

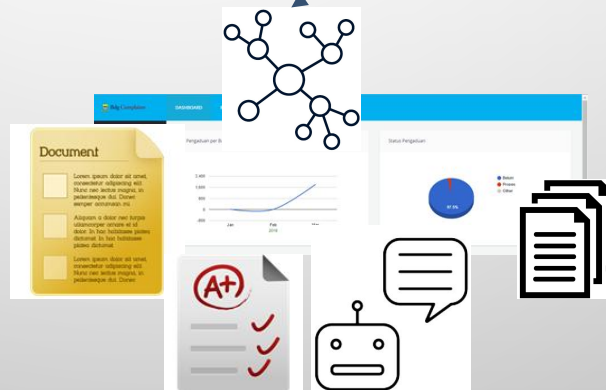


Prosa.ai

AI Powered Text and Speech

# NATURAL LANGUAGE PROCESSING

AYU PURWARIANTI  
AYU@INFORMATIKA.ORG



COMPUTATIONAL  
LINGUISTICS

TEXT MINING

NLU

NLG

**Text is ...**

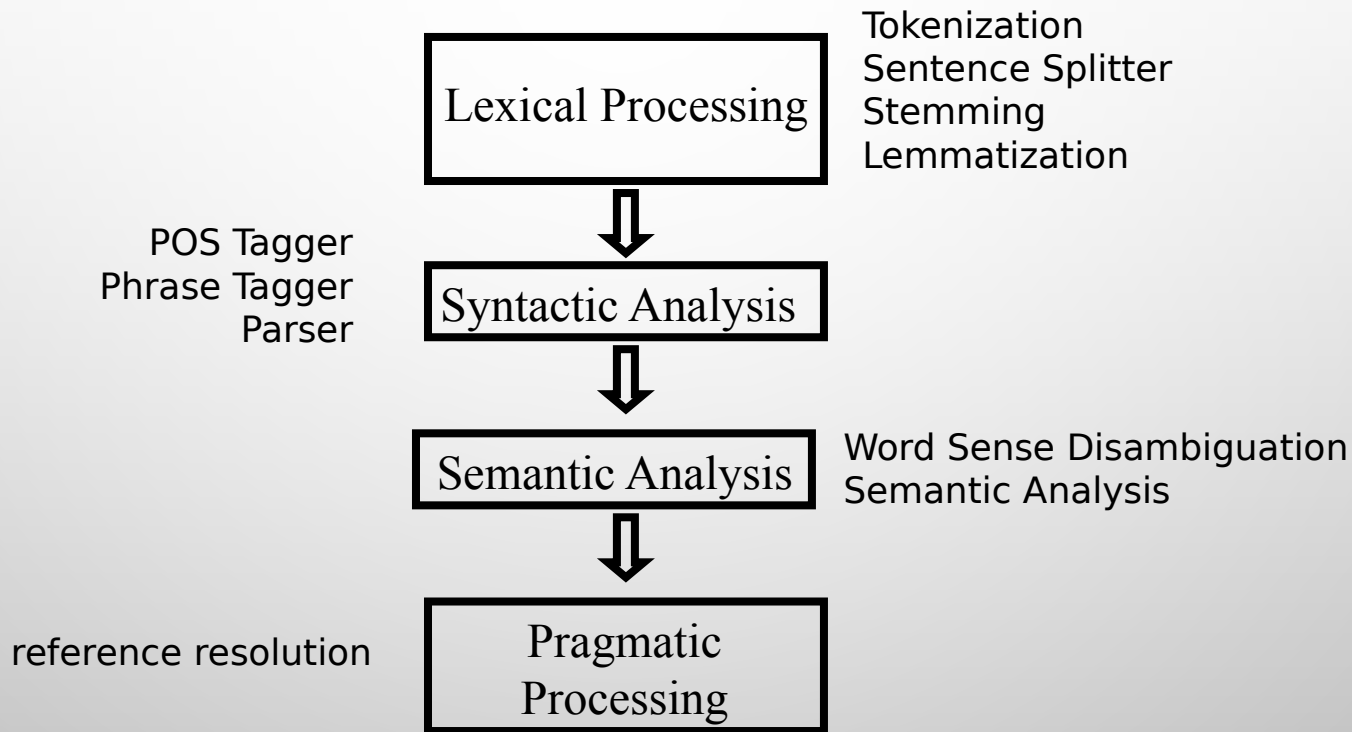
**a string**

**list of tokens**

**having structure**

**represent an intention**

# BASIC NLP TOOLS



# NLP TOOLKIT

- NLTK
  - TOKENIZER, POS TAGGER, NAMED ENTITY TAGGER, PARSER, STEMMING, SEMANTIC ANALYZER, SENTIMENT ANALYZER, MACHINE TRANSLATION, TEXT CLASSIFICATION, CLUSTERING, COLLOCATION
- OPENNLP
  - SENTENCE DETECTION, TOKENIZATION, NAMED ENTITY TAGGER, PART OF SPEECH TAGGER, DOCUMENT CATEGORIZER, CHUNKER, PARSER, COREFERENCE RESOLUTION
- STANFORD CORENLP
  - SENTENCE SPLITTER, TOKENIZATION, POS TAGGER, LEMMATIZATION, NAMED ENTITY TAGGER, PARSER, COREFERENCE RESOLUTION, SENTIMENT ANALYSIS, MENTION DETECTION, OPEN IE

## Prosa NLP API Documentations

Home

Getting Started

APIs Reference ^

[Syntactic Analyzer](#)

Word Normalizer

Named Entities Extractor

News Quotes Extractor

Hate Speech Detector

News Topic Classifier

Document-based Sentiment Analyzer

Concept-based Sentiment Analyzer

Aspect-based Sentiment Analyzer (ABSA)

Dependency Parser

# Syntactic Analyzer

Extracts syntactic information from a given text. The API uses four separate modules:

1. Sentence splitter: Splits the given text by sentences
2. Tokenizer: Splits each split sentence in the given text by tokens
3. POS tagger: Assign part-of-speech tag to each token in the given text
4. Stemmer: Extracts the lemma and the affixes present in each token in the given text

## Request Method

POST

## Request URL

```
1 https://api.prosa.ai/v1/syntax
```

## Table of contents

Request Method

Request URL

Request Header

Request Body

Granularity Enums

Example

Sample Request (JSON)

Sample Response (JSON)

