





Enduring impact of Chomsky's work on theoretical computer science. 15 pages devoted to the Chomsky hierarchy.



No mention of the Chomsky hierarchy. The same with the latest edition of Hopcroft, Motwani, and Ullman (2006).

Semi-Thue System/String Rewriting System					
Production		$\alpha \to \beta \qquad \qquad \alpha, \beta \in (N \cup \Sigma)^*$			
One-s	tep rewriting	$\gamma\alpha\delta \Rightarrow \gamma\beta\delta$			
Langu	age	$L(G) = \{ w \in \Sigma^* \mid S \Rightarrow^* w \}$			
	Туре 0	unrestricted			
	Type 1	$ \alpha \leq \beta $			
	Type 2	$\alpha \in N$			
	Туре З	$\alpha \in N, \beta \in \Sigma N$			

N: finite alphabet of nonterminal symbols

 Σ : finite alphabet of terminal symbols



LBA: linear-bounded automata = NSPACE(n) Turing machines PDA: pushdown automata FA: finite automata



Context-sensitive languages are computationally very complex. Some context-sensitive languages are PSPACE-complete.

- Context-sensitive languages are computationally very complex.
- There is a **huge** gap between context-free and contextsensitive.
- Many intermediate classes of languages have been considered in formal language theory.
- In particular, "mildly context-sensitive" grammars have been investigated by computational linguists.

Indexed Grammars

Hopcroft and Ullman 1979

14.3 INDEXED LANGUAGES

Of the many generalizations of context-free grammars that have been proposed, a class called "indexed" appears the most natural, in that it arises in a wide variety of contexts. We give a grammar definition here. Other definitions of the indexed languages are cited in the bibliographic notes.

An indexed grammar is a 5-tuple (V, T, I, P, S), where V is the set of variables, T the set of terminals, I the set of *indices*, S in V is the start symbol, and P is a finite set of productions of the forms

1) $A \rightarrow \alpha$, 2) $A \rightarrow Bf$, or 3) $Af \rightarrow \alpha$, where A and B are in V, f is in I, and α is in $(V \cup T)^*$. **Example 14.2** Let $G = (\{S, T, A, B, C\}, \{a, b, c\}, \{f, g\}, P, S)$, where P consists of $S \rightarrow Tg$, $Af \rightarrow aA$, $Ag \rightarrow a$, $T \rightarrow T f_{2}$ $Bf \rightarrow bB$, $Bg \rightarrow b$, $T \rightarrow ABC$, $Cf \rightarrow cC$, $Cg \rightarrow c$. An example derivation in this indexed grammar is $S \Rightarrow Tg \Rightarrow Tfg \Rightarrow AfgBfgCfg$ $\Rightarrow aAgBfgCfg \Rightarrow aaBfgCfg \Rightarrow aabBgCfg$ \Rightarrow aabbCfg \Rightarrow aabbcCg \Rightarrow aabbcc. In general, $S \stackrel{*}{\Rightarrow} Tf^{i}g \Rightarrow Af^{l}gBf^{l}gCf^{i}g \stackrel{*}{\Rightarrow} a^{i+1}b^{i+1}c^{i+1}.$ **Theorem 14.7** (a) If L is accepted by a one-way nondeterministic stack automaton, then L is an indexed language. (b) If L is an indexed language, then L is a context-sensitive language.



What type of rewriting system is an indexed grammar??

Туре 0		
Production	$\alpha \rightarrow \beta$	
One-step rewriting	$\gamma \alpha \delta \Rightarrow \gamma \beta \delta$	
Indexed		
Production	$T \rightarrow ABC$	
One-step rewriting	$Tfg \Rightarrow AfgBfgCfg$	
An indexed grammar is not an instance of a type 0 grammar.		



Indexed languages are also sort of complex.

