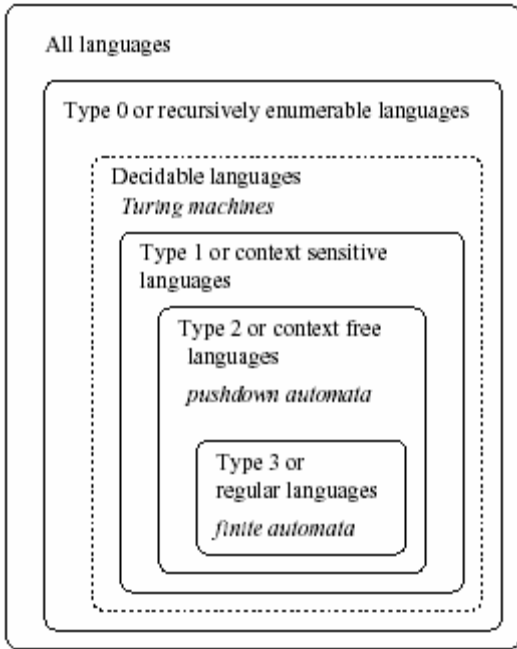


The Chomsky Hierarchy



A language is a set of words.
A word is a sequence of symbols.
The symbols are given by an alphabet, Σ .

The set of all languages is given by $P(\Sigma^*)$

Type 3 (Regular) Languages

Deterministic Finite Automata

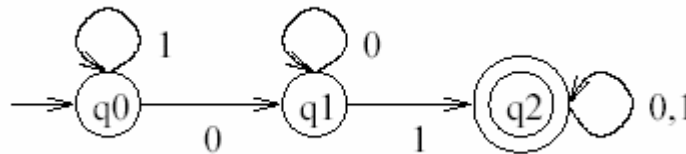
The machine which recognizes Regular languages. Equivalent to a computer with finite (fixed) memory.

A DFA, $A = (Q, \Sigma, \delta, q_0, F)$ where

- Q is a finite set of states
- Σ is the alphabet
- δ is the transition function $\delta \in Q \times \Sigma \rightarrow Q$
- q_0 is the initial state
- F is the finite set of final states

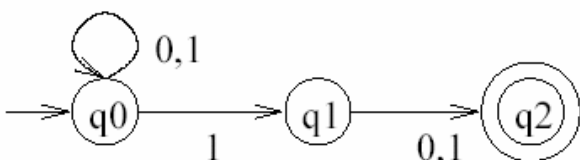
A typical transition sequence looks like this:

$$\begin{aligned} \delta(q_0, 101) &= \delta(\delta(q_0, 1), 01) \\ &= \delta(q_0, 01) \\ &= \delta(\delta(q_0, 0), 1) \\ &= \delta(q_1, 1) \\ &= \delta(\delta(q_1, 1), \epsilon) \\ &= \delta(q_2, \epsilon) \\ &= q_2 \end{aligned}$$



In a deterministic finite automaton, there is always precisely one transition for each alphabet symbol in each state.

Nondeterministic Finite Automata



This is very similar to a DFA, except that there can be multiple (or zero) transitions for a symbol in a state. A

word is accepted if there exists a path through the automaton that ends in a final state.

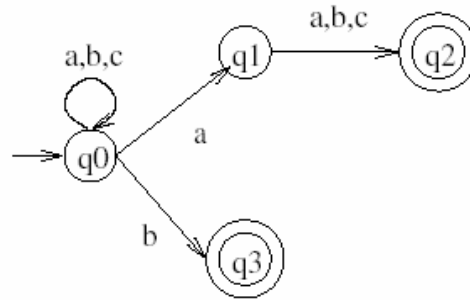
Constructing a DFA from an NFA

To do this, we use the subset construction.

Write down a table, adding any new states to the bottom. Then work through from the initial state – some states will not be reachable.

