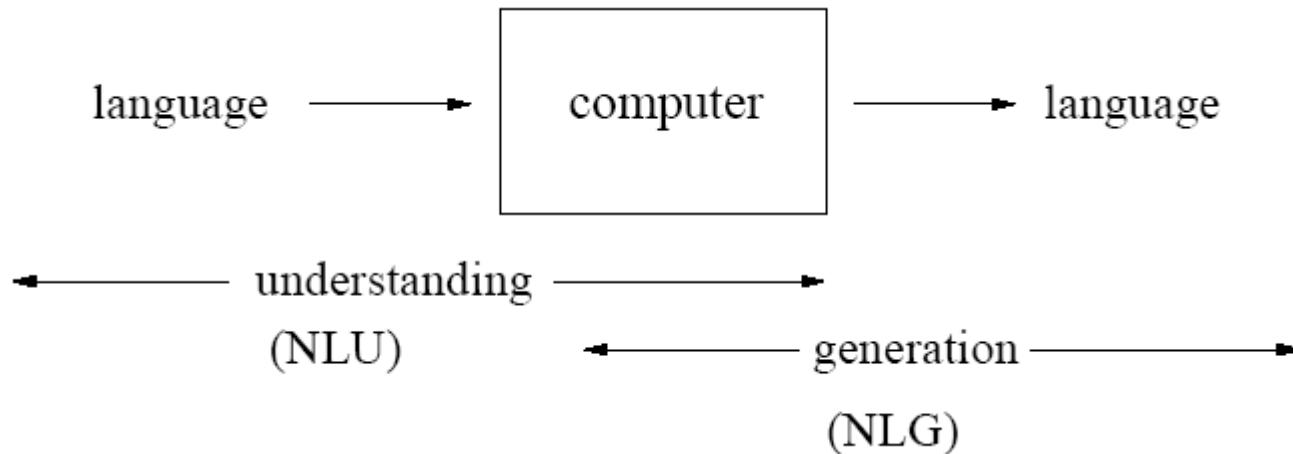


# Introduction to Natural Language Processing

(Some slides adapted from Ralph Grishman at NYU,  
Yejin Choi at UWashington, N. Tomura at Depaul, Jurafsky and Martin,  
CS224N, CS224d at Stanford and other resources on the web)

# What is NLP?

- Natural Language Processing (NLP) is a field in Artificial Intelligence (AI) devoted to creating computers that use natural language as input and/or output.



- An interdisciplinary field with many names corresponding to its many facets, names like speech and language processing, human language technology, natural language processing, computational linguistics, and speech recognition and synthesis.

# What is NLP?

- Natural language processing is a field at the intersection of
  - computer science
  - artificial intelligence
  - and linguistics.
- Goal: for computers to process or “understand” natural language in order to perform tasks that are useful, e.g.
  - Question Answering
- Fully **understanding and representing** the **meaning** of language (or even defining it) is an illusive goal.
- Perfect language understanding is AI-complete



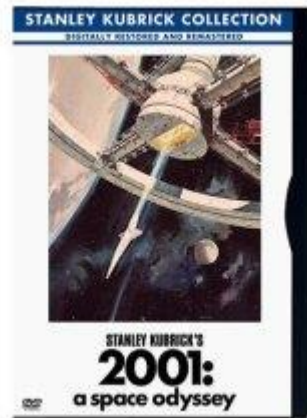
(\*) To call a problem AI-complete reflects an attitude that it would not be solved by a simple specific algorithm.

# Why NLP?

- The goal of this new field is to get computers to perform useful tasks involving human language, tasks like enabling human-machine communication, improving human-human communication, or simply doing useful processing of text or speech.
- To interact with computing devices using human (natural) languages. For example,
  - Building intelligent robots (AI).
  - Enabling voice-controlled operation.
- To access (large amount of) information and knowledge stored in the form of human languages quickly.