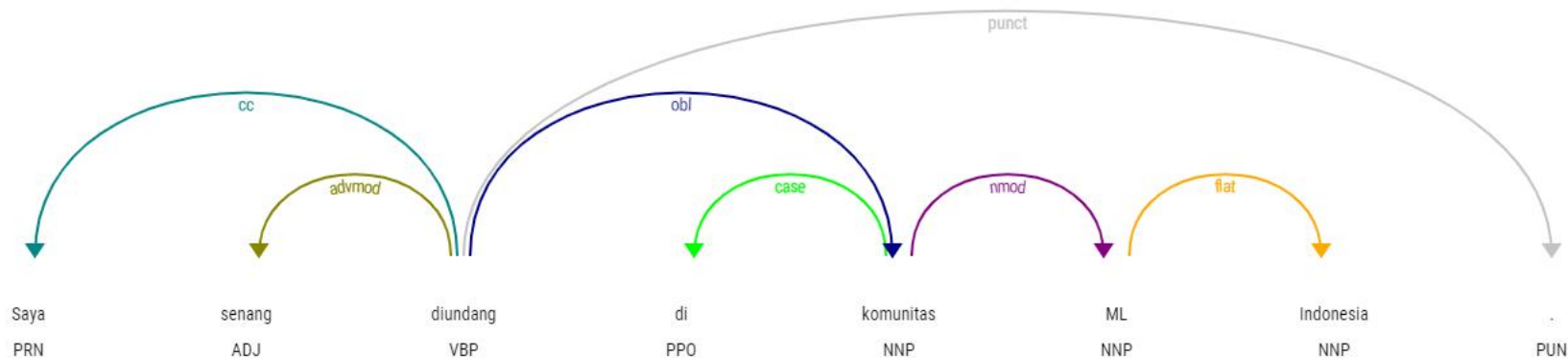
Fine Grained POS Tags

Coarse POS Tags

Version History

Questions?

## Dependency Parser



## Named Entity Recognizer

ORGANIZATION

Saya senang diundang di komunitas **ML Indonesia.**

LEXALYTICS



bing™

Google



**LingvoSoft**  
THE JOY OF UNDERSTANDING



# SEVERAL BENEFITS OF USING NLP

improve user experience

- spelling checker, completer, recommendation

automate support

- chatbot

monitor and analyze feedback

- sentiment analysis

make things faster

- search engine, summarization

make things easier

- virtual assistant, machine translation

# EMAIL



## Spam Filtering

### **Text classification**

Spam vs Not Spam



## Email type

### **Text classification**

Meeting invitation,  
personal question,  
broadcast information



## Scheduler

### **Information Extraction**

Date, time, location,  
agenda





## Summarization

Summary on the  
discussion

# DEFINITION OF TEXT CLASSIFICATION

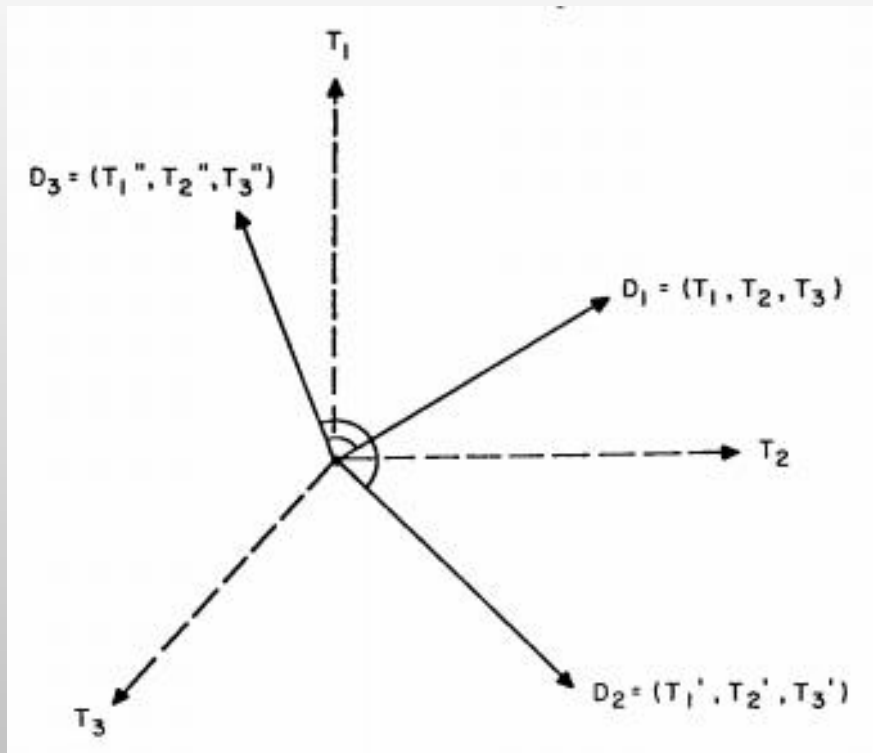
- SYSTEM TO CLASSIFY OR TEXT (OR DOCUMENT) INTO A CERTAIN LABEL
- EXAMPLE:
  - SPAM FILTERING:
    - LABEL THE INPUT (SUCH AS SMS OR EMAIL) INTO SPAM OR NOT SPAM LABELS
  - NEWS CLUSTERING:
    - GROUP THE NEWS ARTICLE INTO CERTAIN CATEGORY SUCH AS POLITICS, SPORTS, ECONOMY, ETC
- = “TEXT CLASSIFICATION”, “DOCUMENT CLASSIFICATION”, “TEXT CATEGORIZATION”, “DOCUMENT CATEGORIZATION”

# SPAM FILTERING

- CLASSIFY TEXT INPUT INTO 2 CLASSES: SPAM OR NOT SPAM
- SIMPLE METHOD:
  - USE A WORD LIST THAT CONSISTS OF THE MOST FREQUENT WORDS IN THE SPAM TEXT ; AND THEN USE THRESHOLD RULE, FOR EXAMPLE:
    - IF  $\geq 75\%$  WORDS IN INPUT TEXT ARE AVAILABLE IN THE WORD LIST THEN THE TEXT IS CLASSIFIED AS SPAM
  - WORD LIST  IS CALLED **FEATURE**
  - RULE WITH  $\geq 75\%$  AS THE THRESHOLD  IS CALLED THE **CLASSIFICATION TECHNIQUE** WHICH CAN BE RULE OR MODEL

# VECTOR SPACE MODEL

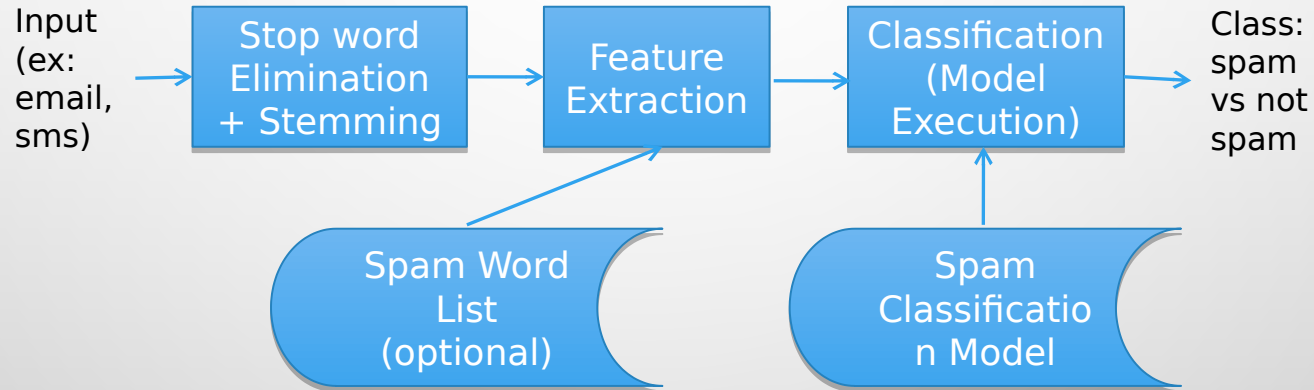
- TEXT INPUT AS VECTOR OF TERMS OR WORDS