E.g.,

	a	b	c
> {q0}	{q0, q1}	{q0, q3}	{q0}
{q1}	{q2}	{q2}	{q2}
{q2}	{}	{}	{}
*{q3}	{}	{}	{}
{q0, q1}	{q0, q1, q2}	{q0, q3, q2}	{q0, q2}
*{q0, q3}	{q0, q1}	{q0, q3}	{q0}
{q0, q1, q2}	{q0, q1, q2}	{q0, q3, q2}	{q0, q2}
*{q0, q3, q2}	{q0, q1}	{q0, q3}	{q0}
{q0, q2}	{q0, q1}	{q0, q3}	{q0}



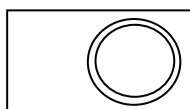
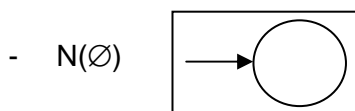
Regular Expressions

Symbol precedence: * (repetition) is strongest, then sequence, then + (or).
I.e., $ab^* = a(b)^*$ and $ab + cd = (ab) + (cd)$

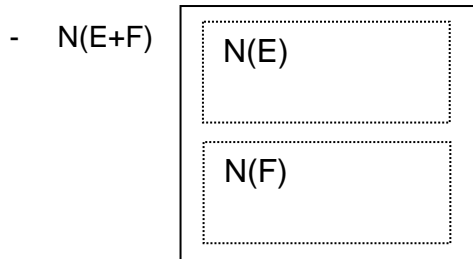
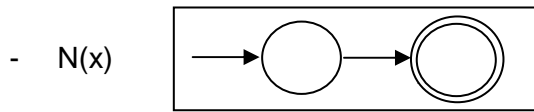
Examples:

- Words which contain a 'b'
 $(a + b + c)^* b (a + b + c)^*$
- Words which do not contain a b.
 $(a + c)^*$
- Words where all as appear before all the cs
 $(a + b)^* (b + c)^*$
- Words s.t. the number of as plus the number of bs is odd.
 $c^*(a + b)c^* ((a + b) c^* (a + b) c^*)^*$
- Words which contain the sequence aa.
 $(a + b + c)^* aa (a + b + c)^*$
- Words which do not contain the sequence aa.
 $(b + c)^* (a (b + c) (b + c)^*)^* (a + \epsilon)$

Deriving an NFA from a Regular Expressions

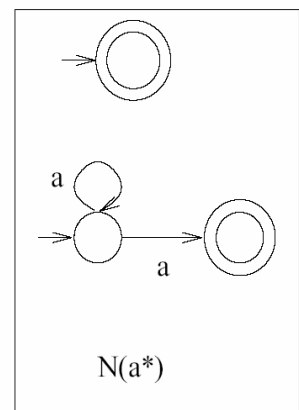
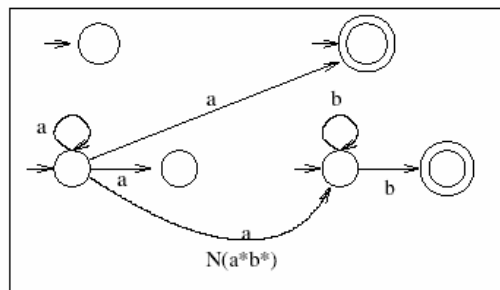
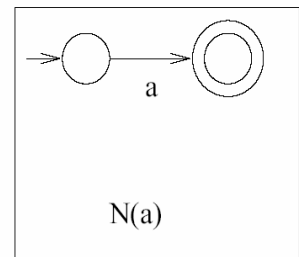
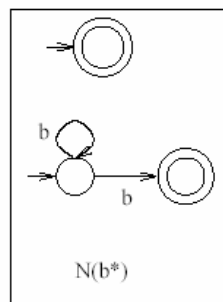
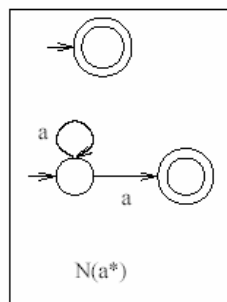


- $N(\epsilon)$ \rightarrow



- $N(EF)$: The initial states of E are the initial states of EF
The final states of F are the final states of EF
Connect the penultimate states of E to the initial states of F.