# Early days of NLP: The Dream: Machines that Can Speak

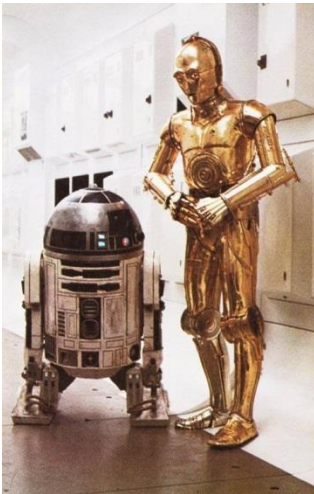
- One example of a useful such task is a **conversational agent**.
- HAL 9000 in “2001: A Space Odyssey” (\*)



(\*) The HAL 9000 computer in Stanley Kubrick's film *2001: A Space Odyssey* is one of the most recognizable characters in 20th century cinema. HAL is an artificial agent capable of such advanced language behavior as speaking and understanding English, and at a crucial moment in the plot, even reading lips.

# Machines that Can Speak (cont.)

- C3PO  
in Star Wars



- KITT  
in Knight Rider



# But Still a Sci-Fi...

- “However, as 2001 approached it became clear that *2001*'s predictions in computer technology were far fetched. Natural language, lip reading, planning and plain common sense in computers were still the stuff of science fiction.”

*HAL 9000 - Wikipedia*

# Test for Intelligence – Turing Test

- Alan Turing (1950) proposed a test of a machine's capability to perform human-like conversation.



- A human judge engages in a natural language conversation with two other parties, one a human and the other a machine; if the judge cannot reliably tell which is which, then the machine is said to pass the test.

# Early Conversational Programs

- ELIZA (by Joseph Weizenbaum), 1966
  - A psychotherapist
  - No real understanding; simple pattern-matching to respond to user input -> canned responses

(my ?x depresses me) (why does your ?x depress you) 5

(life ?x) (why do you say it ?x) 3

(I could ?x) (you could ?x) 2

(because ?x) (that is a good reason) 3

(?x) (tell me more) 0

>> Hello.

How do you do. Please state your problem.

>> Men are all alike.

In what way ?

>> They're always bugging us about something or other.  
I'm not sure I understand you fully.

>> They're always bugging us about something or other.  
Can you think of a specific example ?

>> Well, my boyfriend made me come here.  
Your boyfriend made you come here ?

>> He says I'm depressed much of the time.  
I am sorry to hear that you are depressed.

# Dialogue systems/conversational agents

- ELIZA and HAL examples of dialogue systems
- Require both understanding and generation
  - Dave: Open the pod bay doors, HAL.
  - HAL: I'm sorry Dave, I'm afraid I can't do that.
  - Dave: What's the problem?
  - HAL: I think you know what the problem is just as well as I do.



# Modern NLP: NLP in the Days of Big Data

Three trends:

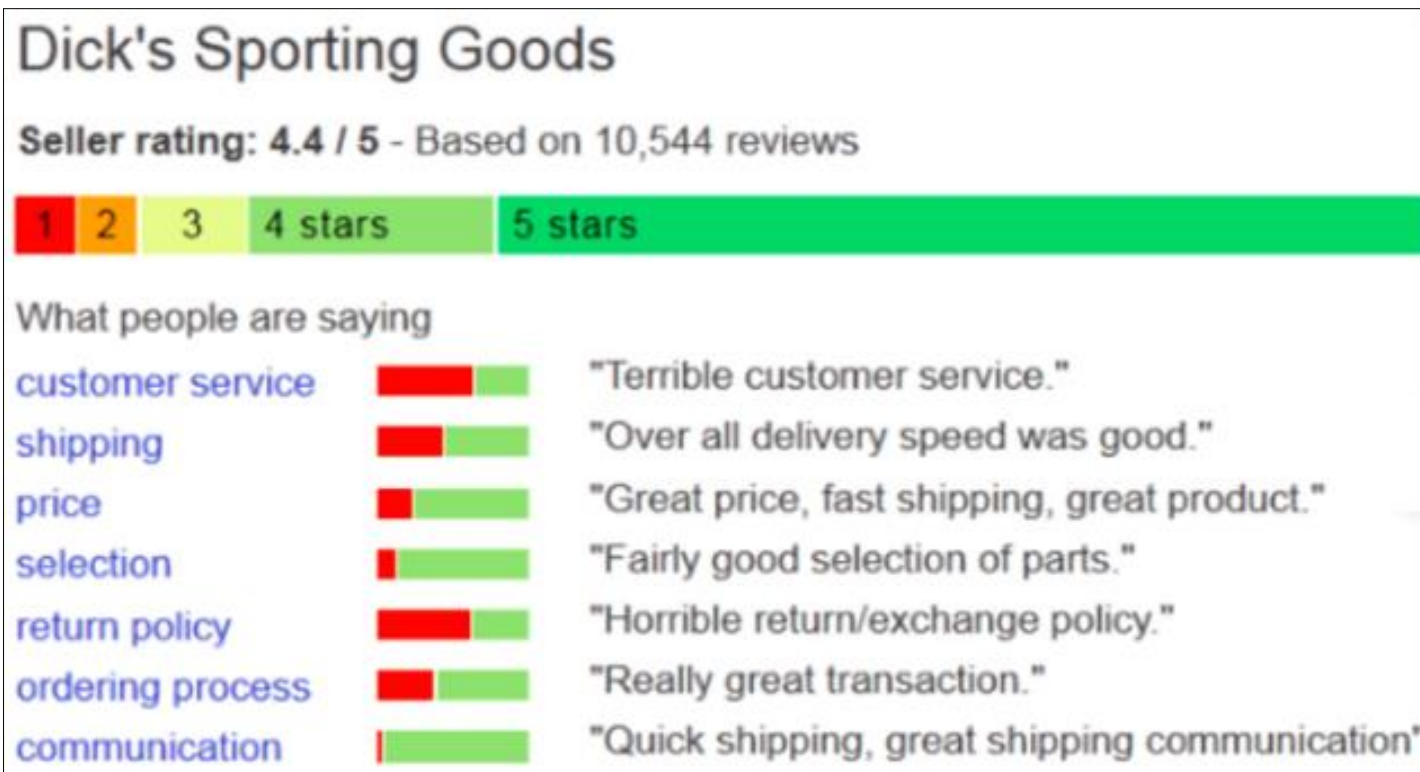
1. An **enormous amount of information** is now available in machine readable form as natural language text (newspapers, web pages, medical records, financial filings, product reviews, discussion forums, etc.)
2. Conversational agents are becoming an important form of human-computer **communication**
3. Much of human-human interaction is now mediated by computers via **social media**

# NLP Applications

- Three prominent application areas:
  - Text analytics/mining (from “***unstructured data***”)
    - Sentiment analysis
    - Topic identification
    - Digital Humanities (“*new ways of doing scholarship that involve collaborative, transdisciplinary, and computationally engaged research, teaching, and publishing.*”)
  - Conversational agents
    - Siri, Cortana, Amazon Alexa, Google Assistant
    - Chatbots
  - Machine translation

# Text Analytics

- Data-mining of weblogs, microblogs, discussion forums, user reviews, and other forms of user-generated media.



# Text Analytics (cont.)

- Typically this involves the extraction of **limited** kinds of semantic and pragmatic information from texts
  - Entity mentions
  - Concept identification
  - Sentiment

The screenshot shows a web interface titled "API TEST TOOL" with a red header. Below the header are three dropdown menus: "English", "Entities", and "Graphical". The main content area displays a paragraph of text with various entities highlighted in colored boxes. To the right of the text is a legend titled "LEGEND color key" with the word "ENTITIES" in red. The legend lists 10 types of entities with corresponding colored icons: Person name (red person icon), Car license plate (pink car icon), Place (purple location pin icon), Phone number (orange phone icon), Email address (orange envelope icon), Date (brown calendar icon), Hour (grey clock icon), Money (black dollar sign icon), Address (purple house icon), and Twitter hashtag (brown hashtag icon).

English    Entities    Graphical

I really enjoyed using the **Canon Ixus** in **Madrid** on **March 4**. The **Panasonic Lumix** is a bit disappointing, but the **Canon** camera is not bad at all. All I want when taking photos is point it and then just press the button. For only **200 dollars**, a really fair price, this camera is perfect for me. Besides, I have had a good customer service experience. **John Faraday** was very nice!