# UNIGRAM APPROACH

- LIST OF WORDS TAKEN FROM ALL TEXT INPUTS AS THE FEATURES
- PROBLEMS: FEATURE SIZE
- DECREASE THE FEATURE SIZE, HOW?
  - STEMMING
  - STOP WORD ELIMINATION
  - TAKE ONLY N-TOP WORDS WITH HIGH WORD WEIGHT
  - SYNONYM

# EXAMPLE OF STATISTICAL BASED SPAM FILTERING



# EXAMPLE OF LEXICAL FEATURE FOR SPAM FILTERING

- INPUT: “INI MAMA ... TOLONG KIRIM PULSA KE NOMOR HP INI”  
(ENGLISH: THIS IS MOM.. PLEASE SEND BALANCE TO THIS PHONE NUMBER”
- WORD LIST: MAMA (MOM), TOLONG (PLEASE), KIRIM (SEND), PULSA (BALANCE), NOMOR (NUMBER), HP (PHONE)
- TRAINING DATA
  - FEATURE: WORDS

Mama (mom)	Tolong (please)	Kirim (send)	Pulsa (balance)	Nomor (number)	Hp (phone)	...	Class
1	1	1	1	1	1	...	Spam
1	0	0	0	0	0	...	Not spam



# LEXICAL FEATURE FOR SPAM FILTERING

- LEXICAL BASED
  - MANUALLY WRITE WORDS OCCUR IN SPAM TEXT
    - BASED ON MANUAL OBSERVATION
    - SMALL DATA SIZE
  - STATISTICAL BASED
    - WORDS IN A WORD LIST HAVE HIGH WORD WEIGHT OR ABOVE CERTAIN THRESHOLD

# WORD LIST (LEXICAL BASED)

- WEAKNESSES:

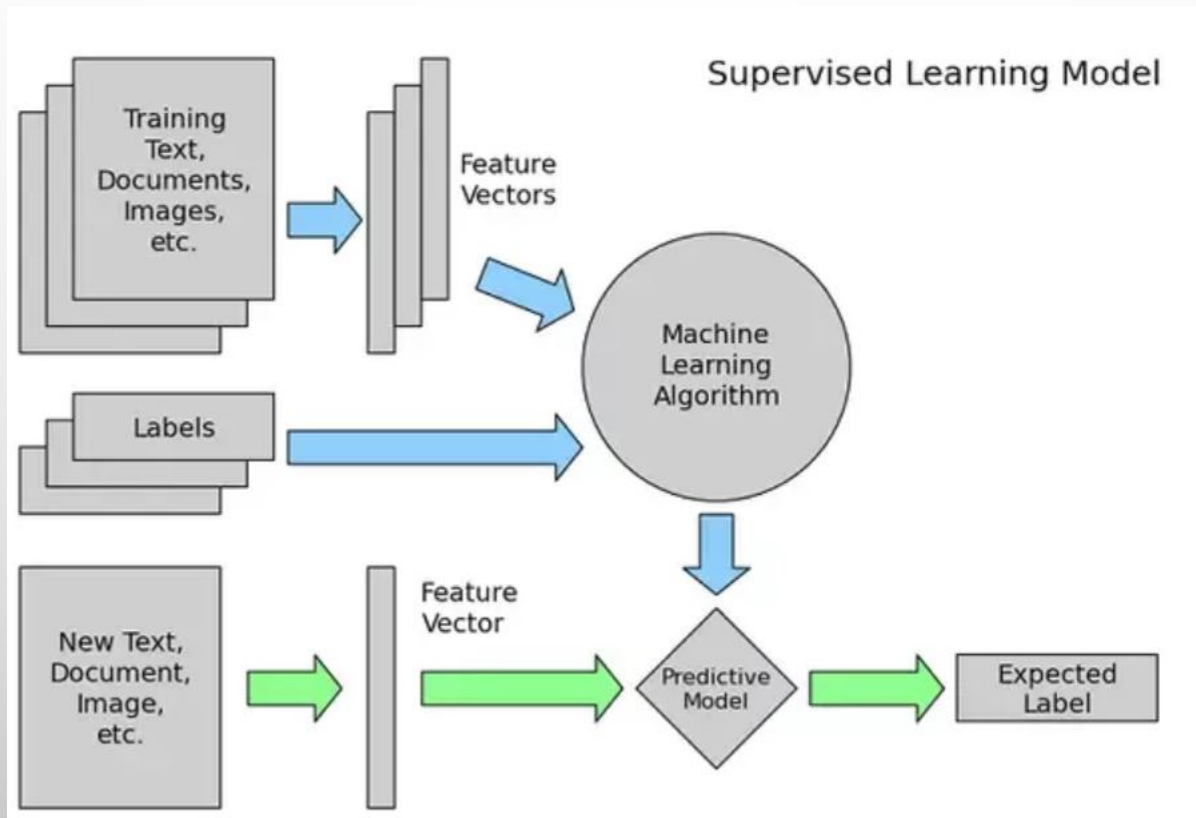
- SPAM WORDS CAN ALSO BE COMMON WORDS WITH HIGH OCCURRENCE FREQUENCY ,  
SUCH AS *INI* (ENG: THIS), *DI* (ENG: AT), ETC
  - STOP WORD ELIMINATION
- SPAM WORDS CAN ALSO OCCUR IN NOT SPAM TEXT, FOR EXAMPLE MAMA
  - WORD WEIGHT
    - $TF \times IDF = \text{TERM FREQUENCY} / \text{DOCUMENT FREQUENCY (CONTAINS TERM)}$ 
      - $IDF = 1/DF$
      - $IDF = \text{LOG} (N/DF)$
    - MUTUAL INFORMATION (MI) =
      - N: NUMBER OF DOCUMENT WITH LABEL (L),  $MI(t, l) = \log_2 \frac{p(t, l)}{p(t) \times p(l)} \approx \log_2 \frac{A \times N}{(A + B) \times (A + C)}$
      - A: FREQUENCY OF WORD (T) WITH LABEL (L);
      - B: FREQUENCY OF WORD (T) WITHOUT LABEL (L);
      - C: NUMBER OF DOCUMENT WITH LABEL (L) WITHOUT WORD (T)

