

Chapter 1

String Acceptors

1.1 Deterministic Finite-state String Acceptors

1.1.1 Orientation

This section is about deterministic finite-state acceptors for strings. The term *finite-state* means that the memory is bounded by a constant, no matter the size of the input to the machine. The term *deterministic* means there is single course of action the machine follows to compute the output from some input. As we will see later, *non-deterministic machines* can be thought of as pursuing multiple computations simultaneously. The term *acceptor* is synonymous with *recognizer*. It means that this machine solves *membership problems*: given a set of objects X and input object x , does x belong to X ? The term *string* means we are considering the membership problem over stringsets. So X is a set of strings (so $X \subseteq \Sigma^*$) and the input x is a string.

1.1.2 Definitions

Definition 1. A deterministic finite-state acceptor (DFA) is a tuple $(Q, \Sigma, q_0, F, \delta)$ where

- Q is a finite set of states;
- Σ is a finite set of symbols (the alphabet);
- $q_0 \in Q$ is the initial state;
- $F \subseteq Q$ is a set of accepting (final) states; and
- δ is a function with domain $Q \times \Sigma$ and co-domain Q . It is called the transition function.

We extend the domain of the transition function to $Q \times \Sigma^*$ as follows. In these notes, the empty string is denoted with λ .

$$\begin{aligned}\delta^*(q, \lambda) &= q \\ \delta^*(q, aw) &= \delta^*(\delta(q, a), w)\end{aligned}\tag{1.1}$$

Consider some DFA $A = (Q, \Sigma, q_0, F, \delta)$ and string $w \in \Sigma^*$. If $\delta^*(q_0, w) \in F$ then we say A *accepts/recognizes/generates* w . Otherwise A *rejects* w .

Definition 2. *The stringset recognized by A is $L(A) = \{w \in \Sigma^* \mid \delta^*(q_0, w) \in F\}$.*

The use of the ‘L’ denotes “Language” as stringsets are traditionally referred to as *formal languages*.

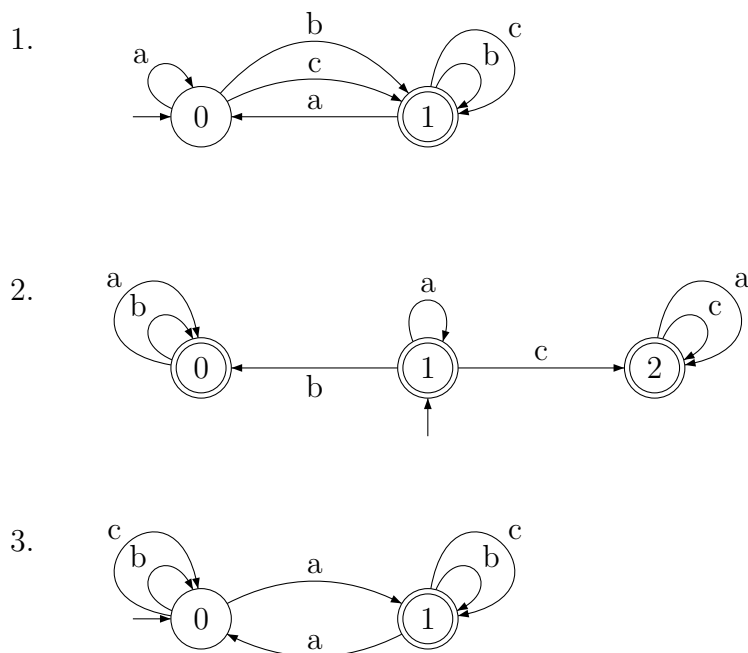
Definition 3. *A stringset is regular if there is a DFA that recognizes it.*

1.1.3 Exercises

Exercise 1. This exercise is about designing DFA. Let $\Sigma = \{a, b, c\}$. Write DFA which express the following generalizations on word well-formedness.

1. All words begin with a consonant, end with a vowel, and alternate consonants and vowels.
2. Words do not contain *aaa* as a substring.
3. If a word begins with *a*, it must end with *c*.
4. Words must contain two *bs*.
5. All words have an even number of vowels.

Exercise 2. This exercise is about reading and interpreting DFA. Provide generalizations in English prose which accurately describe the stringset these DFA describe.



4. Write the DFA in #1-3 in mathematical notation. So what is Q, Σ, q_0, F , and δ ?

1.2 Properties of DFA

Note that for a DFA A , its transition function δ may be partial. That is, there may be some $q \in Q, a \in \Sigma$ such that $\delta(q, a)$ is undefined. If δ is a partial function, δ^* will be also. It is assumed that if $\delta^*(q_0, w)$ is undefined, then A rejects w .