	Conder 1000
GERALD GAZDAR	Gazuar 1966
UNIVERSITY OF SUSSEX	
SCHOOL OF SOCIAL SCIENCES	
BRIGHTON	
APPLICABILITY	
OF INDEXED GRAMMARS	
TO NATURAL LANGUAGES	
In the standard formulations, an indexed grammar can contain rules of three different sorts:	
(1) i. $A[] \rightarrow W[]$ ii. $A[] \rightarrow B[i,]$ iii. $A[i,] \rightarrow W[]$	
Suppose we allow the rule types shown in (2) in addition to those listed in (1):	
(2) i. $A[] \rightarrow W_1[]B[]W_2[]$ ii. $A[] \rightarrow W_1[]B[i]W_2[]$ iii. $A[i] \rightarrow W_1[]B[i]W_2[]$	

Vijayashanker 1987

The productions in an Indexed Grammar can be written as (refer [Gazdar 85a], the original definition of Indexed Grammars appears in [Aho 68]) $X[\cdot\cdot] \rightarrow X_1[i, \cdot\cdot] \dots X_n[i, \cdot\cdot]$ (push) or $X[i, \cdot\cdot] \rightarrow X_1[\cdot\cdot] \dots X_n[\cdot]$ (pop) or $X[\cdot\cdot] \rightarrow a$ In an LIG the productions have the form $X[\cdot\cdot] \rightarrow X_1[\dots X_j[i, \cdot\cdot] \dots X_n[]$ (push) or $X[i, \cdot\cdot] \rightarrow X_1[] \dots X_j[\cdot\cdot] \dots X_n[]$ (pop) or $X[i, \cdot\cdot] \rightarrow a$ *LIG* differs from IG in that the stack is passed on to only one of the children.



recursively enumerable

context-

sensitive

linear

indexed

type 0

type 1

type 1.5?

type 2 type 3 TAG: tree-adjoining grammar; LIG: linear indexed grammar; HG: head grammar

There is a book where the author calls the linear indexed grammars "type 1.5".

Туре 0				
Production	$\alpha ightarrow \beta$			
One-step rewriting	$\gamma \alpha \delta \Rightarrow \gamma \beta \delta$			
Linear indexed				
Production	$T[] \to AB[]C$			
One-step rewriting	$T[fg] \Rightarrow AB[fg]C$			
A linear indexed grammar is not an instance of a type 0 grammar.				



There are actually almost no other grammar formalisms that are instances of Chomsky's type 0 grammars.

Thue→Post→Chomsky

THE JOURNAL OF SYMBOLIC LOGIC Volume 12, Number 1, March 1947

RECURSIVE UNSOLVABILITY OF A PROBLEM OF THUE

EMIL L. POST

Chomsky 1965

Confusion over this matter has been sufficiently persistent to suggest that a terminological change might be in order. Nevertheless, I think that the term "generative grammar" is completely appropriate, and have therefore continued to use it. The term "generate" is familiar in the sense intended here in logic, particularly in Post's theory of combinatorial systems. Furthermore, "generate" seems to be the most appropriate translation for Humboldt's term *erzeugen*, which he frequently uses, it seems, in essentially the sense here intended. Since this use of the term "generate" is well established both in logic and in the tradition of linguistic theory, I can see no reason for a revision of terminology. How did Chomsky come to use semi-Thue systems as the most general type of grammars? I believe Post (1947) coined the term "semi-Thue system".

Post Canonical Systems

FORMAL REDUCTIONS OF THE GENERAL COMBINATORIAL DECISION PROBLEM.*

By EMIL L. Post.

 $\label{eq:theta} The operations of the system are a specified finite set of productions, each of the following form:$

 $g_1P_{i_1}g_2P_{i_2}\cdot\cdot\cdot g_mP_{i_m}g_{m+1}.$

An indexed grammar is an instance of a Post canonical system.

Why Didn't Chomsky Use Post Canonical Systems?

Post 1947

A very special case of the canonical form is what we term the normal form. A system in canonical form will be said to be in *normal form* if it has but one primitive assertion, and, each of its productions is in the form

> g P produces P g'.

The main purpose of the present paper is to demonstrate that every system in canonical form can formally be reduced to a system in normal form. The two forms may therefore in fact be said to be *equipotent*. Perhaps because of this striking result?



Should Context-Free Grammars Be Thought of as String Rewriting Systems?



It is silly to define well-formed sentences in terms of derivations if what we really want is to assign tree structures to sentences.

Bottom-up Interpretation of Context-**Free Rules**

 $B_i \Rightarrow_G^* v_i \ (i=1,\ldots,n) \quad A \rightarrow w_0 \ B_1 \ w_1 \ \ldots \ B_n \ w_n \in P$ $A \Rightarrow_G{}^* W_0 V_1 W_1 \dots V_n W_n$

 $L(G) = \{ w \in \Sigma^* \mid S \Rightarrow_G^* w \}$

All we need is the subrelation of \Rightarrow_{G}^{*}

restricted to N $\times \Sigma^*$. Just a general type of inductive definition.

