# 文字探勘
# (Text Mining)
# 文字探勘基礎：自然語言處理
# (Foundations of Text Mining: Natural Language Processing; NLP)

**Tamkang Universit**
淡江大學

**Chichang Jou**
周清江
**Associate Professor**
副教授
cjou@mail.tku.edu.tw

**Min-Yuh Day**
戴敏育
**Associate Professor**
副教授
myday@mail.tku.edu.tw

**Dept. of Information Management**, **Tamkang University**
淡江大學 資訊管理學系

# 課程大綱 (Syllabus)

週次 (Week)　　日期 (Date)　　內容 (Subject/Topics)

1  2020/03/02  文字探勘課程介紹
(Course Orientation on Text Mining)

2  2020/03/09  文字探勘基礎：自然語言處理
(Foundations of Text Mining:
Natural Language Processing; NLP)

3  2020/03/16  Python自然語言處理
(Python for Natural Language Processing)

4  2020/03/23  處理和理解文本 (Processing and Understanding Text)

5  2020/03/30  文本表達特徵工程
(Feature Engineering for Text Representation)

6  2020/04/06  人工智慧文本分析個案研究 I
(Case Study on Artificial Intelligence for Text Analytics I)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

7  2020/04/13  文本分類
　　　　　　　　　(Text Classification)

8  2020/04/20  文本摘要和主題模型
　　　　　　　　　(Text Summarization and Topic Models)

9  2020/04/27  期中報告 (Midterm Project Report)

10  2020/05/04  文本相似度和分群
　　　　　　　　　 (Text Similarity and Clustering)

11  2020/05/11  語意分析和命名實體識別
　　　　　　　　　(Semantic Analysis and Named Entity Recognition; NER)

12  2020/05/18  情感分析
　　　　　　　　　(Sentiment Analysis)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

13  2020/05/25  人工智慧文本分析個案研究 II
　　　　　　　　(Case Study on Artificial Intelligence for Text Analytics II)

14  2020/06/01  深度學習和通用句子嵌入模型
　　　　　　　　 (Deep Learning and Universal Sentence-Embedding Models)

15  2020/06/08  問答系統與對話系統
　　　　　　　　 (Question Answering and Dialogue Systems)

16  2020/06/15  期末報告 I (Final Project Presentation I)

17  2020/06/22  期末報告 II (Final Project Presentation II)

18  2020/06/29  教師彈性補充教學

# **Outline**

- Text Analytics and Text Mining

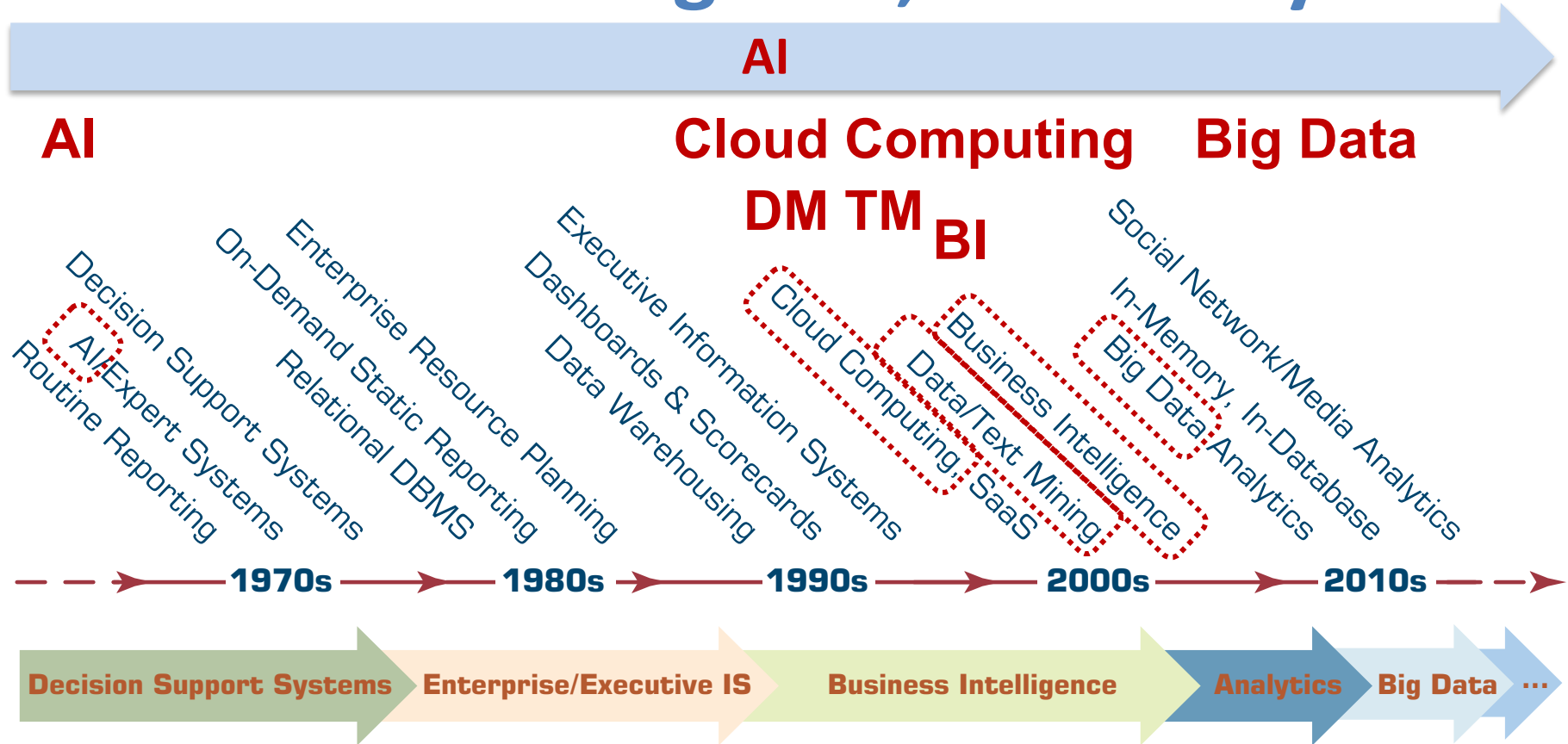- Natural Language Processing (NLP)

