

07.NLP-PARSING

FIRDAUS SOLIHIN
UNIVERSITAS TRUNOJOYO

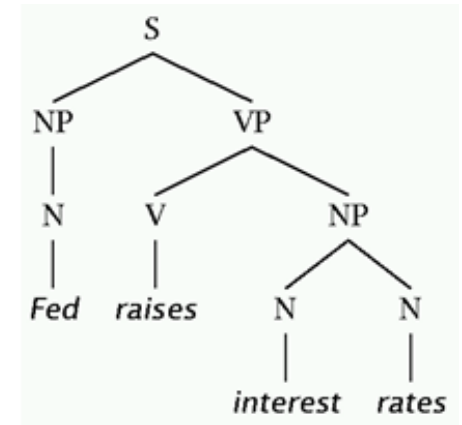


LINGUISTIC STRUCTURE

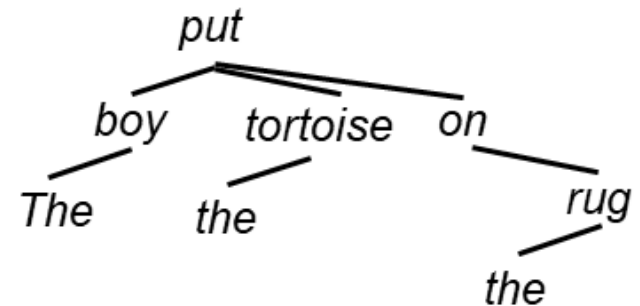
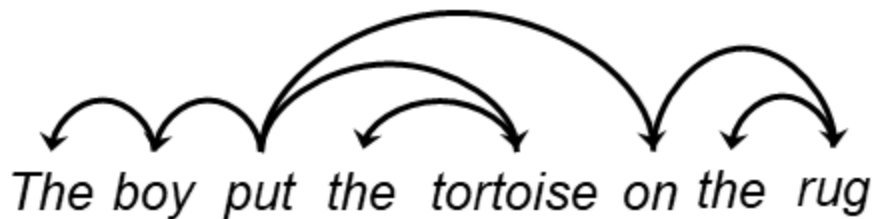


Struktur bahasa memiliki dua representasi

1. Constituency (phrase structure)



2. Dependency structure



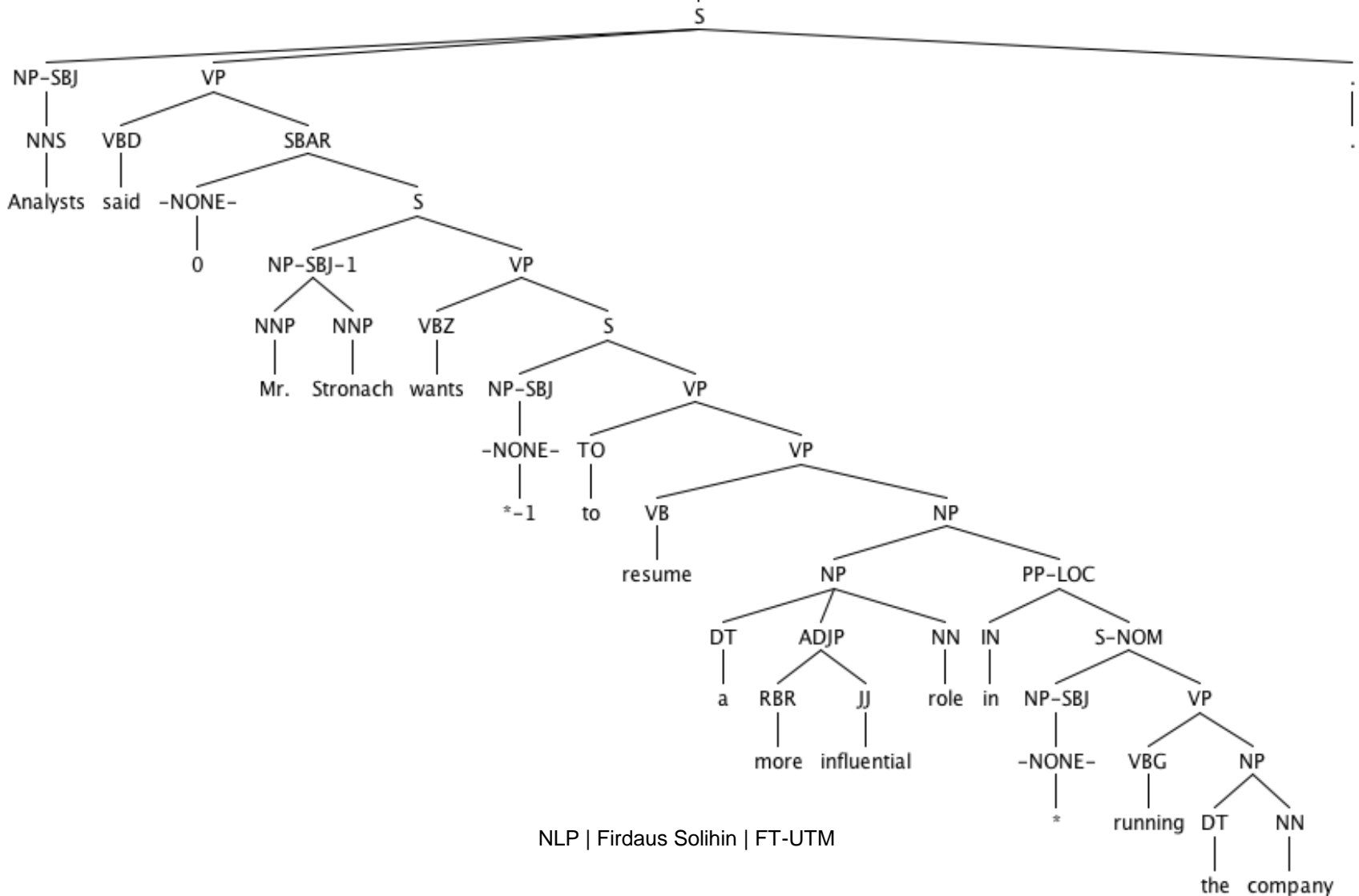
1. CONSTITUENCY (PHRASE STRUCTURE)

