0|1)^* \in L.$ a



NFA \rightarrow DFA.

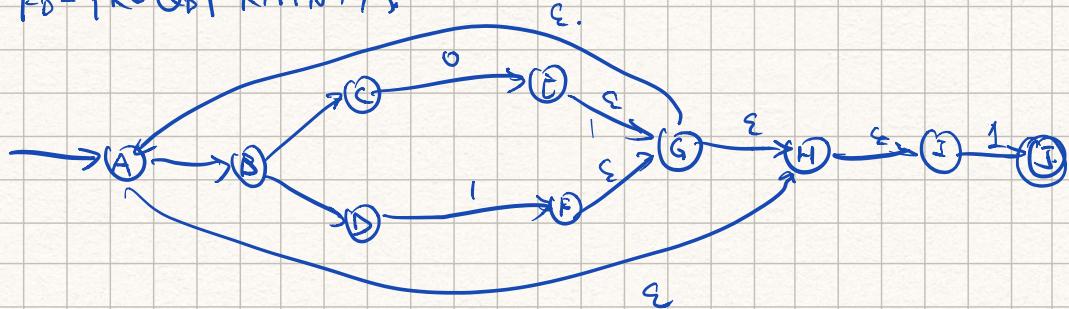
Theorem.

NFA and DFA are equivalent

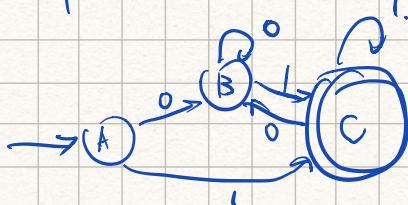
$QD \cong P(Q_N)$.

$\forall R \in \text{BD}, \forall a \in \Sigma_D(r,a) = \bigcup_{r \in R} \epsilon\text{-close}(\delta_N(r,a))$

$F_D = \{R \in Q_D \mid R \cap F_N \neq \emptyset\}$

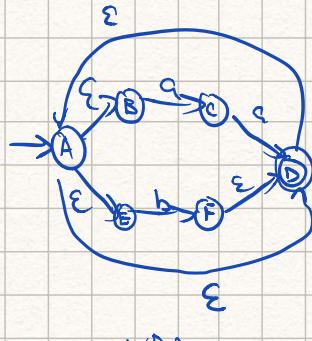


Use epsilon close.

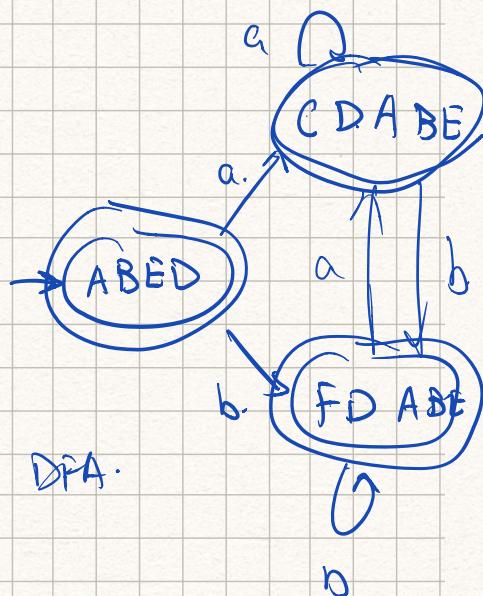


Exercise.

$(a|b)^*$

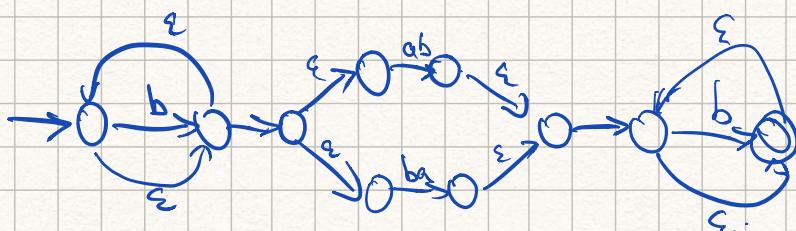


NFA



DFA.