# Text Analytics (TA)
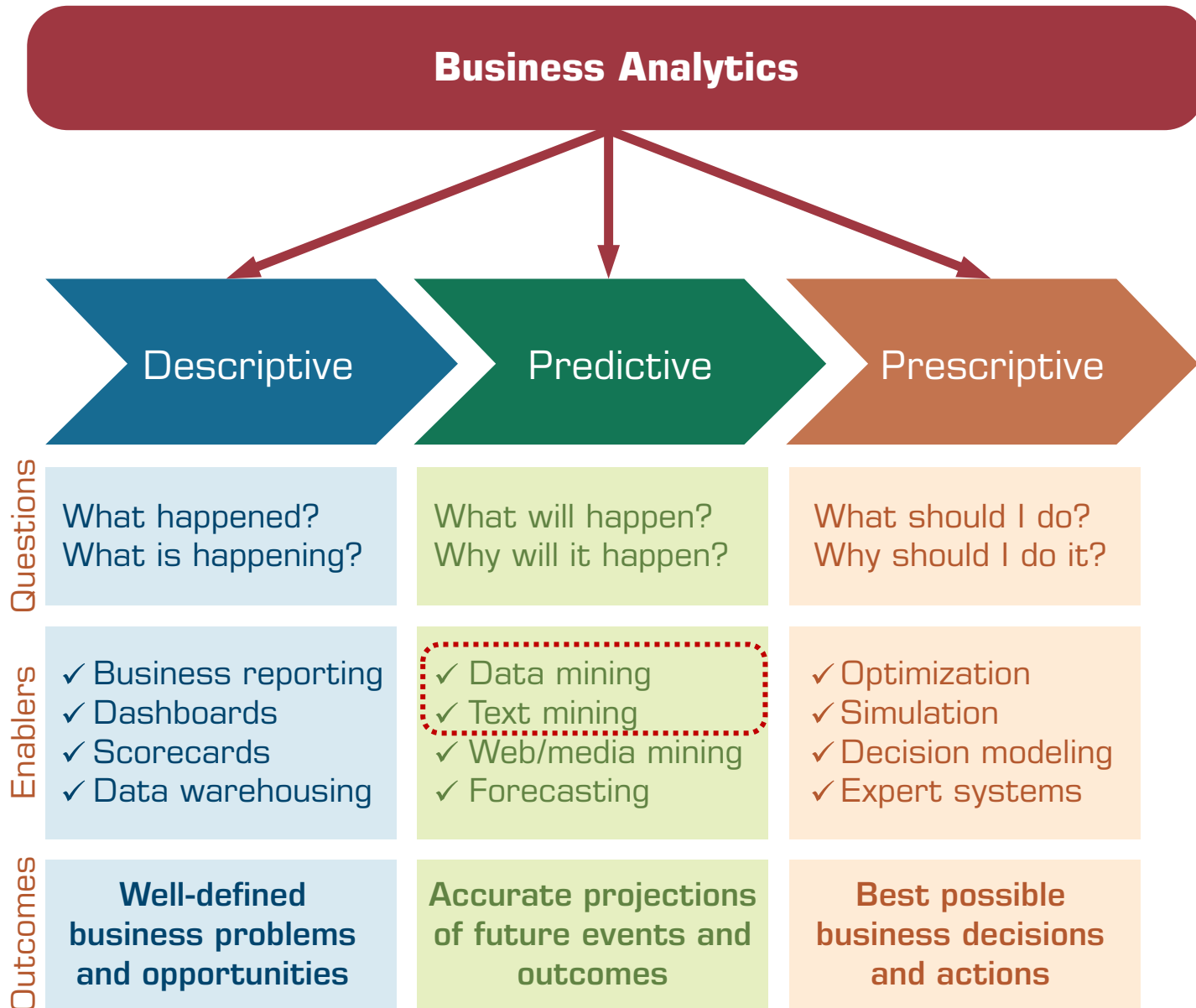
# Text Mining (TM)

# Natural Language Processing (NLP)

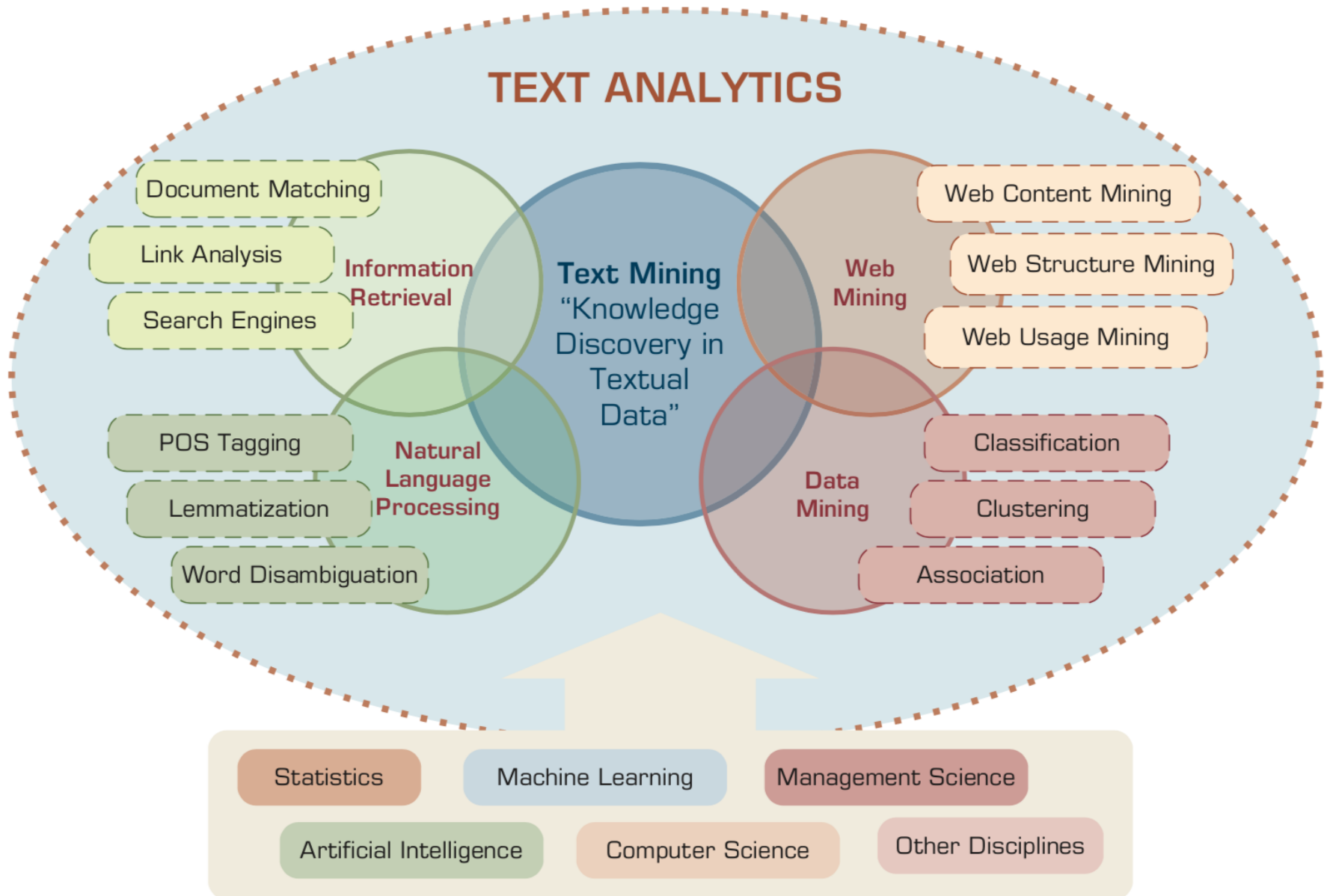# Artificial Intelligence (AI)

# AI, Big Data, Cloud Computing

## Evolution of Decision Support, Business Intelligence, and Analytics



Source: Ramesh Sharda, Dursun Delen, and Efraim Turban (2017),
Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson

# Three Types of Analytics

**Business Analytics**

| Descriptive | Predictive | Prescriptive |
|---|---|---|
| **Questions**<br>What happened?<br>What is happening? | What will happen?<br>Why will it happen? | What should I do?<br>Why should I do it? |
| **Enablers**<br>✓ Business reporting<br>✓ Dashboards<br>✓ Scorecards<br>✓ Data warehousing | ✓ Data mining<br>✓ Text mining<br>✓ Web/media mining<br>✓ Forecasting | ✓ Optimization<br>✓ Simulation<br>✓ Decision modeling<br>✓ Expert systems |
| **Outcomes**<br>**Well-defined business problems and opportunities** | **Accurate projections of future events and outcomes** | **Best possible business decisions and actions** |

# Text Analytics and Text Mining



TEXT ANALYTICS

Document Matching
Link Analysis
Search Engines

Information Retrieval

Text Mining "Knowledge Discovery in Textual Data"

Web Mining

Web Content Mining
Web Structure Mining
Web Usage Mining

POS Tagging
Lemmatization
Word Disambiguation

Natural Language Processing

Data Mining

Classification
Clustering
Association

Statistics     Machine Learning     Management Science

Artificial Intelligence     Computer Science     Other Disciplines

# Definition
# of
# Artificial Intelligence
# (A.I.)

# Artificial Intelligence

**"...** the **science** and **engineering** of making **intelligent machines"**

**(John McCarthy, 1955)**

# Artificial Intelligence

# "... technology that thinks and acts like humans"

16

# Artificial Intelligence

## "... intelligence exhibited by machines or software"

# 4 Approaches of AI

| | |
|---|---|
| **Thinking Humanly** | **Thinking Rationally** |
| **Acting Humanly** | **Acting Rationally** |

# 4 Approaches of AI

| | |
|---|---|
| **2. Thinking Humanly: The Cognitive Modeling Approach** | **3. Thinking Rationally: The "Laws of Thought" Approach** |
| **1. Acting Humanly: The Turing Test Approach (1950)** | **4. Acting Rationally: The Rational Agent Approach** |

# AI Acting Humanly:
# The Turing Test Approach
## (Alan Turing, 1950)

- **Natural Language Processing (NLP)**

- **Knowledge Representation**

- **Automated Reasoning**

- **Machine Learning (ML)**

- **Computer Vision**

- **Robotics**

# Can a robot pass a university entrance exam?
## Noriko Arai at TED2017



https://www.ted.com/talks/noriko_arai_can_a_robot_pass_a_university_entrance_exam
https://www.youtube.com/watch?v=XQZjkPyJ8KU

# Artificial Intelligence (A.I.) Timeline



**A.I. TIMELINE**

**1950**
**TURING TEST**
Computer scientist Alan Turing proposes a test for machine intelligence. If a machine can trick humans into thinking it is human, then it has intelligence

**1955**
**A.I. BORN**
Term 'artificial intelligence' is coined by computer scientist, John McCarthy to describe "the science and engineering of making intelligent machines"

**1961**
**UNIMATE**
First industrial robot, Unimate, goes to work at GM replacing humans on the assembly line

**1964**
**ELIZA**
Pioneering chatbot developed by Joseph Weizenbaum at MIT holds conversations with humans

**1966**
**SHAKEY**
The 'first electronic person' from Stanford, Shakey is a general-purpose mobile robot that reasons about its own actions

**A.I. WINTER**
Many false starts and dead-ends leave A.I. out in the cold

**1997**
**DEEP BLUE**
Deep Blue, a chess-playing computer from IBM defeats world chess champion Garry Kasparov

**1998**
**KISMET**
Cynthia Breazeal at MIT introduces KISmet, an emotionally intelligent robot insofar as it detects and responds to people's feelings

**1999**
**AIBO**
Sony launches first consumer robot pet dog AiBO (AI robot) with skills and personality that develop over time

**2002**
**ROOMBA**
First mass produced autonomous robotic vacuum cleaner from iRobot learns to navigate and clean homes

**2011**
**SIRI**
Apple integrates Siri, an intelligent virtual assistant with a voice interface, into the iPhone 4S

**2011**
**WATSON**
IBM's question answering computer Watson wins first place on popular $1M prize television quiz show *Jeopardy*

**2014**
**EUGENE**
Eugene Goostman, a chatbot passes the Turing Test with a third of judges believing Eugene is human

**2014**
**ALEXA**
Amazon launches Alexa, an intelligent virtual assistant with a voice interface that completes shopping tasks