# FEATURE

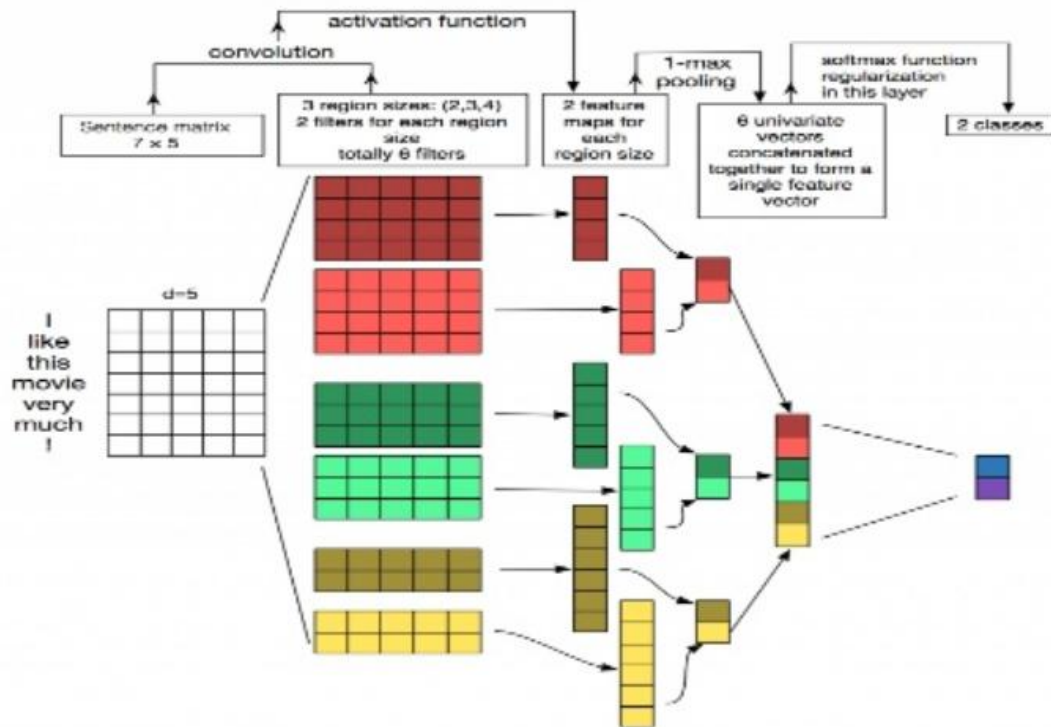
- LEXICAL BASED
  - IDEA: TO CLASSIFY AN INPUT BASED ON ITS WORDS. TO DECREASE THE NUMBER OF WORDS, THE WORDS CAN BE SELECTED FIRST, BY USING:
    - SPECIFIC WORD LIST (MANUALLY SELECTED)
    - NAMED ENTITY
    - TFXIDF
    - MUTUAL INFORMATION
    - POS TAG INFORMATION (FOR EXAMPLE: ONLY NOUN & VERB)
- SYNTACTICAL PARSER
  - IDEA: TO TAKE INTO ACCOUNT, THE WORD ORDER OR GRAMMAR; OR TO DO WORD SELECTION BASED ON ITS SYNTACTICAL INFORMATION
  - SHALLOW PARSER
  - DEEP PARSER
  - N-GRAM
- SEMANTIC

# CLASSIFICATION ALGORITHM FOR SPAM FILTERING

- MANUAL
  - IF “NUMBER OF SPAM WORDS > THRESHOLD” THEN SPAM
- STATISTICAL BASED (MACHINE LEARNING)
  - SPAM RULES ARE DISCOVERED BY MACHINE LEARNING ALGORITHM (EXAMPLE: ID3 OR C4.5)
    - EXAMPLE: IF EMAIL CONTAINS W1 AND W2 THEN THE EMAIL IS SPAM
  - PROBABILITY SCORE OF AN INPUT, DEFINED AS SPAM
    - CONTOH:  $P(W1|SPAM) = 0.3$
  - DEEP LEARNING
    - CNN, RNN (LSTM, GRU)



# CNN for NLP

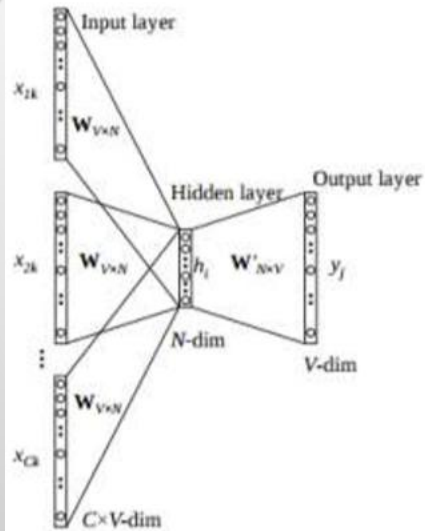


images found in the WildML blog:

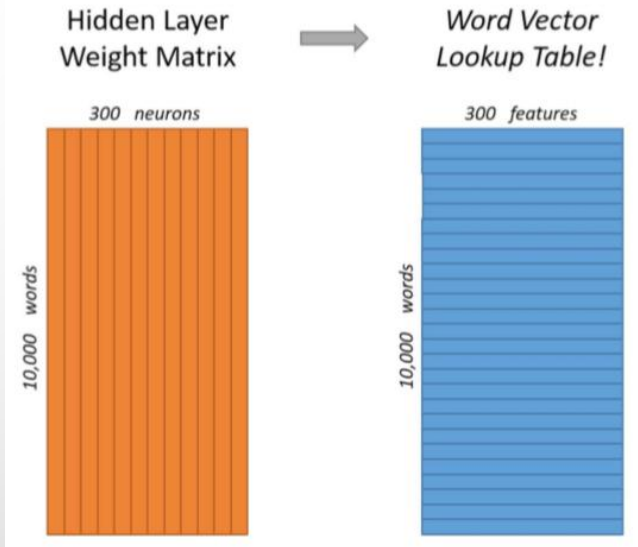
<http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/>

also very good tutorial on CNN for NLP with Tensorflow

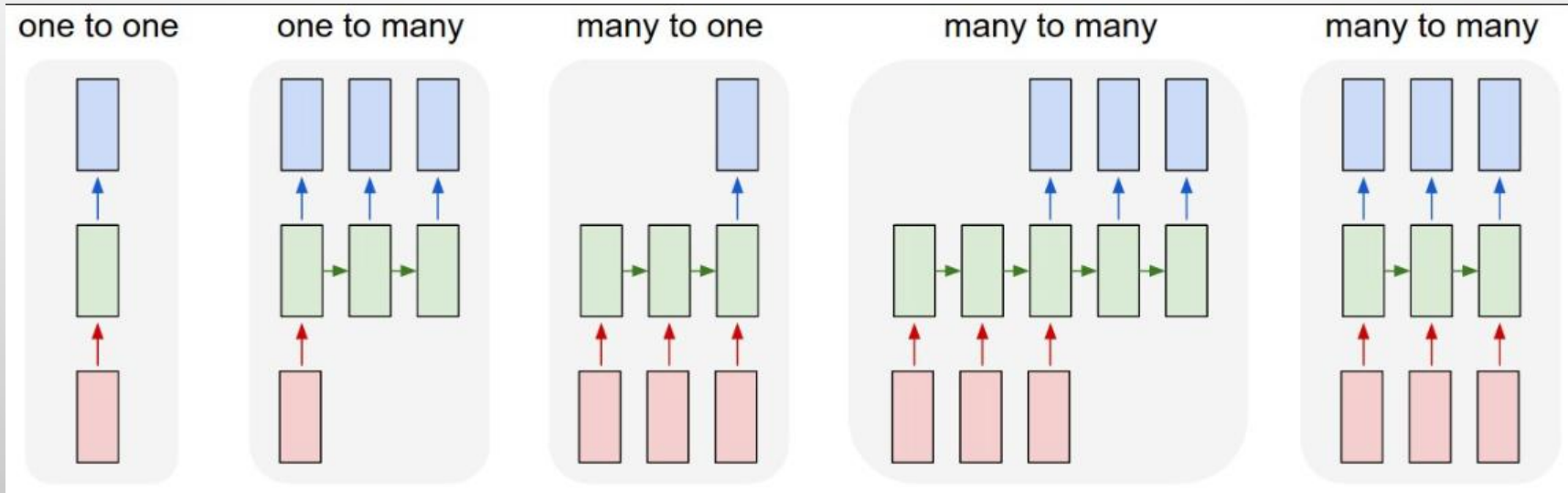
<http://www.wildml.com/2015/12/implementing-a-cnn-for-text-classification-in-tensorflow/>



- The weight between the hidden layer and the output layer is taken as the word vector representation of the word.



