

NLP (Natural language processing)

Dr. Andreas Hadiyono, ST, MMSI

Machine
Learning
Course

23/02/2019





NLP



- “Natural” languages
 - English, Mandarin, French, Swahili, Arabic, Nahuatl,
 - NOT Java, C++, Perl, ...
- Tujuan Utama: Natural human-to-computer communication
- Sub-field of Artificial Intelligence, but very interdisciplinary
 - Computer science, human-computer interaction (HCI), linguistics, cognitive psychology, speech signal processing (EE), ...



Tujuan NLP

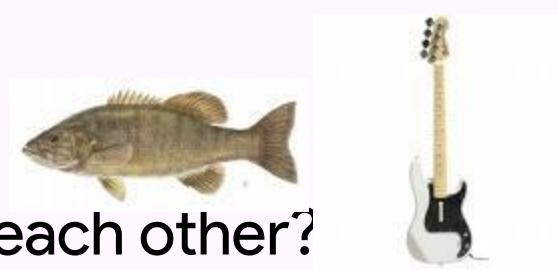
- Understand language analysis & generation
- Communication
- Language is a window to the mind
- Data is in linguistic form
- Data can be in Structured (table form), Semi structured (XML form), Unstructured (sentence form).



Language Processing

- Level 1 – Speech sound (*Phonetics & Phonology*)
- Level 2 – Words & their forms (*Morphology, Lexicon*)
- Level 3 – Structure of sentences (*Syntax, Parsing*)
- Level 4 – Meaning of sentences (*Semantics*)
- Level 5 – Meaning in context & for a purpose (*Pragmatics*)
- Level 6 – Connected sentence processing in a larger body of text (*Discourse*)

- Morphology: What is a word?
- 奧林匹克運動會（希臘語：Ολυμπιακοί Αγώνες，簡稱奧運會或奧運）是國際奧林匹克委員會主辦的包含多種體育運動項目的國際性運動會，每四年舉行一次。
- ای بیوت = “to her houses”
- Lexicography: What does each word mean?
 - He plays bass guitar.
 - That bass was delicious!
- Syntax: How do the words relate to each other?
 - The dog bit **the man**. ≠ **The man** bit **the dog**.
 - But in Russian: **человек** **собаку** **съел** = **человек** **съел** **собаку**





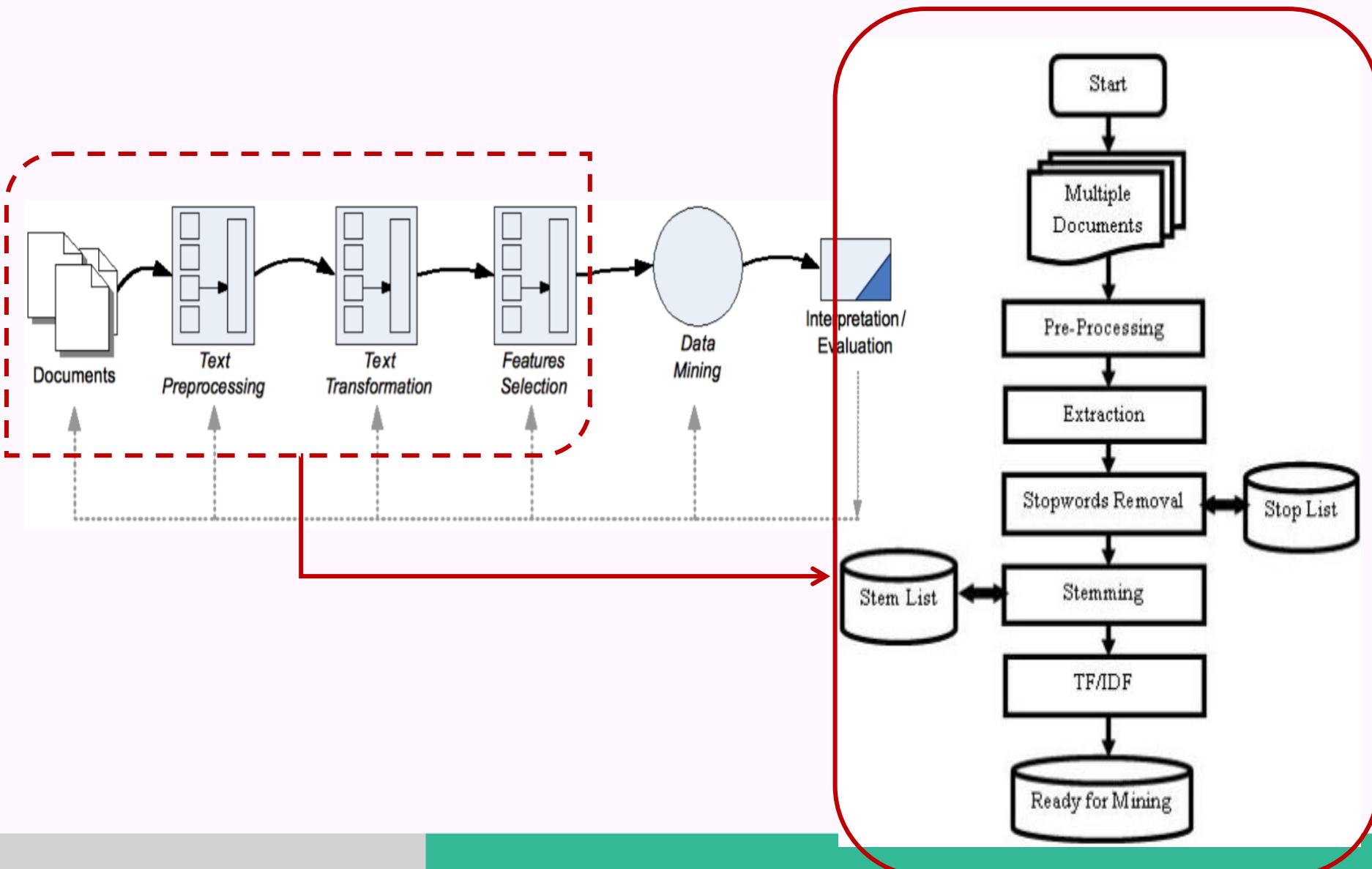
Language Processing

- **Semantics**: How can we infer meaning from sentences?
 - I saw the man on the hill with the telescope.
 - The ipod is so **small!** ☺
 - The monitor is so **small!** ☹
- **Discourse**: How about across many sentences?
 - President Bush met with President-Elect Obama today at the White House. He welcomed him, and showed him around.
 - Who is “he”? Who is “him”? How would a computer figure that out?



Preprocessing

Tahapan (bagian dari data/teks mining)





Pra pemrosesan Teks

- merubah data textual orisinal dalam bentuks struktur yang data-mining-ready
- fase paling kritis dan paling kompleks
- tujuan: mendapatkan fitur kunci dan term kunci dari dokumen teks untuk memperkuat relevansi antara kata dan dokumen , dan relevansi kata dan kategori
- merepresentasikan setiap dokumen sebagai vektor fitur, dengan membagi menjadi unsur terkecil dari dokumen teks --> kata



Aktifitas Pra pemrosesan

- Stopword Removal:
 - tidak semua kata memiliki “informasi” yang dibutuhkan
- Stemming (dan atau lemmatization):
 - teknik untuk mencari akar (root)/stem dari kata,
 - contoh kata user, users, using, used memiliki akar/stem: “USE”
 - stemmer: Porter stemmer (Inggris), Indonesia stemmer (Algoritma Nazief & Adriani)
 - lemmatization: am, is, was, are --> be
- Document indexing:
 - tujuan utama: meningkatkan efisiensi dgn mengekstrak dokumen hasil suatu term terpilih yg digunakan untuk mengindeks
 - dengan memilih himpunan kata kunci yang pantas berdasarkan corpus document, dan memberikan bobot tiap kata kunci tsb ke dokumen tertentu, yang kemudian
 - mentransformasikan setiap dokumen menjadi vektor bobot kata kunci
 - bobotnya berelasi dgn frekuensi kemunculan term dalam dokumen, dan banyak dokumen yg menggunakan term tsb



Aktifitas Pra pemrosesan

- Document Indexing: Term Weighting --> TF-IDF
- Dimensionality Reduction (DF)
 - banyaknya dokumen dimana suatu term muncul
 - cara paling simpel untuk reduksi kosa kata (vocabulary)



Stop word removal

- stop word (SW) perlu dihilangkan dalam dokumen teks karena membuat dokumen lebih berat dan kurang penting utk dianalisis
- penghilangan SW mengurangi dimensi dari ruang term.
- umumnya penghilangan atas preposisi dan pro-noun yang tidak memberikan arti thdp dokumen
- contoh:
 - inggris: *the, in, an, with, dll*
 - indonesia: *ini, itu, kami, di, ke, dari, pada, dengan, si, dll*

stop word removal tidak
sesuai diterapkan untuk
dokumen yg sensitive, misal:
kitab suci

(rio bagus prakoso, 2017)



Stop word removal

Jenis metode SW:

- metode klasik: SW didapat dari daftar pre-compiled
- Zipf-law (Z-Method):
- Mutual Information Methode (MI)
- Term Based Random Sampling (TBRS) --> thn 2005:
 - secara manual mendekripsi stop word dari dokumen web. menggunakan pengukuran divergensi Kull :

$$d_x(t) = P_x(t) \cdot \log_2 \frac{P_x(t)}{p(t)}$$

dimana $P_x(t)$ adalah normalisasi frekuensi term t dalam mass x ,
 $P(t)$ adalah normalisasi frekuensi term t dalam koleksi keseluruhan.

stop list dibentuk dgn mengambil least informative term dalam seluruh chunks, menghapus seluruh duplikasi



Zipf law: didasari observasi/empiris, bukan teori

- Zipf law:

- semakin sering suatu kata yang digunakan pada suatu bahasa tertentu akan menyebabkan kata tersebut semakin sering untuk digunakan kembali, namun semakin jarang suatu kata yang digunakan pada suatu bahasa tertentu akan menyebabkan kata tersebut semakin jarang untuk digunakan kembali.

$$P(r) \approx \frac{1}{r \ln(1.78 R)},$$

where R is the number of different words.

- terdapat suatu rasio penyebab 20:80 dan akibat 80:20.

Contoh 1, 20% pelanggan dari suatu perusahaan akan menghasilkan 80% dari total pemasukan perusahaan tersebut.

Contoh 2, 20% pelanggan dari suatu perusahaan akan mengakibatkan permasalahan sebesar 80% pada perusahaan tersebut. Hal ini menyebabkan fokus untuk menyelesaikan suatu permasalahan dapat diarahkan kepada 20% tersebut, maka permasalahan akan selesai.



Stemming vs Lemmatization

Lemmatization

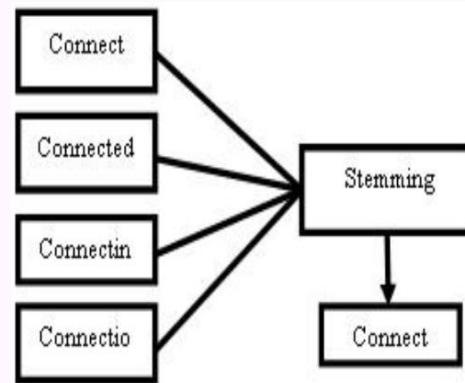
- proses merubah kata dari kalimat ke dalam bentuk (rujukan) kamus
- analisis morphologi
- perlu kamus lengkap
- misal, kata *amusing*, *amusement*, and *amused above*, lemma untuk tiap kata
--> *amuse*

Stemming

- proses merubah kata dari kalimat ke dalam porsi tidak berubah
- biasanya mereduksi berbagai suffixes (akhiran), jika perlu preffix (awalan)
- misal, kata *amusing*, *amusement*, and *amused above*, di stem --> *amus*



Stemming



**sama2 mencari root/stem
dari kata kalimat**



Stemming dalam NLTK

- Misal dgn stemmer Lancaster

```
1 lancaster = nltk.LancasterStemmer()  
2 stems = [lancaster.stem(i) for i in tokens]  
3
```

- Outputnya:

```
1 ['omg', ',', 'nat', 'langu', 'process', 'is', 'so', 'cool', 'and', 'i', "m", '  
2
```

- Misal dgn stemmer Porter

```
1 porter = nltk.PorterStemmer()  
2 stem = [porter.stem(i) for i in tokens]  
3
```

- Outputnya

```
1 ['omg', ',', 'natur', 'languag', 'process', 'is', 'so', 'cool', 'and', 'i', "'m", '  
2
```

pada stemmer Porter, kata “natural”
dipetakan ke “natur” dibanding “nat” pada
lancaster.
dan kata “really” dipetakan ke realli
dibanding “real”



Lemmatization dalam NLP

```
1 from nltk import WordNetLemmatizer  
2  
3 lemma = nltk.WordNetLemmatizer()  
4 text = "Women in technology are amazing at coding"  
5 ex = [i.lower() for i in text.split()]  
6  
7 lemmas = [lemma.lemmatize(i) for i in ex]  
8  
9 ['woman', 'in', 'technology', 'are', 'amazing', 'at', 'coding']  
10
```

- Kata “**women**” diubah/petakan ke “**woman**”



Pendekatan Bag-of-Tokens

Documents

Four score and seven years ago our fathers brought forth on this continent, **a new nation**, conceived in Liberty, and dedicated to the proposition that all men are created equal.