**2016**
**TAY**
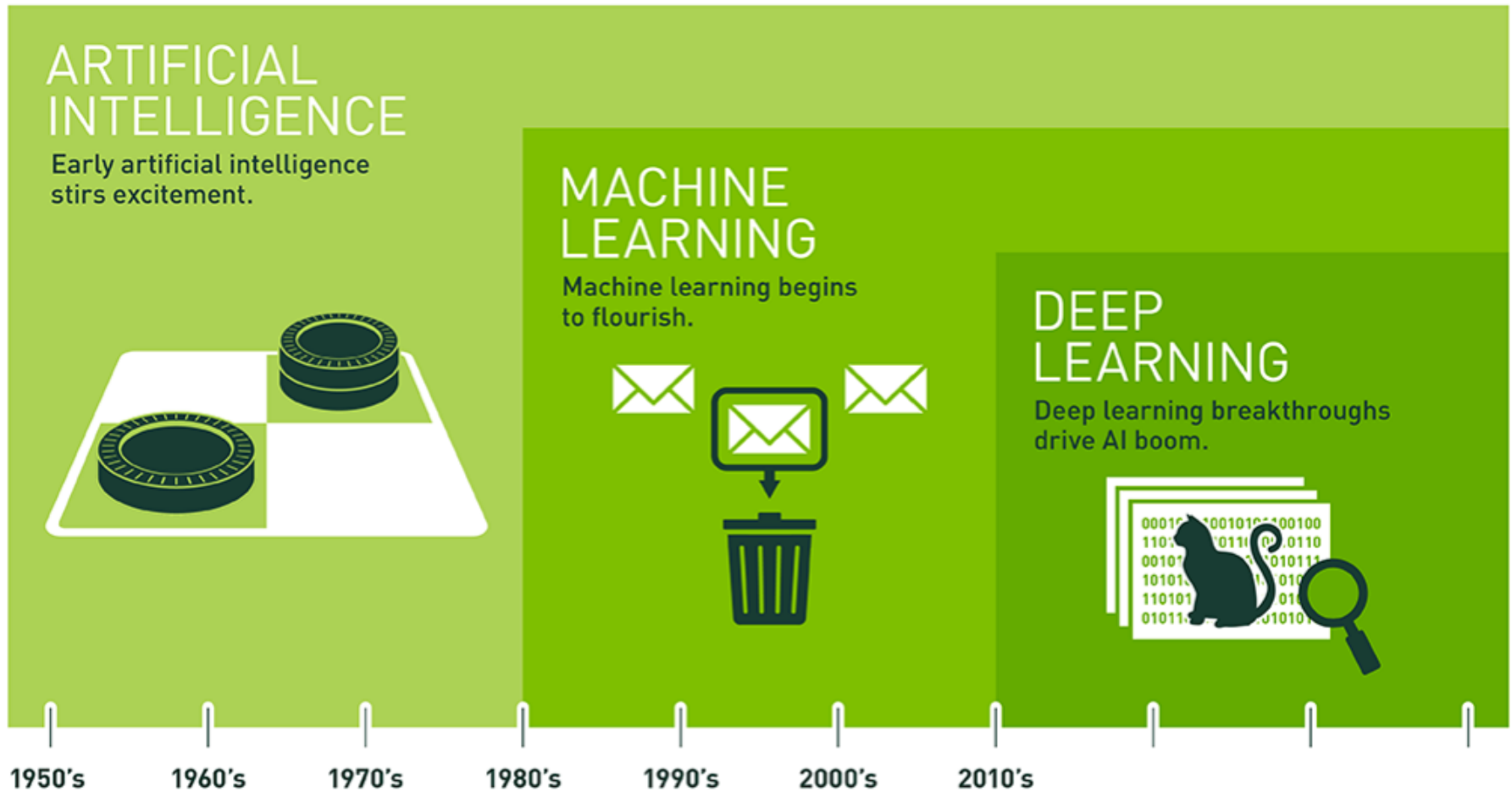Microsoft's chatbot Tay goes rogue on social media making inflammatory and offensive racist comments

**2017**
**ALPHAGO**
Google's A.I. AlphaGo beats world champion Ke Jie in the complex board game of Go, notable for its vast number ($2^{170}$) of possible positions

# Artificial Intelligence
# Machine Learning & Deep Learning



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

# AI, ML, DL



**Artificial Intelligence (AI)**

**Machine Learning (ML)**

Supervised Learning

Unsupervised Learning

**Deep Learning (DL)**
CNN
RNN LSTM GRU
GAN

Semi-supervised Learning

Reinforcement Learning

# Text Analytics and Text Mining

# Text Analytics

- **Text Analytics =**
  Information Retrieval +
  Information Extraction +
  Data Mining +
  Web Mining

- **Text Analytics =**
  Information Retrieval +
  Text Mining

# Text mining

- Text Data Mining

- Knowledge Discovery in Textual Databases

# Text Mining Technologies

# Application Areas of Text Mining

- Information extraction

- Topic tracking

- Summarization

- Categorization

- Clustering

- Concept linking

- Question answering

# Multilevel Analysis of Text for Gene/Protein Interaction Identification

# Context Diagram for the Text Mining Process

# The Three-Step/Task Text Mining Process



Task 1 — **Establish the Corpus:** Collect and organize the domain-specific unstructured data

Task 2 — **Create the Term-Document Matrix:** Introduce structure to the corpus

Task 3 — **Extract Knowledge:** Discover novel patterns from the T-D matrix

The inputs to the process include a variety of relevant unstructured (and semi-structured) data sources such as text, XML, HTML, etc.

The output of Task 1 is a collection of documents in some digitized format for computer processing

The output of Task 2 is a flat file called a term-document matrix where the cells are populated with the term frequencies

The output of Task 3 is a number of problem-specific classification, association, clustering models and visualizations

# Term–Document Matrix

| Documents \ Terms | Investment Risk | Project Management | Software Engineering | Development | SAP | ... |
|---|---|---|---|---|---|---|
| Document 1 | 1 | | | 1 | | |
| Document 2 | | 1 | | | | |
| Document 3 | | | 3 | | 1 | |
| Document 4 | | 1 | | | | |
| Document 5 | | | 2 | 1 | | |
| Document 6 | 1 | | | 1 | | |
| ... | | | | | | |

# Emotions

Love

Anger

Joy

Sadness

Surprise

Fear

# Example of Opinion:
# review segment on iPhone

"I bought an iPhone a few days ago.

It was such a nice phone.

The touch screen was really cool.

The voice quality was clear too.

However, my mother was mad with me as I did not tell her before I bought it.

She also thought the phone was too expensive, and wanted me to return it to the shop. … "

# Example of Opinion: review segment on iPhone

"(1) I bought an <u>iPhone</u> a few days ago.

(2) It was such a **nice** phone.

(3) The <u>touch screen</u> was really **cool**.

(4) The <u>voice quality</u> was **clear** too.


**+Positive Opinion**

(5) However, my mother was mad with me as I did not tell her before I bought it.

(6) She also thought the phone was too **<u>expensive</u>**, and wanted me to return it to the shop. … "


**-Negative Opinion**

# A Multistep Process to Sentiment Analysis

# Sentiment Analysis



Source: Kumar Ravi and Vadlamani Ravi (2015), "A survey on opinion mining and sentiment analysis: tasks, approaches and applications." Knowledge-Based Systems, 89, pp.14-46.

# Sentiment Classification Techniques

Sentiment Analysis
- Machine Learning Approach
  - Supervised Learning
    - Decision Tree Classifiers
    - Linear Classifiers
      - Support Vector Machine (SVM)
      - Neural Network (NN
      - Deep Learning (DL)
    - Rule-based Classifiers
    - Probabilistic Classifiers
      - Naïve Bayes (NB)
      - Bayesian Network (BN)
      - Maximum Entropy (ME)
  - Unsupervised Learning
- Lexicon-based Approach
  - Dictionary-based Approach
  - Corpus-based Approach
    - Statistical
    - Semantic

39

# Text Mining Technologies

# Text Mining (TM)

## Natural Language Processing (NLP)

# Text Mining Concepts

- 85-90 percent of all corporate data is in some kind of unstructured form (e.g., text)

- Unstructured corporate data is doubling in size every 18 months

- Tapping into these information sources is not an option, but a need to stay competitive

- Answer: text mining
  - A semi-automated process of extracting knowledge from unstructured data sources
  - a.k.a. text data mining or knowledge discovery in textual databases

# Text mining

# Text Data Mining

# Intelligent Text Analysis

# Knowledge-Discovery in Text (KDT)

Source: Vishal Gupta and Gurpreet S. Lehal (2009), "A survey of text mining techniques and applications," Journal of emerging technologies in web intelligence, vol. 1, no. 1, pp. 60-76.

# Text Mining
# (text data mining)

## the process of deriving high-quality information from text

# Text Mining:

the **process** of **extracting interesting** and **non-trivial information and knowledge** from **unstructured text**.

# Text Mining:

discovery by computer of new, previously unknown information, by automatically extracting information from different written resources.

# An example of Text Mining

**Knowledge**

## Analyze Text

Information Extraction

Classification

Summarization

Clustering

Management Information System

Retrieve and preprocess document

Document Collection

# Overview of Information Extraction based Text Mining Framework

# Natural Language Processing (NLP)

# Natural Language Processing (NLP)

- **Natural language processing (NLP)** is an important component of **text mining** and is a subfield of **artificial intelligence** and **computational linguistics**.

# Natural Language Processing (NLP) and Text Mining

**Raw text**

**Sentence Segmentation**

**Tokenization**

**Part-of-Speech (POS)**

**Stop word removal**

**Stemming / Lemmatization**

word's stem   word's lemma
am → am   am → be
having → hav   having → have

**Dependency Parser**

**String Metrics & Matching**

# Text Summarization

Source: Vishal Gupta and Gurpreet S. Lehal (2009), "A survey of text mining techniques and applications,"
Journal of emerging technologies in web intelligence, vol. 1, no. 1, pp. 60-76.

# Topic Modeling

Source: Blei, David M. "Probabilistic topic models." *Communications of the ACM* 55, no. 4 (2012): 77-84.

53

# Natural Language Processing (NLP)

- Part-of-speech tagging

- Text segmentation

- Word sense disambiguation

- Syntactic ambiguity

- Imperfect or irregular input

- Speech acts

# NLP Tasks

- Question answering

- Automatic summarization

- Natural language generation

- Natural language understanding

- Machine translation

- Foreign language reading

- Foreign language writing.