We can always make δ total by adding one more state to Q . To see how, call this new state \diamond . Then for each $(q, a) \in Q \times \Sigma$ such that $\delta(q, a)$ is undefined, define $\delta(q, a)$ to equal \diamond . Every string which was formerly undefined w.r.t. to δ^* is now mapped to \diamond , a non-accepting state. This state is sometimes called the *sink* state or the *dead* state.

Definition 4. A DFA is complete if δ is a total function. Otherwise it is incomplete.

It is possible to write DFA which have many useless states. A state can be useless in two ways. First, there may be no string which forces the machine to transition into the state. Second, there may be a state from which no string

Definition 5. A state q in a DFA A is useful if there is a string w such that $\delta^*(q_0, w) = q$ and a string v such that $\delta^*(q, v) \in F$. Otherwise q is useless. If every state in A is useful, then A is called trim.

Not all complete DFAs are trim. If there is a sink state, it is useless in the above sense of the word.

Definition 6. A DFA A is minimal if there is no other DFA A' such that $L(A) = L(A')$ and A' has fewer states than A .

Technically, not all complete DFAs are minimal. If there is a sink state, it is not minimal.

Exercise 3. Consider the DFAs in the exercise 2. Are they complete? Trim? Minimal?

1.3 Some Closure Properties of Regular Languages

A set of objects X is *closed* under an operation \circ if for all objects $x, y \in X$ it is the case that $x \circ y \in X$ too.

We can easily show that the union of any two regular stringsets R and S is also regular. Let $A_R = (Q_R, \Sigma, q_{0R}, F_R, \delta_R)$ be the DFA recognizing R and let $A_S = (Q_S, \Sigma, q_{0S}, F_S, \delta_S)$ be the DFA recognizing S . We can assume A_R and A_S are complete. We assume the same alphabet.

Construct $A = (Q, \Sigma, q_0, F, \delta)$ as follows.

- $Q = Q_R \times Q_S$.
- $q_0 = (q_{0R}, q_{0S})$.
- $F = \{(q_r, q_s) \mid q_r \in F_R \text{ or } q_s \in F_S\}$.
- $\delta((q_r, q_s), a) = (q'_r, q'_s)$ where $\delta_R(q_r, a) = q'_r$ and $\delta_S(q_s, a) = q'_s$.

Theorem 1. $L(A) = R \cup S$.

Similarly, the same kind of construction shows that the intersection of any two regular stringsets is regular. Construct $B = (Q, \Sigma, q_0, F, \delta)$ as follows.

- $Q = Q_R \times Q_S$.
- $q_0 = (q_{0R}, q_{0S})$.
- $F = \{(q_r, q_s) \mid q_r \in F_R \text{ and } q_s \in F_S\}$.
- $\delta((q_r, q_s), a) = (q'_r, q'_s)$ where $\delta_R(q_r, a) = q'_r$ and $\delta_S(q_s, a) = q'_s$.

Theorem 2. $L(B) = R \cap S$.

Here are some additional questions we are interested in for regular stringsets R and S .

1. Is the complement of R (denoted \overline{R}) a regular stringset?
2. Is $R \setminus S$ a regular stringset?
3. Can we decide whether $R \subseteq S$?
4. Is RS a regular stringset? (Note $RS = \{rs \mid r \in R \text{ and } s \in S\}$)
5. Is R^* a regular stringset? (Note $R^0 = \{\lambda\}$, $R^n = R^{n-1}R$, $R^* = \bigcup_{n \in \mathbb{N}^0} R^n$)

The answers to all of these questions is Yes. With a little thought about complete DFA, the answers to first three follow very easily.

Theorem 3. *If R is a regular stringset then the complement of R is regular.*

Proof (Sketch). If R is a regular stringset then there is a complete DFA $A = (Q, \Sigma, q_0, F, \delta)$ which recognizes it. Let $B = (Q, \Sigma, q_0, F', \delta)$ where $F' = Q \setminus F$. We claim $L(B) = \overline{R}$. \square

Corollary 1. *If R, S are regular stringsets then so is $R \setminus S$ since $R \setminus S = R \cap \overline{S}$.*

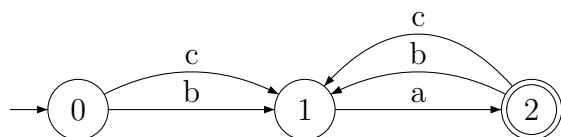
Corollary 2. *If R, S are regular stringsets then it is decidable whether $R \subseteq S$ since $R \subseteq S$ iff $R \setminus S = \emptyset$.*