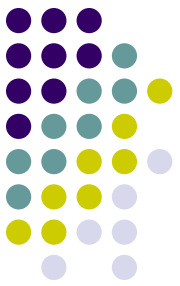


- phrase structure mengatur kata dalam urutan constituency yang bertingkat
- Bagaimana cara mengidentifikasi constituency
 - Distribution: dapat dipindahkan
 - Saya berbicara [kepada mahasiswa] [tentang narkoba]
 - Saya berbicara [tentang narkoba] [kepada mahasiswa]
 - Substitution/Expansion/Pro-form: dapat digantikan
 - Saya duduk [di atas kotak/tepat di atas kotak/di sana].
 - Coordination, regular internal structure, no intrusion, fragments, semantics,...

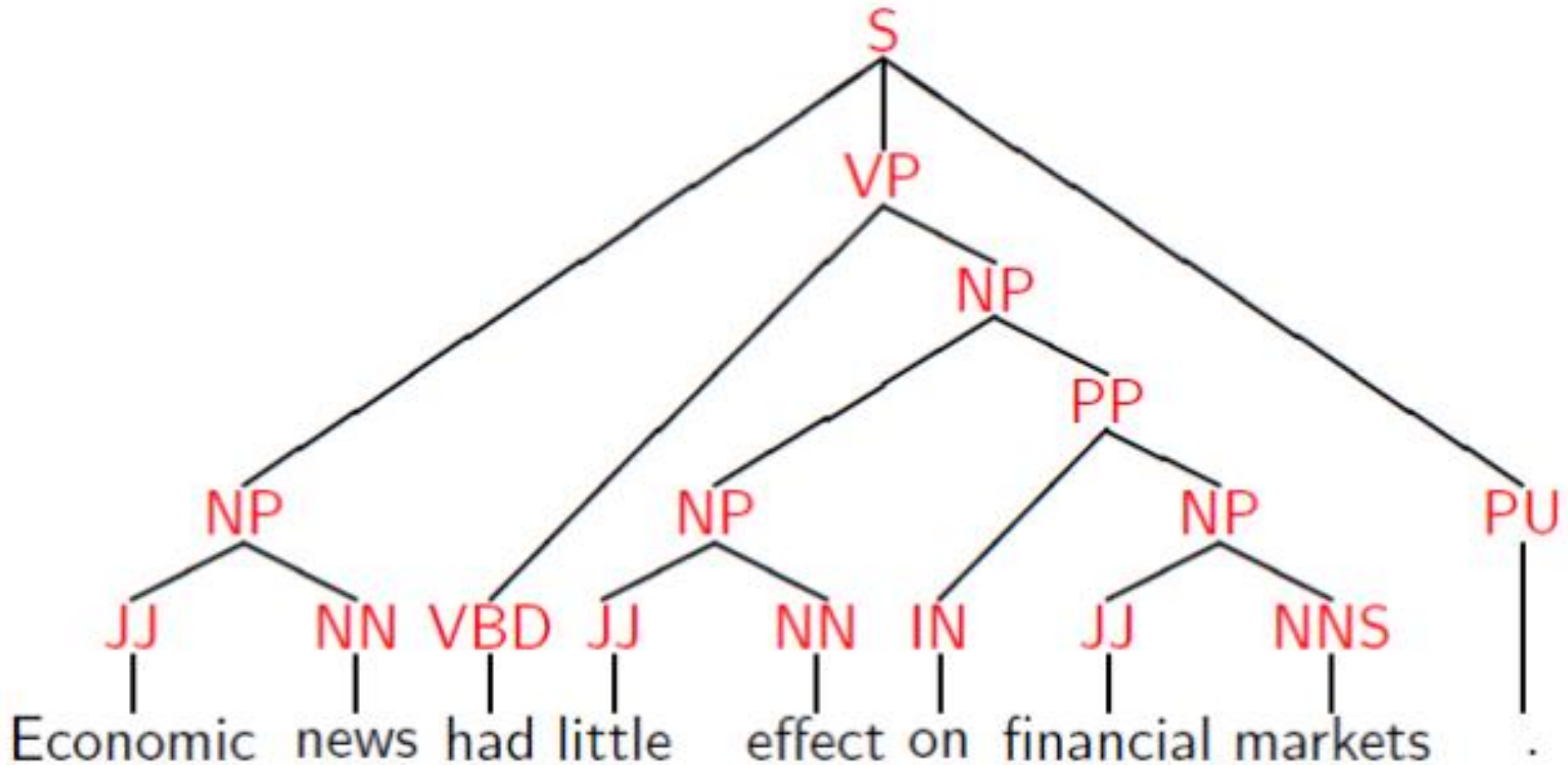


Contoh phrase structure





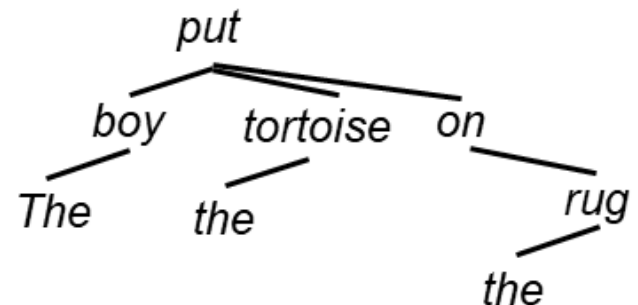
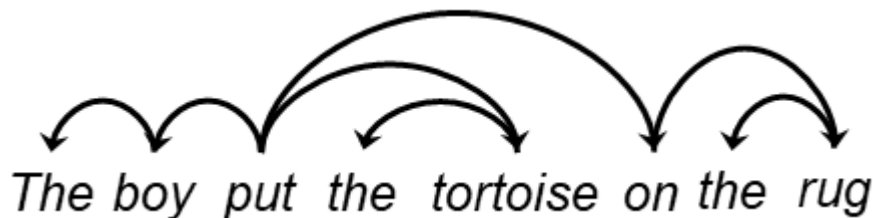
Contoh phrase structure