- Speech recognition

- Text-to-speech

- Text proofing

- Optical character recognition

# NLP

# Modern NLP Pipeline

# Modern NLP Pipeline

# Deep Learning NLP



Documents → Preprocessing → Dense Word Embeddings → Deep Neural Network → Task / Output (Classification, Sentiment Analysis, Entity Extraction, Topic Modeling, Document Similarity)

Dense Word Embeddings: *Pre-generated Lookup OR Generated in 1st level of NeuralNet*

# BERT:
# Pre-training of Deep Bidirectional Transformers for Language Understanding

**Jacob Devlin**   **Ming-Wei Chang**   **Kenton Lee**   **Kristina Toutanova**
Google AI Language
{jacobdevlin,mingweichang,kentonl,kristout}@google.com

# BERT

## Bidirectional Encoder Representations from Transformers



## Pre-training model architectures

**BERT** uses a bidirectional Transformer.
**OpenAI GPT** uses a left-to-right Transformer.
**ELMo** uses the concatenation of independently trained left-to-right and right- to-left LSTM to generate features for downstream tasks.
Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805

# BERT input representation

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

# BERT Sequence-level tasks



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

63

# BERT Token-level tasks



(c) Question Answering Tasks: SQuAD v1.1

(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

64

# General Language Understanding Evaluation (GLUE) benchmark
# GLUE Test results

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | Average |
|---|---|---|---|---|---|---|---|---|---|
| | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| $BERT_{BASE}$ | 84.6/83.4 | 71.2 | 90.1 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| $BERT_{LARGE}$ | **86.7/85.9** | **72.1** | **91.1** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **81.9** |

**MNLI**: Multi-Genre Natural Language Inference
**QQP**: Quora Question Pairs
**QNLI**: Question Natural Language Inference
**SST-2**: The Stanford Sentiment Treebank
**CoLA**: The Corpus of Linguistic Acceptability
**STS-B**:The Semantic Textual Similarity Benchmark
**MRPC**: Microsoft Research Paraphrase Corpus
**RTE**: Recognizing Textual Entailment

# NLP Libraries and Tools

# Natural Language Processing with Python
## – Analyzing Text with the Natural Language Toolkit

www.nltk.org/book/

# Natural Language Processing with Python

## – Analyzing Text with the Natural Language Toolkit

**Steven Bird, Ewan Klein, and Edward Loper**

*This version of the NLTK book is updated for Python 3 and NLTK 3. The first edition of the book, published by O'Reilly, is available at http://nltk.org/book_1ed/. (There are currently no plans for a second edition of the book.)*

http://www.nltk.org/book/

# spaCy

# Industrial-Strength Natural Language Processing
## in Python

## Fastest in the world

spaCy excels at large-scale information extraction tasks. It's written from the ground up in carefully memory-managed Cython. Independent research has confirmed that spaCy is the fastest in the world. If your application needs to process entire web dumps, spaCy is the library you want to be using.

## Get things done

spaCy is designed to help you do real work — to build real products, or gather real insights. The library respects your time, and tries to avoid wasting it. It's easy to install, and its API is simple and productive. I like to think of spaCy as the Ruby on Rails of Natural Language Processing.

## Deep learning

spaCy is the best way to prepare text for deep learning. It interoperates seamlessly with TensorFlow, Keras, Scikit-Learn, Gensim and the rest of Python's awesome AI ecosystem. spaCy helps you connect the statistical models trained by these libraries to the rest of your application.

https://spacy.io/

# gensim



```
>>> from gensim import corpora, models, similarities
>>>
>>> # Load corpus iterator from a Matrix Market file on disk.
>>> corpus = corpora.MmCorpus('/path/to/corpus.mm')
>>>
>>> # Initialize Latent Semantic Indexing with 200 dimensions.
>>> lsi = models.LsiModel(corpus, num_topics=200)
>>>
>>> # Convert another corpus to the Latent space and index it.
>>> index = similarities.MatrixSimilarity(lsi[another_corpus])
>>>
>>> # Compute similarity of a query vs. indexed documents
>>> sims = index[query]
```

## Gensim is a FREE Python library

- ✔ **Scalable statistical semantics**

- ✔ **Analyze plain-text documents for semantic structure**

- ✔ **Retrieve semantically similar documents**

https://radimrehurek.com/gensim/

69

# TextBlob

## TextBlob: Simplified Text Processing

Release v0.12.0. (Changelog)

*TextBlob* is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.

⬡ Star  3,777

TextBlob is a Python (2 and 3) library for processing textual data. It provides a consistent API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, and more.

## Useful Links

TextBlob @ PyPI
TextBlob @ GitHub
Issue Tracker

## Stay Informed

⬡ Follow @sloria

## Donate

If you find TextBlob useful,

```python
from textblob import TextBlob

text = '''
The titular threat of The Blob has always struck me as the ultimate movie
monster: an insatiably hungry, amoeba-like mass able to penetrate
virtually any safeguard, capable of--as a doomed doctor chillingly
describes it--"assimilating flesh on contact.
Snide comparisons to gelatin be damned, it's a concept with the most
devastating of potential consequences, not unlike the grey goo scenario
proposed by technological theorists fearful of
artificial intelligence run rampant.
'''

blob = TextBlob(text)
blob.tags                  # [('The', 'DT'), ('titular', 'JJ'),
                           #  ('threat', 'NN'), ('of', 'IN'), ...]

blob.noun_phrases          # WordList(['titular threat', 'blob',
                           #            'ultimate movie monster',
                           #            'amoeba-like mass', ...])

for sentence in blob.sentences:
    print(sentence.sentiment.polarity)
# 0.060
```

https://textblob.readthedocs.io

70

# Polyglot

Search docs

Installation

Language Detection

Tokenization

Command Line Interface

Downloading Models

Word Embeddings

Part of Speech Tagging

Named Entity Extraction

Morphological Analysis

Transliteration

Sentiment

polyglot

Docs » Welcome to polyglot's documentation!                    ⬡ Edit on GitHub

## Welcome to polyglot's documentation!

## polyglot

`downloads 17k/month`  `pypi package 16.7.4`  `build passing`  `docs passing`

Polyglot is a natural language pipeline that supports massive multilingual applications.

- Free software: GPLv3 license
- Documentation: http://polyglot.readthedocs.org.

## Features

- Tokenization (165 Languages)
- Language detection (196 Languages)
- Named Entity Recognition (40 Languages)
- Part of Speech Tagging (16 Languages)
- Sentiment Analysis (136 Languages)
- Word Embeddings (137 Languages)
- Morphological analysis (135 Languages)
- Transliteration (69 Languages)

https://polyglot.readthedocs.io/

71

# scikit-learn



## Classification

Identifying to which category an object belongs to.

**Applications**: Spam detection, Image recognition.
**Algorithms**: SVM, nearest neighbors, random forest, ...
— Examples

## Regression

Predicting a continuous-valued attribute associated with an object.

**Applications**: Drug response, Stock prices.
**Algorithms**: SVR, ridge regression, Lasso, ...
— Examples

## Clustering

Automatic grouping of similar objects into sets.

**Applications**: Customer segmentation, Grouping experiment outcomes
**Algorithms**: k-Means, spectral clustering, mean-shift, ...
— Examples

## Dimensionality reduction

Reducing the number of random variables to consider.

**Applications**: Visualization, Increased efficiency

## Model selection

Comparing, validating and choosing parameters and models.

**Goal**: Improved accuracy via parameter tuning

## Preprocessing

Feature extraction and normalization.

**Application**: Transforming input data such as text for use with machine learning algorithms.
**Modules**: preprocessing, feature extraction.

http://scikit-learn.org/

# The Stanford Natural Language Processing Group

home · people · teaching · research · publications · software · events · local

The Stanford NLP Group makes parts of our Natural Language Processing software available to everyone. These are statistical NLP toolkits for various major computational linguistics problems. They can be incorporated into applications with human language technology needs.

All the software we distribute here is written in Java. All recent distributions require Oracle Java 6+ or OpenJDK 7+. Distribution packages include components for command-line invocation, jar files, a Java API, and source code. A number of helpful people have extended our work with bindings or translations for other languages. As a result, much of this software can also easily be used from Python (or Jython), Ruby, Perl, Javascript, and F# or other .NET languages.

## Supported software distributions

This code is being developed, and we try to answer questions and fix bugs on a best-effort basis.

All these software distributions are open source, **licensed under the GNU General Public License** (v2 or later). Note that this is the *full* GPL, which allows many free uses, but *does not allow* its incorporation into any type of distributed proprietary software, even in part or in translation. **Commercial licensing** is also available; please contact us if you are interested.

Stanford CoreNLP
> An integrated suite of natural language processing tools for English and (mainland) Chinese in Java, including tokenization, part-of-speech tagging, named entity recognition, parsing, and coreference. See also: Stanford Deterministic Coreference Resolution, and the online CoreNLP demo, and the CoreNLP FAQ.

Stanford Parser
> Implementations of probabilistic natural language parsers in Java: highly optimized PCFG and dependency parsers, a lexicalized PCFG parser, and a deep learning reranker. See also: Online parser demo, the Stanford Dependencies page, and Parser FAQ.

Stanford POS Tagger
> A maximum-entropy (CMM) part-of-speech (POS) tagger for English,

# Stanford NLP Software

73

# Stanford CoreNLP   http://nlp.stanford.edu:8080/corenlp/process
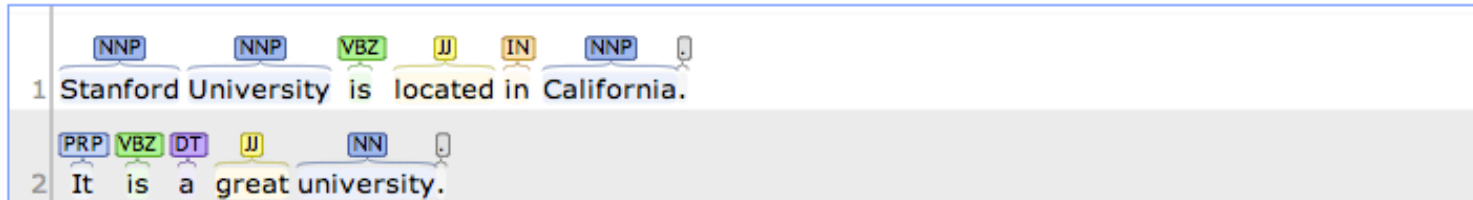
## Stanford CoreNLP

Output format:  Visualise  ⬍

Please enter your text here:

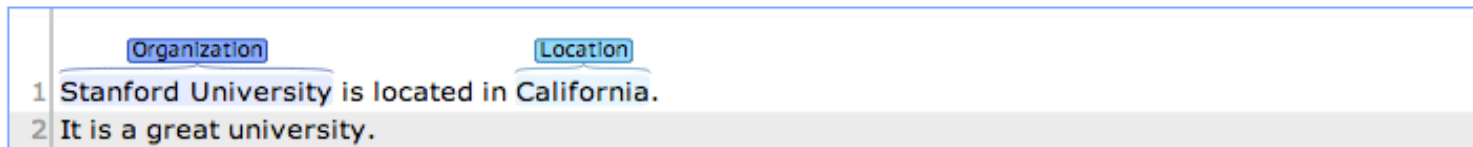Stanford University is located in California. It is a great university.
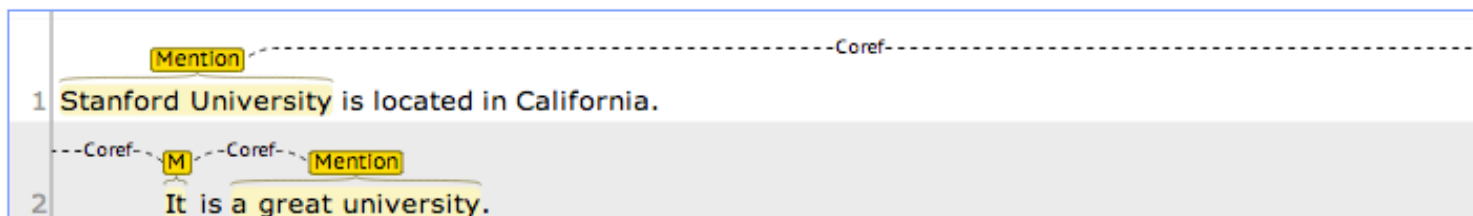
Submit    Clear

### Part-of-Speech:

```
        NNP         NNP      VBZ    JJ    IN    NNP    .
1  Stanford University  is  located in  California.

   PRP VBZ DT   JJ         NN      .
2  It   is   a  great university.
```

### Named Entity Recognition:

```
        Organization                         Location
1  Stanford University is located in California.
2  It is a great university.
```

### Coreference:

```
        Mention -------------------------------Coref-------------------------------
1  Stanford University is located in California.

   ---Coref--  M  --Coref--  Mention
2           It  is a great university.
```

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

Stanford University is located in California.
It is a great university.

**Part-of-Speech:**

| | NNP | NNP | VBZ | JJ | IN | NNP | . |
|---|---|---|---|---|---|---|---|
| 1 | Stanford | University | is | located | in | California | . |

| | PRP | VBZ | DT | JJ | NN | . |
|---|---|---|---|---|---|---|
| 2 | It | is | a | great | university | . |

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

Stanford University is located in California.
It is a great university.

**Named Entity Recognition:**

| | Organization | | | Location | |
|---|---|---|---|---|---|
| 1 | Stanford University | is located in | California | . |
| 2 | It is a great university . | | | | |

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

Stanford University is located in California.
It is a great university.



Coreference:

Mention ··········································· -Coref--

1 Stanford University is located in California .

--Coref-- M --Coref-- Mention

2 It is a great university .

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

> Stanford University is located in California.
> It is a great university.



**Basic dependencies:**

1  Stanford University is located in California.

2  It is a great university.

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process



**Collapsed dependencies:**

1  Stanford University is located in California.

2  It is a great university.

**Collapsed CC-processed dependencies:**

1  Stanford University is located in California.

2  It is a great university.

Visualisation provided using the brat visualisation/annotation software.
Copyright © 2011, Stanford University, All Rights Reserved.

# Stanford CoreNLP

Output format: Pretty print ⬍

Please enter your text here:

Stanford University is located in California. It is a great university.

Submit    Clear

## Stanford CoreNLP XML Output

**Document**

**Document Info**

**Sentences**

*Sentence #1*

*Tokens*

| Id | Word | Lemma | Char begin | Char end | POS | NER | Normalized NER | Speaker |
|----|------|-------|------------|----------|-----|-----|----------------|---------|
| 1 | Stanford | Stanford | 0 | 8 | NNP | ORGANIZATION | | PER0 |
| 2 | University | University | 9 | 19 | NNP | ORGANIZATION | | PER0 |
| 3 | is | be | 20 | 22 | VBZ | O | | PER0 |
| 4 | located | located | 23 | 30 | JJ | O | | PER0 |
| 5 | in | in | 31 | 33 | IN | O | | PER0 |
| 6 | California | California | 34 | 44 | NNP | LOCATION | | PER0 |
| 7 | . | . | 44 | 45 | . | O | | PER0 |

*Parse tree*
(ROOT (S (NP (NNP Stanford) (NNP University)) (VP (VBZ is) (ADJP (JJ located) (PP (IN in) (NP (NNP California)))))) (. .)))

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

Stanford University is located in California.
It is a great university.

*Sentence #1*

*Tokens*

| Id | Word | Lemma | Char begin | Char end | POS | NER | Normalized NER | Speaker |
|---|---|---|---|---|---|---|---|---|
| 1 | Stanford | Stanford | 0 | 8 | NNP | ORGANIZATION | | PER0 |
| 2 | University | University | 9 | 19 | NNP | ORGANIZATION | | PER0 |
| 3 | is | be | 20 | 22 | VBZ | O | | PER0 |
| 4 | located | located | 23 | 30 | JJ | O | | PER0 |
| 5 | in | in | 31 | 33 | IN | O | | PER0 |
| 6 | California | California | 34 | 44 | NNP | LOCATION | | PER0 |
| 7 | . | . | 44 | 45 | . | O | | PER0 |

*Parse tree*
(ROOT (S (NP (NNP Stanford) (NNP University)) (VP (VBZ is) (ADJP (JJ located) (PP (IN in) (NP (NNP California)))))) (. .)))

# Stanford CoreNLP

Stanford University is located in California.
It is a great university.

*Sentence #2*

*Tokens*

| Id | Word | Lemma | Char begin | Char end | POS | NER | Normalized NER | Speaker |
|----|------|-------|-----------|----------|-----|-----|----------------|---------|
| 1 | It | it | 46 | 48 | PRP | O | | PER0 |
| 2 | is | be | 49 | 51 | VBZ | O | | PER0 |
| 3 | a | a | 52 | 53 | DT | O | | PER0 |
| 4 | great | great | 54 | 59 | JJ | O | | PER0 |
| 5 | university | university | 60 | 70 | NN | O | | PER0 |
| 6 | . | . | 70 | 71 | . | O | | PER0 |

*Parse tree*
(ROOT (S (NP (PRP It)) (VP (VBZ is) (NP (DT a) (JJ great) (NN university))) (. .)))

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

Stanford University is located in California.
It is a great university.

## Coreference resolution graph

1.

| Sentence | Head | Text | Context |
|---|---|---|---|
| 1 | 2 (gov) | Stanford University | |
| 2 | 1 | It | |
| 2 | 5 | a great university | |

Tokens

| Id | Word | Lemma | Char begin | Char end | POS | NER | Normalized NER | Speaker |
|---|---|---|---|---|---|---|---|---|
| 1 | Stanford | Stanford | 0 | 8 | NNP | ORGANIZATION | | PER0 |
| 2 | University | University | 9 | 19 | NNP | ORGANIZATION | | PER0 |
| 3 | is | be | 20 | 22 | VBZ | O | | PER0 |
| 4 | located | located | 23 | 30 | JJ | O | | PER0 |
| 5 | in | in | 31 | 33 | IN | O | | PER0 |
| 6 | California | California | 34 | 44 | NNP | LOCATION | | PER0 |
| 7 | . | . | 44 | 45 | . | O | | PER0 |

Parse tree
(ROOT (S (NP (NNP Stanford) (NNP University)) (VP (VBZ is) (ADJP (JJ located) (PP (IN in) (NP (NNP California)))))) (. .)))

Uncollapsed dependencies

root ( ROOT-0 , located-4 )
nn ( University-2 , Stanford-1 )
nsubj ( located-4 , University-2 )
cop ( located-4 , is-3 )
prep ( located-4 , in-5 )
pobj ( in-5 , California-6 )
Collapsed dependencies

root ( ROOT-0 , located-4 )
nn ( University-2 , Stanford-1 )
nsubj ( located-4 , University-2 )
cop ( located-4 , is-3 )
prep_in ( located-4 , California-6 )
Collapsed dependencies with CC processed

root ( ROOT-0 , located-4 )
nn ( University-2 , Stanford-1 )
nsubj ( located-4 , University-2 )
cop ( located-4 , is-3 )
prep_in ( located-4 , California-6 )

# Stanford CoreNLP

http://nlp.stanford.edu:8080/corenlp/process

Stanford University is located in California.
It is a great university.

# Stanford CoreNLP

Output format: [ XML ▲▼ ]

Please enter your text here:

```
Stanford University is located in California. It is a great university.
```

[ Submit ]  [ Clear ]

```xml
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet href="CoreNLP-to-HTML.xsl" type="text/xsl"?>
<root>
  <document>
    <sentences>
      <sentence id="1">
        <tokens>
          <token id="1">
            <word>Stanford</word>
            <lemma>Stanford</lemma>
            <CharacterOffsetBegin>0</CharacterOffsetBegin>
            <CharacterOffsetEnd>8</CharacterOffsetEnd>
            <POS>NNP</POS>
            <NER>ORGANIZATION</NER>
            <Speaker>PER0</Speaker>
          </token>
          <token id="2">
            <word>University</word>
            <lemma>University</lemma>
            <CharacterOffsetBegin>9</CharacterOffsetBegin>
            <CharacterOffsetEnd>19</CharacterOffsetEnd>
            <POS>NNP</POS>
            <NER>ORGANIZATION</NER>
            <Speaker>PER0</Speaker>
          </token>
```

# NER for News Article

money.cnn.com/2014/05/02/technology/gates-microsoft-stock-sale/index.html

2K
TOTAL SHARES

461

1K

74

25

## Bill Gates no longer Microsoft's biggest shareholder

CNNMoney

By Patrick M. Sheridan   @CNNTech May 2, 2014: 5:46 PM ET

Recommend 1.2k

PHOTO: CHIP SOMODEVILLA/GETTY IMAGES

Bill Gates sold nearly 8 million shares of Microsoft over the past two days.

2K
TOTAL SHARES

461    1K    74    25

NEW YORK (CNNMoney)

For the first time in Microsoft's history, founder Bill Gates is no longer its largest individual shareholder.

In the past two days, Gates has sold nearly 8 million shares of Microsoft ( MSFT, Fortune

---

Bill Gates no longer Microsoft's biggest shareholder
By Patrick M. Sheridan  @CNNTech May 2, 2014: 5:46 PM ET

Bill Gates sold nearly 8 million shares of Microsoft over the past two days.