$b^* (ab|ba)b^*$



Finite State Transducers

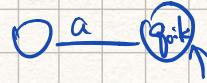
Mealy Machine. $\delta = Q \times \Sigma_E \rightarrow T_\Sigma \times Q$

input alphabets

output alphabets.

input output

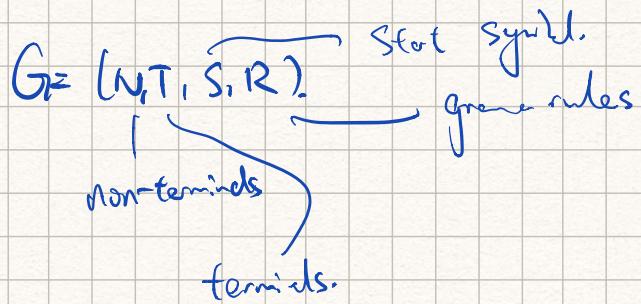
Moore Machine. $S = Q \rightarrow \Gamma \times Q$



output.

FST in lexer label transitions to error states / report lexical errors

Context Free Grammar.



$\text{EXPR} \rightarrow \text{if } \text{EXPR} \text{ then } \text{EXPR} \text{ else } \text{EXPR. fi}$
 | while EXPR loop EXPR pool.
 | id.

For a grammar G_1 , $L(G_1)$ is $\{a_1 \dots a_n \mid S \xrightarrow{*} a_1 \dots a_n \text{ where } a_i \in T\}$
 all set of strings a gram can generate

In CFG,

$\text{str} \in L(G_1)$ is a g^s-on.

also need to produce a parse tree.

handle errors gracefully

need an implementation CUP

The form of a grammar is important. (Ambiguities)

The Chomsky Hierarchy

a right-regular grammar is the grammar $G = (N, T, S, R)$

R s all obey the constraints

$$A \rightarrow a$$

$$A \rightarrow aB$$

$$A \rightarrow \epsilon.$$

$A \in N$ and $a \in T$.

for left-regular.

$$A \rightarrow Ba.$$

if left and right are mixed \rightarrow CFG.

Grammer.	Production Constraint	Automata.
Regular (right)	$A \rightarrow a \mid aB \mid \epsilon.$	FSA. NFA / DFA.
Context free.	$A \rightarrow \alpha.$	PDA.
Context sensitive	$\alpha A \beta \rightarrow \alpha \delta \beta$	LBA Linear Bounded Automata
Universal.	$\alpha \rightarrow \beta.$	TM.

Derivations

a seq. of sentential forms connected by the application of a single production from S .

$$S \rightarrow X_1 \dots X_n \rightarrow \dots \rightarrow Y_1 \dots Y_m.$$

$$\text{if } x_i \rightarrow y_j \text{ GR.}$$

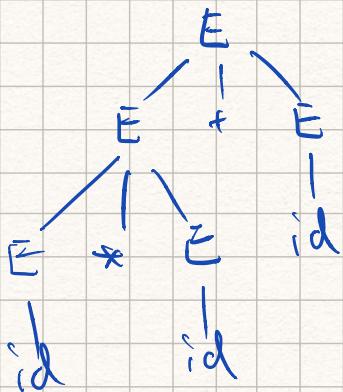
Parse Trees

S is the root.

$X \rightarrow Y_1 \dots Y_n. \quad Y_1 \dots Y_n$ are children of X

$E \rightarrow E+E$
 $\rightarrow E^*E$
 $\rightarrow (E)$
 $\rightarrow id$.

String: id+id+id



A parse tree has terminals at leaves
non-terminals as interior nodes

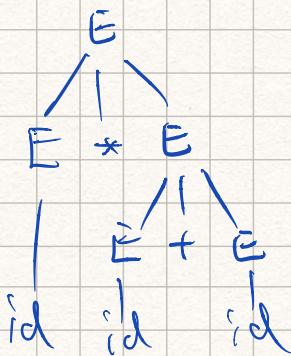
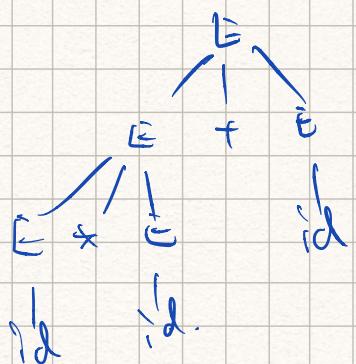
in-order gives the original input.