# RNN FOR NLP



- SOURCE: [HTTP://KARPATHY.GITHUB.IO/2015/05/21/RNN-EFFECTIVENESS/](http://karpathy.github.io/2015/05/21/rnn-effectiveness/)



# ARTICLES

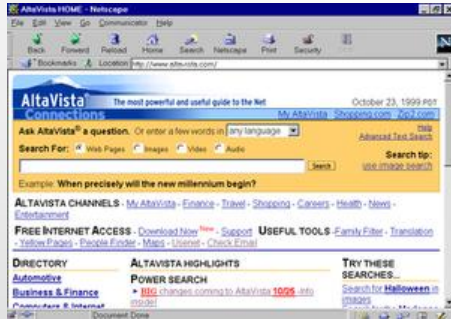
# MILLION OF ARTICLES

- SEARCH ENGINE
  - SEARCHING RELEVANT ARTICLES USING GIVEN QUERY
- RECOMMENDATION
  - RECOMMEND A RELATED ARTICLES BASED ON USER HISTORY
- SUMMARIZATION
  - SUMMARIZE ONE OR SEVERAL ARTICLES
- NEWS AGGREGATOR
  - NEWS SUMMARIZATION AND TOPIC DETECTION
- PLAGIARISM DETECTION
  - CHECK THE SIMILARITY OF NEW ARTICLE TO EXISTING ARTICLES

# SEARCH ENGINE

MAKE SEARCHING EASIER

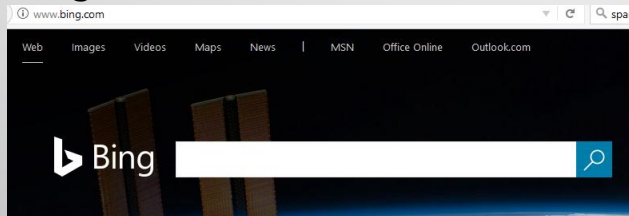
Altavista – 1995, 2003, 2013      Yahoo – 1994



Google - 1998



Bing – 2009



Baidu 百度 2000

# SEARCH ENGINE



# NEWS AGGREGATOR

The screenshot shows a news aggregator website with a sidebar on the left and a main content area. The sidebar lists various sections: Top Stories, World, U.S., Business, Technology, Entertainment, Sports, Science, Health, and a 'Manage sections' option at the bottom. The main content area is divided into three columns: 'Top Stories', 'In the News', and 'Recent'. The 'Top Stories' column features two news items. The first item, 'New GE Chief Shakes Up Management Team', includes a photo of a man, the source 'New York Times', and a timestamp '44m ago'. Below the headline, there is a 'RELATED COVERAGE' section with a link to 'General Electric' and a 'MORE ABOUT' section with a link to 'Jeffrey S. Bornstein'. The second item, 'What storm, Mr. President? Trump puts world on edge with cryptic cliffhanger', includes a photo of Donald Trump and Melania Trump, the source 'Washington Post', and a timestamp '1h ago'. Below this headline, there is a 'RELATED COVERAGE' section with a link to 'Trump plans to declare that Iran nuclear deal is not in the national interest' and a 'Highly Cited' section with a link to 'Washington Post'. The 'In the News' column lists several topics: Hurricane Nate, Harvey Weinstein, Donald Trump, National Hurricane Center, Iran, AOL Instant Messenger, Ralphie May, Elon Musk, Tesla, Inc., and Jefferson Sessions. The 'Recent' column lists a news item: 'Fourth ex-governor from Mexico's PRI arrested on corruption charges' from Reuters, dated '35m ago'. Three blue arrows point from the interface to labels at the bottom: one from the 'Manage sections' option to 'Text Classification', one from the 'RELATED COVERAGE' section of the first news item to 'Summarization', and one from the 'In the News' column to 'Topic Detection'.

Headlines Local For You U.S. ▼

SECTIONS

- Top Stories
- World
- U.S.
- Business
- Technology
- Entertainment
- Sports
- Science
- Health
- Manage sections

### Top Stories

**New GE Chief Shakes Up Management Team**  
New York Times · 44m ago

RELATED COVERAGE  
General Electric announces slew of executive changes, including new CFO  
Highly Cited · CNBC · 3h ago

MORE ABOUT  
General Electric  
Jeffrey S. Bornstein

**New GE CEO Shakes Up Leadership Team**  
Fortune · 1h ago

GE's Shakeup Is Better Sooner Rather Than Later  
Opinion · Bloomberg · 54m ago

View full coverage →

**'What storm, Mr. President?' Trump puts world on edge with cryptic cliffhanger**  
Washington Post · 1h ago

RELATED COVERAGE  
Trump plans to declare that Iran nuclear deal is not in the national interest  
Highly Cited · Washington Post · Oct 6, 2017

### In the News

- Hurricane Nate
- Harvey Weinstein
- Donald Trump
- National Hurricane Center
- Iran
- AOL Instant Messenger
- Ralphie May
- Elon Musk
- Tesla, Inc.
- Jefferson Sessions

### Recent

Fourth ex-governor from Mexico's PRI arrested on corruption charges  
Reuters · 35m ago

