# Course 2BA1: Hilary Term 2003
# Section 5: Formal Languages

David R. Wilkins

Copyright c David R. Wilkins 2001{03

# Contents

.and Wor7.559 -39.073 Tc
$f0;1g$ of binary digits, or the set $f0;1;2;3;4;5;6;7;8;9g$ of decimal digits.)

For any natural number $n$, we de ne a *word* of lend5F15 5,9.539.955 Tf 127.528 0 Td[(n)]

Note that $|w_1 \cdot w_2| = |w_1| + |w_2|$ for all words $w_1$ and $w_2$ over some alphabet.

The operation $\cdot$ of concatenation on the set of words over some alphabet is not commutative if that alphabet has more than one element. Indeed if $a$ and $b$ are distinct elements of this alphabet then $a \cdot b$ is the string $ab$, and $b \cdot a$ is the string $ba$, and therefore $a \cdot b \neq b \cdot a$.

Let $w_1$, $w_2$ and $w_3$

$\langle T \rangle$ represents the definite article `the' (which will subsequently replace it);

$\langle N \rangle$ represents a noun chosen from the set the set $\{dog, cat\}$;

$\langle Adj \rangle$ represents an adjective chosen from the set $\{black, old, small\}$;

$\langle V \rangle$ represents a verb chosen from the set $\{saw, chased\}$;

Each of these productions specifies that the entity on the left hand side of the arrow may be replaced by the string on the right hand side of the arrow. We may apply these productions successively, of

are applied, one at a time, to transform strings made up of terminals and

Let $p$ and $q$ be Boolean variables. The *conjunction $p \wedge q$* of $p$ and $q$ is true if and only if both $p$ and $q$ are true. The *disjunction $p \vee q$* of $p$ and $q$

| $p$ | $q$ | $r$ | $p \wedge q$ | $q \_ r$ | $(p \wedge q) \_ r$ | $p \wedge (q \_ r)$ |
|-----|-----|-----|--------------|----------|---------------------|---------------------|
|     |     |     |              |          |                     |                     |

determine in the usual fashion the order in which the binary operations are to be performed and the subformulae to which they are to be applied. We could introduce characters **T** and **F** to denote the Boolean constants `true' and `false' respectively. It remains to consider how Boolean variables are to be represented. We could certainly use single letters $p$, $q$, $r$, $s$. But this would only enable us to write down formulae with at most four distinct propositional variables. Were we to use single letters from the English alphabet in both upper and lower case to denote propositional variables, this would restrict us to formulae with at most fty-two distinct Boolean variables. But there should be no limit to the number of distinct Boolean variables that we could introduce into a well-formed formula. We therefore need a scheme for representing unlimited quantities of Boolean variables within our formula. We would do 32(therr)-g1(w)6,

Our grammar now has productions to generate any atomic formula. We

ambiguity in the cases when $G$ is atomic, the negation of a well-formed formula or a conjunction of well-formed formulae. (Our rules of precedence ensure that $: G$

We can then apply further productions in order to obtain formulae such as $p \wedge \neg q$ and $p' \wedge p'' \wedge p'''$.

Finally we have to specify the productions for handling disjunction. These are analogous to those for conjunctions, and are the following: |

$$\langle\text{disjunction}\rangle \quad \rightarrow \quad \langle\text{df}\rangle \vee \langle\text{df}\rangle$$
$$\langle\text{df}\rangle \quad \rightarrow \quad \langle\text{atom}\rangle$$
$$\langle\text{df}\rangle \quad \rightarrow \quad \langle\text{negation}\rangle$$
$$\langle\text{df}\rangle \quad \rightarrow \quad \langle\text{disjunction}\rangle$$
$$\langle\text{df}\rangle \quad \rightarrow \quad (\langle \vee \rangle$$

*h*conjunction*i ! h*

$$\langle variable \rangle \ \rightarrow \ \langle letter \rangle$$
$$\langle letter \rangle \ \rightarrow \ p$$
$$\langle letter \rangle \ \rightarrow \ q$$
$$\langle letter \rangle \ \rightarrow \ r$$
$$\langle letter \rangle \ \rightarrow \ s$$

The well-formed formulae of the Propositional Calculus are those words $F$ over the alphabet

) $(p \wedge r') \_ (^h\text{conjunction}\, i)$

) $(p \wedge r') \_ (^h\text{cf}\, i \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((^h\text{w}\; i) \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((^h\text{compound}\, i) \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((^h\text{negation}\, i) \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: {}^h\text{nf}\, i) \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: {}^h\text{atom}\, i) \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: {}^h\text{variable}\, i) \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: {}^h\text{variable}\, i') \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: {}^h\text{variable}\, i''') \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: {}^h\text{letter}\, i''') \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge {}^h\text{conjunction}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge {}^h\text{cf}\, i \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge {}^h\text{atom}\, i \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge {}^h\text{variable}\, i \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge {}^h\text{letter}\, i \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge {}^h\text{cf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge {}^h\text{negation}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{nf}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{atom}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{variable}\, i)$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{variable}\, i')$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{variable}\, i'')$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{variable}\, i''')$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : {}^h\text{letter}\, i''')$