left-most derivations
at each step / replace left non-terminal.
right - - - - - right - - - - -

They produce SAME result.

it's just the producer order differences

Two Parse Trees



Ambiguity.

a grammar is ambiguous if
it generates ≥ 1 left-most
parse tree for some string

it's BAD can leave
programs ill-defined.

b*4+2

$E+E$ first gives 2b

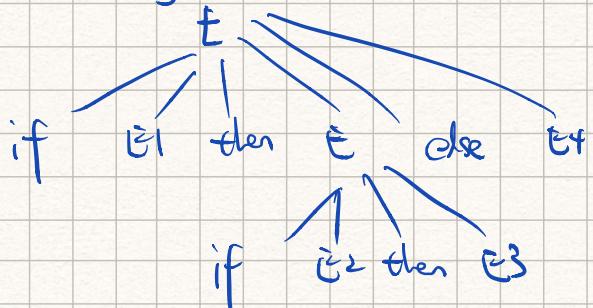
E^*E first gives 3b

Handling ambiguities

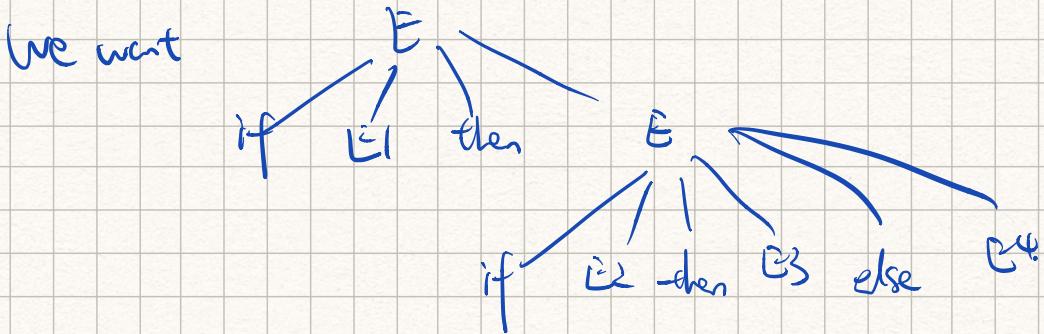
1. $E \rightarrow E' + E | E'$ ← too complicated
 $E' \rightarrow id * E' | id | (E) * E' | (E)$

2. Or forces precedence of * over +

Example Dangling Else



← not what we want.



Ambiguous definition can simplify a grammar.

Instead of we hacky grammar.

use natural grammar / with disambiguating annotations

% left + %left *:
precedence decl.

Top Down Parsing

parse tree constructed from Top \rightarrow down.
left \rightarrow right.

Recursive Descent Parsing

Def parsing. find parse tree for tokens $t_1 \dots t_n$.

Def Recursive Descent Parsing.

From start symbol, try all productions exhaustively

- The fringe of the parse tree is $t_1 \dots t_k A$.

- Try $A \rightarrow BC / CD \dots$

\downarrow
 $t_1 \dots t_k BC \dots$

- Backtrack when fringe doesn't match prefix of string

- Stop when all terminals

Consider $S \rightarrow Sa$.

$S \rightarrow Sa \rightarrow S Sa \dots$



left recursive grammar

No fringe to match.

1. Eliminate left recursion

$S \rightarrow Sa / \beta$

$S \rightarrow I / So$

$S \rightarrow \beta S'$

$S \rightarrow IS'$

$S' \rightarrow \alpha S / \epsilon$

$S' \rightarrow OS' / \epsilon$

In general:

$S \rightarrow S\alpha_1 \dots S\alpha_n | \beta_1 | \beta_2 | \dots | \beta_m$

$S \rightarrow \beta_1 S' | \beta_2 S' | \dots | \beta_m S'$

$S' \rightarrow \alpha_1 S' | \alpha_2 S' | \dots | \epsilon$

Indirect Left Recursion

$S \rightarrow A\alpha | S$

$A \rightarrow S\beta$

If put A back in $S \rightarrow S\beta\alpha | S$

$$\downarrow$$

$$S \rightarrow S S'$$

$$S' \rightarrow \beta \alpha S' \epsilon.$$

ϵ -market Recursion.