Corollary 3. *If R, S are regular stringsets then it is decidable if $R = S$ since $R = S$ iff $R \subseteq S$ and $S \subseteq R$.*

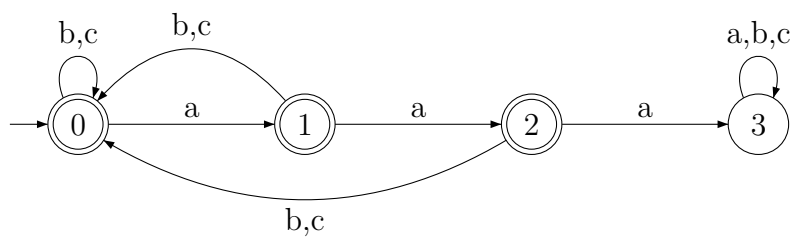
We postpone explaining how and why for the last two questions.

Some answers to Exercise 1

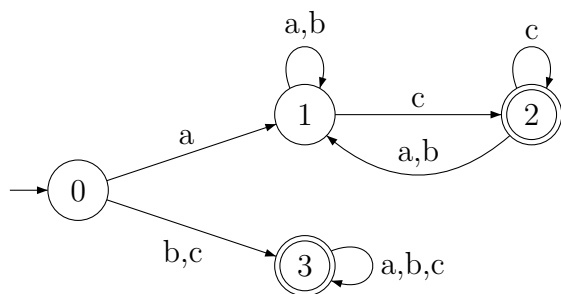
1.



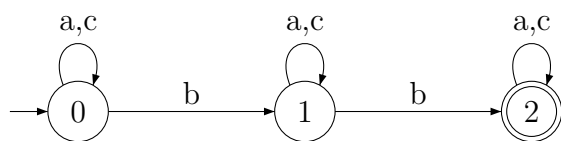
2.



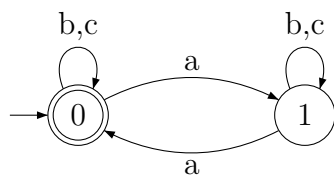
3.



4.



5.



Some answers to Exercise 2

1. Words must end in a consonant.
2. Words cannot contain both “b” and “c”. So words permit only one type of consonant.
3. Words must have an odd number of vowels.

Chapter 2

Tree Acceptors

2.1 Deterministic Bottom-up Finite-state Tree Acceptors

2.1.1 Orientation

This section is about deterministic bottom-up finite-state tree acceptors. The term *finite-state* means that the memory is bounded by a constant, no matter the size of the input to the machine. The term *deterministic* means there is single course of action the machine follows to compute its output. The term *acceptor* means this machine solves *membership problem*: given a set of objects X and input object x , does x belong to X ? The term *tree* means we are considering the membership problem over *treesets*. The term *bottom-up* means that for each node a in a tree, the computation solves the problem by assigning states to the children of a before assigning states to a itself. This contrasts with *top-down* machines which assign states to a first and then the children of a . Visually, these terms make sense provided the root of the tree is at the top and branches of the tree move downward.

Acceptor is synonymous with *recognizer*. *Treeset* is synonymous with *tree language*.

A definitive reference for finite-state automata for trees is freely available online. It is “Tree Automata Techniques and Applications” (TATA) (Comon *et al.*, 2007). The presentation here differs from the one there, as mentioned below.

2.1.2 Definitions

We will use the following definition of trees.

Definition 7 (Trees). *We assume an alphabet Σ and symbols $[]$ not belonging to Σ .*

Base Cases: *For each $a \in \Sigma$, $a[]$ is a tree.*

Inductive Case: *If $a \in \Sigma$ and $t_1 t_2 \dots t_n$ is a string of trees of length n then $a[t_1 t_2 \dots t_n]$ is a tree.*

Let Σ^T denote the set of all trees of finite size using Σ . We also write $a[\lambda]$ for $a[]$.

It will be helpful to review the following concepts related to functions: *domain*, *co-domain*, *image*, and *pre-image*. A function $f : X \rightarrow Y$ is said to have domain X and co-domain Y . This means that if $f(x)$ is defined, we know $x \in X$ and $f(x) \in Y$. However, f may not be defined for all $x \in X$. Also, f may not be *onto* Y , there may be some elements in Y that are never “reached” by f .

This is where the concepts *image* and *pre-image* come into play. The image of f is the set $\{f(x) \in Y \mid x \in X, f(x) \text{ is defined}\}$. The pre-image of f is the set $\{x \in X \mid f(x) \text{ is defined}\}$. So the pre-image of f is the subset of the domain of f where f is defined. The image of f is the corresponding subset of the co-domain of f .