NEW YORK (CNNMoney)

For the first time in Microsoft's history, founder Bill Gates is no longer its largest individual shareholder.
In the past two days, Gates has sold nearly 8 million shares of Microsoft (MSFT, Fortune 500), bringing down his total to roughly 330 million.

That puts him behind Microsoft's former CEO Steve Ballmer who owns 333 million shares.
Related: Gates reclaims title of world's richest billionaire
Ballmer, who was Microsoft's CEO until earlier this year, was one of Gates' first hires.
It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the Bill & Melinda Gates foundation.
The foundation has spent $28.3 billion fighting hunger and poverty since its inception back in 1997.

# Stanford Named Entity Tagger (NER)

# Stanford Named Entity Tagger (NER)

http://nlp.stanford.edu:8080/ner/process

## Stanford Named Entity Tagger

Classifier:  english.muc.7class.distsim.crf.ser.gz

Output Format:  inlineXML

Preserve Spacing:  yes

Please enter your text here:

Bill Gates no longer Microsoft's biggest shareholder
By Patrick M. Sheridan  @CNNTech May 2, 2014: 5:46 PM ET

Bill Gates sold nearly 8 million shares of Microsoft over the past two
days.

NEW YORK (CNNMoney)

[Submit]  [Clear]

Bill Gates no longer <ORGANIZATION>Microsoft</ORGANIZATION>'s biggest shareholder By <PERSON>Patrick M. Sheridan</PERSON> @CNNTech <DATE>May 2, 2014</DATE>: 5:46 PM ET Bill Gates sold nearly 8 million shares of <ORGANIZATION>Microsoft</ORGANIZATION> over the past two days. <LOCATION>NEW YORK</LOCATION> (CNNMoney) For the first time in <ORGANIZATION>Microsoft</ORGANIZATION>'s history, founder <PERSON>Bill Gates</PERSON> is no longer its largest individual shareholder. In the <DATE>past two days</DATE>, Gates has sold nearly 8 million shares of <ORGANIZATION>Microsoft</ORGANIZATION> (<ORGANIZATION>MSFT</ORGANIZATION>, Fortune 500), bringing down his total to roughly 330 million. That puts him behind <ORGANIZATION>Microsoft</ORGANIZATION>'s former CEO <PERSON>Steve Ballmer</PERSON> who owns 333 million shares. Related: Gates reclaims title of world's richest billionaire <PERSON>Ballmer</PERSON>, who was <ORGANIZATION>Microsoft</ORGANIZATION>'s CEO until <DATE>earlier this year</DATE>, was one of Gates' first hires. It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the <ORGANIZATION>Bill & Melinda Gates</ORGANIZATION> foundation. The foundation has spent <MONEY>$28.3 billion</MONEY> fighting hunger and poverty since its inception back in <DATE>1997</DATE>.

# Stanford Named Entity Tagger (NER)

http://nlp.stanford.edu:8080/ner/process

# Stanford Named Entity Tagger (NER)
## http://nlp.stanford.edu:8080/ner/process

**Stanford Named Entity Tagger**

Classifier: [ english.muc.7class.distsim.crf.ser.gz ÷ ]

Output Format: [ slashTags ÷ ]

Preserve Spacing: [ yes ÷ ]

Please enter your text here:

```
Bill Gates no longer Microsoft's biggest shareholder
By Patrick M. Sheridan  @CNNTech May 2, 2014: 5:46 PM ET

Bill Gates sold nearly 8 million shares of Microsoft over the past two
days.

NEW VORK (CNNMoney)
```

[ Submit ]  [ Clear ]

Bill/O Gates/O no/O longer/O Microsoft/ORGANIZATION's/O biggest/O shareholder/O By/O Patrick/PERSON M./PERSON Sheridan/PERSON @CNNTech/O May/DATE 2/DATE,/DATE 2014/DATE:/O 5:46/O PM/O ET/O Bill/O Gates/O sold/O nearly/O 8/O million/O shares/O of/O Microsoft/ORGANIZATION over/O the/O past/O two/O days/O./O NEW/LOCATION YORK/LOCATION -LRB-/OCNNMoney/O-RRB-/O For/O the/O first/O time/O in/O Microsoft/ORGANIZATION's/O history/O,/O founder/O Bill/PERSON Gates/PERSON is/O no/O longer/O its/O largest/O individual/O shareholder/O./O In/O the/O past/DATE two/DATE days/DATE,/O Gates/O has/O sold/O nearly/O 8/O million/O shares/O of/O Microsoft/ORGANIZATION -LRB-/OMSFT/ORGANIZATION,/O Fortune/O 500/O-RRB-/O,/O bringing/O down/O his/O total/O to/O roughly/O 330/O million/O./O That/O puts/O him/O behind/O Microsoft/ORGANIZATION's/O former/O CEO/O Steve/PERSON Ballmer/PERSON who/O owns/O 333/O million/O shares/O./O Related/O:/O Gates/O reclaims/O title/O of/O world/O's/O richest/O billionaire/O Ballmer/PERSON,/O who/O was/O Microsoft/ORGANIZATION's/O CEO/O until/O earlier/DATE this/DATE year/DATE,/O was/O one/O of/O Gates/O'/O first/O hires/O./O It/O's/O a/O passing/O of/O the/O torch/O for/O Gates/O who/O has/O always/O been/O the/O largest/O single/O owner/O of/O his/O company/O's/O stock/O./O Gates/O now/O spends/O his/O time/O and/O personal/O fortune/O helping/O run/O the/O Bill/ORGANIZATION &/ORGANIZATION Melinda/ORGANIZATION Gates/ORGANIZATION foundation/O./O The/O foundation/O has/O spent/O $/MONEY28.3/MONEY billion/MONEY fighting/O hunger/O and/O poverty/O since/O its/O inception/O back/O in/O 1997/DATE./O

# Stanford Named Entity Tagger (NER)

## http://nlp.stanford.edu:8080/ner/process

**Stanford Named Entity Tagger**

Classifier: [ english.conll.4class.distsim.crf.ser.gz ÷ ]

Output Format: [ highlighted ÷ ]

Preserve Spacing: [ yes ÷ ]

Please enter your text here:

> Bill Gates no longer Microsoft's biggest shareholder
> By Patrick M. Sheridan  @CNNTech May 2, 2014: 5:46 PM ET
>
> Bill Gates sold nearly 8 million shares of Microsoft over the past two days.
>
> NEW YORK (CNNMoney)

[ Submit ]  [ Clear ]

Bill Gates no longer Microsoft's biggest shareholder By Patrick M. Sheridan @CNNTech May 2, 2014: 5:46 PM ET Bill Gates sold nearly 8 million shares of Microsoft over the past two days. NEW YORK (CNNMoney) For the first time in Microsoft's history, founder Bill Gates is no longer its largest individual shareholder. In the past two days, Gates has sold nearly 8 million shares of Microsoft (MSFT, Fortune 500), bringing down his total to roughly 330 million. That puts him behind Microsoft's former CEO Steve Ballmer who owns 333 million shares. Related: Gates reclaims title of world's richest billionaire Ballmer, who was Microsoft's CEO until earlier this year, was one of Gates' first hires. It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the Bill & Melinda Gates foundation. The foundation has spent $28.3 billion fighting hunger and poverty since its inception back in 1997.

Potential tags:
LOCATION
ORGANIZATION
PERSON
MISC

# Stanford Named Entity Tagger (NER)

http://nlp.stanford.edu:8080/ner/process

## Stanford Named Entity Tagger

Classifier: `english.all.3class.distsim.crf.ser.gz` ▼

Output Format: `highlighted` ▼

Preserve Spacing: `yes` ▼

Please enter your text here:

> Bill Gates no longer Microsoft's biggest shareholder
> By Patrick M. Sheridan  @CNNTech May 2, 2014: 5:46 PM ET
>
> Bill Gates sold nearly 8 million shares of Microsoft over the past two days.
> NEW YORK (CNNMoney)

[ Submit ]  [ Clear ]

Bill Gates no longer Microsoft's biggest shareholder By Patrick M. Sheridan @CNNTech May 2, 2014: 5:46 PM ET Bill Gates sold nearly 8 million shares of Microsoft over the past two days. NEW YORK (CNNMoney) For the first time in Microsoft's history, founder Bill Gates is no longer its largest individual shareholder. In the past two days, Gates has sold nearly 8 million shares of Microsoft (MSFT, Fortune 500), bringing down his total to roughly 330 million. That puts him behind Microsoft's former CEO Steve Ballmer who owns 333 million shares. Related: Gates reclaims title of world's richest billionaire Ballmer, who was Microsoft's CEO until earlier this year, was one of Gates' first hires. It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the Bill & Melinda Gates foundation. The foundation has spent $28.3 billion fighting hunger and poverty since its inception back in 1997.

Potential tags:
 LOCATION
 ORGANIZATION
 PERSON

Copyright © 2011, Stanford University, All Rights Reserved.

92

## Classifier: english.muc.**7class**.distsim.crf.ser.gz

Bill Gates no longer Microsoft's biggest shareholder By Patrick M. Sheridan @CNNTech May 2, 2014: 5:46 PM ET Bill Gates sold nearly 8 million shares of Microsoft over the past two days. NEW YORK (CNNMoney) For the first time in Microsoft's history, founder Bill Gates is no longer its largest individual shareholder. In the past two days, Gates has sold nearly 8 million shares of Microsoft (MSFT, Fortune 500), bringing down his total to roughly 330 million. That puts him behind Microsoft's former CEO Steve Ballmer who owns 333 million shares. Related: Gates reclaims title of world's richest billionaire Ballmer, who was Microsoft's CEO until earlier this year, was one of Gates' first hires. It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the Bill & Melinda Gates foundation. The foundation has spent $28.3 billion fighting hunger and poverty since its inception back in 1997.