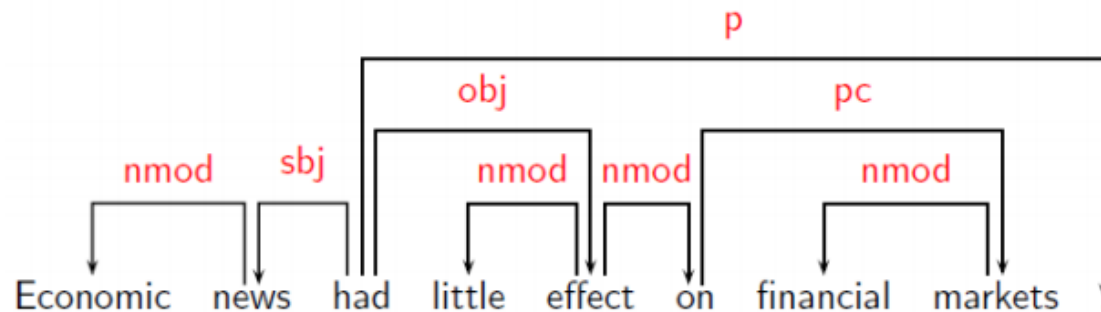
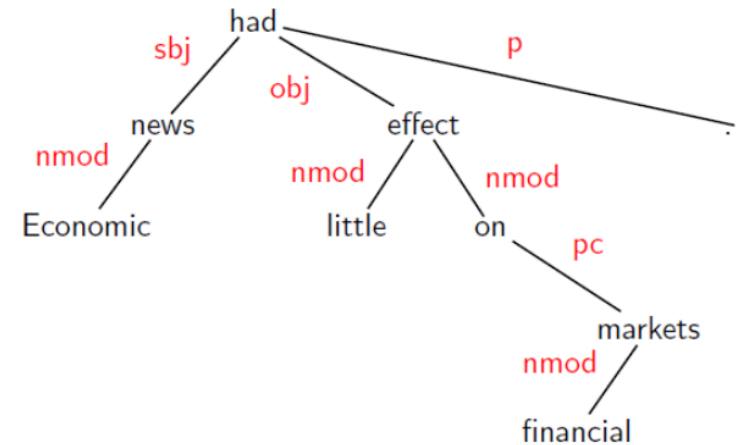
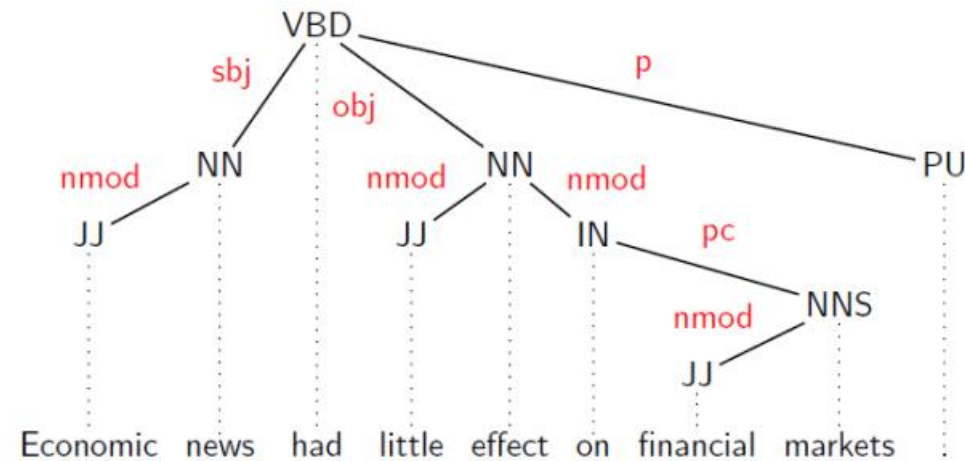
2. DEPENDENCY STRUCTURE



- Dependency structure atau struktur ketergantungan yaitu kata mana yang saling bergantung dengan kata lainnya



Contoh dependency structure





CONTEXT FREE GRAMMAR (CFG)



Context-Free Grammar

- Context-Free Grammar (CFG/Phrase-Structure Grammar/Backus-Naur Form) adalah suatu notasi matematis yang menyatakan aturan sebuah bahasa berdasarkan constituency.
- CFG terdiri dari dua bagian:
 1. Lexicon: daftar symbol/kata
 2. Sehimpunan (rewrite) rule atau production, yang menyatakan bagaimana symbol dalam bahasa dikelompokkan



Komponen dan Symbol CFG

- Komponen & symbol CFG terbagi 2 kelas:
 1. **Terminal symbols:** symbol yang merepresentasikan kata dalam bahasa (muncul dalam string/kalimat). Daftar terminal symbol disebut juga Lexicon
 2. **Non-terminal symbols:** symbol yang merepresentasikan kelompok/aggregate/generalisasi terminal.
- Non-terminal yang diasosiasikan dengan terminal kadang disebut **preterminal symbol**. Pada NLP, preterminal = part of speech