With this in place, we can define our first tree acceptor.

Definition 8 (DBFTA). A Deterministic Bottom-up Finite-state Acceptor (DBFTA) is a tuple (Q, Σ, F, δ) where

- Q is a finite set of states;
- Σ is a finite alphabet;
- $F \subseteq Q$ is a set of accepting (final) states; and
- $\delta : Q^* \times \Sigma \rightarrow Q$ is the transition function. The pre-image of δ must be finite. This means we can write it down—for example, as a list.

We use the transition function δ to define a new function $\delta^* : \Sigma^T \rightarrow Q$ as follows.

$$\begin{aligned} \delta^*(a[\lambda]) &= \delta(\lambda, a) \\ \delta^*(a[t_1 \cdots t_n]) &= \delta(\delta^*(t_1) \cdots \delta^*(t_n), a) \end{aligned} \tag{2.1}$$

There are some important consequences to the formulation of δ^* shown here. One is that δ^* is undefined on tree $a[t_1 \cdots t_n]$ if no transition $\delta(q_1 \cdots q_n, a)$ is defined.

(Also, I am abusing notation since δ^* is strictly speaking not the transitive closure of δ .)

Definition 9 (Treeset of a DBFTA). Consider some DBFTA $A = (Q, \Sigma, F, \delta)$ and tree $t \in \Sigma^T$. If $\delta^*(t)$ is defined and belongs to F then we say A accepts/recognizes t . Otherwise A rejects t . The treeset recognized by A is $L(A) = \{t \in \Sigma^T \mid \delta^*(t) \in F\}$.

The use of the ‘L’ denotes “Language” as treesets are traditionally referred to as *formal tree languages*.

We observe that since the pre-image of δ is finite, there is some n such that for all $(\vec{q}, a, q') \in \delta$, it is the case that $|\vec{q}| \leq n$. In other words, δ is effectively a function from $Q^n \times \Sigma$ to Q . We say δ is n -wide and any DBFTA with an n -wide δ is also called n -wide.

Definition 10 (n -wide Recognizable Treesets). A treeset is n -wide recognizable if there is a n -wide DBFTA that recognizes it.

Definition 11 (Recognizable Treesets). A treeset is recognizable if there exists $n \in \mathbb{N}$, such that an n -wide DBFTA that recognizes it.

2.1.3 Notes on Definitions

We have departed a bit from standard definitions. In particular, most introductions to tree automata make use of a particular kind of alphabet called a *ranked alphabet*. A ranked alphabet Σ_r is a finite alphabet Σ with an arity function $ar : \Sigma \rightarrow \mathbb{N}$. We write $\Sigma_r = (\Sigma, ar)$. The idea is that each symbol comes pre-equipped with a number which indicates how many children it has in trees. This is reasonable provided a node's label determines how many children it can have.

Strictly speaking, a ranked alphabet is not a necessary feature of tree automata. There are two substantive reasons to adopt it. First, using it helps ensure that the transition function is finite. (So it can accomplish the same thing as our requirement that the pre-image of δ be finite.). Second, it can help ensure our transition function is total; that is, defined for every element of the alphabet and the possible states of its children.

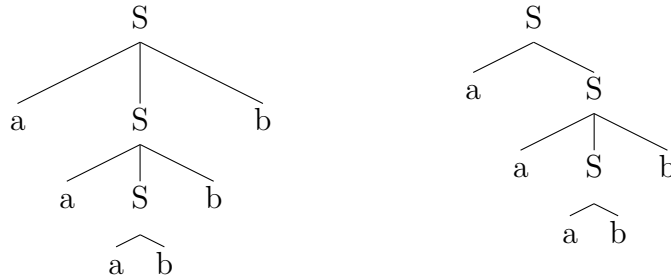
A ranked alphabet effectively defines a maximum width n for trees where $n = \max\{ar(a) \mid a \in \Sigma\}$. Thus a DBFTA defined with a ranked alphabet will always recognize a treeset which is a subset of $\Sigma^{T,n}$.

2.1.4 Examples

Example 1. Let $A = (Q, \Sigma, F, \delta)$ with its parts defined as follows.