Now we are engaged in a great civil war, testing whether **that nation**, or ...

Feature Extraction

Token Sets

nation – 5
civil - 1
war – 2
men – 2
died – 4
people – 5
Liberty – 1
God – 1
...

Loses all order-specific information!
Severely limits context!



Chunking Teks

- Biasa disebut shallow/parsial parsing
- Identifikasi part-of-speech dan frase pendek (noun phrases) kata benda
- merupakan struktur non-rekursif yang dapat ditangani dengan metode stata hingga



- Karena full parsing expensive dan tidak terlalu robust
- Contoh kita ingin tahu entitas nam, misal President Obama dalam dokumen berikut;



President Barack Obama criticized insurance companies and banks as he urged supporters to pressure Congress to back his moves to revamp the health-care system and overhaul financial regulations. ([source](#))

[NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to] [NP only \$1.8 billion] [PP in] [NP September]



Chunking Teks

- Chunking dapat dilakukan pada saat ingin mengetahui entitas bernama (named-entity), seperti:
 - nama orang
 - organisasi
 - ekspresi waktu (tahun)
 - lokasi
 - nilai uang
 - persentase
 - dll
- misal, dokumen tak beranotasi/berlabel:

Jim bought 300 shares of Acme Corp. in 2006.

- setelah di-chunking:

[Jim]_{Person} bought 300 shares of [Acme Corp.]_{Organization} in [2006]_{Time}.

Proses ini disebut: **Named-Entity Recognition**



Chunking Teks

- Deteksi Lokasi

KEEP UP ON YOUR READING WITH AUDIO BOOKS

Vietnam UK Louisiana, USA

Audio books are highly popular with library patrons in the town

Louisiana, USA S.Carolina, USA Pennsylvania, USA Mass., USA

of Springfield, Greene County, MO. "People are mobile

Turkey Virginia, USA Maine, USA Norway Alabama, USA

and busier, and audio books fit into that lifestyle" says Gary

Louisiana, USA Indiana, USA

Sanchez, who oversees the library's \$2 million budget...

Dominican Republic Pennsylvania, USA Kentucky, USA

- NER

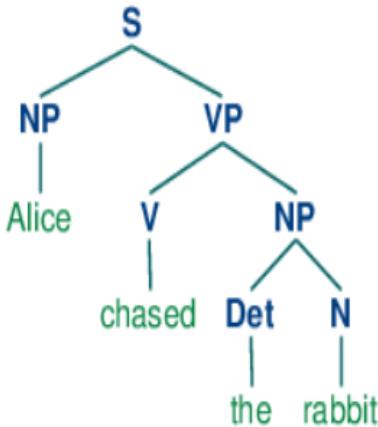
NE Type	Examples
ORGANIZATION	<i>Georgia-Pacific Corp., WHO</i>
PERSON	<i>Eddy Bonte, President Obama</i>
LOCATION	<i>Murray River, Mount Everest</i>
DATE	<i>June, 2008-06-29</i>
TIME	<i>two fifty a m, 1:30 p.m.</i>
MONEY	<i>175 million Canadian Dollars, GBP 10.40</i>
PERCENT	<i>twenty pct, 18.75 %</i>
FACILITY	<i>Washington Monument, Stonehenge</i>
GPE	<i>South East Asia, Midlothian</i>



(Parsing) Tree

- Tree dalam NLTK

- Kumpulan titik berlabel yang terhubung



- format teks dari Tree

```
(S  
  (NP Alice)  
  (VP  
    (V chased)  
    (NP  
      (Det the)  
      (N rabbit))))
```

```
>>> tree1 = nltk.Tree('NP', ['Alice'])  
>>> print(tree1)  
(NP Alice)  
>>> tree2 = nltk.Tree('NP', ['the', 'rabbit'])  
>>> print(tree2)  
(NP the rabbit)
```

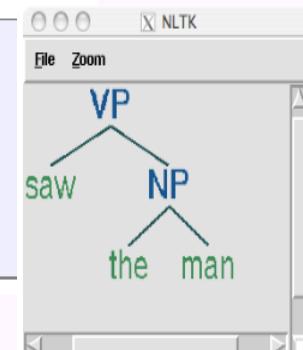
We can incorporate these into successively larger trees as follows:

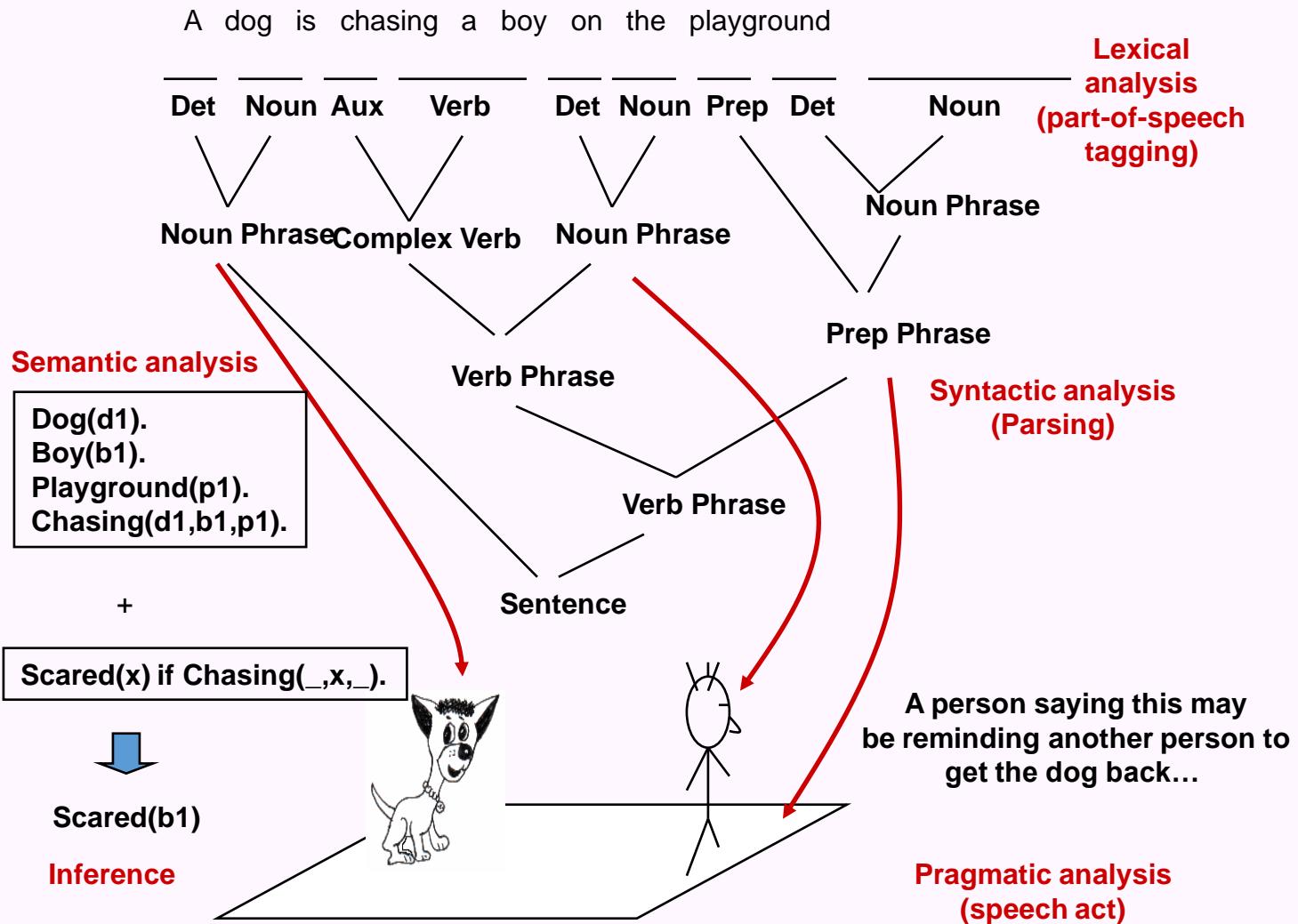
```
>>> tree3 = nltk.Tree('VP', ['chased', tree2])  
>>> tree4 = nltk.Tree('S', [tree1, tree3])  
>>> print(tree4)  
(S (NP Alice) (VP chased (NP the rabbit)))
```

Here are some of the methods available for tree objects:

```
>>> print(tree4[1])  
(VP chased (NP the rabbit))  
>>> tree4[1].label()  
'VP'  
>>> tree4.leaves()  
['Alice', 'chased', 'the', 'rabbit']  
>>> tree4[1][1][1]  
'rabbit'
```

```
>>> tree3.draw()
```





(Taken from ChengXiang Zhai, CS 397cxz – Fall 2003)



Parsing Struktur Frase

Interpreting Language is Hard!