KOMPONEN CFG

- **Terminal** = ditulis huruf kecil {a, λ , .. }
- **Non Terminal** = ditulis huruf besar {A, S=Start Here, ...}
- Himpunan String / Kata yang terbentuk



Contoh CFG

- Himp = $\{\lambda, ab, aabb, aaabbb, \dots\}$

- CFG

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

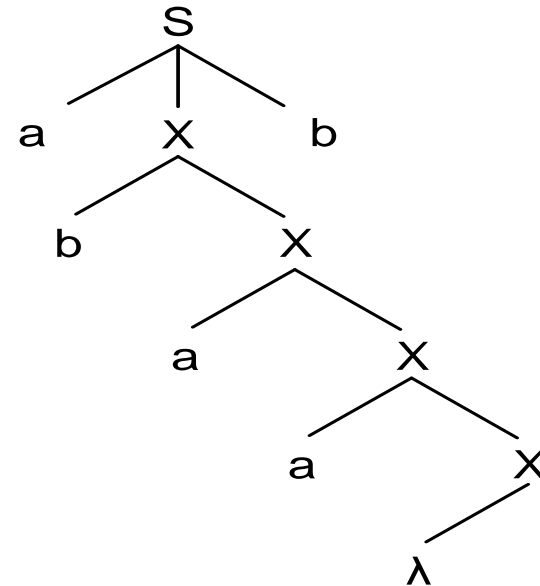
Atau

$$S \rightarrow aSb \mid \lambda$$

Struktur pohon (TREE)



- Struktur pohon (tree) dapat digunakan untuk menggambarkan proses pembentukan kata/string dari CFG
- $S \rightarrow a X b \mid b X a$
 $X \rightarrow a X \mid b X \mid \lambda$





Contoh CFG dlm NLP

Contoh *rule*:

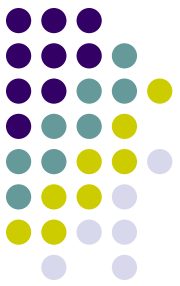
Dalam bhs. Inggris, sebuah NP bisa terdiri dari ProperNoun atau Determiner diikuti Nominal. Sebuah Nominal bisa terdiri dari satu atau lebih Noun.

NP → Det Nominal NP → ProperNoun Nominal → Noun | Noun Nominal

Contoh *lexicon*:

Lexicon bisa saja dinyatakan sebagai aturan context-free sebagai berikut:

Det → a Det → the Noun → flight



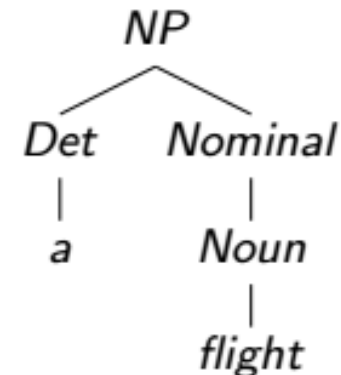
CFG sebagai pseudo-NLG

Sebagai generator:

CFG sebagai mesin yang menghasilkan kalimat sah sebagai output.

Sebagai generator, bacalah tanda “ \rightarrow ” sebagai: “*rewrite symbol on the left with the string of symbols on the right*”.

Mis, mulai dengan *NP*,
kita ganti *NP* dengan *Det Nominal*
lalu ganti *Nominal* dengan *Noun*,
lalu ganti *Det* dengan *a*,
dan ganti *Noun* dengan *flight*.



CFG bisa digunakan untuk menghasilkan semua kalimat yang bisa di-*derive* dari sebuah non-terminal, misalnya *NP*. Urutan pengaplikasian rule disebut sebagai *derivation* sebuah kalimat – biasanya direpresentasikan sebagai sebuah *parse tree*.



CFG sebagai pseudo-NLG

- Bayangkan kita membuat sebuah program yang inputnya sebuah CFG, dan secara random menghasilkan beberapa kalimat yang sah. Ini bisa saja dianggap contoh sistem NLG (Natural Language Generation) yang paling sederhana.
 - Random nonsense: <http://www.elsewhere.org/pomo/> (The Postmodernism Generator)
 - Computer science paper generator :-)
<http://pdos.csail.mit.edu/scigen/> (SCIgen)
 - Ketika digabung dengan parsing, bisa membuat “conversational agent”, mis. ELIZA (cf. Turing Test)
- Program-program ini mengetahui syntax, tetapi tidak semantics (makna).
- Sistem NLG “sesungguhnya” secara eksplisit menyatakan makna.

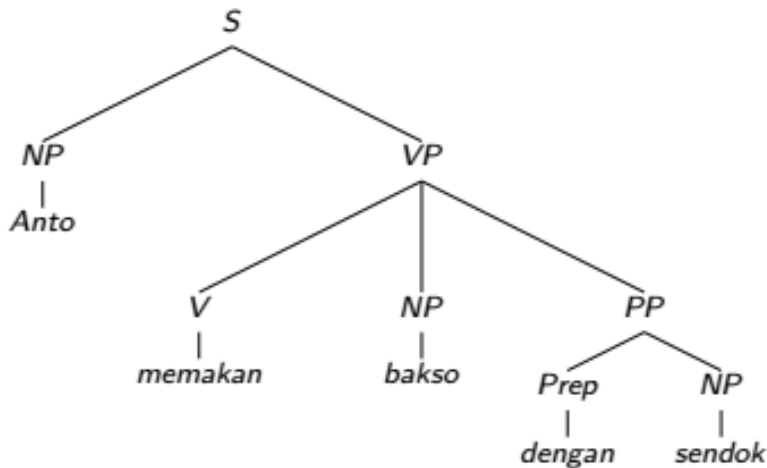


CFG sebagai parser

Sebagai parser:

CFG sebagai mesin yang menghasilkan struktur untuk sebuah kalimat input.