Def. non-terminal is nullable if $N \xrightarrow{*} \epsilon$.

$$S \rightarrow Y S \mid q$$

$$Y \rightarrow A \times B \mid Cc \mid cd \mid \epsilon.$$

Y is nullable. $\rightarrow S$ is left recursive.

Put Y in S

$$S \rightarrow A \times B S \mid Cc S \mid cd S \mid q.$$

Don't need $S \rightarrow S$ from $Y \rightarrow \epsilon$. ← useless

Repeat.

Eliminate,

Unreachable rules and $N \xrightarrow{*} N$

Direct left recursion via transformation

Indirect via substitution

ϵ -market recursion

Until no elimination applies

Survey of Recursive Descent (Top-Down)

- Simple / general.
- (left recursion elimination.
can be done automatically)
- Unpopular. because of backtracking
 - inefficient.
 - outdated (LOL)

We can avoid backtracking using Predictive Parsers

LL(k) grammars.
 $\uparrow \downarrow \leftarrow$

LL(1) Languages.

Left
to
Right look ahead

LL(1) means that for each non-terminal and token,
only one production leads to
a successful parse.

In a 2D lookUp Table.

Usually LL(1).

Consider grammar:

$$\begin{aligned} E &\rightarrow T + E \mid T \\ T &\rightarrow \text{int} \mid \text{int} * T \mid (E) \end{aligned}$$

for T two productions start with int.

E not clear what to do.

S.t. left factoring is
needed for LL(1)

$S \rightarrow \text{ID} \mid \text{ID} [E] \mid$ is not LL(1) two alternatives start with ID

In general.

$S \rightarrow AB_1 \mid AB_2 \mid \dots \mid AB_n$ < rules not start by A>

S.t.

$S \rightarrow AB' \mid$ < rules not prefix A >

$A' \rightarrow B_1 \mid B_2 \mid \dots \mid B_n$ where B_i can be ϵ

$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$

$$E \rightarrow TE'$$

$$E' \rightarrow +E \mid \epsilon.$$

$$T \rightarrow \text{int } T' \mid (E)$$

$$T' \rightarrow *T \mid \epsilon$$

\$ EOF.

	int	*	+	()	\$
E	TE'			TE'		
E'			TE		ϵ	ϵ
T	int'			(E)		
T'		*T	ϵ		ϵ .	ϵ

Blank entries in table error.

No way to get a string start with "*" from E

Method similar to recursive descent, except.

For each nonterminal S , look at next - a.

choose $[S, a]$

Use stack to keep track of non-terminals

reject on error state.

Accept on EOF.

LL(1) Example.

int * int.

Constructing LL(1) Parsing Table

LL(1) language are those defined
by LL(1) parsing table.

- No multiple definition for one table
entry.

- If there is a table,
parsing single and fast.
no backtracking necessary

LL(1) parse table is a matrix $|N| \times (|T| + 1)$

$T[A, b] = \rightarrow p$: $A \rightarrow p$

$T[A, b] = \rightarrow \epsilon$ discard A

Consider State $S \rightarrow p A$

↓
consider A: b is the lookahead.

1. $b \in \text{prefix } A : b \in \text{First}(A)$

2. $b \notin \text{an expansion of } A$ but. $A \rightarrow \epsilon$ b belongs to another
expansion of γ .