I saw a girl with a telescope

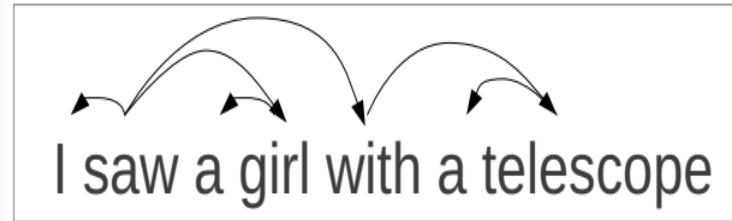


- “Parsing” resolves structural ambiguity in a formal way

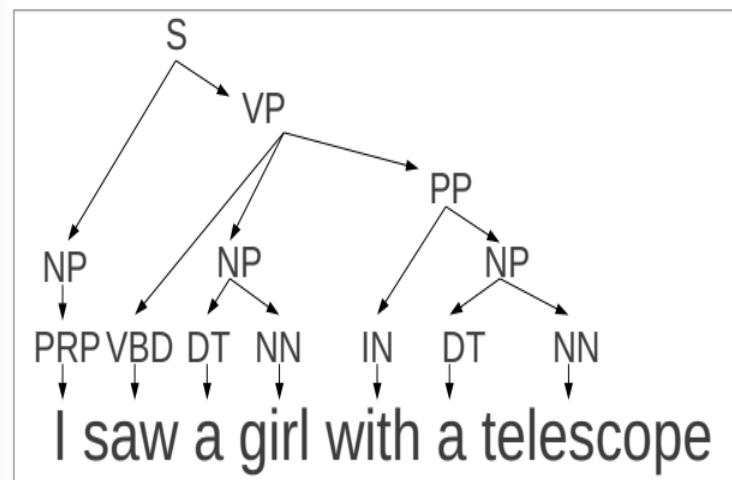
Parsing Struktur Frase: 2

Jenis Parsing

- Dependensi: berfokus pada relasi antara dua kata

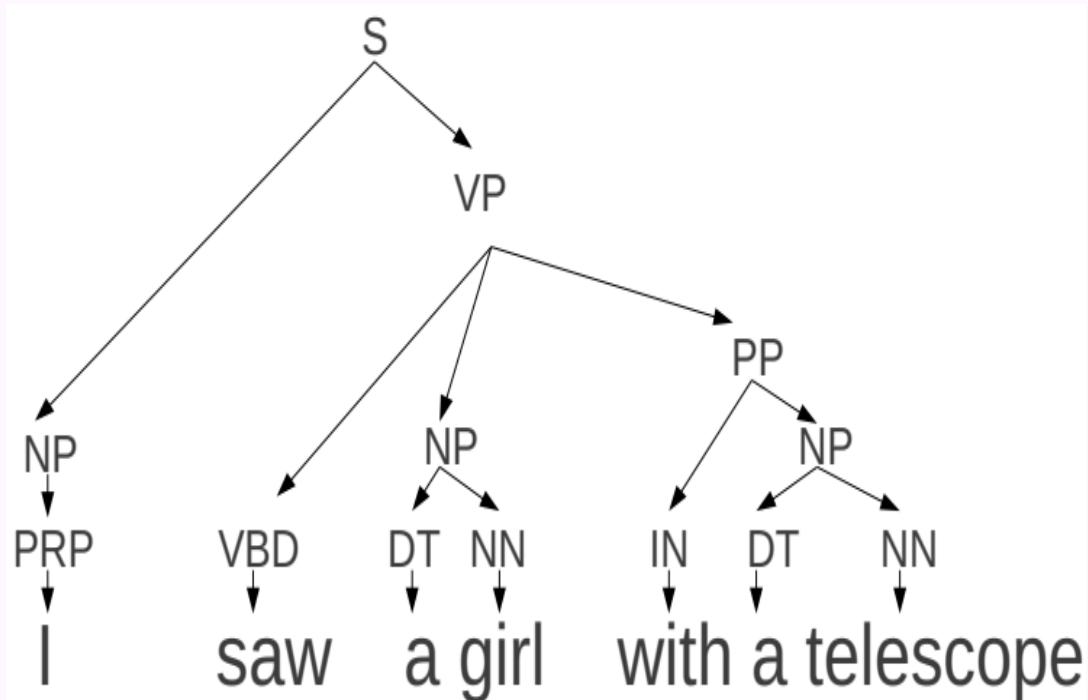


- Struktur frase: berfokus pada identifikasi frase dan struktur rekursifnya



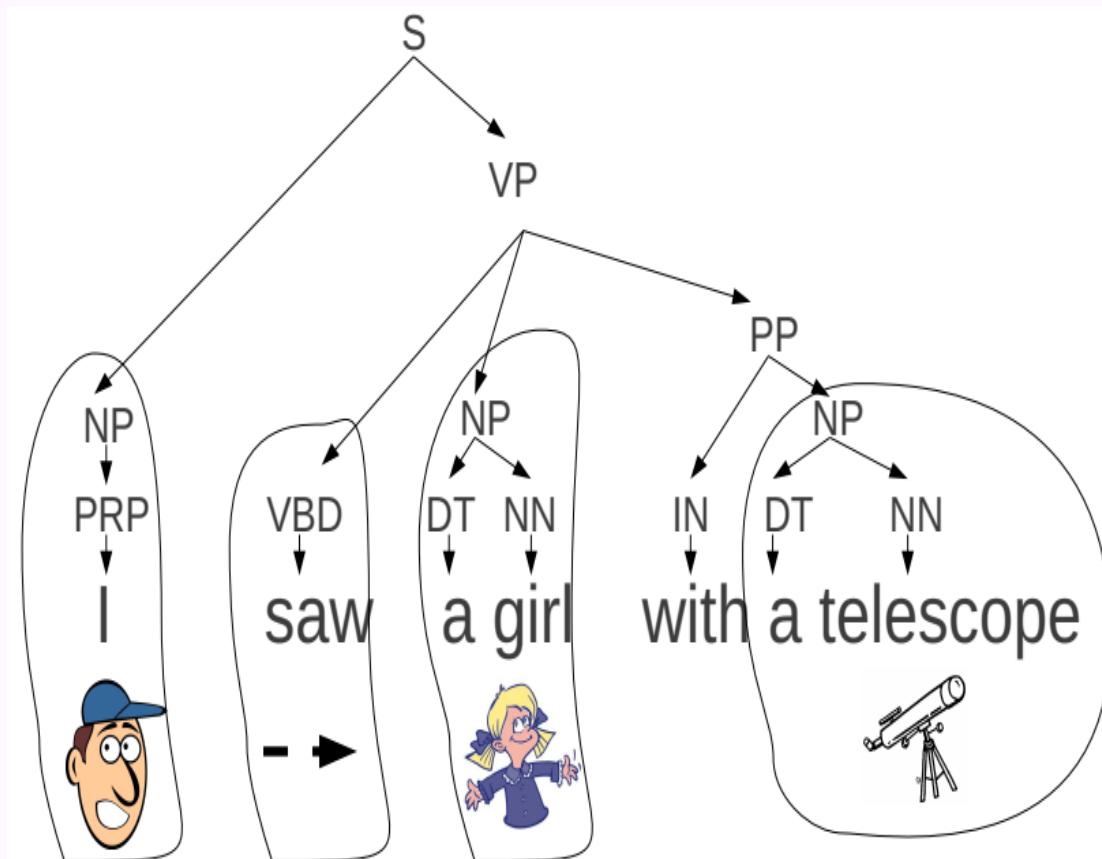


Parsing Struktur Frase: Struktur Rekursif



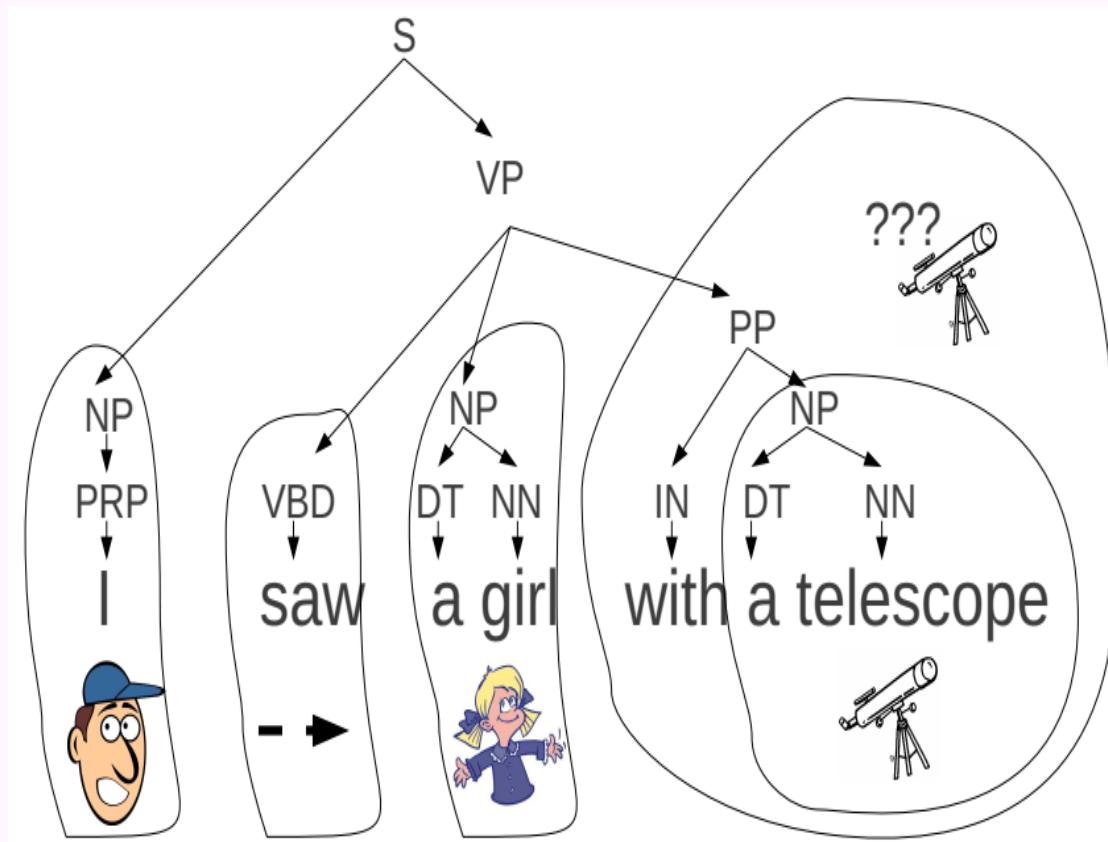


Parsing Struktur Frase: Struktur Rekursif



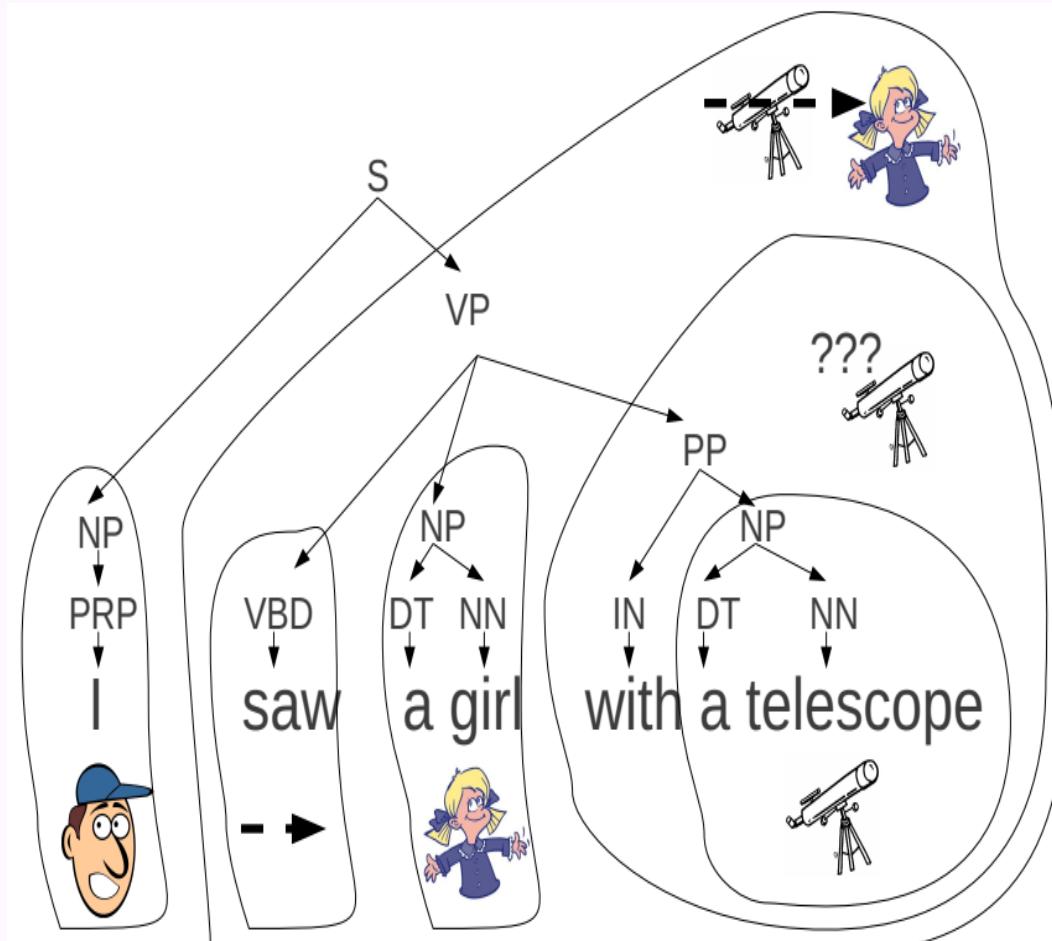


Parsing Struktur Frase: Struktur Rekursif

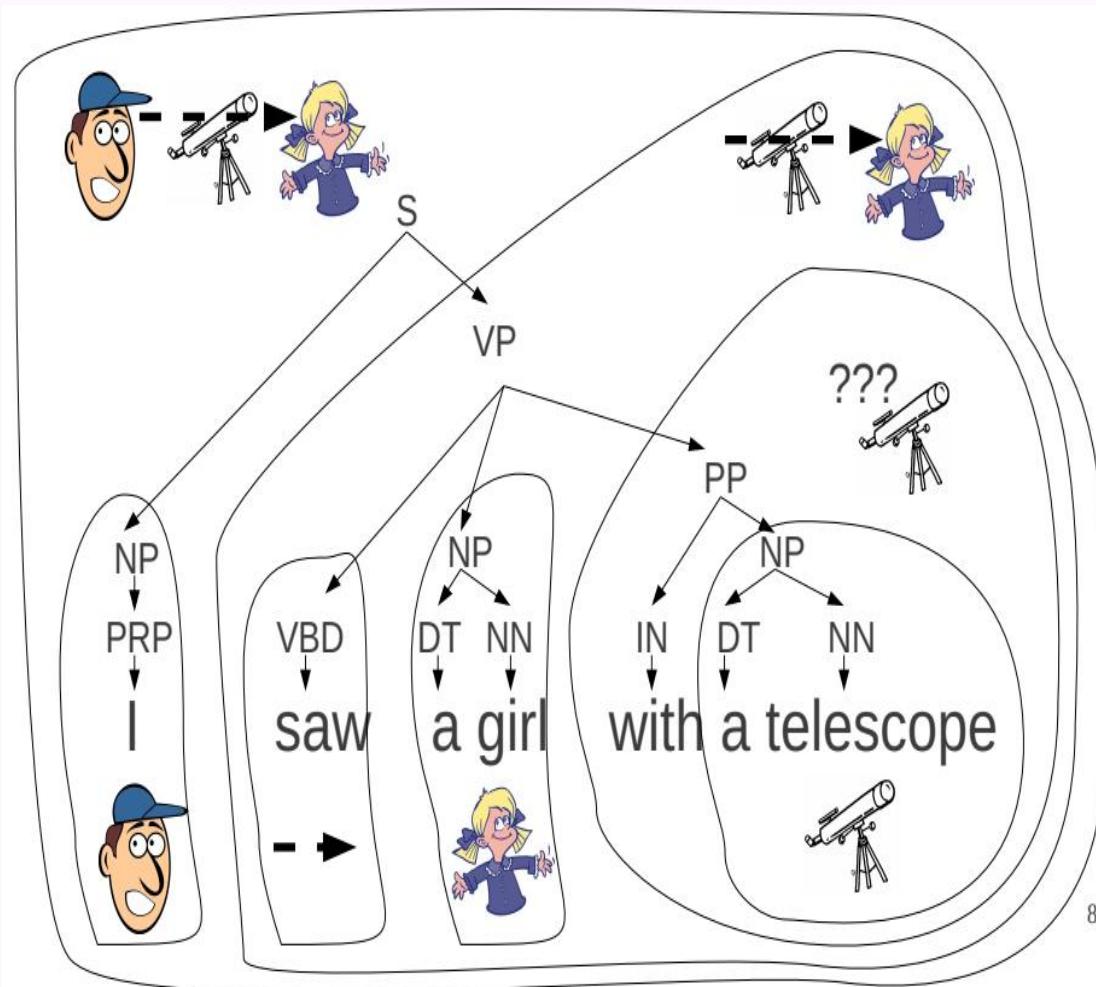




Parsing Struktur Frase: Struktur Rekursif

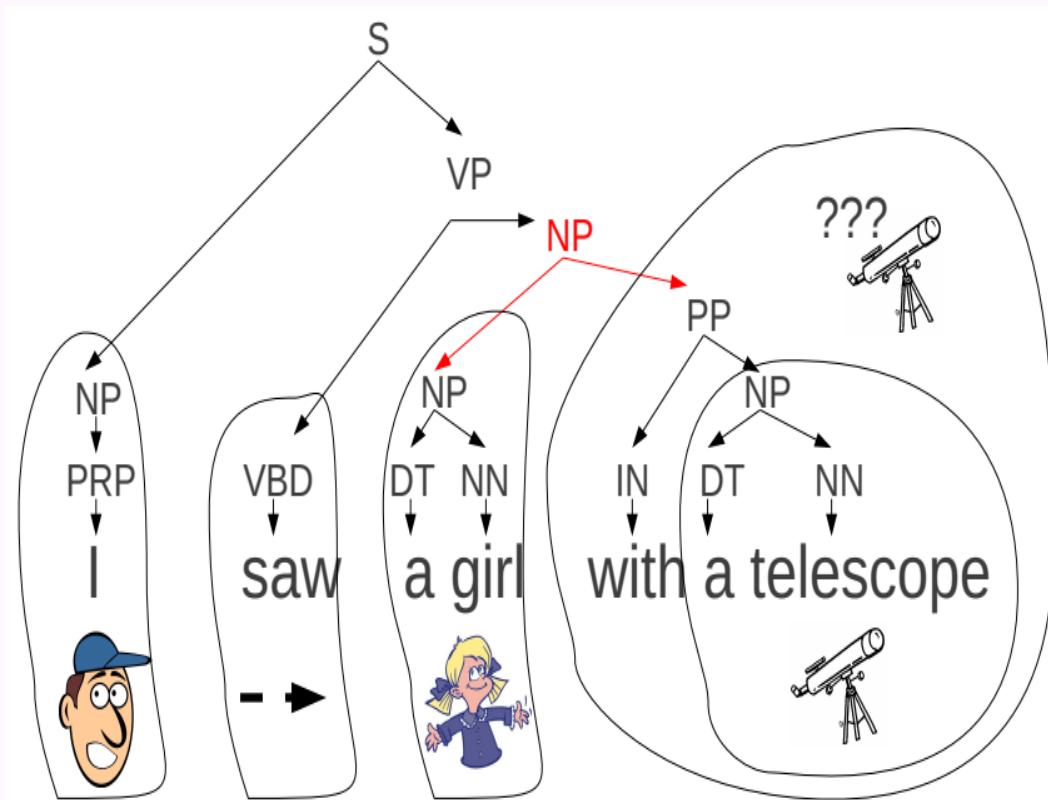


Parsing Struktur Frase: Struktur Rekursif



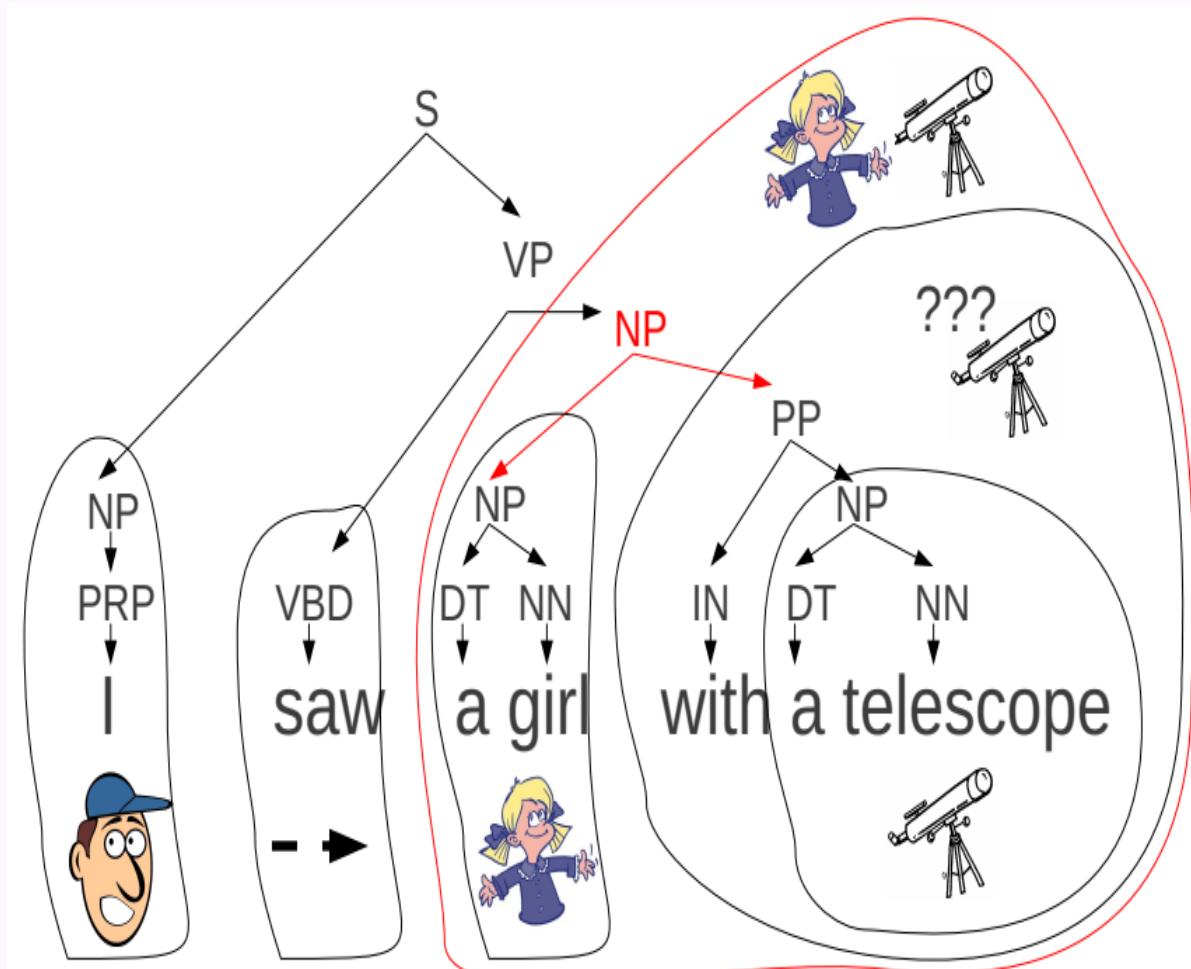


Strukur Beda, Interpretasi Beda