Sebagai parser, bacalah tanda “ \rightarrow ” sebagai: *“if you see the symbols on the right, rewrite with the symbol on the left”*.



- Masalah bahasa manusia: satu kalimat, banyak struktur (=ambiguity/kerancuan)
- Tujuan akhir *parsing*: agar program dapat “memahami” semantics/makna dari kalimat.

Notasi alternatif, *bracketed notation*:

[S [NP Anto] [VP [V memakan] [NP bakso] [PP [Prep dengan] [NP sendok]]]]

Grammar singkat untuk bahasa Inggris (ATIS)



Lexicon untuk \mathcal{L}_0

<i>Noun</i>	→	<i>flights</i> <i>breeze</i> <i>trip</i> <i>morning</i> ...
<i>Verb</i>	→	<i>is</i> <i>prefer</i> <i>like</i> <i>need</i> ...
<i>Adjective</i>	→	<i>cheapest</i> <i>first</i> <i>other</i> <i>direct</i> ...
<i>Pronoun</i>	→	<i>me</i> <i>I</i> <i>you</i> <i>it</i> ...
<i>ProperNoun</i>	→	<i>Alaska</i> <i>Baltimore</i> <i>Chicago</i> <i>Garuda</i> ...
<i>Determiner</i>	→	<i>the</i> <i>a</i> <i>an</i> <i>this</i> ...
<i>Preposition</i>	→	<i>from</i> <i>to</i> <i>on</i> <i>near</i> ...

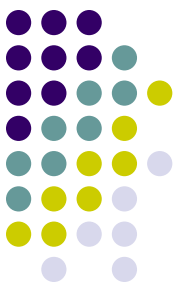
Grammar untuk \mathcal{L}_0

<i>S</i>	→	<i>NP VP</i>	I + prefer a morning flight
<i>NP</i>	→	<i>Pronoun</i>	I
		<i>ProperNoun</i>	Los Angeles
		<i>Det Nominal</i>	a + flight
<i>Nominal</i>	→	<i>Noun Nominal</i>	morning + flight
		<i>Noun</i>	flights
<i>VP</i>	→	<i>Verb</i>	do
		<i>Verb NP</i>	prefer a morning flight
		<i>Verb NP PP</i>	leave Boston in the morning
		<i>Verb PP</i>	leave in the morning
<i>PP</i>	→	<i>Prep NP</i>	from Los Angeles



Grammaticality

- Kalimat-kalimat yang bisa di-derive dari S dikatakan grammatical.
- Kalimat-kalimat yang **TIDAK** bisa di-derive dari S dikatakan **ungrammatical**.
- Perbedaan yang sangat “tajam” untuk bahasa formal seperti ini terkadang kurang cocok untuk bahasa natural/manusia . . .
- Dalam bidang linguistics, pemodelan ini disebut generative grammar (Chomsky), akhir '60-an.

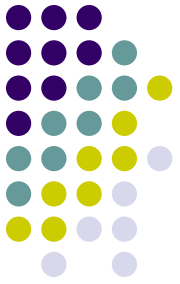


Denisi formal CFG

Ingat, sebuah CFG adalah 4-tuple:

- ① himpunan simbol non-terminal N
- ② himpunan simbol terminal Σ (di mana $N \cap \Sigma = \emptyset$)
- ③ himpunan **production rule** P , masing-masing berbentuk $A \rightarrow \alpha$ di mana
 - $A \in N$, dan
 - $\alpha \in (\Sigma \cup N)^*$ (dkl. α adalah string simbol terminal/nonterminal)
- ④ sebuah **start symbol** $S \in N$

Bahasa formal yang dinyatakan oleh sebuah CFG adalah himpunan string yang bisa di-*derive* dari symbol khusus: *start symbol* (S). Dalam NLP, S sering diartikan sebagai “*sentence*”.



PARSING

Syntatic Parsing



- Memeriksa makna yang terkandung dalam text dan membandingkannya dengan aturan tata bahasa formal
- *Grammar* : aturan pembentuk suatu kalimat dalam sebuah bahasa dan dikenali dengan tata bahasa
- Proses searching:
 - *Context Free Grammar (CFG)*