E.g.,



- $N(E^*)$:
- Link penultimate states to initial states
- Add an initial state which is final

Showing a Language is not Regular



Revision Notes

Regular languages can be recognized by a computer with finite (fixed) memory. Such a computer corresponds to a DFA.

The language 0^n1^n cannot be recognized using a finite memory – n can be arbitrarily large.

The *Pumping Lemma* can prove a language is/is not regular.

- Given a regular language, L , there is a number $n \in \mathbb{N}$ such that all words, $w \in L$ which are longer than n ($|w| \geq n$) can be split into three words $w = xyz$ such that:
 - $y \neq \epsilon$
 - $|xy| \leq n$
 - $\forall k \in \mathbb{N}. xy^kz \in L$

E.g.,

$$L_{\text{pali}} = \{wxw^R \mid w \in \Sigma^*, x \in \Sigma\} \cup \{ww^R \mid w \in \Sigma^*\}$$

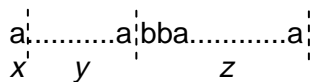
where w^R is the word w backwards.

Proposition: L_{pali} is not regular

Proof: Assume L_{pali} is regular and derive a contradiction.

Let $w = a^n b b a^n \in L$

$|w| \geq n$ – True.



Remember: $|xy| \leq n$.

Consider $xyyz = aa^{n-1}a^{n-1}a^{n-1}bba^n$ which is plainly not in L_{pali} .

Contradiction – hence, L_{pali} is not regular.



Type 2 (Context-free) Languages

Context Free Grammars

A context-free grammar, $G = (V, \Sigma, S, P)$ where:

V is the finite set of variables or nonterminal symbols

Σ is the set of terminal symbols (we assume V and Σ are disjoint)

S is the start symbol $\in V$

P is the finite set of productions: $P \subseteq V \times (V \cup T)^*$

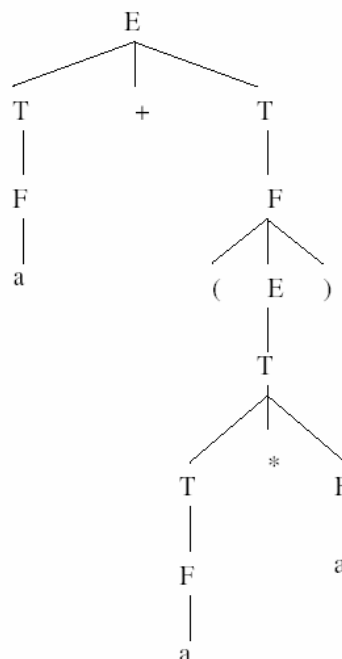
E.g., $G = (V = \{E, T, F\}$
 $\Sigma = \{a, +, *, (,)\}$
 $S = E$
 $P = \{ \begin{array}{l} E \rightarrow T \mid E + T \\ T \rightarrow F \mid T * F \\ F \rightarrow a \mid (E) \end{array} \}$

Let's try and form $a + (a * a)$:

$E \rightarrow E + T$
 $\rightarrow T + T$
 $\rightarrow F + T$
 $\rightarrow a + T$
 $\rightarrow a + F$
 $\rightarrow a + (E)$
 $\rightarrow a + (T)$
 $\rightarrow a + (T * F)$
 $\rightarrow a + (F * F)$
 $\rightarrow a + (a * F)$
 $\rightarrow a + (a * a)$

Parse Trees

A graphical representation of the derivation above:



A grammar is said to be ambiguous a word has more than one parse tree.



Adding symbol precedence to a grammar

Example:

$G = (\{E, A\}, \{p, q, r, (,), \vee, \wedge, \neg, \text{true}, \text{false}\})$ with P given by:

$E \rightarrow A \mid \neg E \mid E \wedge E \mid E \vee E \mid (E)$

$A \rightarrow p \mid q \mid r \mid \text{true} \mid \text{false}$

This is an ambiguous grammar. Suggest an alternative grammar for the same language which is not ambiguous. The grammar should reflect the following conventions on how to read boolean expressions:

- \neg binds stronger than \wedge and \vee
- \wedge binds stronger than \vee

Add in two more intermediary nonterminals so we have one for each level of precedence. The first (lowest) level of precedence is \vee , the next is \wedge , the next is \neg , the next is $(, p, q, r, \text{true}, \text{false}$.