Inductive Definition = CFG Interpreted Bottom-up

By ME_a is understood the set of meaningful expressions of type a; this notion has the following recursive definition:

(1) Every variable and constant of type a is in ME_a.

(2) If $\alpha \in ME_a$ and u is a variable of type b, then $\lambda u\alpha \in ME_{(b,a)}$.

(3) If $\alpha \in ME_{\langle a,b \rangle}$ and $\beta \in ME_a$, then $\alpha(\beta) \in ME_b$. (4) If $\alpha, \beta \in ME_a$, then $\alpha = \beta \in ME_t$.

- (5) If $\phi, \psi \in ME_t$ and u is a variable, then $\neg \phi$, $[\phi \land \psi]$, $[\phi \lor \psi]$, $[\phi \rightarrow \psi], [\phi \leftrightarrow \psi], \forall u\phi, \land u\phi, \Box \phi, W\phi, H\phi \in \mathsf{ME}_t.$
- (6) If $\alpha \in ME_a$, then $[\alpha] \in ME_{\langle s,a \rangle}$.

CFGs as Logic Programs on Strings

 $A \rightarrow w_0 B_1 w_1 \dots B_n w_n$

 $A(w_0 \mathbf{x}_1 w_1 \dots \mathbf{x}_n w_n) \leftarrow B_1(\mathbf{x}_1), \dots, B_n(\mathbf{x}_n)$

Horn clause

 $L(G) = \{ w \in \Sigma^* \mid G \vdash S(w) \}$



It's a small step to use n-ary predicates for $n \ge 2$.



It's best to think of an MCFG as a kind of logic program. Each rule is a definite clause. Nonterminals are predicates on strings.











Elementary Formal Systems

Smullyan 1961

"Elementary Formal Systems can be looked at as variants of the canonical languages of Post The reason for our choice of elementary formal systems, in lieu of Post's canonical systems, is that their structure is more simply described, and their techniques are more easily applied. The general notion of "production" used in Post systems is replaced by the simpler logistic rules of substitution and detachment (modus ponens), which are more easily formalized."

Chomsky HierarchyRewriting
SystemsMachinesLogic Programs on StringsLanguagesType 0TuringElementary Formal Systems
(Smullyan 1961)r.e.TuringLength-Bounded EFS (Arikawa
CSL =CSL =

Type 1	LBA	et al. 1989)	NSPACE(n)
	Poly-time Turing	Simple LMG (Groenink 1997) / Hereditary EFS (Ikeda and Arimura 1997)	Р
		MCFG	MCFL
Type 2	PDA	Simple EFS (Arikawa 1970)	CFL
Туре З	FA		Reg



Chomsky on Mathematical Linguistics (1979)

... what happened is that there was a period of fairly fruitful contact between automata theory and linguistics in

the late fifties and the early sixties. It mostly had to do with the properties of context-free grammars. Things turned up which are quite interesting.

... Certainly context-free grammars represent some of the properties of languages. This seems to me what one would expect from applied mathematics, to see if you can find systems that capture some of the properties of the complex system that you are working with, and to ask whether those systems have any intrinsic mathematical interest, and whether they are worth studying in abstraction. And that has happened at exactly one level, the level of context-free grammar.

Chomsky (2004): Generative Enterprise Revisited

A long quote, from an interview held in the year when Hopcroft and Ullman's textbook was published.

Chomsky on Mathematical Linguistics (1979)

... At any other level it has not hap-

pened. The systems that capture other properties of language, for example transformational grammar, hold no interest for mathematics. But I do not think that that is a necessary truth. It could turn out that there would be richer or more appropriate mathematical ideas that would capture other, maybe deeper properties of language than context-free grammars do. In that case you have another branch of applied mathematics which might have linguistic consequences. That could be exciting.

Chomsky (2004): Generative Enterprise Revisited

Summary

- The four levels of the Chomsky hierarchy are not of equal importance.
- The choice of semi-Thue systems as the most general form of grammar hampered investigation of some important language classes.
- Smullyan's elementary formal system provides the right level of generality and simplicity.
- Bottom-up semantics of rules is basic.
- Natural generalization of "inductive definitions".

What is the Significance of Grammar Formalisms?

- Abstract conception of grammars.
- Makes possible investigation of algorithmic properties of grammars.
- Should have good formal properties (like familiar closure properties).
- The choice of a grammar formalism should **not** be thought of as a tight characterization of possible human grammars.
- A grammar formalism does not have explanatory power.