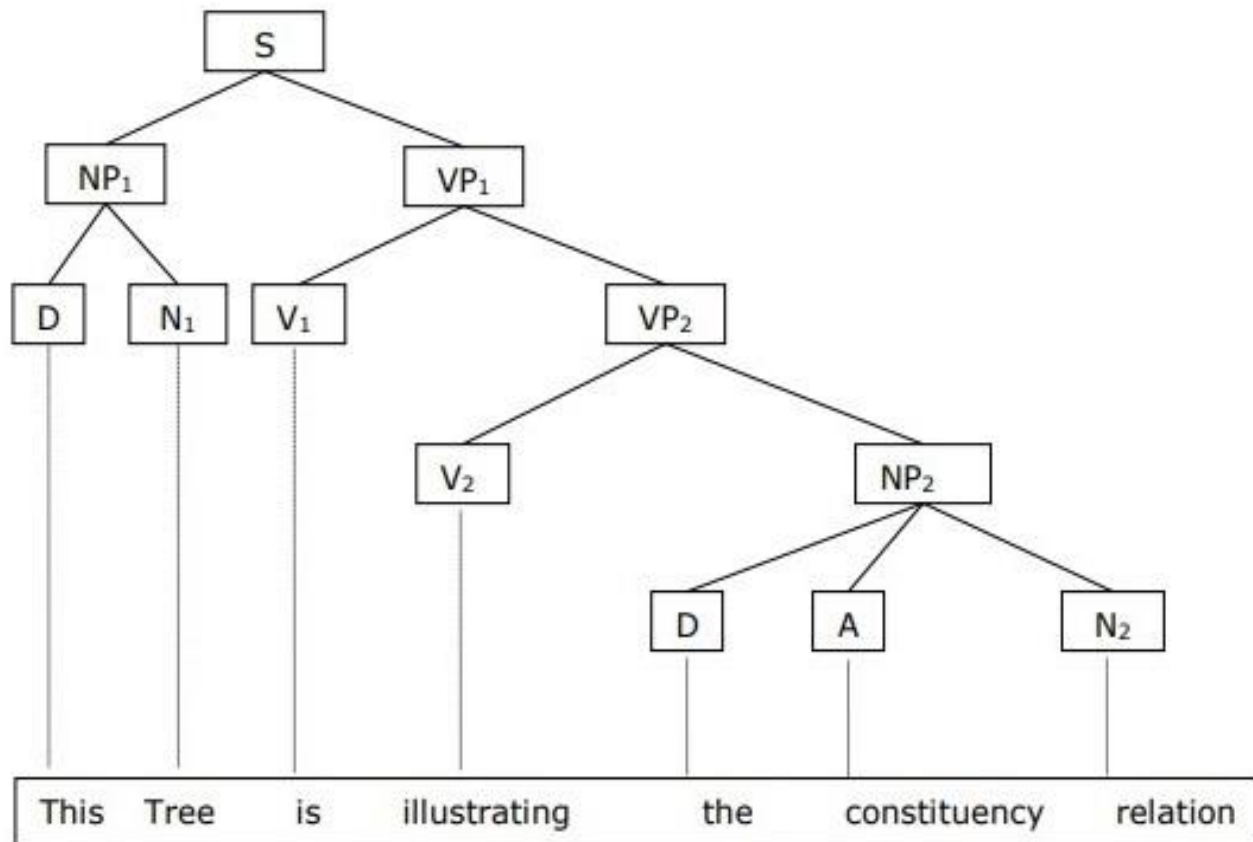
Grammar

$S \rightarrow NP VP$
 $S \rightarrow Aux NP VP$
 $S \rightarrow VP$
 $NP \rightarrow Pronoun$
 $NP \rightarrow Proper-Noun$
 $NP \rightarrow Det Nominal$
 $Nominal \rightarrow Noun$
 $Nominal \rightarrow Nominal Noun$
 $Nominal \rightarrow Nominal PP$
 $VP \rightarrow Verb$
 $VP \rightarrow Verb NP$
 $VP \rightarrow VP PP$
 $PP \rightarrow Prep NP$

Lexicon

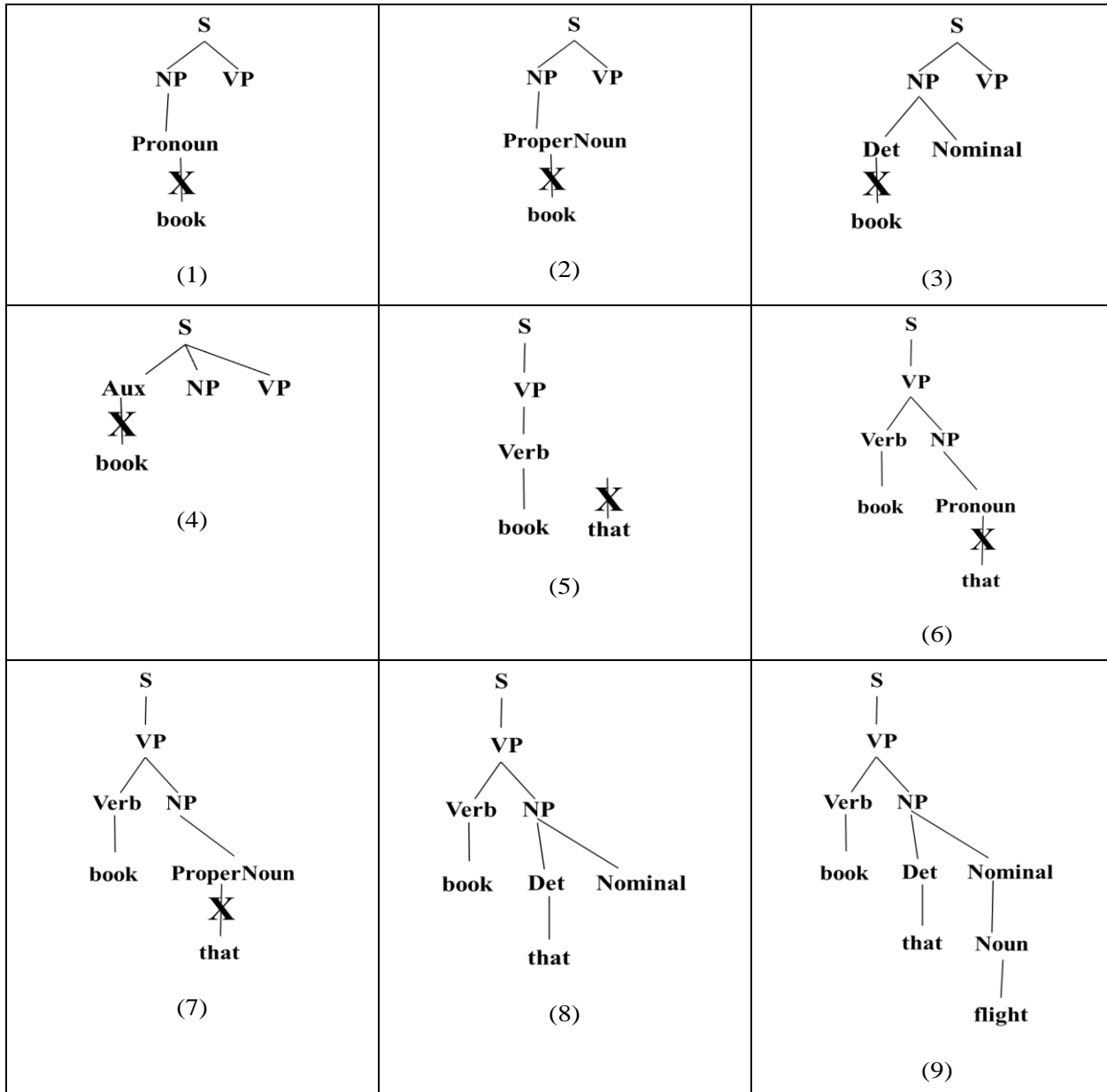
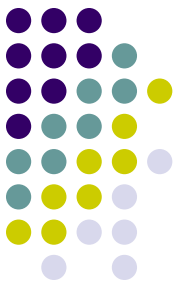
$Det \rightarrow the | a | that | this$
 $Noun \rightarrow book | flight | meal | money$
 $Verb \rightarrow book | include | prefer$
 $Pronoun \rightarrow I | he | she | me$
 $Proper-Noun \rightarrow Houston | NWA$
 $Aux \rightarrow does$
 $Prep \rightarrow from | to | on | near | through$