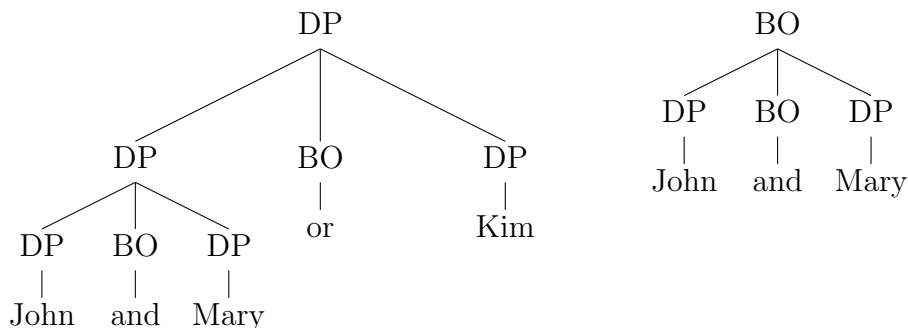
- $Q = \{q_a, q_b, q_S\}$
- $\Sigma = \{a, b, S\}$
- $F = \{q_S\}$
- $\delta(\lambda, a) = q_a$
- $\delta(\lambda, b) = q_b$
- $\delta(q_a q_b, S) = q_S$
- $\delta(q_a q_S q_b, S) = q_S$

Let us see how the acceptor A processes the two trees below as inputs.



Example 2. Let $A = (Q, \Sigma, F, \delta)$ with its parts defined as follows.

- $Q = \{q_{DP}, q_{BO}\}$
- $\Sigma = \{\text{and, or, Kim, John, Mary, DP, BO}\}$
- $F = \{q_{DP}\}$
- $\delta(q_{DP}, DP) = q_{DP}$
- $\delta(q_{BO}, BO) = q_{BO}$
- $\delta(q_{DP} q_{BO} q_{DP}, DP) = q_{DP}$
- $\delta(\lambda, \text{Kim}) = q_{DP}$
- $\delta(\lambda, \text{John}) = q_{DP}$
- $\delta(\lambda, \text{Mary}) = q_{DP}$
- $\delta(\lambda, \text{and}) = q_{BO}$
- $\delta(\lambda, \text{or}) = q_{BO}$



2.1.5 Observations

- For every symbol $a \in \Sigma$ which can be leaf in a tree, you will need to define a transition $\delta(\lambda, a)$.
- For every symbol $a \in \Sigma$ which can have n children, you will need to define a transition $\delta(q_1 \cdots q_n, a)$.

2.1.6 Connection to Context-Free Languages

Recognizable treesets are closely related to the derivation trees of context-free languages.

Theorem 4.

- Let G be a context-free word grammar, then the set of derivation trees of $L(G)$ is a recognizable tree language.
- Let L be a recognizable tree language then $\text{Yield}(L)$ is a context-free word language.
- There exists a recognizable tree language not equal to the set of derivation trees of any context-free language. Thus the class of derivation treesets of context-free word languages is a proper subset of the class of recognizable treesets.

2.2 Deterministic Top-down Finite-state Tree Acceptors

2.2.1 Orientation

This section is about deterministic top-down finite-state tree acceptors. The term *finite-state* means that the memory is bounded by a constant, no matter the size of the input to the machine. The term *deterministic* means there is single course of action the machine follows to compute its output. The term *acceptor* means this machine solves *membership problem*: given a set of objects X and input object x , does x belong to X ? The term *tree* means we are considering the membership problem over *treesets*. The term *top-down* means that for each node a in a tree, the computation solves the problem by assigning a state to the parent

of a before assigning a state to a itself. This contrasts with *bottom-up* machines which assign states to the children of a first and then a . Visually, these terms make sense provided the root of the tree is at the top and branches of the tree move downward.

Acceptor is synonymous with *recognizer*. *Treeset* is synonymous with *tree language*.

A definitive reference for finite-state automata for trees is freely available online. It is “Tree Automata Techniques and Applications” (TATA) (Comon *et al.*, 2007). The presentation here differs from the one there, as mentioned below.

2.2.2 Definition

Definition 12 (DTFTA). *A Deterministic Top-down Finite-state Acceptor (DTFTA) is a tuple (Q, Σ, F, δ) where*

- Q is a finite set of states;
- q_0 is a initial state;
- Σ is a finite alphabet;
- $\delta : Q \times \Sigma \times \mathbb{N} \rightarrow Q^*$ is the transition function. Note the pre-image of δ is necessarily finite.

The transition function takes a state, a letter, and a number n and returns a string of states. The idea is that the length of this output string should be n . Basically, when moving top-down, the states of the child sub-trees depend on these three things: the state of the parent, the label of the parent, and the number of children the parent has.

We use the transition function δ to define a new function $\delta^* : Q \times \Sigma^T \rightarrow Q^*$ as follows.