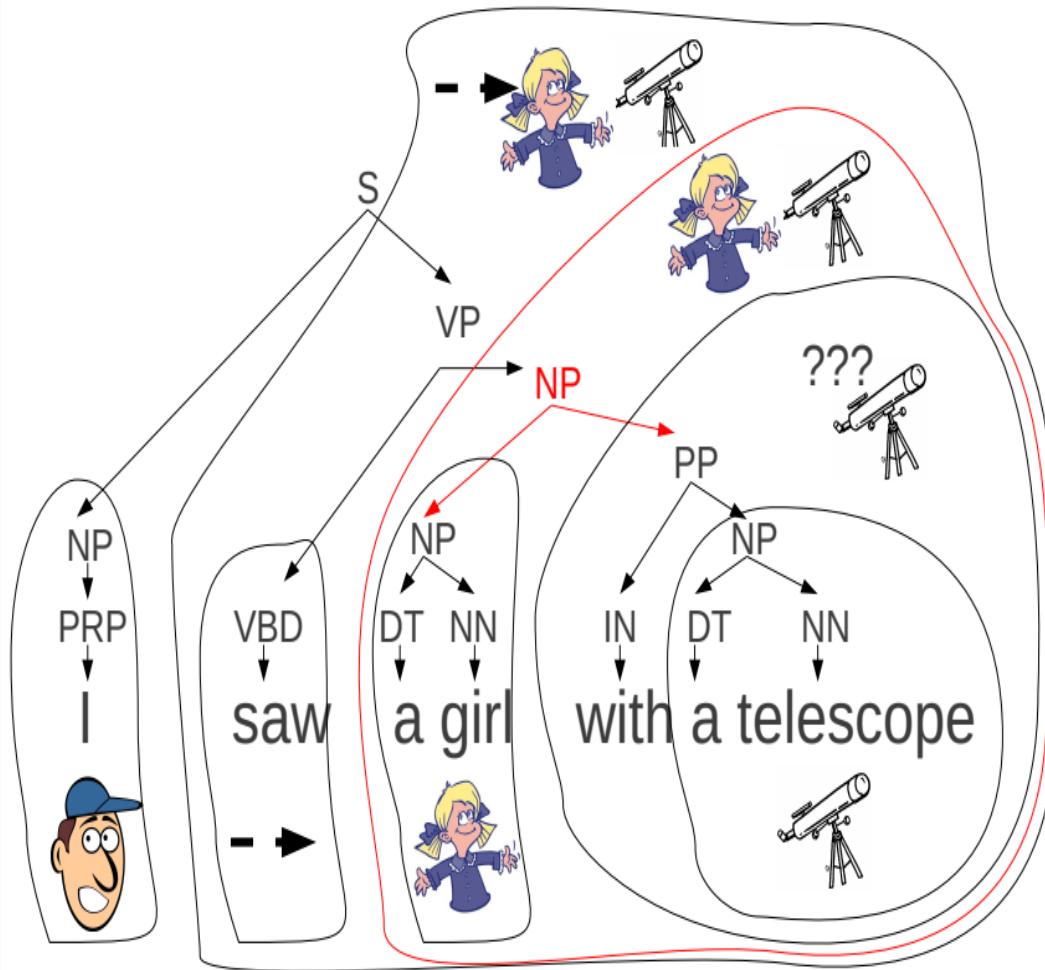




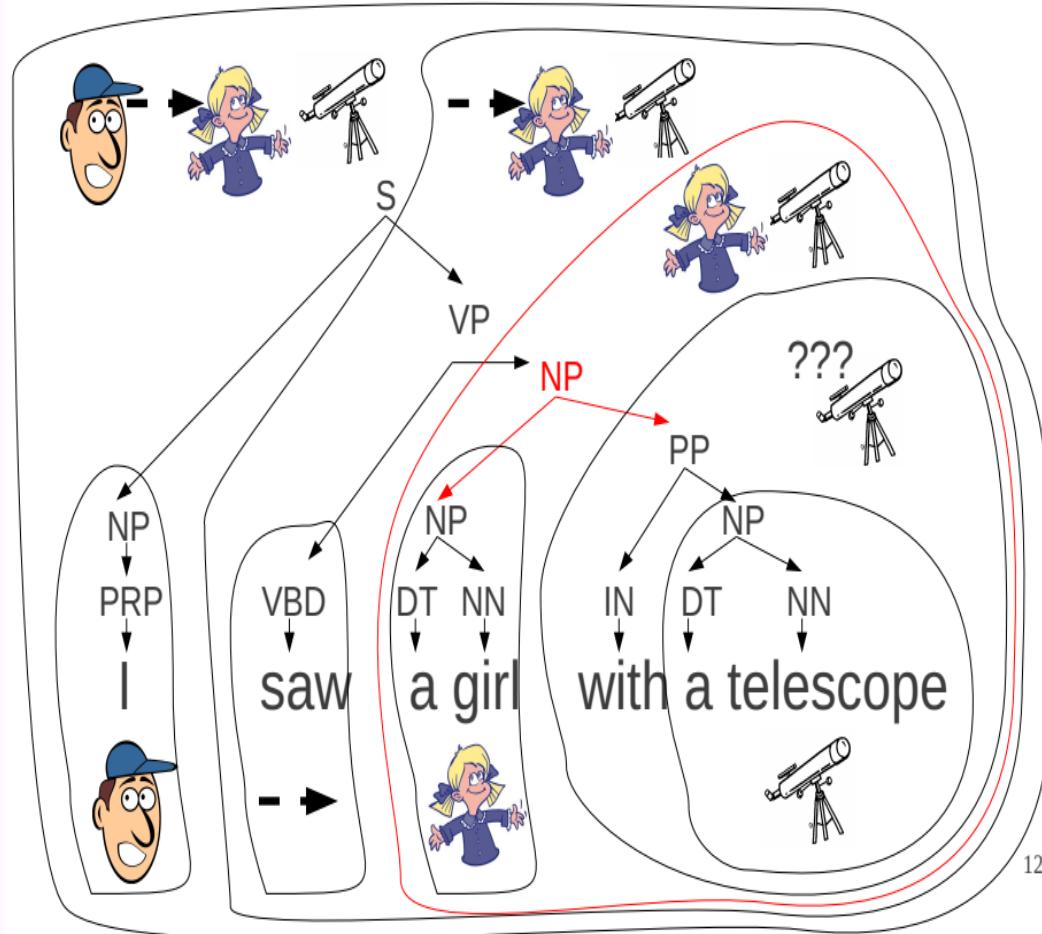
Strukur Beda, Interpretasi Beda



Strukur Beda, Interpretasi Beda

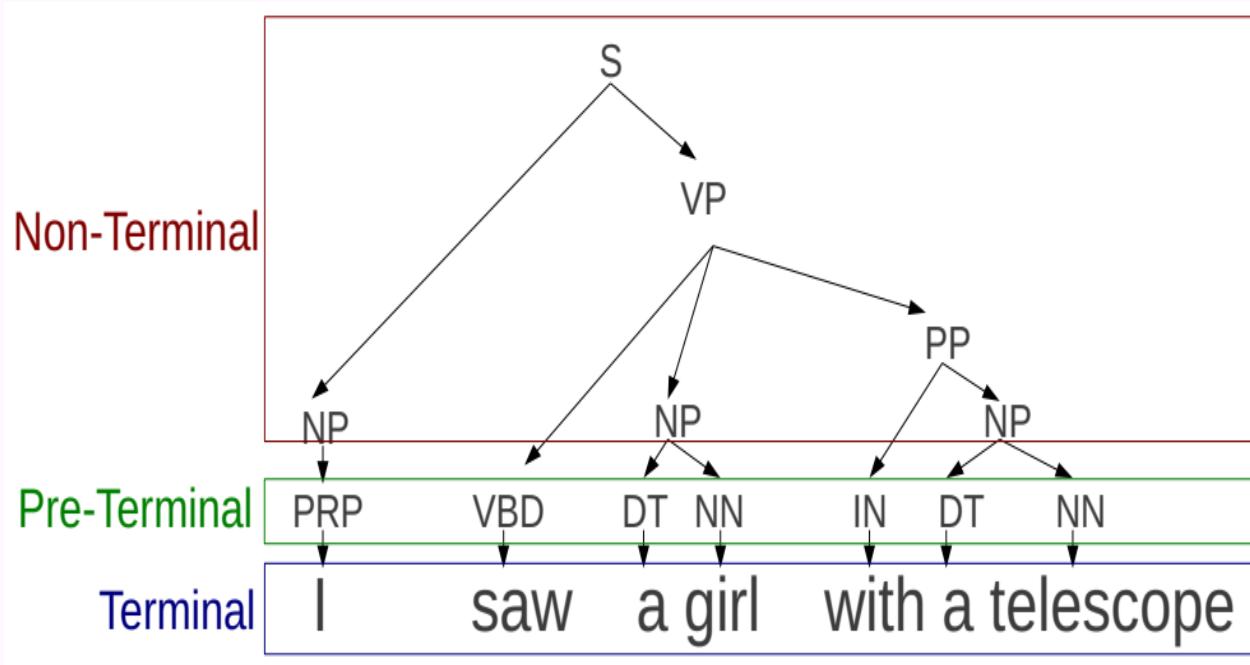


Strukur Beda, Interpretasi Beda



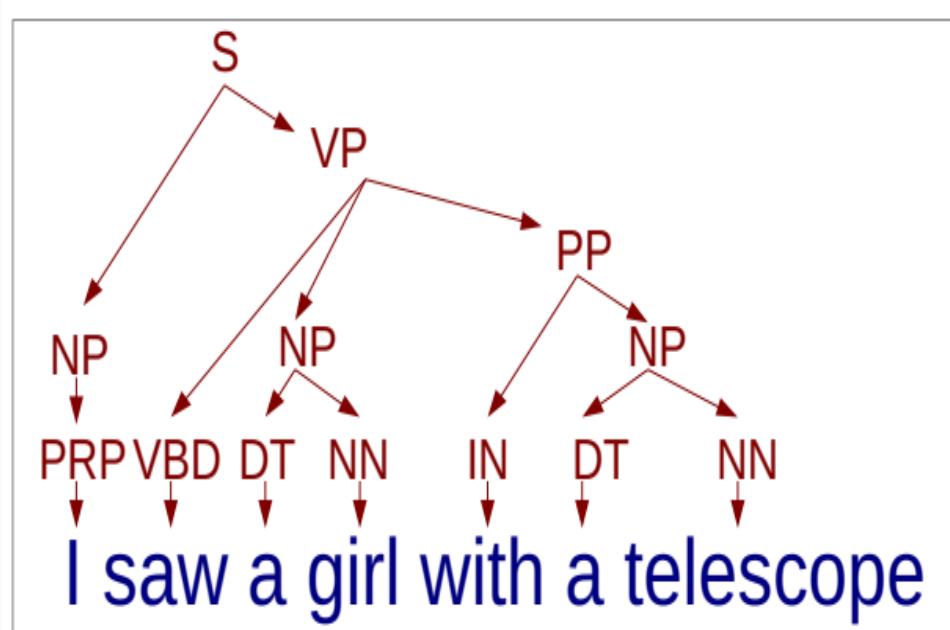


Non Terminal, Pre Terminal dan Terminal



Parsing sebagai suatu Prediksi

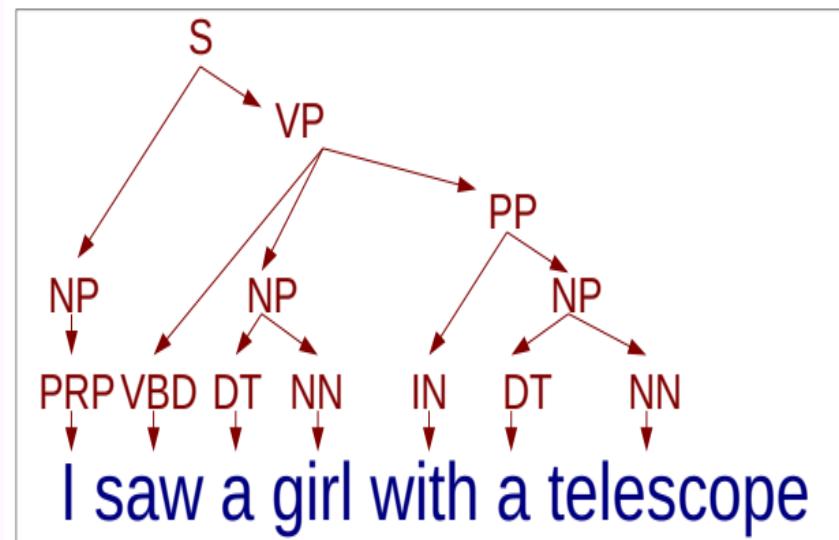
- Diberikan kalimat X, prediksi parse tree Y





Parsing Model Probabilistik

- Diberikan/diketahui kalimat X, prediksi probabilitas terbesar dari parse tree Y

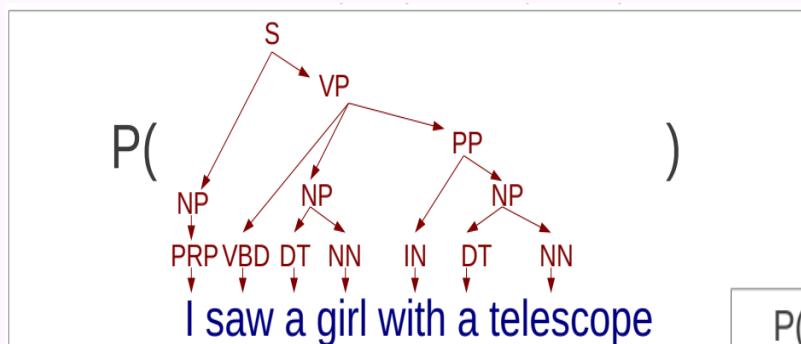


$$\underset{Y}{\operatorname{argmax}} P(Y|X) = \underset{Y}{\operatorname{argmax}} P(Y, X)$$

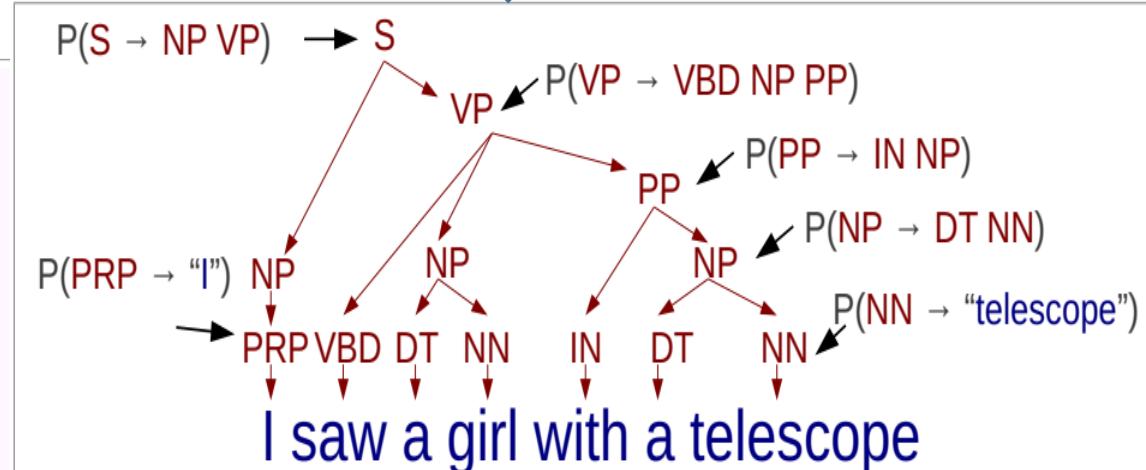
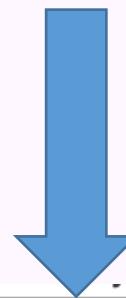
$$\underset{Y}{\operatorname{argmax}} P(Y|X)$$

Probabilistic Context Free Grammar (PCFG)