LEGEND color key    ENTITIES

Type of entities:

Person name	Date
Car license plate	Hour
Place	Money
Phone number	Address
Email address	Twitter hashtag

# Information retrieval

Topic: Advantages and disadvantages of using potassium hydroxide in any aspect of organic farming, especially...

information need

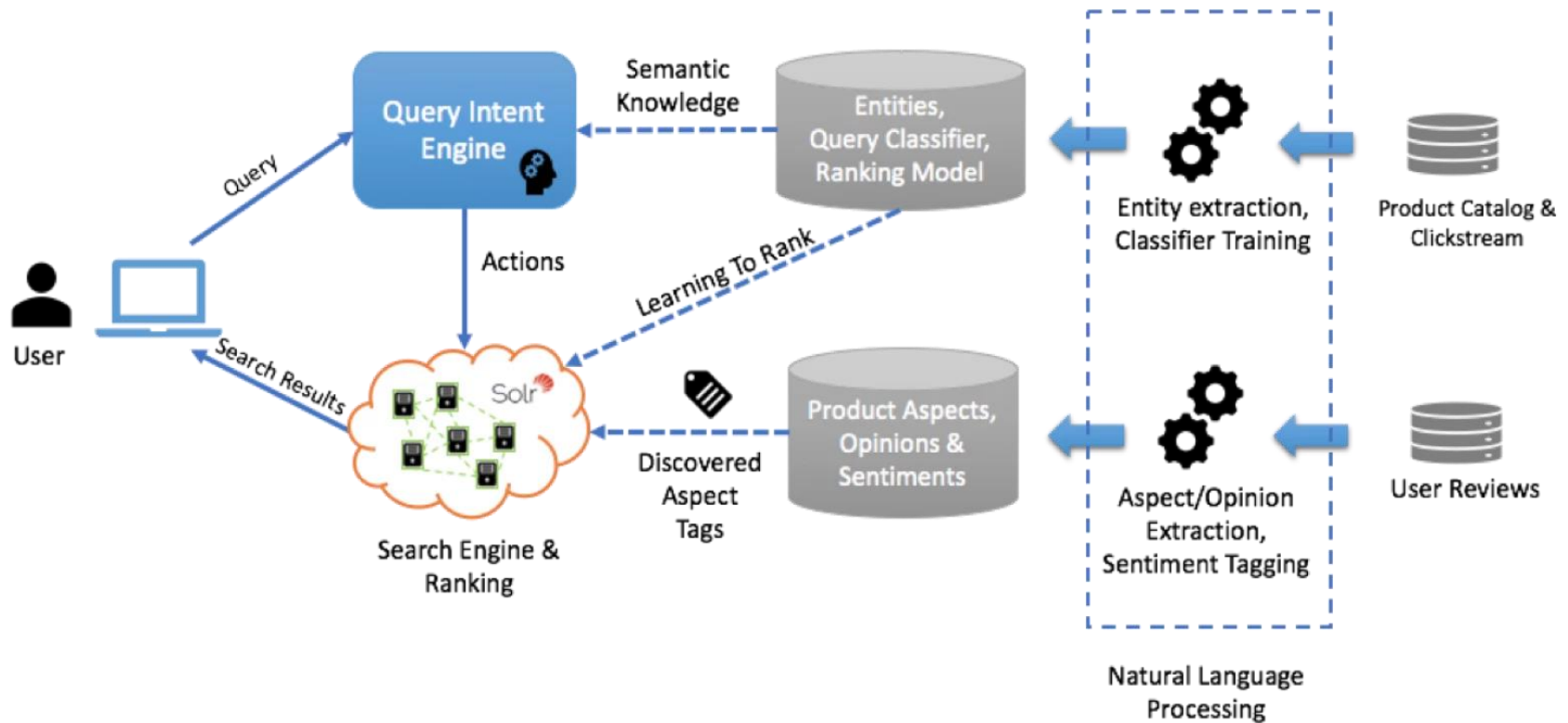


<i>doc 1</i>	<i>score</i>
<i>doc 2</i>	<i>score</i>
<i>doc 3</i>	<i>score</i>
...	
<i>doc n</i>	<i>score</i>

relevant documents  
(ranked)

IR system

# Use Machine Learning and NLP to Empower Search



# Search, and way beyond search

The image shows a Google search interface for the term "sars". The search bar at the top contains "sars" and shows a search icon. Below the search bar, there are navigation tabs for "All", "News", "Images", "Books", "Videos", and "More". The search results show "About 554,000,000 results (0.49 seconds)". The first result is from WHO: "Severe Acute Respiratory Syndrome (SARS) - WHO | World ...". The second result is from CDC: "Severe Acute Respiratory Syndrome | SARS-CoV Disease". A "COVID-19" tag is visible. A "Common questions" section lists two questions: "What is the difference between SARS-CoV-2 and COVID-19?" and "How are COVID-19 and SARS-CoV-2 related?". On the right, a knowledge panel titled "Severe acute respiratory syndrome" provides an overview, symptoms, and treatments. It states that SARS is a contagious and sometimes fatal respiratory illness caused by a coronavirus, first identified in 2002 in China. Symptoms include fever, dry cough, headache, muscle aches, and difficulty breathing. It notes that no treatment exists except supportive care and that the disease is extremely rare, with fewer than 1,000 US cases per year.