Potential tags:
LOCATION
TIME
PERSON
ORGANIZATION
MONEY
PERCENT
DATE

## Classifier: english.all.**3class**.distsim.crf.ser.gz

Bill Gates no longer Microsoft's biggest shareholder By Patrick M. Sheridan @CNNTech May 2, 2014: 5:46 PM ET Bill Gates sold nearly 8 million shares of Microsoft over the past two days. NEW YORK (CNNMoney) For the first time in Microsoft's history, founder Bill Gates is no longer its largest individual shareholder. In the past two days, Gates has sold nearly 8 million shares of Microsoft (MSFT, Fortune 500), bringing down his total to roughly 330 million. That puts him behind Microsoft's former CEO Steve Ballmer who owns 333 million shares. Related: Gates reclaims title of world's richest billionaire Ballmer, who was Microsoft's CEO until earlier this year, was one of Gates' first hires. It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the Bill & Melinda Gates foundation. The foundation has spent $28.3 billion fighting hunger and poverty since its inception back in 1997.

Potential tags:
LOCATION
ORGANIZATION
PERSON

# Stanford Named Entity Tagger (NER)

## http://nlp.stanford.edu:8080/ner/process

## Stanford NER Output Format: inlineXML

Bill Gates no longer <ORGANIZATION>Microsoft</ORGANIZATION>'s biggest shareholder By <PERSON>Patrick M. Sheridan</PERSON> @CNNTech <DATE>May 2, 2014</DATE>: 5:46 PM ET Bill Gates sold nearly 8 million shares of <ORGANIZATION>Microsoft</ORGANIZATION> over the past two days. <LOCATION>NEW YORK</LOCATION> (CNNMoney) For the first time in <ORGANIZATION>Microsoft</ORGANIZATION>'s history, founder <PERSON>Bill Gates</PERSON> is no longer its largest individual shareholder. In the <DATE>past two days</DATE>, Gates has sold nearly 8 million shares of <ORGANIZATION>Microsoft</ORGANIZATION> (<ORGANIZATION>MSFT</ORGANIZATION>, Fortune 500), bringing down his total to roughly 330 million. That puts him behind <ORGANIZATION>Microsoft</ORGANIZATION>'s former CEO <PERSON>Steve Ballmer</PERSON> who owns 333 million shares. Related: Gates reclaims title of world's richest billionaire <PERSON>Ballmer</PERSON>, who was <ORGANIZATION>Microsoft</ORGANIZATION>'s CEO until <DATE>earlier this year</DATE>, was one of Gates' first hires. It's a passing of the torch for Gates who has always been the largest single owner of his company's stock. Gates now spends his time and personal fortune helping run the <ORGANIZATION>Bill & Melinda Gates</ORGANIZATION> foundation. The foundation has spent <MONEY>$28.3 billion</MONEY> fighting hunger and poverty since its inception back in <DATE>1997</DATE>.

# Stanford Named Entity Tagger (NER)

## http://nlp.stanford.edu:8080/ner/process

## Stanford NER Output Format: slashTags

Bill/O Gates/O no/O longer/O Microsoft/ORGANIZATION's/O biggest/O shareholder/O By/O Patrick/PERSON M./PERSON Sheridan/PERSON @CNNTech/O May/DATE 2/DATE,/DATE 2014/DATE:/O 5:46/O PM/O ET/O Bill/O Gates/O sold/O nearly/O 8/O million/O shares/O of/O Microsoft/ORGANIZATION over/O the/O past/O two/O days/O./O NEW/LOCATION YORK/LOCATION -LRB-/OCNNMoney/O-RRB-/O For/O the/O first/O time/O in/O Microsoft/ORGANIZATION's/O history/O,/O founder/O Bill/PERSON Gates/PERSON is/O no/O longer/O its/O largest/O individual/O shareholder/O./O In/O the/O past/DATE two/DATE days/DATE,/O Gates/O has/O sold/O nearly/O 8/O million/O shares/O of/O Microsoft/ORGANIZATION -LRB-/OMSFT/ORGANIZATION,/O Fortune/O 500/O-RRB-/O,/O bringing/O down/O his/O total/O to/O roughly/O 330/O million/O./O That/O puts/O him/O behind/O Microsoft/ORGANIZATION's/O former/O CEO/O Steve/PERSON Ballmer/PERSON who/O owns/O 333/O million/O shares/O./O Related/O:/O Gates/O reclaims/O title/O of/O world/O's/O richest/O billionaire/O Ballmer/PERSON,/O who/O was/O Microsoft/ORGANIZATION's/O CEO/O until/O earlier/DATE this/DATE year/DATE,/O was/O one/O of/O Gates/O'/O first/O hires/O./O It/O's/O a/O passing/O of/O the/O torch/O for/O Gates/O who/O has/O always/O been/O the/O largest/O single/O owner/O of/O his/O company/O's/O stock/O./O Gates/O now/O spends/O his/O time/O and/O personal/O fortune/O helping/O run/O the/O Bill/ORGANIZATION &/ORGANIZATION Melinda/ORGANIZATION Gates/ORGANIZATION foundation/O./O The/O foundation/O has/O spent/O $/MONEY28.3/MONEY billion/MONEY fighting/O hunger/O and/O poverty/O since/O its/O inception/O back/O in/O 1997/DATE./O

# Vector Representations of Words

# Word Embeddings

# Word2Vec

# GloVe

# Modern NLP Pipeline

# Facebook Research FastText

Pre-trained word vectors
Word2Vec
wiki.zh.vec (861MB)
332647 word
300 vec

Pre-trained word vectors for 90 languages, trained on Wikipedia using fastText.

These vectors in dimension 300 were obtained using the skip-gram model with default parameters.

https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md

Source: Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov. "Enriching word vectors with subword information." *arXiv preprint arXiv:1607.04606* (2016).

# Facebook Research FastText Word2Vec: wiki.zh.vec

## (861MB) (332647 word 300 vec)



https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md

# Word Embeddings in LSTM RNN

Time Expanded LSTM Network



LSTM Internal States

Word Embeddings

Input Question    Is    this    person    dancing    ?

Fixed length question vector encoded by the LSTM

# NLP Tools: spaCy vs. NLTK

| | SPACY | SYNTAXNET | NLTK | CORENLP |
|---|---|---|---|---|
| Easy installation | + | − | + | + |
| Python API | + | − | + | − |
| Multi-language support | ○ | + | + | + |
| Tokenization | + | + | + | + |
| Part-of-speech tagging | + | + | + | + |
| Sentence segmentation | + | + | + | + |
| Dependency parsing | + | + | − | + |
| Entity Recognition | + | − | + | + |
| Integrated word vectors | + | − | − | − |
| Sentiment analysis | + | − | + | + |
| Coreference resolution | − | − | − | + |

Source: https://spacy.io/docs/api/

101

# Natural Language Processing (NLP) spaCy

1. Tokenization
2. Part-of-speech tagging
3. Sentence segmentation
4. Dependency parsing
5. Entity Recognition
6. Integrated word vectors
7. Sentiment analysis
8. Coreference resolution

# spaCy:
# Fastest Syntactic Parser

| SYSTEM | LANGUAGE | ACCURACY | SPEED (WPS) |
|---|---|---|---|
| **spaCy** | **Cython** | **91.8** | **13,963** |
| ClearNLP | Java | 91.7 | 10,271 |
| CoreNLP | Java | 89.6 | 8,602 |
| MATE | Java | **92.5** | 550 |
| Turbo | C++ | 92.4 | 349 |

Source: https://spacy.io/docs/api/

# Processing Speed of NLP libraries

| SYSTEM | ABSOLUTE (MS PER DOC) | | | RELATIVE (TO SPACY) | | |
|---|---|---|---|---|---|---|
| | TOKENIZE | TAG | PARSE | TOKENIZE | TAG | PARSE |
| **spaCy** | **0.2ms** | **1ms** | **19ms** | **1x** | **1x** | **1x** |
| CoreNLP | 2ms | 10ms | 49ms | 10x | 10x | 2.6x |
| ZPar | 1ms | 8ms | 850ms | 5x | 8x | 44.7x |
| NLTK | 4ms | 443ms | n/a | 20x | 443x | n/a |

Source: https://spacy.io/docs/api/

# Google SyntaxNet (2016):
# Best Syntactic Dependency Parsing Accuracy

| SYSTEM | NEWS | WEB | QUESTIONS |
|---|---|---|---|
| spaCy | 92.8 | n/a | n/a |
| Parsey McParseface | 94.15 | 89.08 | 94.77 |
| Martins et al. (2013) | 93.10 | 88.23 | 94.21 |
| Zhang and McDonald (2014) | 93.32 | 88.65 | 93.37 |
| Weiss et al. (2015) | 93.91 | 89.29 | 94.17 |
| Andor et al. (2016) | **94.44** | **90.17** | **95.40** |

# Named Entity Recognition (NER)

| SYSTEM | PRECISION | RECALL | F-MEASURE |
|---|---|---|---|
| spaCy | 0.7240 | 0.6514 | 0.6858 |
| **CoreNLP** | **0.7914** | **0.7327** | **0.7609** |
| NLTK | 0.5136 | 0.6532 | 0.5750 |
| LingPipe | 0.5412 | 0.5357 | 0.5384 |

# Text Classification

Source: https://developers.google.com/machine-learning/guides/text-classification/

# Text Classification Workflow

- Step 1: Gather Data

- Step 2: Explore Your Data

- Step 2.5: Choose a Model*

- Step 3: Prepare Your Data

- Step 4: Build, Train, and Evaluate Your Model

- Step 5: Tune Hyperparameters

- Step 6: Deploy Your Model

# Text Classification Flowchart

# Text Classification S/W<1500: N-gram

# Text Classification S/W>=1500: Sequence

Source: https://developers.google.com/machine-learning/guides/text-classification/step-2-5

# Step 2.5: Choose a Model
## Samples/Words < 1500
## 150,000/100 = 1500



IMDb review dataset,
the samples/words-per-sample ratio is ~ 144

Source: https://developers.google.com/machine-learning/guides/text-classification/step-2-5

# Step 2.5: Choose a Model
## Samples/Words < 15,000
## 1,500,000/100 = 15,000

Source: https://developers.google.com/machine-learning/guides/text-classification/step-2-5