Ilustrasi Parse Tree

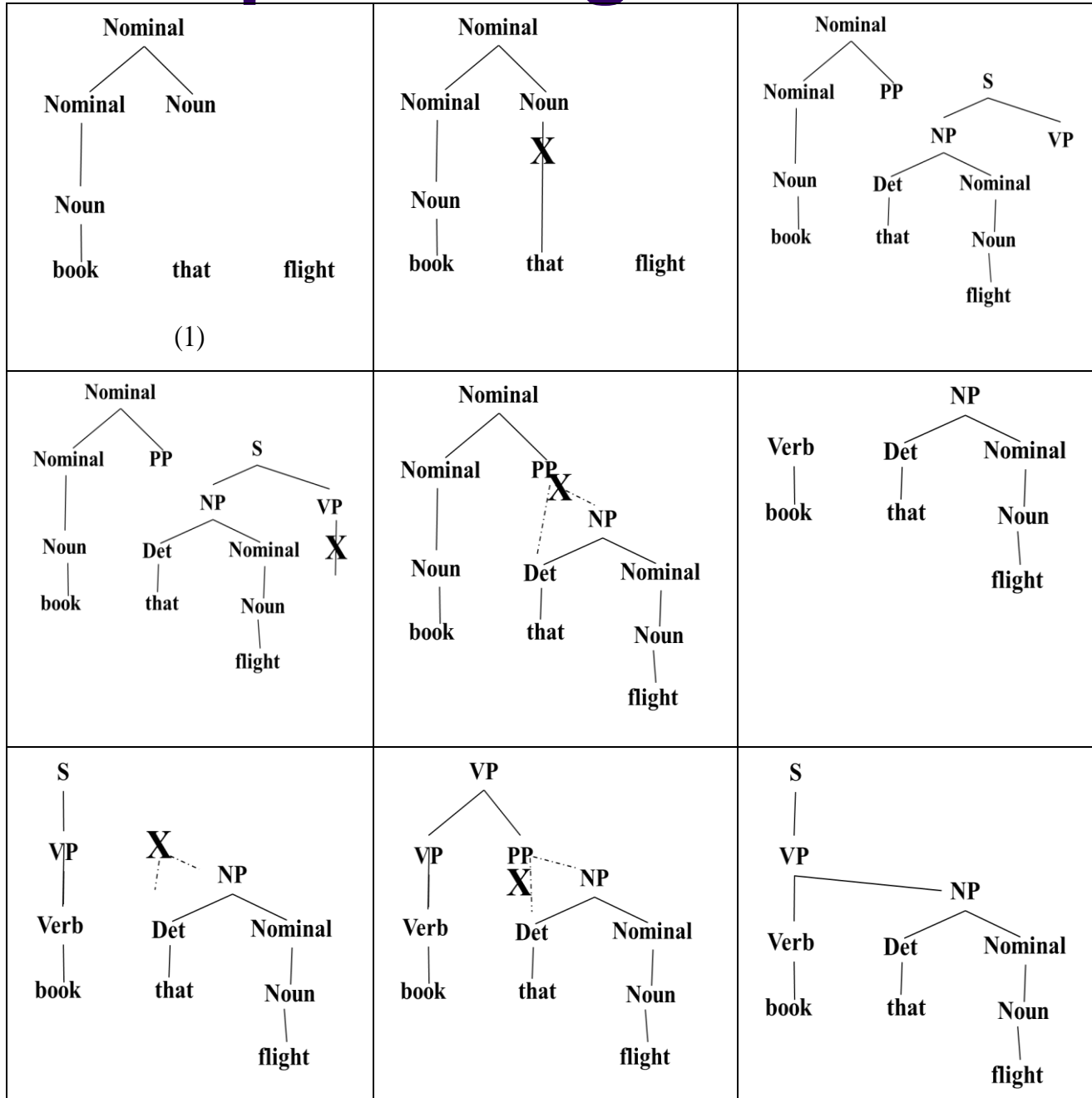


Top Down Parsing

- Kalimat : 'book that flight'