- Bagaimana mendefinisikan probabilitas gabungan untuk parse tree brkt:



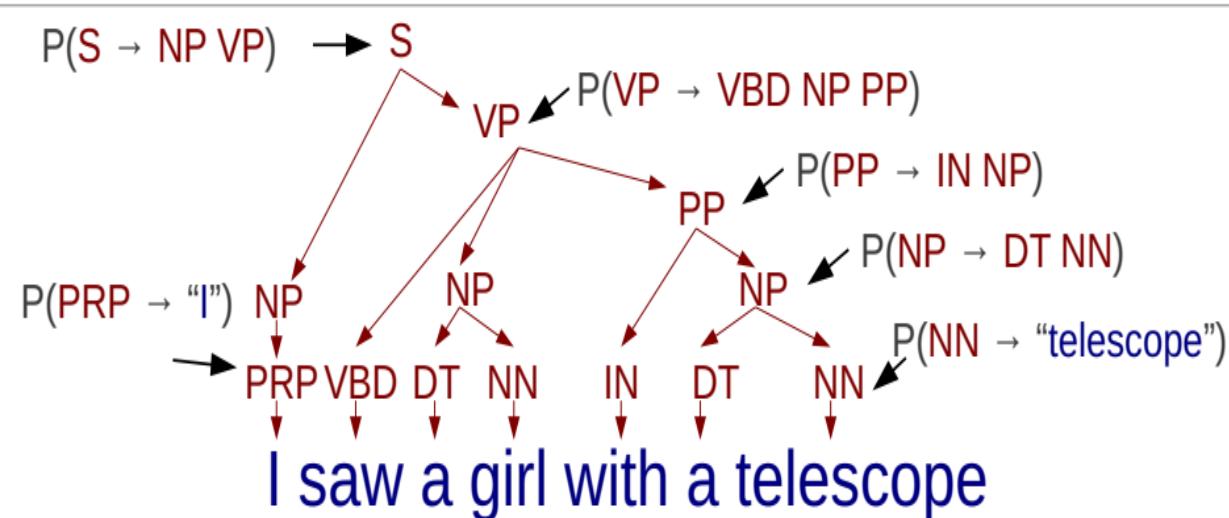
- PCFG: Definisikan probabilitas utk tiap node





Probabilistic Context Free Grammar (PCFG)

- PCFG: Define probability for each node



- Parse tree probability is product of node probabilities

$$\begin{aligned} & P(S \rightarrow NP VP) * P(NP \rightarrow PRP) * P(PRP \rightarrow "I") \\ & * P(VP \rightarrow VBD NP PP) * P(VBD \rightarrow "saw") * P(NP \rightarrow DT NN) \\ & * P(DT \rightarrow "a") * P(NN \rightarrow "girl") * P(PP \rightarrow IN NP) * P(IN \rightarrow "with") \\ & * P(NP \rightarrow DT NN) * P(DT \rightarrow "a") * P(NN \rightarrow "telescope") \end{aligned}$$

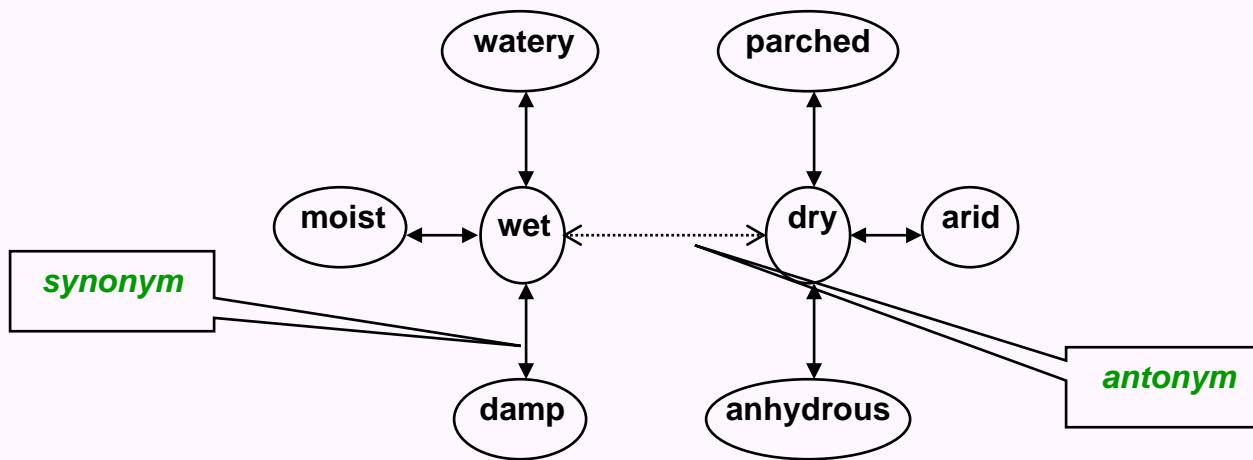


Linguistik

- Progress on **Useful Sub-Goals:**
 - English Lexicon
 - Part-of-Speech Tagging
 - Word Sense Disambiguation
 - Phrase Detection / Parsing