Text Classification

Summarization

Topic Detection

# PLAGIARISM DETECTION

105.000.000

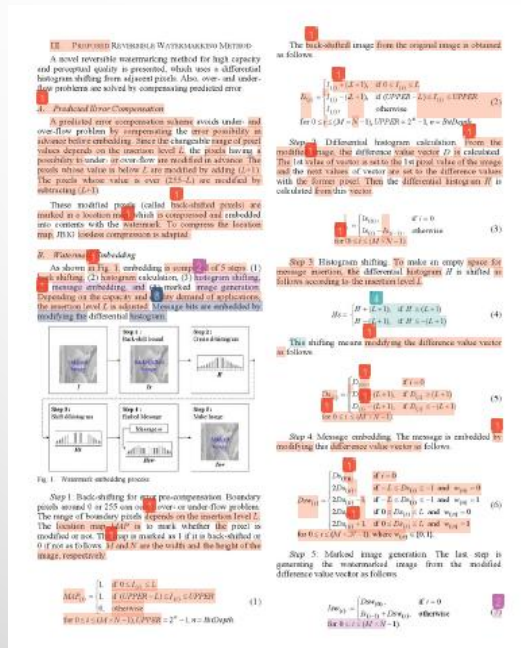
Scientific articles

60.000.000.000

Web articles

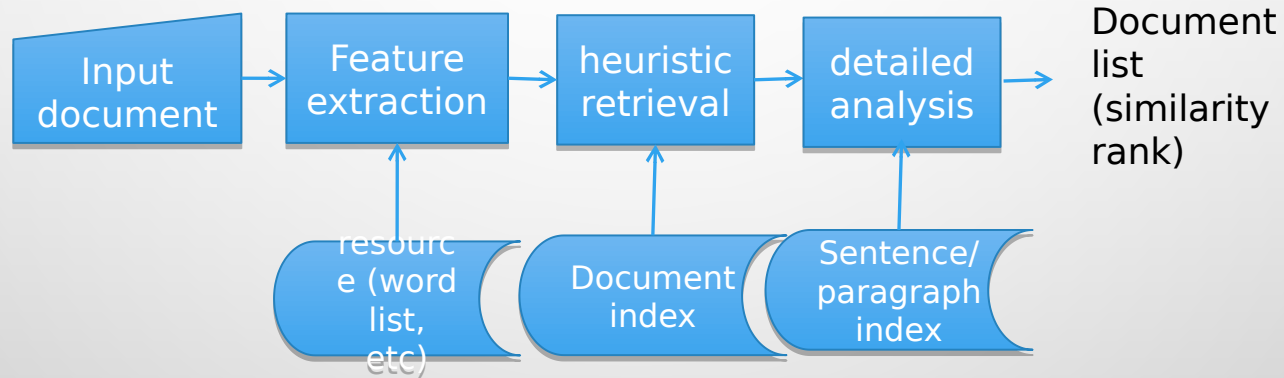
Source: <http://www.ithenticate.com/>

• ITHENTICATE  
• CROSSCHECK  
• K  
• TURNITIN



65%	
EMILY Y. Y. CHOI	
44 VOLUME 20 NUMBER 2	
1	Yoo, Dong-Gyu, Hae-Yoon Lee, and Byoung Man Kim. "Differential Histogram Modification-Based Reversible Watermarking with Predicted Error Compensation", 2010 Sixth International Conference on Intelligent Information Hiding and Multimedia Signal Processing, 2010. <span>1432 words — 57%</span>
2	Yoo, Dong-Gyu. "High capacity reversible watermarking using differential histogram shifting and predicted error compensation", Journal of Electronic Imaging, 2011. <span>58 words — 2%</span>
3	www.koroscience.or.kr <span>22 words — 1%</span>
4	hkko.kaist.ac.kr <span>21 words — 1%</span>
5	Kim, K.S. "Region-based tampering detection and recovery using homogeneity analysis in quality-sensitive imaging", Computer Vision and Image Understanding, 201109 <span>17 words — 1%</span>
6	Dong-Gyu Yoo. "High capacity reversible watermarking using differential histogram shifting and predicted error compensation", Journal of Electronic Imaging, 2011 <span>17 words — 1%</span>
7	scholar.ndsl.kr <span>16 words — 1%</span>
8	Lecture Notes in Computer Science, 2015.

# SYSTEM ARCHITECTURE




- COMPARE INPUT DOCUMENT WITH ALL DOCUMENTS IN THE COLLECTION (CORPUS)



# SOCIAL MEDIA



# MILLION OF USER GENERATED CONTENT

- CUSTOMER FEEDBACK
    - OPINION MINING
    - TRENDING TOPIC
    - USER MONITORING
  - CUSTOMER SERVICE
    - CHATBOT
  - FAKE INFORMATION
    - HOAX CLASSIFIER
  - REAL INFORMATION
    - SUMMARIZE EXTRACTED INFORMATION
- 

# OPINION MINING

- SENTIMENT CLASSIFICATION



## Sentiment analysis

Label: positive, negative, neutral

- ASPECT BASED

**Aspect/Feature Identification**  
Identify salient topics

*battery life,  
sound quality,  
ease of use,...*

**Sentiment Prediction**  
Determine polarity of text containing topics

*battery life is great → +ve  
long battery life → +ve*

**Summary Presentation**  
• Aggregate polarity ratings  
• Present opinion summaries

*Battery Life:* ★★★★★  
*Sound Quality:* ★★★★★

(1) *I bought a Samsung camera and my friends brought a Canon camera yesterday.* (2) *In the past week, we both used the cameras a lot.* (3) *The photos from my Samy are not that great, and the battery life is short too.* (4) *My friend was very happy with his camera and loves its picture quality.* (5) *I want a camera that can take good photos.* (6) *I am going to return it tomorrow.*

- Entity: samsung, samy, canon. Samsung and Samy are grouped together
- Aspect: picture, photo, battery life. Picture & photo are grouped for the camera
- Holder: the blog author (sentence 3), bigJohn's friend (sentence 4)
- Time; Sept-15-2011
- Sentiment: negative on picture quality, negative on battery life, position on camera as a whole

# PROBLEMS IN OPINION MINING

- NE RECOGNITION
  - BARACK OBAMA BERHASIL MENGURANGI ...
- ANAPHORA RESOLUTION
  - KAMI BERDUA MENONTON FILM "...". ITU BENAR-BENAR PENGALAMAN BURUK
- SARCASM (IRONI)
  - PEMERINTAH MAKIN HEBAT AJA, NGAMBIL MENTRI DARI PARTAI YG TERKENAL KORUPSI
- NON-FORMAL LANGUAGE
- FINDING MAIN OBJECT AND SENTIMENT ASPECT
  - MENTARI SINYALNYA MAKIN JELEK! MENDING PAKE XL, LEBIH BAGUS
- RELATIVE SENTIMENT
  - HARGA MAHASISWA
- ONE INPUT WITH BOTH POSITIVE AND NEGATIVE SENTIMENT
- COMPARISON
  - XL LEBIH MURAH DARI TELKOMSEL TAPI LEBIH SERING PUTUS-PUTUS



Germany



Brazil



Indonesia



Nigeria

# Events on Twitter not in the News



**Aditya Danang R.N #8**

@Daraditya\_



Follow

Besok aliansi BEM se-Indonesia demo di Istana Negara.. @BemFPsiUHT min gak ada aksi??