# Step 3: Prepare Your Data

```
Texts:
T1: 'The mouse ran up the clock'
T2: 'The mouse ran down'

Token Index:
{'the': 1, 'mouse': 2, 'ran': 3, 'up': 4, 'clock': 5, 'down': 6,}.
   NOTE: 'the' occurs most frequently,
         so the index value of 1 is assigned to it.
         Some libraries reserve index 0 for unknown tokens,
         as is the case here.

Sequence of token indexes:
```
T1: 'The mouse ran up the clock' =
     [1, 2, 3, 4, 1, 5]
T1: 'The mouse ran down' =
     [1, 2, 3, 6]

# One-hot encoding

'The mouse ran up the clock' =

| | |
|---|---|
| The | 1 |
| mouse | 2 |
| ran | 3 |
| up | 4 |
| the | 1 |
| clock | 5 |

```
[ [0, 1, 0, 0, 0, 0, 0],
  [0, 0, 1, 0, 0, 0, 0],
  [0, 0, 0, 1, 0, 0, 0],
  [0, 0, 0, 0, 1, 0, 0],
  [0, 1, 0, 0, 0, 0, 0],
  [0, 0, 0, 0, 0, 1, 0] ]

  [0, 1, 2, 3, 4, 5, 6]
```

Source: https://developers.google.com/machine-learning/guides/text-classification/step-3

# Word embeddings



Male-Female     Verb Tense     Country-Capital

# Word embeddings

Source: https://developers.google.com/machine-learning/guides/text-classification/step-3

# Sequence to Sequence (Seq2Seq)

# Transformer (Attention is All You Need)
## (Vaswani et al., 2017)

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT (Bidirectional Encoder Representations from Transformers)

## Overall pre-training and fine-tuning procedures for BERT

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT (Bidirectional Encoder Representations from Transformers)

### BERT input representation

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

# BERT, OpenAI GPT, ELMo

# Fine-tuning BERT on Different Tasks



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

(c) Question Answering Tasks:
SQuAD v1.1

(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

# Pre-trained Language Model (PLM)



**Semi-supervised Sequence Learning**
**context2Vec**
**Pre-trained seq2seq**

ULMFiT — ELMo — GPT

Multi-lingual — Transformer — Bidirectional LM

MultiFiT

Cross-lingual — BERT — Larger model / More data — GPT-2 — Defense — Grover

Multi-task

XLM UDify

MT-DNN — + Generation

Knowledge distillation — MASS UniLM

MT-DNN$_{KD}$

Span prediction / Remove NSP

Longer time / Remove NSP / More data

Permutation LM / Transformer-XL / More data

+Knowledge Graph

Cross-modal

Whole Word Masking

SpanBERT

RoBERTa

XLNet

ERNIE (Tsinghua)

Neural entity linker

KnowBert

VideoBERT / CBT / ViLBERT / VisualBERT / B2T2 / Unicoder-VL / LXMERT / VL-BERT / UNITER

ERNIE (Baidu) / BERT-wwm

By Xiaozhi Wang & Zhengyan Zhang @THUNLP

Source: https://github.com/thunlp/PLMpapers

124

# Turing Natural Language Generation (T-NLG)



Source: https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/

125

# Transformers

## State-of-the-art Natural Language Processing for TensorFlow 2.0 and PyTorch

- Transformers
  - pytorch-transformers
  - pytorch-pretrained-bert
- provides state-of-the-art general-purpose architectures
  - (BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet, CTRL...)
  - for Natural Language Understanding (NLU) and
    Natural Language Generation (NLG)
    with over 32+ pretrained models
    in 100+ languages
    and deep interoperability between
    TensorFlow 2.0 and
    PyTorch.

# Transfer Learning
# in Natural Language Processing

# NLP Benchmark Datasets

| Task | Dataset | Link |
|---|---|---|
| Machine Translation | WMT 2014 EN-DE<br>WMT 2014 EN-FR | http://www-lium.univ-lemans.fr/~schwenk/cslm_joint_paper/ |
| Text Summarization | CNN/DM<br>Newsroom<br>DUC<br>Gigaword | https://cs.nyu.edu/~kcho/DMQA/<br>https://summari.es/<br>https://www-nlpir.nist.gov/projects/duc/data.html<br>https://catalog.ldc.upenn.edu/LDC2012T21 |
| Reading Comprehension<br>Question Answering<br>Question Generation | ARC<br>CliCR<br>CNN/DM<br>NewsQA<br>RACE<br>SQuAD<br>Story Cloze Test<br>NarativeQA<br>Quasar<br>SearchQA | http://data.allenai.org/arc/<br>http://aclweb.org/anthology/N18-1140<br>https://cs.nyu.edu/~kcho/DMQA/<br>https://datasets.maluuba.com/NewsQA<br>http://www.qizhexie.com/data/RACE_leaderboard<br>https://rajpurkar.github.io/SQuAD-explorer/<br>http://aclweb.org/anthology/W17-0906.pdf<br>https://github.com/deepmind/narrativeqa<br>https://github.com/bdhingra/quasar<br>https://github.com/nyu-dl/SearchQA |
| Semantic Parsing | AMR parsing<br>ATIS (SQL Parsing)<br>WikiSQL (SQL Parsing) | https://amr.isi.edu/index.html<br>https://github.com/jkkummerfeld/text2sql-data/tree/master/data<br>https://github.com/salesforce/WikiSQL |
| Sentiment Analysis | IMDB Reviews<br>SST<br>Yelp Reviews<br>Subjectivity Dataset | http://ai.stanford.edu/~amaas/data/sentiment/<br>https://nlp.stanford.edu/sentiment/index.html<br>https://www.yelp.com/dataset/challenge<br>http://www.cs.cornell.edu/people/pabo/movie-review-data/ |
| Text Classification | AG News<br>DBpedia<br>TREC<br>20 NewsGroup | http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html<br>https://wiki.dbpedia.org/Datasets<br>https://trec.nist.gov/data.html<br>http://qwone.com/~jason/20Newsgroups/ |
| Natural Language Inference | SNLI Corpus<br>MultiNLI<br>SciTail | https://nlp.stanford.edu/projects/snli/<br>https://www.nyu.edu/projects/bowman/multinli/<br>http://data.allenai.org/scitail/ |
| Semantic Role Labeling | Proposition Bank<br>OneNotes | http://propbank.github.io/<br>https://catalog.ldc.upenn.edu/LDC2013T19 |

# A High-Level Depiction of DeepQA Architecture

# Chatbots
# Bot Maturity Model

Customers want to have simpler means to interact with businesses and get faster response to a question or complaint.

Source: https://www.capgemini.com/2017/04/how-can-chatbots-meet-expectations-introducing-the-bot-maturity/

# Dialogue
# on
# Airline Travel Information System (ATIS)

# The ATIS
# (Airline Travel Information System) Dataset

https://www.kaggle.com/siddhadev/atis-dataset-from-ms-cntk

| Sentence | what | flights | leave | from | phoenix |
|----------|------|---------|-------|------|---------|
| Slots | O | O | O | O | B-fromloc |
| Intent | atis_flight | | | | |

Training samples: 4978
Testing samples: 893
Vocab size: 943
Slot count: 129
Intent count: 26

# SF-ID Network (E et al., 2019)
# Slot Filling (SF)
# Intent Detection (ID)

**A Novel Bi-directional Interrelated Model for Joint Intent Detection and Slot Filling**

133

# Intent Detection on ATIS State-of-the-art

## Intent Detection on ATIS



| RANK | METHOD | ACCURACY | PAPER TITLE | YEAR | PAPER | CODE |
|---|---|---|---|---|---|---|
| 1 | SF-ID | 0.9776 | A Novel Bi-directional Interrelated Model for Joint Intent Detection and Slot Filling | 2019 | 📄 | ⭕ |
| 2 | Capsule-NLU | 0.950 | Joint Slot Filling and Intent Detection via Capsule Neural Networks | 2018 | 📄 | ⭕ |

# Slot Filling on ATIS State-of-the-art

## Slot Filling on ATIS



| RANK | METHOD | F1 | PAPER TITLE | YEAR | PAPER | CODE |
|------|--------|-----|-------------|------|-------|------|
| 1 | SF-ID | 0.958 | A Novel Bi-directional Interrelated Model for Joint Intent Detection and Slot Filling | 2019 | 📄 | 🔵 |
| 2 | Capsule-NLU | 0.952 | Joint Slot Filling and Intent Detection via Capsule Neural Networks | 2018 | 📄 | 🔵 |

Source: https://paperswithcode.com/sota/slot-filling-on-atis

135

# TensorFlow NLP Examples

- ## Basic Text Classification (Text Classification) (46 Seconds)
  - https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/keras/basic_text_classification.ipynb

- ## NMT with Attention (20-30 minutes)
  - https://colab.research.google.com/github/tensorflow/tensorflow/blob/master/tensorflow/contrib/eager/python/examples/nmt_with_attention/nmt_with_attention.ipynb

# Text Classification
# IMDB Movie Reviews

https://colab.research.google.com/drive/1x16h1GhHsLIrLYtPCvCHaoO1W-i_gror

Source: https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/keras/basic_text_classification.ipynb

# Summary

- Text Analytics and Text Mining

- Natural Language Processing (NLP)

# References

- Dipanjan Sarkar (2019), Text Analytics with Python: A Practitioner's Guide to Natural Language Processing, Second Edition. APress.

- Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python: Enabling Language-Aware Data Products with Machine Learning, O'Reilly.

- Charu C. Aggarwal (2018), Machine Learning for Text, Springer.

- Gabe Ignatow and Rada F. Mihalcea (2017), An Introduction to Text Mining: Research Design, Data Collection, and Analysis, SAGE Publications.

- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

- Sebastian Ruder, Matthew E. Peters, Swabha Swayamdipta, and Thomas Wolf (2019). "Transfer learning in natural language processing." In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Tutorials, pp. 15-18.

- Amirsina Torfi, Rouzbeh A. Shirvani, Yaser Keneshloo, Nader Tavvaf, and Edward A. Fox (2020). "Natural Language Processing Advancements By Deep Learning: A Survey." arXiv preprint arXiv:2003.01200.