$b \in \text{Follow}(A)$

$\epsilon \in \text{First}(A)$

Example

$E \rightarrow T X$

$X \rightarrow + E / \epsilon$

Stack	Input	Action.
$E \$$	int * int \$	$T E'$
$T E' \$$	int * int \$	in T' .
$\text{int} T' E' \$$	int * int \$	terminal
$T' E' \$$	+ int \$	$\neq T$.
$* T E' \$$	* int \$	terminal.
$T E' \$$	int \$	int'
$\text{int} T E' \$$	int \$	terminal
$T E' \$$	\$	ϵ .
$E' \$$	\$	ϵ
\$	\$	ACCEPT.

$$T \rightarrow (E) \mid \text{int } Y \quad Y \rightarrow *T \mid \epsilon.$$

$$\begin{aligned}\text{First}(T) &= \{\text{int}\} \\ \text{First}(E) &= \text{First}(T) \\ \text{First}(X) &= \{E, \epsilon\} \\ \text{First}(Y) &= \{*, \epsilon\}\end{aligned}$$

Follow Sets

$$S \rightarrow * \beta A Y$$

$$\text{Follow}(A) = \text{First}(Y) - \epsilon \cup \text{First}(Y_1) - \epsilon$$

↑
if $y \rightarrow \epsilon$.

$$\epsilon \notin \text{Follow}(X)$$

$$\begin{aligned}F_w(E) &= \{\$, \}, \} \cup F_w(X) \\ &= F_w(E) \cup \{\$, \} \\ &= \{\$, \}\}\end{aligned}$$

Construct a parse table.

$$\begin{array}{c} A \rightarrow a. \\ b \in \text{First}(a) - \{\epsilon\}. \quad T(a, b) = a. \\ \hline a \rightarrow \epsilon \quad b \in \text{Follow}(a) \quad \text{do.} \\ \quad \quad \quad T(a, b) = \epsilon. \end{array}$$

$$F_w(X) = F_w(E) = \{+, \$, \}$$

$$\begin{aligned}F_w(T) &= \{\text{First}(X) - \{\epsilon\}\} \cup F_w(Y) \\ &\quad \cup F_w(E) \\ &= \{+, \} \cup F_w(T) \cup \{+, \$, \} \\ &= \{+, \}, \$\}\end{aligned}$$

$$F_w(Y) = F_w(T) = \{+, \}, \$\}$$

Notes on LL(1).

- G is not LL(1) if
 - G is ambiguous
 - G is left recursive
 - G is not left-factored

Bottom Up Parsing

Handles same grammar just as efficient.

Consider $E \rightarrow E + (E) \mid \text{int}$ ← not LL(1) / due to left recursion.

Idea.

str ← input.

repeat.

Identify β has been excised in str. (rightmost) $A \rightarrow \beta$ is a production.

reduce $\beta \rightarrow A$. → str becomes $\alpha A \gamma$

until $\text{str} = \gamma$

Example

↓

Bottom-up parsing

int + (int) + (int)
 $\rightarrow \underline{E} + (\underline{E}) + (\underline{E})$
 $\rightarrow \underline{E} + (\underline{E})$
 $\rightarrow E$

reversely
 on right most
 derivation.

Idea:

Split the string into two substrings.

- right substring \rightarrow a string of terminals
- left substring \rightarrow NT and Ts.

divide point marked by |

Initially $| x_1 x_2 \dots x_n$.

Two things to do: Shift Reduce.

Shift: move | \rightarrow |

Reduce: Apply inverse. $\underline{E + (E)}| \rightarrow E|$

Implement with a stack



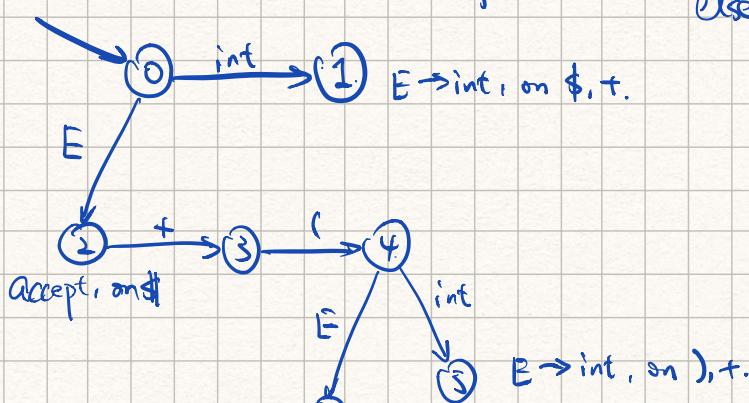
char left of l. Shift pushes a new elem in.
 Reduces pops ts and ntS out, push a NT in.