View translation

10:50 AM - 20 May 2015



**Organizing strikes**



**PIDI 林忠龙**

@da\_vidicode

gelodok dan pecinan seluruh indonesia mogok kerja. kalo tuntutan kami g d jg harga kabel dan elektronik kami n gilaan

Translated from Indonesian by bing

it indonesia mogok. Kalo our demands g jg sup ian



Follow



**Benhard Sitorus**

@BenhardSitorus

Ayo Kita dukung aksi sosial #Viking dan #Bomber utk Pengalengan..  
[fb.me/7b3AdGJBq](https://fb.me/7b3AdGJBq)

View translation

4:51 PM



**Asking for solidarity**

Aksi demo mhswa/massa yg besar, masif dan maraton selalu diback up oleh isu, dana, backing yg kuat. Jika tdk, ya pasti tdk efektif

View translation

RETWEETS

10

12:51 AM - 20 May 2015



**General Strike discussion**

**Threatening with Strike**



**Info Lalulintas**

@masialin



Follow

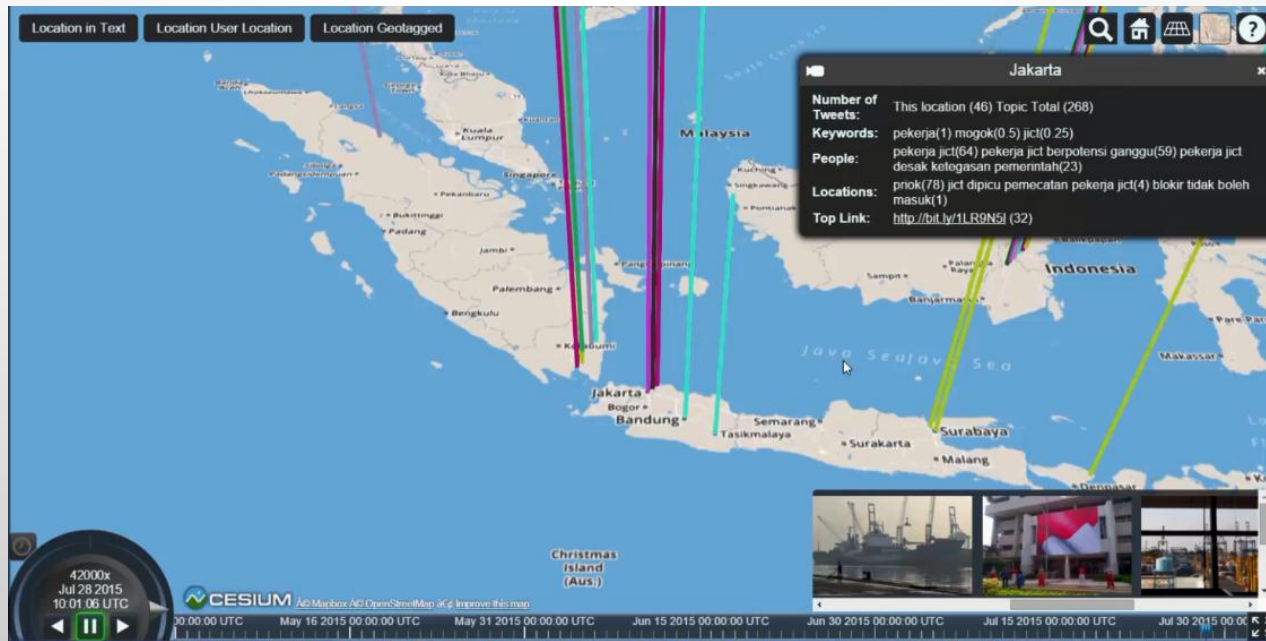
@PuspitaFM demo aksi donasi untuk muslim rohingya di veteran

View translation

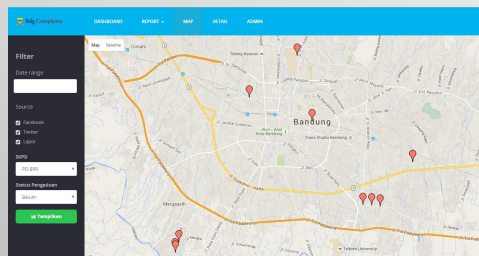
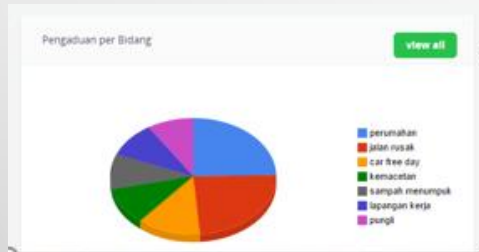
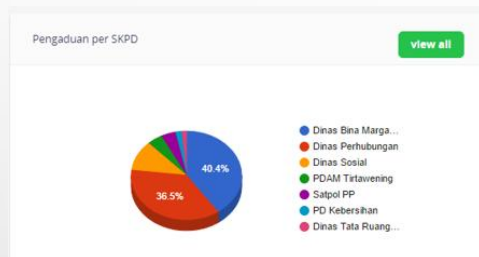


**Asking for Donation**









# COMPLAINT MANAGEMENT SYSTEM



## Complaint Classification

Doc classification:  
complaint vs  
response; nama  
dinas



## Topic Classification

Doc classification



## Information Extraction

Lokasi, waktu,  
kondisi, penyebab

# CLASSIFICATION: TEXT/DOCUMENT VS SENTENCE VS TOKEN

- TEXT/DOCUMENT CLASSIFICATION: ENTIRE DOCUMENT AS INPUT FEATURE
- SENTENCE CLASSIFICATION:
  - SENTENCE SPLITTER
  - RELATION AMONG SENTENCE ➡ REFERENCE RESOLUTION, SENTENCE POSITION
  - SENTENCE AS INPUT FEATURE
- TOKEN CLASSIFICATION:
  - SENTENCE SPLITTER
  - TOKENIZATION
  - TOKEN/KATA/TERM AS INPUT FEATURE

# EXAMPLE ON TOKEN CLASSIFICATION

- POS (PART OF SPEECH) TAGGER
  - CLASS: POS TAG, EX: NN (NOUN), VB (VERB), ADJ (ADJECTIVE), ADV (ADVERB)
- NAMED ENTITY TAGGER
  - CLASS: NE TAG, EX: PERSON, ORGANIZATION, LOCATION
- KEYWORD EXTRACTION
  - CLASS: YES/NO

# INFORMATION EXTRACTION (NER + RELATION EXT)



Subject: **US-TN**-SOFTWARE PROGRAMMER  
Date: **17 Nov 1996** 17:37:29 GMT  
Organization: Reference.Com Posting Service  
Message-ID:  
<**56nigp\$mrs@bilbo.reference.com**>

## **SOFTWARE PROGRAMMER**

Position available for Software Programmer experienced in generating software for PC-Based **Voice Mail** systems. Experienced in **C** Programming. Must be familiar with communicating with and controlling voice cards; preferable Dialogic, however, experience with others such as Rhetorix and

computer\_science\_job  
id: **56nigp\$mrs@bilbo.reference.com**  
title: **SOFTWARE PROGRAMMER**  
company:  
recruiter:  
state: **TN**  
city:  
country: **US**  
language: **C**  
platform: **PC \ DOS \ OS-2 \ UNIX**  
rea: **Voice Mail**  
req\_years\_experience: **2**  
desired\_years\_experience: **5**

For years, [Microsoft Corporation CEO Bill Gates](#) was against open source. But today he appears to have changed his mind. "We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#)

Select Name  
From PEOPLE  
Where Organization =  
'Microsoft'

## PEOPLE

<u>Name</u>	<u>Title</u>	<u>Organization</u>
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft..

Bill Gates  
Bill Veghte

# IE=SEGMENTATION+CLASSIFICATION+ASSOCIATION+CLUSTERING

For years, [Microsoft Corporation](#) [CEO Bill Gates](#) was against open source.

But today [Microsoft](#) appears to have changed his mind. [Gates](#) himself says [Microsoft](#) will gladly disclose its crown jewels - the coveted code behind the Windows operating system - to select customers.