Wordnet



An extensive **lexical network** for the English language

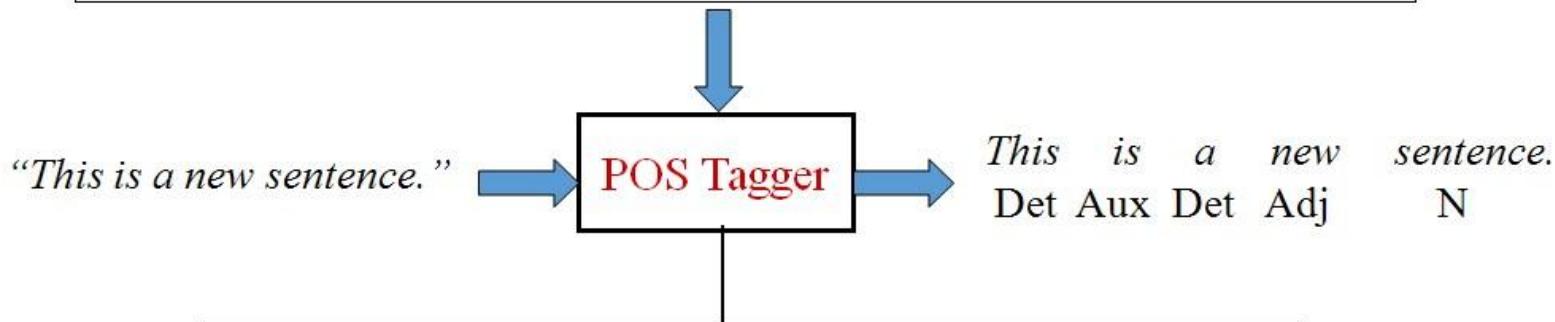
- Contains over **138,838** words.
- Several graphs, one for each **part-of-speech**.
- **Synsets** (synonym sets), each defining a semantic sense.
- **Relationship** information (antonym, hyponym, meronym ...)
- Downloadable for **free** (UNIX, Windows)
- Expanding to **other languages** (Global WordNet Association)
- Funded **>\$3 million**, mainly government (translation interest)
- Founder **George Miller**, National Medal of Science, 1991.



Part of Speech Tagging

Training data (Annotated text)

This	sentence	serves	as	an	example	of	annotated	text...
Det	N	V1	P	Det	N	P	V2	N



Pick the **most likely** tag sequence.

$$p(w_1, \dots, w_k, t_1, \dots, t_k) = \frac{p(t_1 | w_1) \dots p(t_k | w_k) p(w_1) \dots p(w_k)}{\prod_{i=1}^k p(w_i | t_i) p(t_i | t_{i-1})}$$

Independent assignment
Most common tag

Partial dependency
(HMM)



Parsing

Choose most likely parse tree... Probability of this tree=0.000015

Probabilistic CFG

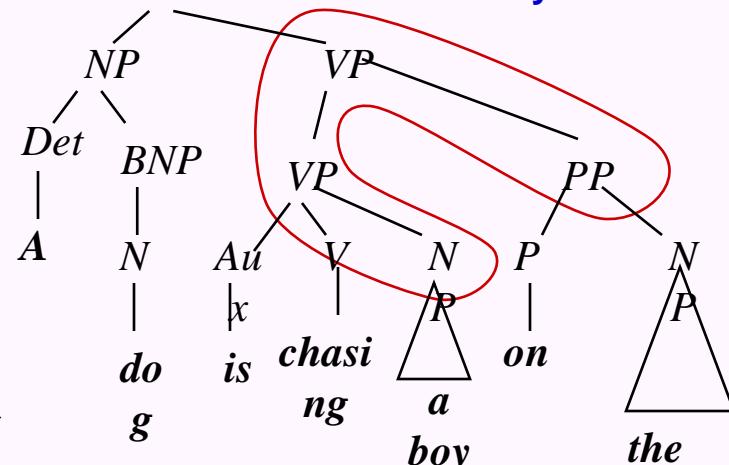
Grammar	$S \rightarrow NP\ VP$ 1.0 $NP \rightarrow Det$ 0.3 BNP 0.4 $NP \rightarrow BNP$ 0.3 $NP \rightarrow NP\ PP$... $BNP \rightarrow N$ $VP \rightarrow V$ $VP \rightarrow Aux\ V$... NP 1.0 $VP \rightarrow VP\ PP$
----------------	---

Lexicon	$PP \rightarrow P\ NP$ 0.01
----------------	-----------------------------

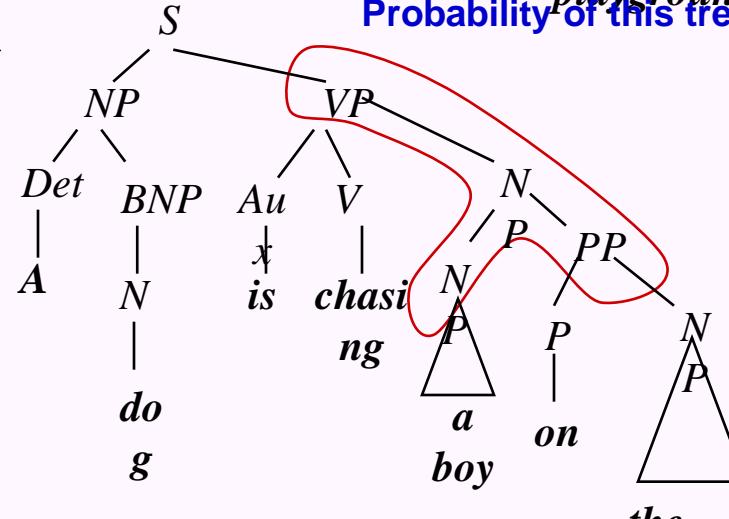
Lexicon	$V \rightarrow chasing$ 0.003 $Aux \rightarrow is$ $N \rightarrow dog$... $N \rightarrow boy$ $N \rightarrow$... $playground$
----------------	--

$Det \rightarrow the$

$Det \rightarrow a$



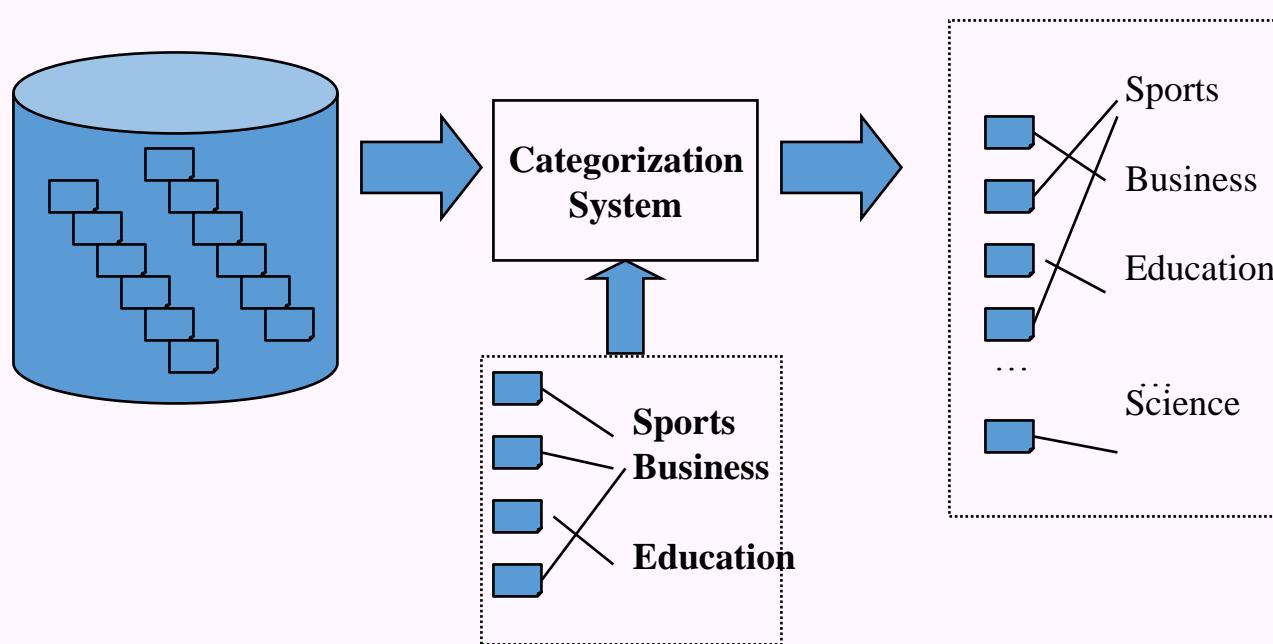
Probability of this tree=0.000011



playground

Kategorisasi Tekst

- Pre-given categories dan contoh dokument berlabel
(Categories may form hierarchy)
- Mengklasifikasikan dokumen baru
- Masalah klasifikasi standar (supervised learning)



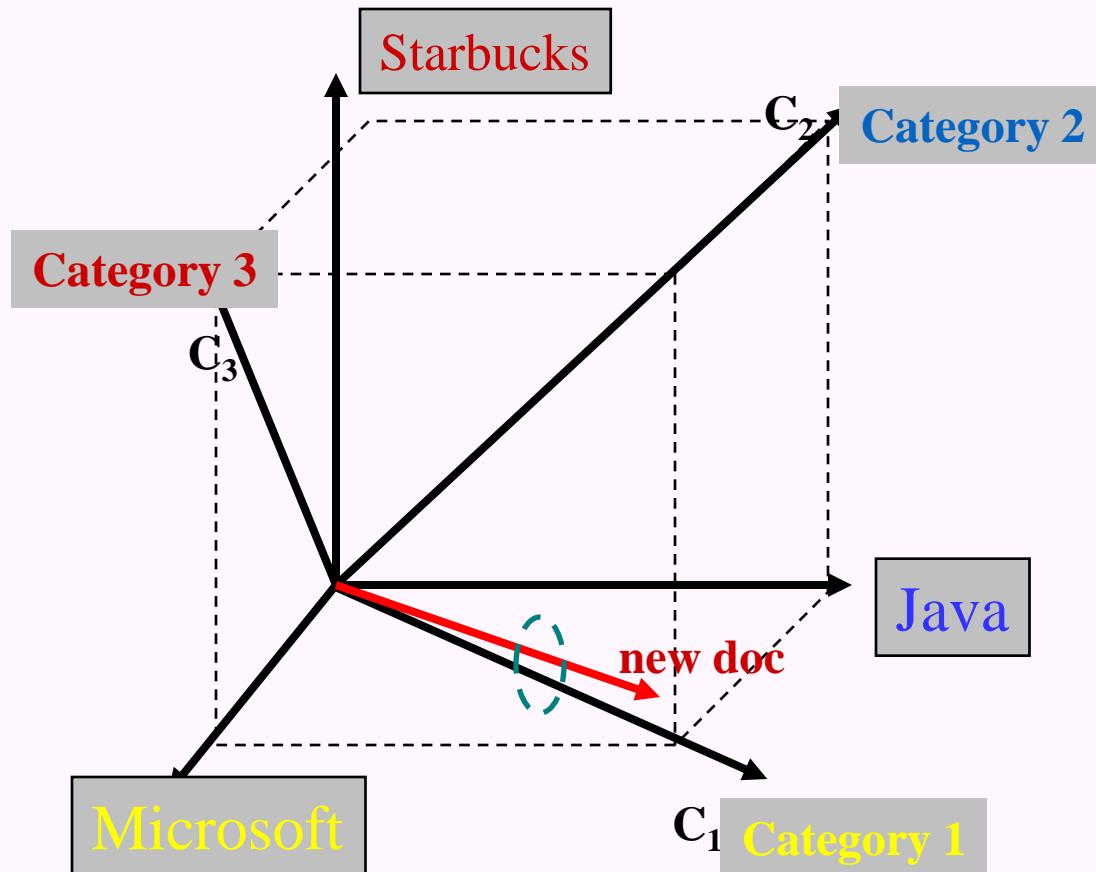


Model Ruang Vector (Vector Space Model)

- MerRepresentasikan suatu dokumen dengan vektor term
 - Term: basic concept, e.g., word or phrase
 - Setiap term mendefinisikan satu dimensi
 - N terms mendefinisikan ruang berdimensi N (N-dimensional space)
 - Element dari vector dipetakan ke bobot term (term weight)
 - E.g., $d = (x_1, \dots, x_N)$, x_i is “importance” of term i
- New document is assigned to the most likely category based on vector similarity.