The modified grammar is now:

$E \rightarrow F \mid E \vee F$

$F \rightarrow L \mid F \wedge L$

$L \rightarrow A \mid \neg L$

$A \rightarrow (E) \mid p \mid q \mid r \mid \text{true} \mid \text{false}$

N.B. Each nonterminal either goes to the same symbol connected to the next symbol | just the next symbol.

Pushdown Automata

The machine which recognizes Context-free languages. Uses a stack memory (LIFO).

A Pushdown automata, $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$ where

Q is a finite set of states

Σ is the alphabet

Γ is the stack symbols

δ is the transition function $Q, \Sigma, \Gamma, \delta \in Q \times (E \cup \{\epsilon\}) \times \Gamma \rightarrow P_{fin}(Q \times \Gamma^*)$

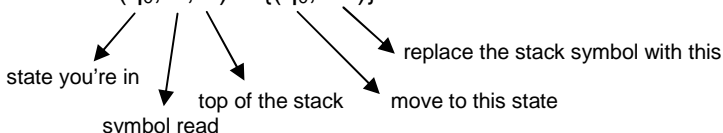
q_0 is the initial state

Z_0 is the initial stack symbol

F is the finite set of final states

A typical transition looks like this:

$$\delta(q_0, 0, \#) = \{(q_0, 0\#\}$$



PDAs can accept by final state *or* by empty stack (in which case F is omitted from the definition).

Determinism – a PDA is deterministic if there is never a choice of transition. It may get stuck. NPDAs are more powerful than DPDAs.



Turning a grammar into a PDA

Remember this grammar: $G = (\{E, T, F\}$
 $\Sigma = \{a, +, *, (,)\}$
 $S = E$
 $P = \{ E \rightarrow T \mid E + T$
 $T \rightarrow F \mid T * F$
 $F \rightarrow a \mid (E) \}$

We can define the PDA, $P(G) = (\{q_0\}, \{(,), a, +, *\}, \{E, T, F, (,), a, +, *\}, \delta, q_0, E)$ with:

$\delta(q_0, \epsilon, E) = \{(q_0, T), (q_0, E + T)\}$

$\delta(q_0, \epsilon, T) = \{(q_0, F), (q_0, T * F)\}$

$\delta(q_0, \epsilon, F) = \{(q_0, a), (q_0, (E))\}$

$\delta(q_0, (,)) = \{(q_0, \epsilon)\}$

$\delta(q_0, a, a) = \{(q_0, \epsilon)\}$

$\delta(q_0, +, +) = \{(q_0, \epsilon)\}$

$\delta(q_0, *, *) = \{(q_0, \epsilon)\}$

$\delta(q_0, \epsilon, \epsilon) = \{(q_0, \epsilon)\}$

LL(1) Parsers

In order to be unambiguous, a parser must have only one choice at each transition stage. The class of grammars where this is possible is called LL(1).

For each non-terminal symbol, A, we define First(A) which is the set of terminals that A can start with, and Follow(A) which is the set of terminals that can come immediately after A.

We then define a Lookahead set for each production. A grammar is LL(1) if each lookahead set is different.

E.g.