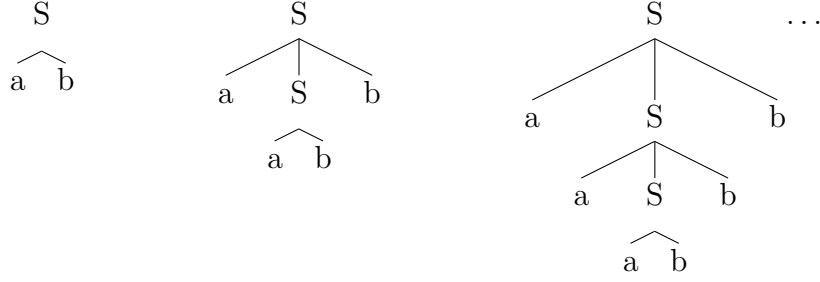
$$\begin{aligned} \delta^*(q, a[\lambda]) &= \delta(q, a, 0) \\ \delta^*(q, a[t_1 \cdots t_n]) &= \delta^*(q_1, t_1) \cdots \delta^*(q_n, t_n) \text{ where } \delta(q, a, n) = q_1 \cdots q_n \end{aligned} \quad (2.2)$$

As before, there are some important consequences to the formulation of δ^* . One is that δ^* is undefined on tree $a[t_1 \cdots t_n]$ if transition $\delta(q, a, n)$ does not return a string from Q^* of length n . (Also, I am abusing notation since δ^* is strictly speaking not the transitive closure of δ .)

Definition 13 (Treeset of a DTFTA). *Consider some DTFTA $A = (Q, \Sigma, F, \delta)$ and tree $t \in \Sigma^T$. If $\delta^*(q_0, t)$ is defined and equals λ then we say A accepts/recognizes t . Otherwise A rejects t . Formally, the treeset recognized by A is $L(A) = \{t \in \Sigma^T \mid \delta^*(t) = \lambda\}$.*

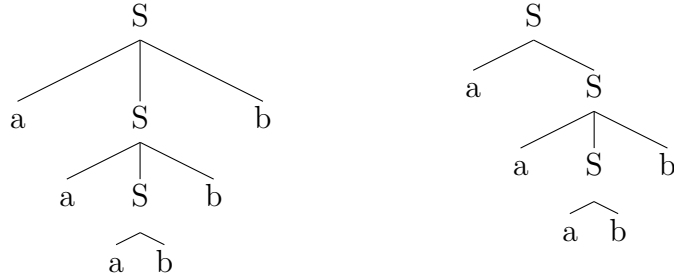
The use of the ‘L’ denotes “Language” as treesets are traditionally referred to as *formal tree languages*.

Example 3. Recall the example from last week which generates trees like



- $Q = \{q_a, q_b, q_S\}$
- $\Sigma = \{a, b, S\}$
- $q_0 = q_S$
- $\delta(q_a, a, 0) = \lambda$
- $\delta(q_b, b, 0) = \lambda$
- $\delta(q_S, S, 3) = q_a q_S q_b$
- $\delta(q_S, S, 2) = q_a q_b$

Let us see how the acceptor A processes the two trees below as inputs.



Theorem 5. *Every treeset recognizable by a DTFTA is recognizable, but there are recognizable treesets which cannot be recognized by a DTFTA.*

The following example helps show why this is the case. Consider the treeset T containing only the two trees shown below.



This is a recognisable treeset because the DBFTA below recognizes exactly these two trees and no others.

- $Q = \{q, q_S\}$
- $\Sigma = \{a, b, S\}$
- $q_0 = q_S$
- $\delta(\lambda, a) = q_a$
- $\delta(\lambda, b) = q_b$
- $\delta(q_a q_b, S) = q_S$
- $\delta(q_b q_a, S) = q_S$

Notice that this DBFTA fails on these two trees.



A DTFTA cannot recognize the trees in T without also recognizing the trees shown immediately above. This is because moving top down there can only be one value for $\delta(q_S, S, 2)$. Suppose it equals $q_1 q_2$. To recognize the first tree, we would also have to make sure that $\delta(q_1, a, 0)$ and $\delta(q_2, b, 0)$ are defined. Similarly, to recognize the second tree, we would have to make sure that $\delta(q_1, b, 0)$ and $\delta(q_2, a, 0)$ are defined. But it follows then that the aforementioned trees above are also recognized by this DTFTA. For instance the tree with two a leaves is recognized because both $\delta(q_1, a, 0)$ and $\delta(q_2, a, 0)$ are defined. Thus no DTFTA recognizes T .

2.2.3 Observations

- For every symbol $a \in \Sigma$ which can be leaf in a tree, you will need to define a transition $\delta(q, a, 0) = \lambda$.
- For every symbol $a \in \Sigma$ which can have n children, you will need to define a transition $\delta(q, a, n) = q_1 \cdots q_n$.