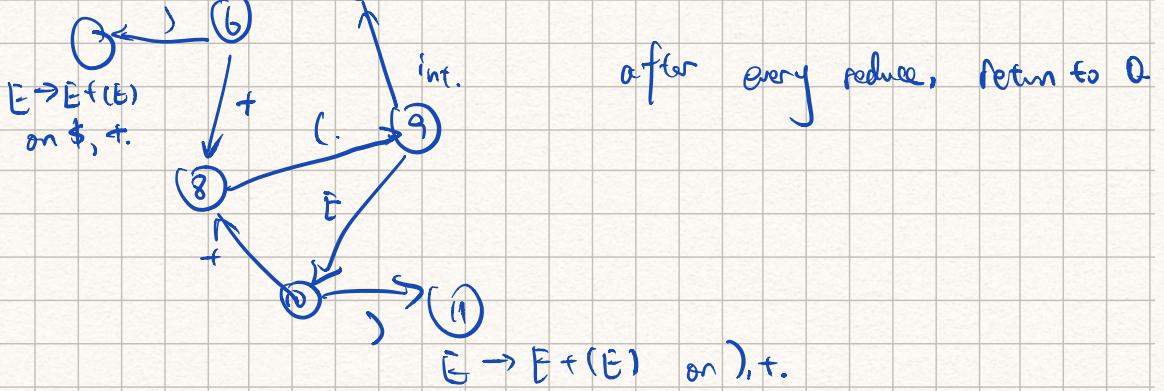
Key Issue: When to shift/reduce?

- Decide based on stack's content and a lookahead.
- Idea: use FSA. $S = NVT$.

run FSA on stack / check state X and lookahead S

if X has a transition label s. shift.
 else reduce





FSA represented as a 2D table.

FSA states are rows

Terminals / NTs are columns

After every shift/reduce, we run FSA again / (inefficient)
for each stack element., remember which state FSA on.

$$Q, \Sigma = NFT.$$

$$\sigma_i \in Q$$

LR parser maintains a stack $\langle \sigma_1, \alpha_1 \rangle - \langle \sigma_n, \alpha_n \rangle$

We jump to state σ_n after shift/reduce

Let $I\$$ Be input.

$j = 0$ be index into I , σ_0 is FSA's start state.

stack = push $(\langle \sigma_0, \epsilon \rangle)$

repeat

$\langle \sigma, \alpha \rangle :=$ top stack

shift: action $(\sigma, I[j])$;

shift $\sigma_t :=$ push $(\langle \sigma_t, I[j+1] \rangle)$

reduce $X \rightarrow RHS$.

| RHS. pop (stack).

push $(\langle X, goto, (\sigma, \alpha) \rangle)$

accept: halt.

error: halt & err.

Constructing FSA Table

LR(1) item. a pair $X \rightarrow \alpha \cdot \beta, q$.

$X \rightarrow \alpha \beta$ is a production marked by what is on the stack. α .
Pencils to be handled β .

a is a terminal.

LR(1) means 1 lookahead.

$[X \rightarrow \alpha \cdot \beta, a]$ describes a context.

X followed by a .

top of stack has α .

Need to find a prefix derived from β a next.

Convention new start symbol $S^1 \rightarrow S, \#$ empty stack

LR(1) items

In $E \rightarrow E^+ \cdot (E), +$

we can shift if next is (

$E \rightarrow E^+ (\cdot E), +$

In $E \rightarrow E^+ (E)^*, +$

we can reduce $E \rightarrow E^+ (E)$ if a + follows

Consider $E \rightarrow E^+ (\cdot E), +$

↑
expect has a prefix from E and suffix
is) +.

E has two productions:

$E \rightarrow \cdot \text{int},)$

$E \rightarrow \cdot E^+ (E),)$

The closure operation.

For $G = \langle N, T, R, S \rangle$

Closure (items) =

repeat

A $[X \rightarrow A \cdot Y \beta, a] \in \text{items}$

A $Y \rightarrow y \in R$.

A be first (βa)

items = item $\cup [Y \rightarrow \cdot y, b]$

until items unchanged.