Ilustrasi Model Ruang Vektor





Yang tidak dispesifikasikan oleh Model Ruang Vektor

- How to select terms to capture “basic concepts”
 - Word stopping
 - e.g. “a”, “the”, “always”, “along”
 - Word stemming
 - e.g. “computer”, “computing”, “computerize” => “compute”
 - Latent semantic indexing
- How to assign weights
 - Not all words are equally important: Some are more indicative than others
 - e.g. “algebra” vs. “science”
- How to measure the similarity



Cara Mengukur Similaritas

- Given two document

$$D_i = (w_{i1}, w_{i2}, \dots, w_{iN}) \quad D_j = (w_{j1}, w_{j2}, \dots, w_{jN})$$

- Similarity definition

- dot product

$$\text{Sim}(D_i, D_j) = \sum_{t=i}^N w_{it} * w_{jt}$$

- normalized dot product (or cosine)

$$\text{Sim}(D_i, D_j) = \frac{\sum_{t=i}^N w_{it} * w_{jt}}{\sqrt{\sum_{t=1}^N (w_{it})^2 * \sum_{t=1}^N (w_{jt})^2}}$$



Contoh Ilustrasi

doc1

text
mining
search
engine
text

$$\text{Sim(newdoc, doc1)} = 4.8 * 2.4 + 4.5 * 4.5$$

doc2

travel
text
map
travel

$$\text{Sim(newdoc, doc3)} = 0$$

doc3

government
president
congress

text mining travel map search engine govern president congress

IDF(faked) 2.4 4.5 2.8 3.3 2.1 5.4 2.2 3.2 4.3

doc1 2(4.8) 1(4.5) 1(2.1) 1(5.4)

doc2 1(2.4) 2 (5.6) 1(3.3)

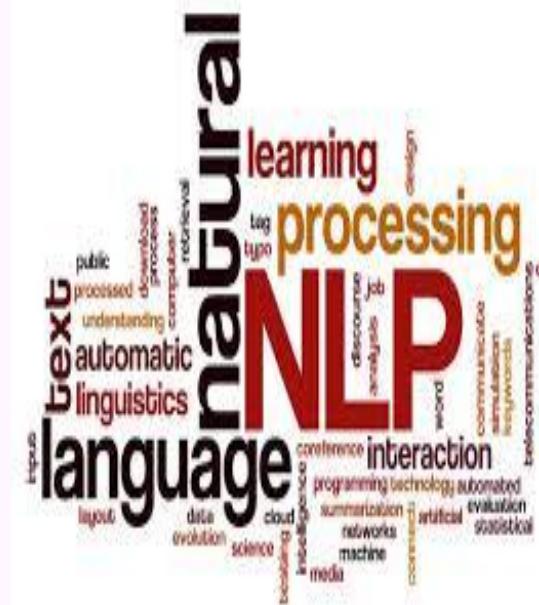
doc3 1 (2.2) 1(3.2) 1(4.3)

.....

newdoc 1(2.4) 1(4.5)

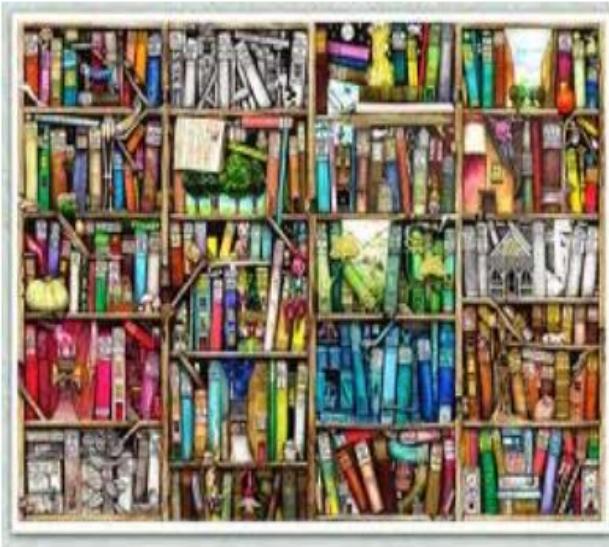
To whom is
newdoc more
similar?

Analisis sintak, topic modelling & sosmed mining





Topic Model



- Mencari pola topik tersembunyi dalam koleksi dokumen menggunakan metode statistika
- Menganotasi dokumen dgn topik-topik yang ditemukan tsb
- Memanfaatkan anotasi topik utk mengelola, menyimpulkan , mencari teks ...



Contoh Topic

Topic 247

word	prob.
DRUGS	.069
DRUG	.060
MEDICINE	.027
EFFECTS	.026
BODY	.023
MEDICINES	.019
PAIN	.016
PERSON	.016
MARIJUANA	.014
LABEL	.012
ALCOHOL	.012
DANGEROUS	.011
ABUSE	.009
EFFECT	.009
KNOWN	.008
PILLS	.008

Topic 5

word	prob.
RED	.202
BLUE	.099
GREEN	.096
YELLOW	.073
WHITE	.048
COLOR	.048
BRIGHT	.030
COLORS	.029
ORANGE	.027
BROWN	.027
PINK	.017
LOOK	.017
BLACK	.016
PURPLE	.015
CROSS	.011
COLORED	.009

Topic 43

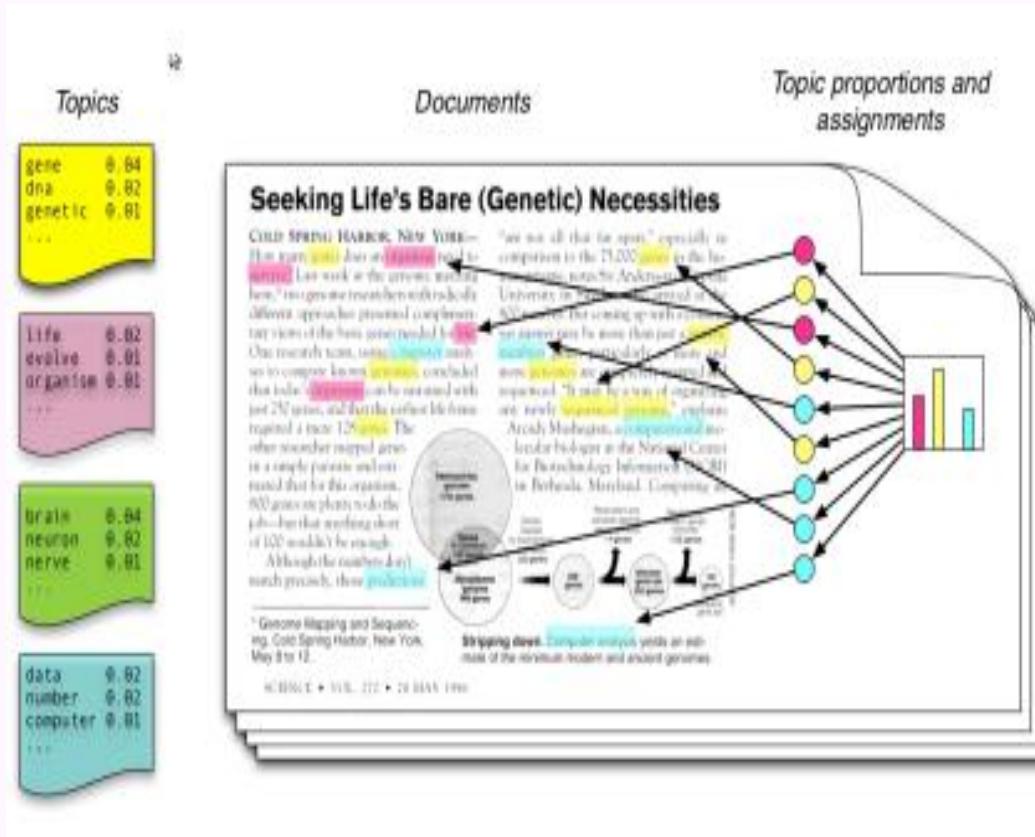
word	prob.
MIND	.081
THOUGHT	.066
REMEMBER	.064
MEMORY	.037
THINKING	.030
PROFESSOR	.028
FELT	.025
REMEMBERED	.022
THOUGHTS	.020
FORGOTTEN	.020
MOMENT	.020
THINK	.019
THING	.016
WONDER	.014
FORGET	.012
RECALL	.012

Topic 56

word	prob.
DOCTOR	.074
DR.	.063
PATIENT	.061
HOSPITAL	.049
CARE	.046
MEDICAL	.042
NURSE	.031
PATIENTS	.029
DOCTORS	.028
HEALTH	.025
MEDICINE	.017
NURSING	.017
DENTAL	.015
NURSES	.013
PHYSICIAN	.012
HOSPITALS	.011

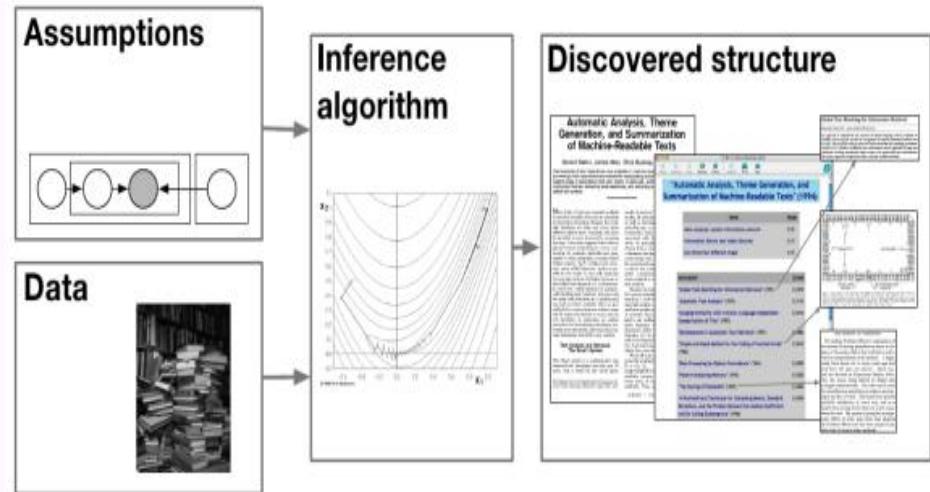
Figure 1. An illustration of four (out of 300) topics extracted from the TASA corpus.

Topic Model atas Dokumen



Topic Model: Struktur

- Observasi: koleksi teks
- Asumsi: teks telah digenerate menurut ke suatu/beberapa model
- Output: model yang telah menGenerate teks





Contoh Topic



SKY WATER TREE
MOUNTAIN PEOPLE



SCOTLAND WATER
FLOWER HILLS TREE



SKY WATER BUILDING
PEOPLE WATER



FISH WATER OCEAN
TREE CORAL



PEOPLE MARKET PATTERN
TEXTILE DISPLAY



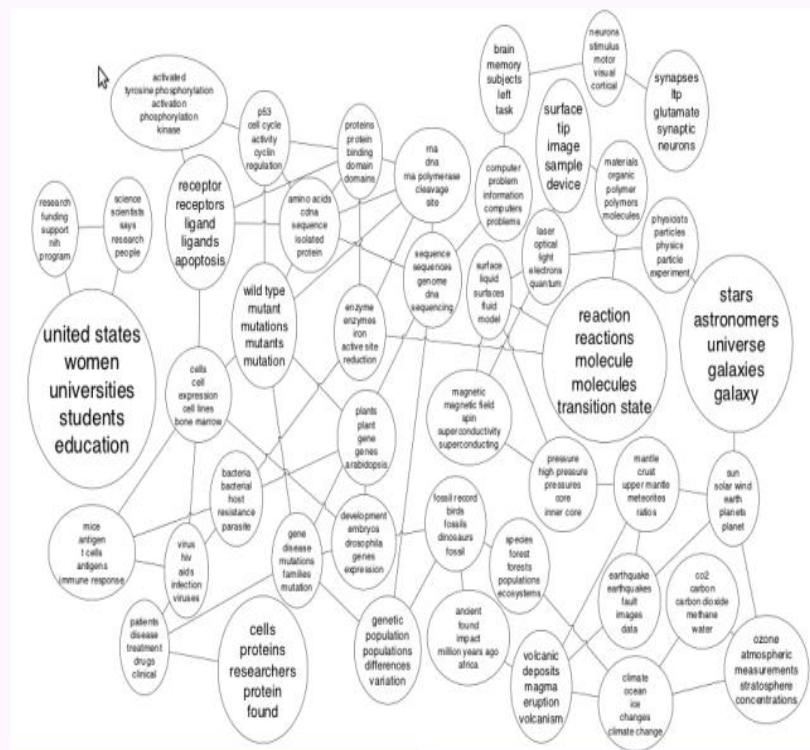
BIRDS NEST TREE
BRANCH LEAVES

Object \equiv bag of words with labels

Contoh Topic

- Komponen dasar:

- himpunan entitas (mis. dokumen, image, individual, gen)
- himpunan relasi (mis. sitasi/kutipan/referensi, co-tag, pertemanan)





Topic Discovery

Topic : “Religion”

Post body							
romney	huckabee	muslim	political	hagee	cabinet	mitt	
consider	true.	anti	problem	course	views	life	
real	speech	moral	answer	jobs	difference	muslims	
hardly	going	christianity					
people	just	American	church	believe	god	black	
jesus	mormon	faith	jews	right	religious	point	
say	mormons						
religion	think	know	really	christian	obama	white	
wright	way	said	good	world	science	time	
dawkins	human	man	things	fact	years	mean	
atheists	blacks	christians					
Post comments							