2.3 Properties of recognizable tree languages

Theorem 6 (Closure under Boolean operations). *The class of recognizable tree languages are closed under union and intersection.*

The class of n -wide recognizable tree languages are closed under union, intersection, and complementation with respect to the $\Sigma^{T,n}$.

The proofs of these cases are very similar to the ones for finite-state acceptors over strings. For every DBFTA A recognizing a treeset T , it can be made *complete* by adding a sink state and transitions to it. Then product constructions can be used to establish closure under union and intersection. Closure under complement is established the same as before too: everything is the same except the final states are now the non-final states of A .

Theorem 7 (Minimal, deterministic, canonical form). *For every recognizable tree language T , there is a unique, smallest DBFTA A which recognizes T . That is, if DBFTA A' also recognizes T then there are at least as many states in A' as there are in A .*

2.4 Connection to Context-Free Languages

Context Free Grammars (CFGs) are studied in detail in a number of textbooks including Harrison (1978); Davis and Weyuker (1983); Hopcroft *et al.* (2001) and Sipser (1997).

Definition 14. *A rewrite grammar is a tuple $\langle T, N, S, \mathcal{R} \rangle$ where*

- \mathcal{T} is a nonempty finite alphabet of symbols. These symbols are also called the terminal symbols, and we usually write them with lowercase letters like a, b, c, \dots

- \mathcal{N} is a nonempty finite set of non-terminal symbols, which are distinct from elements of \mathcal{T} . These symbols are also called *category symbols*, and we usually write them with uppercase letters like A, B, C, \dots
- S is the start category, which is an element of \mathcal{N} .
- A finite set of production rules \mathcal{R} . A production rule has the form $\alpha \rightarrow \beta$ where α, β belong to $(\mathcal{T} \cup \mathcal{N})^*$. In other words, α and β are strings of non-terminal and terminal symbols. While β may be the empty string we require that α include at least one symbol.

Rewrite grammars are also called *phrase structure grammars*.

Definition 15. A CFG is a rewrite grammar with the following properties.

- For all rules $\alpha \rightarrow \beta$, α is an element of \mathcal{N} . So the left-hand-side of each rule is a single non-terminal.
- For all rules, $\alpha \rightarrow \beta$, β is an element of $(\mathcal{N} \cup \mathcal{T})^*$. So the right-hand-side of each rule is a sequence of non-terminal symbols or a single terminal symbol.

Example 4. Define a CFG G_{NE} as follows. Let $\mathcal{N} = \{A, B, S\}$ and $\mathcal{T} = \{a, b\}$. Let \mathcal{R} include the rules $S \rightarrow a S b$ and $S \rightarrow a b$.

Next we define the languages of the rewrite grammars in addition to the derivation treeset of context free grammars.

The language of a rewrite grammar is defined recursively below.

Definition 16. The (partial) derivations of a rewrite grammar $G = \langle \mathcal{T}, \mathcal{N}, S, \mathcal{R} \rangle$ is written $D(G)$ and is defined recursively as follows.

1. The base case: S belongs to $D(G)$.
2. The recursive case: For all $\alpha \rightarrow \beta \in \mathcal{R}$ and for all $\gamma_1, \gamma_2 \in (\mathcal{T} \cup \mathcal{N})^*$, if $\gamma_1 \alpha \gamma_2 \in D(G)$ then $\gamma_1 \beta \gamma_2 \in D(G)$.
3. Nothing else is in $D(G)$.

Then the language of the grammar $L(G) = \{w \in \mathcal{T}^* \mid w \in D(G)\}$.

Exercise 4. Using the definition above, explain why $aaabbb$ belongs to $L(G_{NE})$.

The derivation treeset of a context free grammar is defined recursively below. The *derivation treeset* of a CFG $G = \langle \mathcal{T}, \mathcal{N}, S, \mathcal{R} \rangle$ is written $D_T(G)$.

It is defined as all and only those trees t such that

1. $\text{yield}(t) \in L(G)$ and
2. for all non-leaf nodes $a[a_1[ts_1]a_2[ts_2] \cdots a_n[ts_n]]$ in t , the rule $a \rightarrow a_1 a_2 \cdots a_n$ belongs to \mathcal{R} .

Recognizable treesets are closely related to the derivation trees of context-free languages.