Example

$$S \rightarrow E$$

$$E \rightarrow \text{int} \mid E + (E)$$

Closure ($S \rightarrow^* E, \$$)

$$S \rightarrow E, \$$$

$$E \rightarrow \cdot E + (E), \$$$

$$E \rightarrow \cdot \text{int}, \$$$

$$E \rightarrow E + (E), +$$

$$E \rightarrow \text{int}, +$$

$$S \rightarrow E, \$$$

$$E \rightarrow \cdot E + (E), \$$$

$$E \rightarrow \text{int}, \$$$

FSA state is a closed set of

LR(0) items

Start state contains $[S \rightarrow E, \$]$

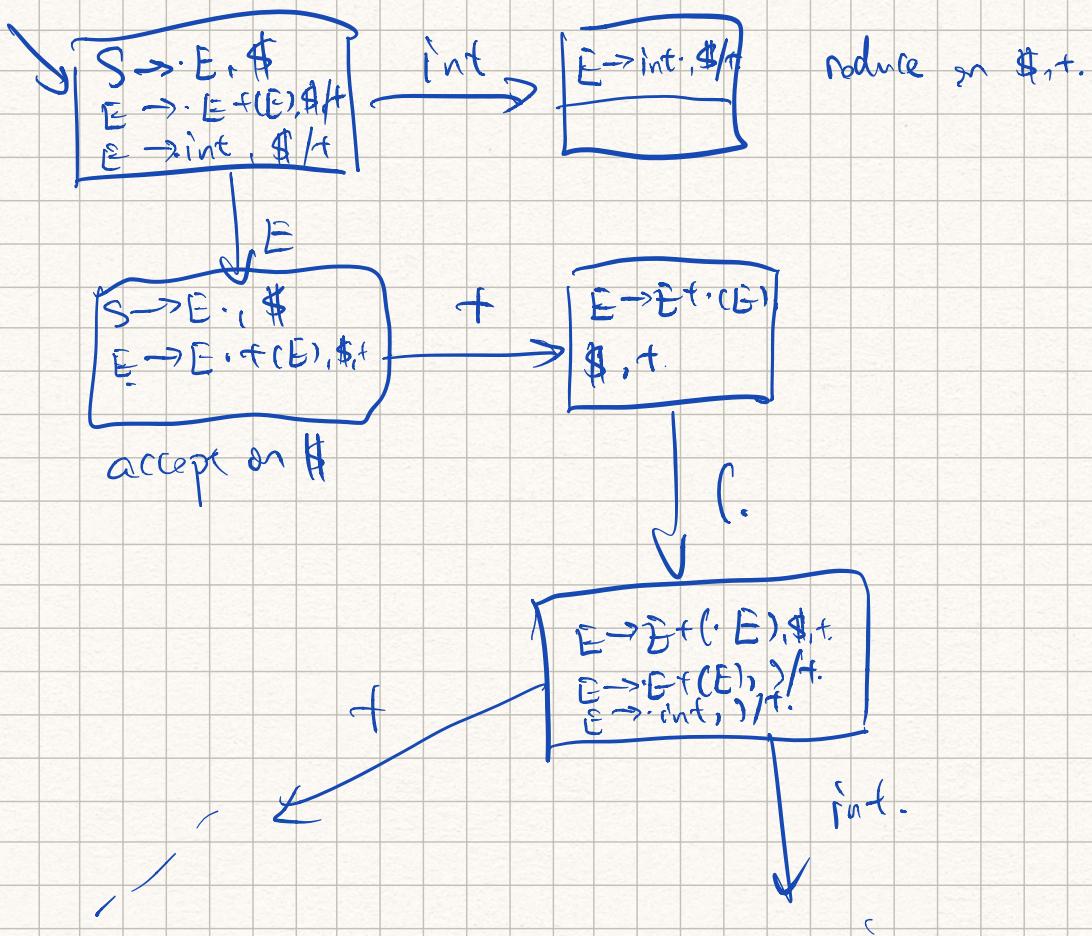
a state that contains $[X \rightarrow a, b]$

reduce $X \rightarrow a$ on b
backtrack

Building LR(0) Transitions.

$[X \rightarrow A \cdot \gamma \beta, b]$ has

a transition labeled γ .