"We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#), countered saying...

[MICROSOFT CORPORATION](#)

[CEO](#)

[BILL GATES](#)

[MICROSOFT](#)

[GATES](#)

[MICROSOFT](#)

[BILL VEGHTE](#)

[MICROSOFT](#)

[VP](#)

[RICHARD STALLMAN](#)

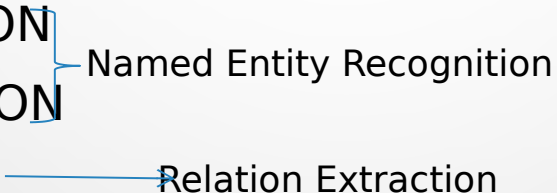
[FOUNDER](#)

[FREE SOFTWARE](#)

[FOUNDATION](#)

# INFORMATION EXTRACTION

- IE

- SEGMENTATION
  - CLASSIFICATION
  - ASSOCIATION
  - CLUSTERING
- Named Entity Recognition
- Relation Extraction
- 

# NAMED ENTITY RECOGNITION


- TYPES:
  - SUPERVISED
  - UNSUPERVISED
- FEATURES:
  - WORD LEVEL FEATURES
  - LIST LOOKUP FEATURES
  - DOCUMENT & CORPUS FEATURES
- SENTENCE:  $S = W_1 W_2 \dots W_{N-1} W_N$
- NAMED ENTITY CLASSIFICATION:
  - SEGMENTATION + CLASSIFICATION
  - RAW SEQUENCE LABELING
    - CLASS: NAMED ENTITY TYPE + “BEGIN/IN/END” INFORMATION
    - $W_1$  , WORD BEFORE, WORD AFTER, CLASS



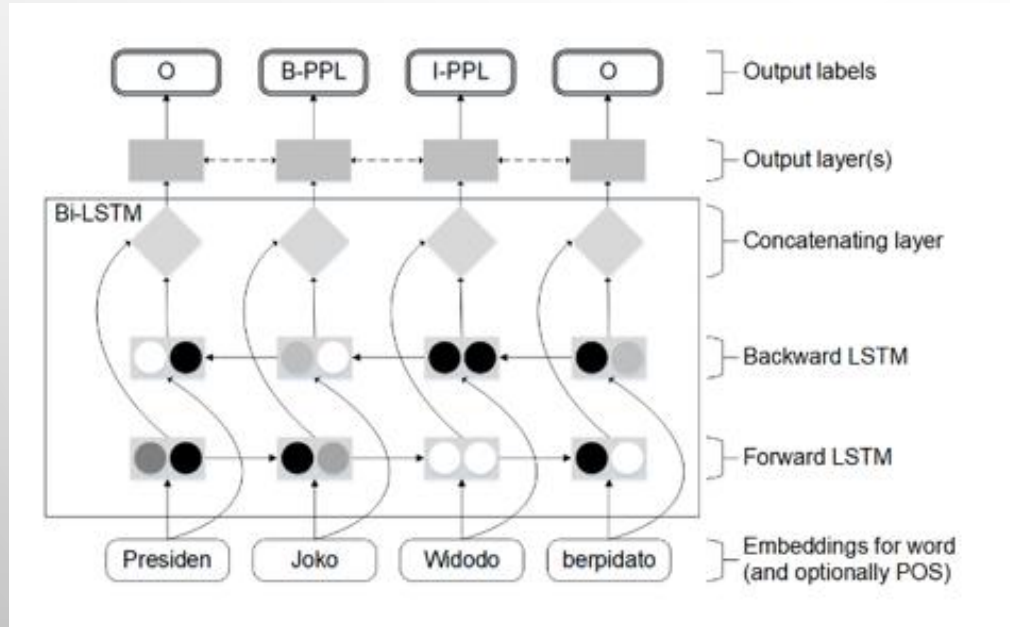
# RELATION EXTRACTION

- SENTENCE:  $S = W_1 W_2 \dots E_1 \dots W_1 \dots E_2 \dots W_{N-1} W_N$
- RELATION CLASSIFICATION:
  - +1 IF  $E_1$  AND  $E_2$  ARE RELATED BY A RELATION  $R$ 
    - $R$  BISA BERUPA NAMA RELASI ATAU HANYA MENYATAKAN ADA RELASI ATAU TIDAK
  - -1 OTHERWISE
- FEATURE EXAMPLE:
  - $E_1, E_2$ , WORD BETWEEN, WORD BEFORE, WORD AFTER, RELATION CLASS

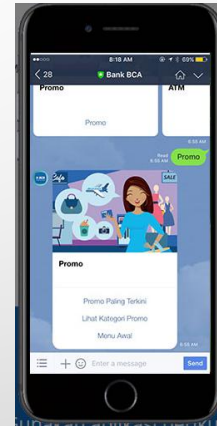
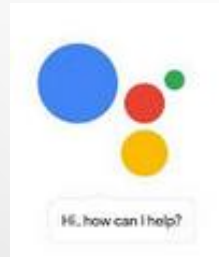
# NER WITH MACHINE LEARNING

- FEATURE:
  - IMPORTANT:
    - CURRENT WORD
    - PRECEDING NE TAG
    - POS TAG
  - OPTIONAL:
    - WORD WINDOW: PRECEDING WORDS, SUCCEEDING WORDS
    - WORD LIST  CAN BE USED AS SINGLE FEATURE OR TO LIMIT FEATURE (WORD VOCABULARY)
- CLASS:
  - NE TAGS: PERSON-B, PERSON-I, ORG-B, ORG-I, OTHER, ETC

# NER WITH DEEP LEARNING

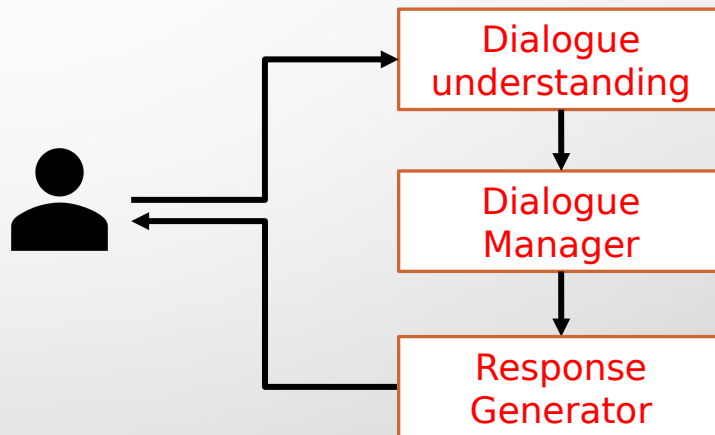


# CHATBOT & VIRTUAL ASSISTANT



# CHATBOT

- DIALOGUE UNDERSTANDING
  - INTENTION CLASSIFICATION
  - ENTITY EXTRACTION
- DIALOGUE MANAGER
  - EXECUTE DIALOGUE SCENARIO
- RESPONSE GENERATOR
  - TEMPLATE / GENERATOR



# CHALLENGES



COMPLETE  
LANGUAGE  
PROBLEMS



LOW  
RESOURCE



DEEP  
LEARNING

**TERIMA KASIH**