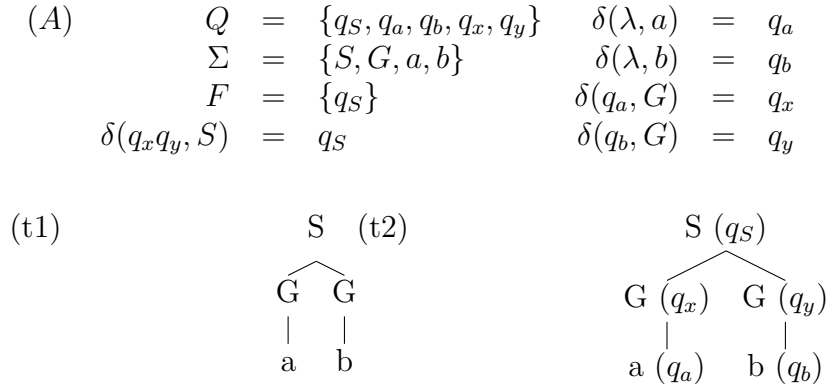
Theorem 8.

- Let G be a context-free word grammar, then the set of derivation trees $D_T(G)$ is a recognizable tree language.
- There exists a recognizable tree language not equal to the set of derivation trees of any context-free language. Thus the class of derivation treesets of context-free word languages is a proper subset of the class of recognizable treesets.
- Let L be a recognizable tree language then $\mathbf{yield}(L)$ is a context-free word language.

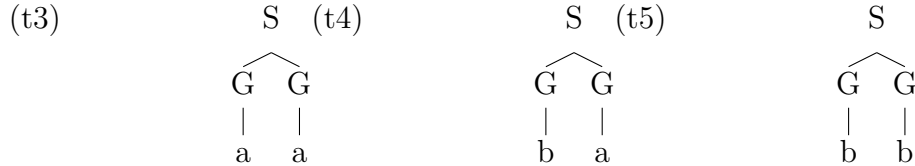
Each of these has a straightforward explanation.

For (1), the recognizable tree language which recognizes $D_T(G)$ for a CFG G can be constructed based on the rules of G . For each symbol a in N , the DBFTA should include $\delta(\lambda, a) = q_a$. And, for each rule $A \rightarrow B_1 \cdots B_n$ in R , the DBFTA should include $\delta(q_{B_1} \cdots q_{B_n}, A) = q_A$. That's it!

For (2), consider the DBFTA A shown below. The claim is that there is no CFG whose derivation language is exactly this recognizable treeset. The only tree in $L(A)$ is (t1), which is shown below A at left. Tree (t2) shows (t1) with the states A assigns to its subtrees.



That no CFG can recognize this treeset follows from the fact that such any CFG G which includes the tree above will need to have the following rules: $S \rightarrow GG$, $G \rightarrow a$, $G \rightarrow b$. But then this G will not only generate (t1) above but also the derivation trees shown below.



Thus $L(A) \neq D_T(G)$.

This example shows that the states of the DBFTA are more abstract than the labels on the nodes. The reason recognizable tree languages are more expressive than the derivation treesets of CFGs follows from this. The DBFTA uses states q_x and q_y to distinguish the subtrees bearing the label G . But the CFG cannot distinguish these trees in this way.

For (3), observe that we can write a CFG that generates the same stringset as the one above. For instance, we could write a CFG with the rule $S \rightarrow ab$.

More generally, though, we can always write a CFG that puts the state information into the nodes themselves. The trees in the derivation treeset for this CFG would be “structurally the same” as the trees in the recognizable treeset, but the labels on the nodes would be different. So they are not the same trees. In the example above for instance we can write a CFG with rules $S \rightarrow G_x G_y, G_x \rightarrow a, G_x \rightarrow b$. There is only one tree in this CFG’s derivation treeset shown below.



Importantly, (t6) is not the same as (t1). The inner nodes are labeled differently! So they are different trees. However, they only differ with respect to how the inner nodes are labeled, so it follows that the stringsets obtained by taking the **yield** of these trees are the same. This is the kind of argument used to show that the yield of any recognizable treeset is a context-free language.

Bibliography

- Comon, H., M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. 2007. Tree automata techniques and applications. Available on: <http://tata.gforge.inria.fr/>. Release October, 12th 2007.
- Davis, Martin D., and Elaine J. Weyuker. 1983. *Computability, Complexity and Languages*. Academic Press.
- Harrison, Michael A. 1978. *Introduction to Formal Language Theory*. Addison-Wesley Publishing Company.
- Hopcroft, John, Rajeev Motwani, and Jeffrey Ullman. 2001. *Introduction to Automata Theory, Languages, and Computation*. Boston, MA: Addison-Wesley.
- Sipser, Michael. 1997. *Introduction to the Theory of Computation*. PWS Publishing Company.