Also write up.

workSet := $S_0 = \text{done}(S^1 \rightarrow \cdot S, \$)$; $Q = \emptyset$ $\delta = \emptyset$

while (workSet ≠ ∅)

$S_S = \text{choose}(\text{workSet})$ removes an item set.

$Q := Q \cup \{S_S\}$.

$\forall [x \rightarrow \alpha \cdot y \beta, t] \in \text{items}(S_S)$

$S_t = \text{transition}(S_S, y)$

$\delta = \delta \cup \{S_S \xrightarrow{y} S_t\}$.

workSet = workSet ∪ { S_t }

$F := \{q \in Q \mid \exists i \in G \text{ items}(q) \quad i = [Y \rightarrow \alpha \cdot, \$]\}$

return $\langle Q, \delta = NFT, \delta, F, S_0 \rangle$

Labels are implicitly items in a node

For $G = \langle N, T, R, S \rangle$

1. $\forall n \in N$, compute $F(n)$

2. $G' = \langle N \cup \{S'\}, T, R \cup \{S' \rightarrow S\}, S' \rangle$ add Rule $S' \rightarrow S, \$$.

3. Run FSA.

Create item sets.

Add transitions

4. Turn FSA to pos. NFA

The algo builds a NFA

δ doesn't have all transitions to all states

No - ε transitions / ambiguities

It's trivially a DFA if we connect all unwritten transitions to trash.

Shift reduce conflict

$X \rightarrow [A \cdot \alpha \beta, b]$.

$Y \rightarrow [\Gamma \cdot, \alpha]$?? Shift or reduce on a .

$[S \rightarrow \text{if } E \text{ then } S_1, \text{else } S_2]$

$[S \rightarrow \text{if } E \text{ then } S_1 \text{ else } S_2, x]$

Default is not shift. (correct in this case)

Example

$$E \rightarrow E+E \mid E+E \mid \text{int}$$

$$\begin{cases} E \rightarrow E+E \cdot, + \\ E \rightarrow E \cdot+E, + \end{cases}$$

reduce. + have same precedence.
 $E \rightarrow E+E \cdot$ rule
 + is left assoc.
L. reduce.

$$\begin{cases} E \rightarrow E+E \cdot, + \\ E \rightarrow E \cdot+E, + \end{cases}$$

reduce $E \rightarrow E+E$
 then + higher precedence
 termnd.

Dangling else

else precedence > then.

$\leftarrow S \rightarrow \text{if } E \text{ then } S_1 \text{ else } S_2$
 if E then S lower then. else

$S \rightarrow \text{if } E \text{ then } S \text{ lower then. }$ else

shift

Best avoid use of too much precedence.

Reduce - Reduce Conflicts

Example. Given Grammer.

$$S \rightarrow \epsilon \mid id \mid id S$$

$$S^1 \rightarrow \cdot S, \#$$

$$S \rightarrow \cdot id, \#$$

$$S \rightarrow \cdot id S, \#$$

$$S \rightarrow \cdot \cdot, \#$$

Σ

id.

$$\begin{array}{l} S \rightarrow id \cdot, \$ \\ S \rightarrow id \cdot S, \$ \end{array}$$

$$S \rightarrow \cdot \cdot, \$$$

$$S \rightarrow \cdot id S, \$$$

Reduce/Reduce Conflict $S \rightarrow \cdot id, \$$

Solve this by rewriting $S \rightarrow id\$ | \epsilon$

Using parser generators

procedure declarations

best parser gen do not construct PFA (PFA/NFA)

LR(1) Tables are big.

Many states are similar. → merge together.

The core of a set of LR items.

$\{x \rightarrow a \cdot \beta, b], t \cdot y \rightarrow y \cdot \gamma, d]\} \rightarrow \{x \rightarrow a \cdot \beta, t \cdot y \rightarrow y \cdot \gamma\}$

Some core / can be merged.



LALR(1) States.

Look-ahead LR states (statistically 10 times fewer)

LALR can have more reduce-reduce conflicts
(in practice rare)

LALR are efficiency hack.

Reasonable programming languages have a LALR(1) grammar.