Finding related, but not obvious words
unique to this corpus



Topic Discovery

Topic : “Iraq War”

Post body

kind
foreign countries forces international presence political states
shiite john role need making course problem
main

american iran just iraq people support point
country nuclear world power military really government
war army right iraqi think

israel sadr bush state way oil years
time going good weapons saddam know maliki
want say policy fact said shia troops

Post comments

The blogger discusses strategy, comments focus on “tangibles”



Topic Model dgn LDA



Pengantar LDA

- Ilustrasi LDA:

Misal ada himpunan kalimat:

- (1) *I like to eat broccoli and bananas.*
- (2) *I ate a banana and spinach smoothie for breakfast.*
- (3) *Chinchillas and kittens are cute.*
- (4) *My sister adopted a kitten yesterday.*
- (5) *Look at this cute hamster munching on a piece of broccoli.*

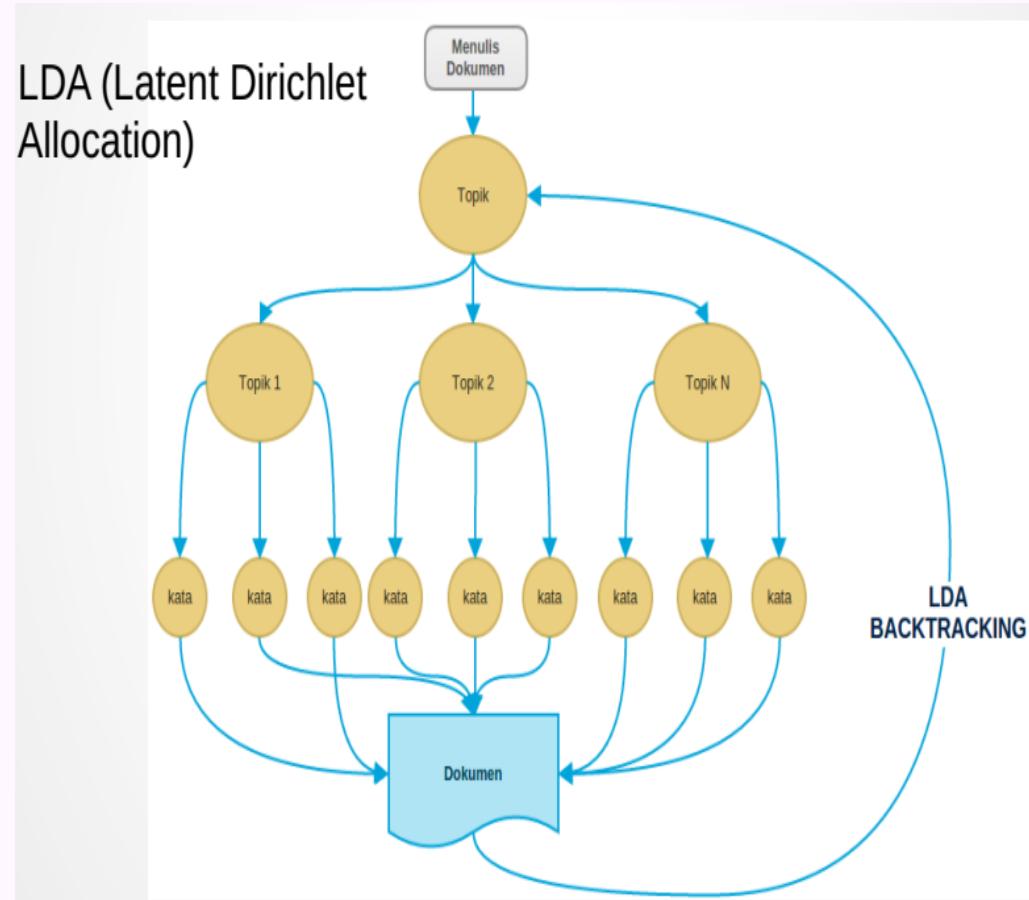
Apa yang dimaksud Latent Dirichlet allocation? Yakni cara untuk menemukan topik secara otomatis yang dimiliki kalimat. Sebagai contoh, diberikan kalimat tsbt dan ditanya tentang 2 topik, LDA akan menghasilkan luaran sbb:

- Kalimat 1 and 2: 100% Topic A
- Kalimat 3 and 4: 100% Topic B
- Kalimat 5: 60% Topic A, 40% Topic B
- Topic A: 30% broccoli, 15% bananas, 10% breakfast, 10% munching, ... (dapat diinterpretasikan bahwa topic a tentang **makanan**)
- Topic B: 20% chinchillas, 20% kittens, 20% cute, 15% hamster, ... (dapat diinterpretasikan bahwa topic B to tentang **hewan imut-imut**)



Pengantar LDA

- LDA merepresentasi dokumen sebagai gabungan topik-topik dari kumpulan kata-kata dengan probabilitas tertentu. Jika kita menulis dokumen:





Perhitungan dan Notasi LDA

Notasi:

- Kata: unit dasar suatu data diskrit, didefinisikan sebagai item dari kosakata berindeks $\{1, \dots, V\}$
- Dokumen: sekuen kata N , $w = (w_1, w_2, \dots, w_N)$, kata ke- n dalam sekuen
- Korpus: koleksi dokumen M , $D = (w_1, w_2, \dots, w_M)$

LDA berasumsi proses generatif utk setiap dokumen w dalam korpus D , sbb:

- Pilih $N \sim \text{Poisson}(\xi)$
- Pilih $\sim \text{Dir}(\alpha)$
- Utk setiap N dari kata w_n :
 - Pilih topik $z_n \sim \text{Multinomial}(\theta)$
 - Pilih sebuah kata w_n dari $p(w_n | z_n, \beta)$, probabilitas multinomial bersyarat dalam topic z_n .



Komparasi Topic

- Topik Model yang digenerasi selanjutnya perlu dibandingkan dengan penilaian pakar
- Hasil penilaian pakar menjadi training set baru sebagai data baru utk proses machine learning

Contoh output perhitungan topic model dgn LDA

NO. TEMA	VERBATIM CCD	TEMA/LABEL/ANOTASI	TOPIK MODEL
1	<p>"Tapi sebelum kita mulai wawancaranya mbak, bisa tolong di perkenalkan dulu nama mbak?"</p> <p>"perkenalkan nama saya <u>Sheba</u>, umur saya <u>27</u> tahun, saya <u>ibu rumah tangga</u>." (mengatur posisi duduk)</p>	Latar belakang significant other	topik 1: 1.000*perkenalkan
2	<p>"subjek kalau dirumah seperti apa mbak?"</p> <p>"dia kalau dirumah <u>biasa aja</u>, dia <u>berangkat kerja pagi pulang malam</u>, dia juga <u>di rumah blom lama</u>, karena <u>sebelumnya dia ngekost</u>, kalau kerumah paling cuma satu minggu sekali atau ngak kalau lagi ada tanggal merah aja, kalau malem pulang kerja saya <u>jarang ketemu dia</u>, karena saya pas jam dia pulang juga udah masuk kamar."</p> <p>"jadi mbak jarang ketemu subjek kalau udah malam?"</p> <p>"<u>Iya jarang</u>, itu juga kalau sekali-<u>sekali aja kalau saya</u> kebetulan <u>tidur agak malem</u>, kalau ketemu lagi paling pagi bareng suami saya kerja"</p>	Pendapat significant other mengenai subjek	<p>topik 1: 0.156*ketemu + 0.151*pulang + 0.145*kerja + 0.142*malem + 0.139*pagi + 0.138*malam + 0.130*dirumah</p> <p>topik 2: 0.174*kerja + 0.172*pulang + 0.170*ketemu + 0.123*dirumah + 0.121*malam + 0.120*pagi + 0.119*malem</p>

Contoh output perhitungan topic model dgn LDA

NO. TEMA	VERBATIM CCD	TEMA/LABEL/ANOTASI	TOPIK MODEL
3	<p>"subjek pernah cerita ngak soal kerjaan dia di kantor sama mbak?"</p> <p><u>"Pernah sekali-sekali dia cerita tentang kerjaannya,</u> biasanya tentang personal <u>karyawan disitu</u>, misalnya ada sesuatu yang lucu tentang <u>temen kantornya yang kalau ngomong pake bahasa inggris kurang</u> ngak <u>jelas</u> tapi kalau masalah kerjaan, paling masalah <u>kerjaannya</u> yang <u>banyak aja.</u>"</p>	Subjek bercerita mengenai pekerjaannya significant other	topik 1: 0.250*masalah + 0.250*kerjaannya + 0.250*kerjaan + 0.250*cerita



NO. TEMA	VERBATIM CCD	TEMA/LABEL/ANOTASI	TOPIK MODEL
4	<p>"kalau di rumah subjek pernah ngak menunjukkan gejala stres kerja?</p> <p>"Stres kerja ya, paling kayaknya dia <u>males berdesak-desakan aja</u> <u>kalau berangkat dan pulang kerja di kereta</u>, jadi stres kerja yang keliatan dari dia dia keliatan kecapeaan."</p>	Stres kerja yang subjek tunjukkan di depan significant other	topik 1: 0.333*stres + 0.333*kerja + 0.333*keliatan





Analisis Sentimen



Analisis Sentimen

- Is the attitude of this text positive or negative?
- More complex:
 - Rank the attitude of this text from 1 to 5
- Advanced:
 - Detect the target, source, or complex attitude types

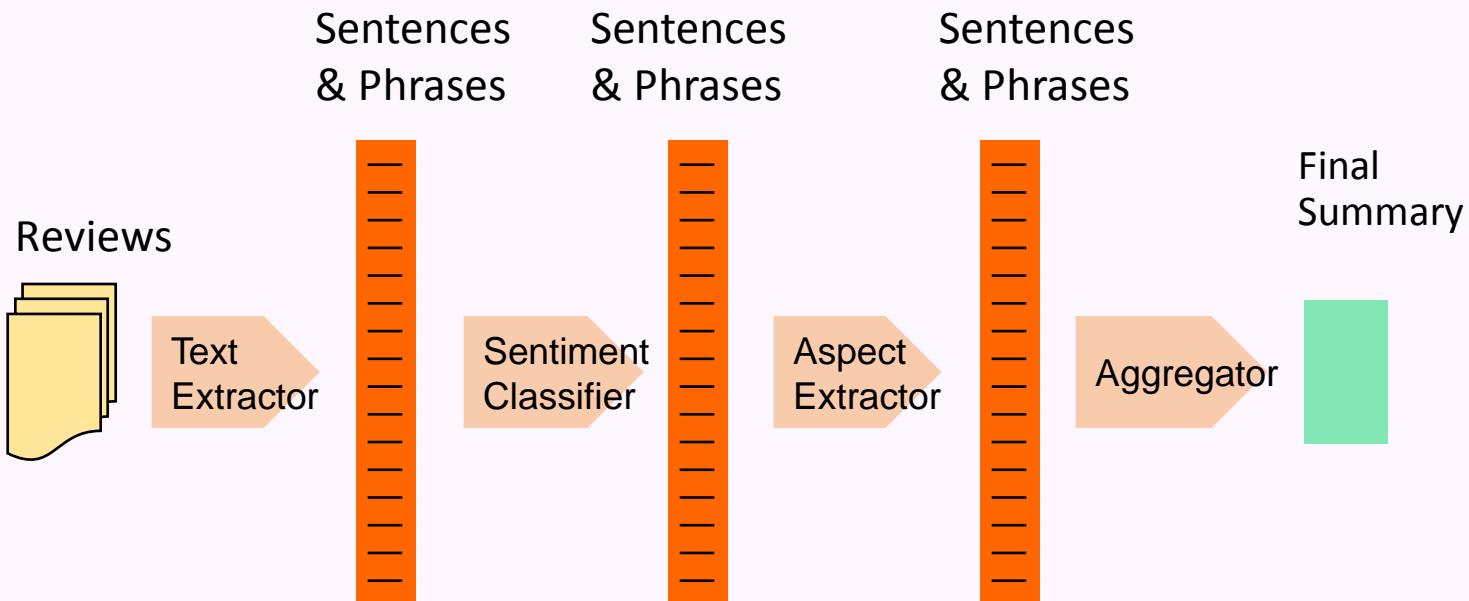


Tipologi Keadaan/Kondis Afektif Scherer

- **Emotion:** brief organically synchronized ... evaluation of a major event
 - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood:** diffuse non-caused low-intensity long-duration change in subjective feeling
 - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances:** affective stance toward another person in a specific interaction
 - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes: enduring, affectively colored beliefs, dispositions towards objects or persons**
 - *liking, loving, hating, valuing, desiring*
- **Personality traits:** stable personality dispositions and typical behavior tendencies
 - *nervous, anxious, reckless, morose, hostile, jealous*



Tahapan Analisis Sentimen





Computational work on other affective states

- **Emotion:**
 - Detecting annoyed callers to dialogue system
 - Detecting confused/frustrated versus confident students
- **Mood:**
 - Finding traumatized or depressed writers
- **Interpersonal stances:**
 - Detection of flirtation or friendliness in conversations
- **Personality traits:**
 - Detection of extroverts



Detection of Friendliness

- Friendly speakers use collaborative conversational style
 - Laughter
 - Less use of negative emotional words
 - More sympathy
 - More agreement
 - Less hedges



TERIMA KASIH