The grammar P is given by:

$Prog \rightarrow \{Stmts\}$

$Stmts \rightarrow \epsilon \mid Stmt Stmts$

$Stmt \rightarrow Name = Expr;$

 | if (Expr) Stmt

 | while (Expr) Stmt

 | print Expr ;

 | ;

 | Prog

$Expr \rightarrow Name \mid Num \mid (Expr Op Expr)$

$Name \rightarrow x \mid y \mid z$

$Op \rightarrow + \mid * \mid - \mid div$

$Num \rightarrow 0 \mid 1$

$First(Prog) = \{ \{ \}$

$First(Stmts) = First(Stmt) = \{x, y, z, if, while, print, ;, \{ \}$

$First(Expr) = \{x, y, z, 0, 1, (\}$

$First(Name) = \{x, y, z\}$

$First(Op) = \{+, *, -, div\}$

$First(Num) = \{0, 1\}$

$Follow(Prog) = \{ \$, \}, x, y, z, if, while, print, ;, \{ \}$

$Follow(Stmts) = \{ \}$

$Follow(Stmt) = First(Stmts) = \{ \}, x, y, z, if, while, print, ;, \{ \}$

$Follow(Expr) = \{ ;,), +, *, -, div \}$

$Follow(Name) = \{ ;,), +, *, -, div, = \}$

$Follow(Op) = First(Expr) = \{x, y, z, 0, 1, (\}$

$Follow(Num) = First(Expr) = \{x, y, z, 0, 1,) \}$



Revision Notes

Lookahead(Prog \rightarrow {Stmts}) = {}
Lookahead(Stmts \rightarrow ϵ) = {}
Lookahead(Stmts \rightarrow Stmt Stmts) = {x, y, z, if, while, print, ;, }
Lookahead(Stmt \rightarrow Name = Expr;) = {x, y, z}
Lookahead(Stmt \rightarrow if (Expr) Stmt) = {if}
Lookahead(Stmt \rightarrow while (Expr) Stmt) = {while}
Lookahead(Stmt \rightarrow print Expr ;) = {print}
Lookahead(Stmt \rightarrow ;) = {;}
Lookahead(Stmt \rightarrow Prog) = {}
Lookahead(Expr \rightarrow Name) = {x, y, z}
Lookahead(Expr \rightarrow Num) = {0, 1}
Lookahead(Expr \rightarrow (Expr Op Expr)) = {}
Lookahead(Name \rightarrow x) = {x}
Lookahead(Name \rightarrow y) = {y}
Lookahead(Name \rightarrow z) = {z}
Lookahead(Op \rightarrow +) = {+}
Lookahead(Op \rightarrow *) = {*}
Lookahead(Op \rightarrow -) = {-}
Lookahead(Op \rightarrow div) = {div}
Lookahead(Num \rightarrow 0) = {0}
Lookahead(Num \rightarrow 1) = {1}

Yes, the grammar is LL(1).

Type 0 (Recursively Enumerable) Languages

Turing Machines

Like a PDA, but uses an infinite *tape* not a *stack*.

A language is *Recursive* or *Decidable* if there is a Turing Machine which will always stop.

A Turing Machine $M = (Q, \Sigma, \Gamma, \delta, q_0, B, F)$ where

Q is the finite set of states

Σ is the alphabet

Γ is the set of tape symbols

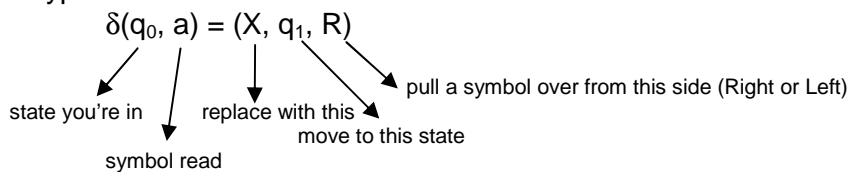
δ is the transition function: $\delta \in Q \times \Gamma \rightarrow \{\text{stop}\} \cup Q \times \Gamma \times \{L, R\}$

q_0 is the initial state

$B \in \Gamma$ is the blank tape symbol

F is the set of final states

A typical transition looks like this:



For example - $(X, q_0, aYbZc)$

You're in q_0 , and read an 'a' (always read from the right).

Replace this a with an X: $(X, q_0, XYbZc)$

Go into q_1 : $(X, q_1, XYbZc)$

Then pull a symbol over from the right: $(XX, q_0, YbZc)$

A grammar is said to be context-sensitive (type 1) if for any production, the left hand side is at least as long as the right hand side.