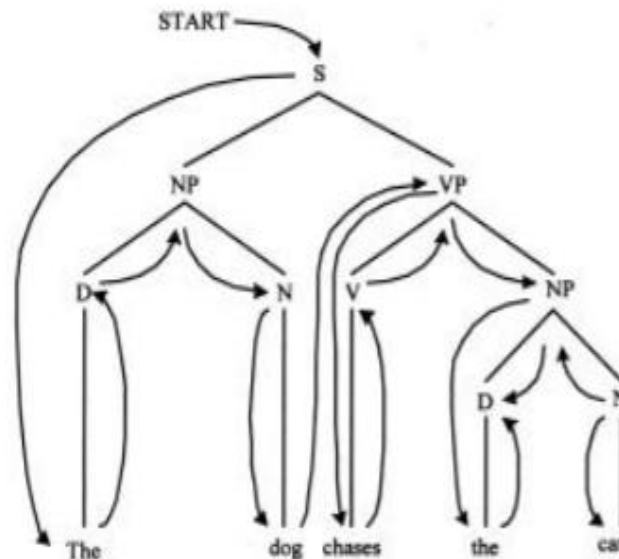
Bottom Up Parsing





Kelemahan

- Top Down parsing → tidak dapat menangani masalah *left-recursion*
- Bottom Up parsing → tidak dapat menangani grammar dengan *empty production*.
- Gabungan : model *Left-Corner* dan *Earley's Parsing*



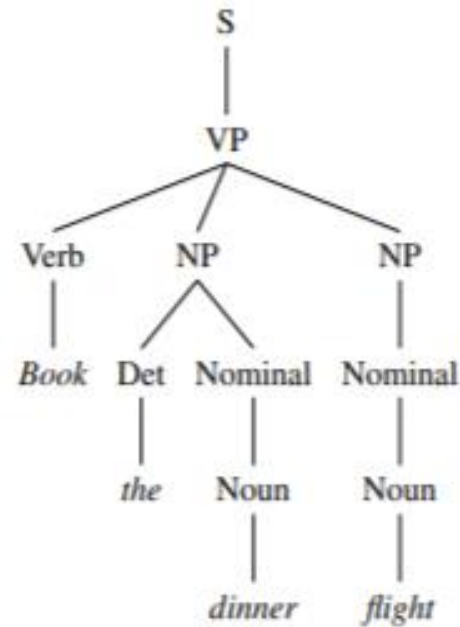
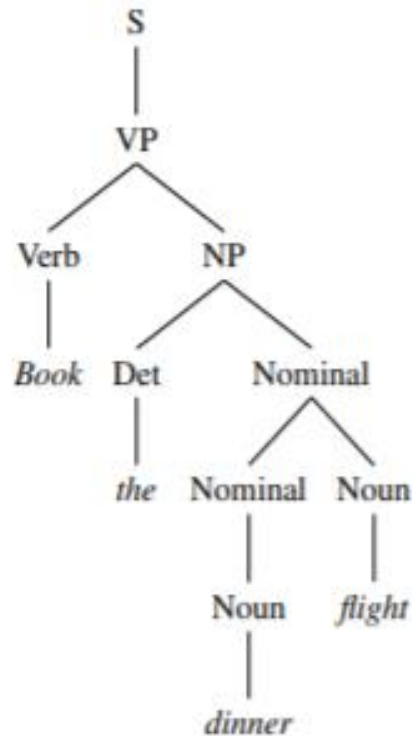
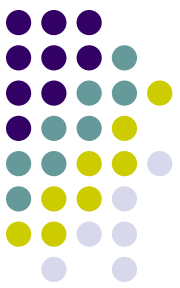


Statistical Parsing

- Model sintaktis parsing tetapi dalam grammarnya terdapat nilai probabilitas
- Probabilistic Context Free Grammar (PCFG)

Grammar	Prob	Lexicon
S → NP VP	0.8	Det → the a that this
S → Aux NP VP	0.1	0.6 0.2 0.1 0.1
S → VP	0.1	Noun → book flight meal money
NP → Pronoun	0.2	0.1 0.5 0.2 0.2
NP → Proper-Noun	0.2	Verb → book include prefer
NP → Det Nominal	0.6	0.5 0.2 0.3
Nominal → Noun	0.3	Pronoun → I he she me
Nominal → Nominal Noun	0.2	0.5 0.1 0.1 0.3
Nominal → Nominal PP	0.5	Proper-Noun → Houston NWA
VP → Verb	0.2	0.8 0.2
VP → Verb NP	0.5	Aux → does
VP → VP PP	0.3	1.0
PP → Prep NP	1.0	Prep → from to on near through
		0.25 0.25 0.1 0.2 0.2

Contoh



Rules			Rules		
	Rules	P		Rules	P
S	→ VP	.05	S	→ VP	.05
VP	→ Verb NP	.20	VP	→ Verb NP NP	.10
NP	→ Det Nominal	.20	NP	→ Det Nominal	.20
Nominal	→ Nominal Noun	.20	NP	→ Nominal	.15
Nominal	→ Noun	.75	Nominal	→ Noun	.75
Verb	→ book	.30	Nominal	→ Noun	.75
Det	→ the	.60	Verb	→ book	.30
Noun	→ dinner	.10	Det	→ the	.60
Noun	→ flight	.40	Noun	→ dinner	.10
			Noun	→ flight	.40



Dependency Parsing

- Model parsing dengan memperhatikan konsep keterkaitan atau hubungan antar kata dalam kalimat

- Contohnya:

Kalimat 1 : I've never been afraid

Dependency Parser : nsubj (afraid-5, I-1)
aux (afraid-5, 've-2)
neg (afraid-5, never-3)
cop (afraid-5, been-4)
root (ROOT-0, afraid-5)

POS Tagging : I've|NNP never|RB been|VBN afraid|JJ



Referensi

- Buku Rerensi Utama Kuliah
- Natural Language Processing, Lecture Slides by Dan Jurafsky and Christopher Manning, Stanford Coursera course
- Komputasi Bahasa Alami, Dr. Fika Hastarita Rachman, S.T., M.Eng. Jurusan Teknik Informatik, Universitas Trunojoyo Madura 2020
- CSC4602354 Pengolahan Bahasa Manusia, Rahmad Mahendra, Fakultas Ilmu Komputer, Universitas Indonesia