) $(p \wedge r') \_ ((: r''') \wedge q \wedge : r''')$

**Remark** ɪThe grammar we have constructed to describe the well-formed for-

parentheses, and that no subformula is enclosed by itself within two or more sets of parentheses. This modified grammar is expressed in Backus-Naur form as follows:

$$
\begin{aligned}
\langle\text{wff}\rangle &\to \langle\text{atom}\rangle \mid \langle\text{compound}\rangle \\
\langle\text{atom}\rangle &\to \langle\text{constant}\rangle \mid \langle\text{variable}\rangle \\
\langle\text{compound}\rangle &\to \langle\text{negation}\rangle \mid \langle\text{conjunction}\rangle \mid \langle\text{disjunction}\rangle \\
\langle\text{negation}\rangle &\to \lnot\langle\text{nf}\rangle \\
\langle\text{nf}\rangle &\to \langle\text{atom}\rangle \mid \langle\text{negation}\rangle \mid (\langle\text{compound}\rangle) \\
\langle\text{conjunction}\rangle &\to \langle\text{cf}\rangle \land \langle\text{cf}\rangle \\
\langle\text{cf}\rangle &\to \langle\text{atom}\rangle \mid \langle\text{negation}\rangle \mid \langle\text{conjunction}\rangle \mid (\langle\text{compound}\rangle) \\
\langle\text{disjunction}\rangle &\to \langle\text{df}\rangle \lor \langle\text{df}\rangle \\
\langle\text{df}\rangle &\to \langle\text{atom}\rangle \mid \langle\text{negation}\rangle \mid \langle\text{disjunction}\rangle \mid (\langle\text{compound}\rangle) \\
\langle\text{constant}\rangle &\to T \mid F \\
\langle\text{variable}\rangle &\to \langle\text{letter}\rangle \mid \langle\text{variable}\rangle' \\
\langle\text{letter}\rangle &\to p \mid q \mid r \mid s
\end{aligned}
$$

## 5.4 Context-Free Grammars

We have discussed examples of context-free grammars. We now present and discuss a formac definition of such grammars.

**Definition** A *context-free grammar* $(V, A, \langle S\rangle, P)$ consists of a finite set $V$, a subset $A$ of $V$, an element $\langle S\rangle$ of $V \setminus A$, and a finite subset $P \subseteq (V \setminus A) \times V^*$.

Let $(V, A, \langle S\rangle, P)$ be a context-free grammar. The elements of $A$ are referred to as *terminals*. Let $N = V \setminus A$. The elements of $N$ are referred to as *nonterminals*. The nonterminal $\langle S\rangle$ is the *start symbol*. The set $N$ of nonterminals is non-empty since $\langle S\rangle \in N$.

The finite set $P$ specifies the *productions*

production $\langle T_i \rangle ! w$ of the grammar such that $w$

As in the case of context-free grammars, the elements of $A$ are referred to as *terminals*, the elements of $V \setminus A$ are referred to as *nonterminals*, the nonterminal $\langle S \rangle$ is the *start symbol* and the elements of $P$ specify the *productions* of the grammar. The pro speci ed bp a1 element ( $r; w$) of $P$ is denoted bp $r \to w$. wever the left hand side $r$ of a pro $r \to w$ in a phrase strF15ure grammar need not consist solely of a single nonterminal, but map be a nite string $r$ of elements of $V$, provided that this string $r$ contains at least one nonterminal. (Note that $V$ denotes the set of all nite words over the alphabet $V$ whose elements are terminals and nonterminals, $A$ denotes the set of all nite words consisting entirely of terminals, and thus $V \setminus A$ denotes the set of all nite words belonging to $V$ which contain

atterminal.) Det mittantl $w^{00}$ $r5\ w$ wthe

**Example** Let $L$

where *⟨A⟩* and *⟨B⟩* represent nonterminals, *b* represents a terminal, and *"*
denotes the empty word. A regular grammar is said to be in *normal form* if

$hBi$ ! $2hBi$

$\vdots$

$hBi$ ! $9hBi$

$hBi$ ! "

**Definition** Let $(S, A, i, t, F)$ be a finite state acceptor. A language $L$ over the alphabet $A$ is said to be *recognized* or *accepted* by the finite state acceptor if $L$