Google

sars

All News Images Books Videos More Tools

About 554,000,000 results (0.49 seconds)

<https://www.who.int> > Health topics >

### Severe Acute Respiratory Syndrome (SARS) - WHO | World ...

Severe acute respiratory syndrome (SARS) is a viral respiratory disease caused by a SARS-associated coronavirus. It was first identified at the end of ...

<https://www.cdc.gov> > sars >

### Severe Acute Respiratory Syndrome | SARS-CoV Disease

Severe acute respiratory syndrome (SARS) is a viral respiratory illness caused by a coronavirus called SARS-associated coronavirus (SARS-CoV).

Basics Fact Sheet · About SARS · Frequently Asked Questions · SARS

COVID-19 >

### Common questions

What is the difference between SARS-CoV-2 and COVID-19? ▾

How are COVID-19 and SARS-CoV-2 related? ▾

### Severe acute respiratory syndrome

Also called: SARS

OVERVIEW SYMPTOMS TREATMENTS SPEC

A contagious and sometimes fatal respiratory illness caused by a coronavirus.

SARS appeared in 2002 in China. It spread worldwide within a few months, though it was quickly contained. SARS is a virus transmitted through droplets that enter the air when someone with the disease coughs, sneezes, or talks. No known transmission has occurred since 2004.

Fever, dry cough, headache, muscle aches, and difficulty breathing are symptoms.

No treatment exists except supportive care.

### Extremely rare

Fewer than 1,000 US cases per year

# Information Extraction

- Information extraction is the process of extracting information from unstructured textual sources to enable finding entities as well as classifying and storing them in a database.
- Ex: Convert unstructured text to database entries

New York Times Co. named Russell T. Lewis, 45, president and general manager of its flagship New York Times newspaper, responsible for all business-side activities. He was executive vice president and deputy general manager. He succeeds Lance R. Primis, who in September was named president and chief operating officer of the parent.

Person	Company	Post	State
Russell T. Lewis	New York Times newspaper	president and general manager	start
Russell T. Lewis	New York Times newspaper	executive vice president	end
Lance R. Primis	New York Times Co.	president and CEO	start

# Text Summarization using NLP

## Text Summarization using NLP

### Natural Language Processing

Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data. The result is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves.

Summary

`summarize(text, 0.6)`

### Natural Language Processing

Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data.

# Image captioning

A person riding a motorcycle on a dirt road.



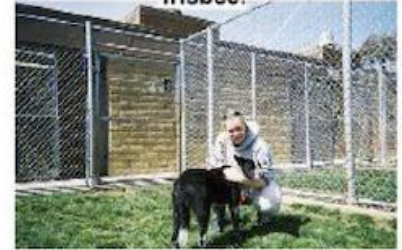
Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Describes without errors

Describes with minor errors

Somewhat related to the image

Unrelated to the image

Sutskever et al. 2014

# Demo

- Sentiment Analysis with Python NLTK Text Classification
  - <http://text-processing.com/demo/sentiment/>
- Tweet Sentiment Visualization Tool
  - [https://www.csc2.ncsu.edu/faculty/healey/tweet\\_viz/tweet\\_app/](https://www.csc2.ncsu.edu/faculty/healey/tweet_viz/tweet_app/)
- Concept Extraction
  - <http://aylien.com/>
  - <https://www.textrazor.com/>
- Text Summarization
  - <https://quillbot.com/summarize>
  - <https://www.paraphraser.io/text-summarizer>
  - <https://www.summarizer.org/>

# Demo: Text generation

**INPUT**

PRODUCT NAME

Glossier

DESCRIPTION OF YOUR PRODUCT

We're creating the new beauty essentials: easy-to-use skincare and makeup that form the backbone to your routine. Try it out this Christmas.

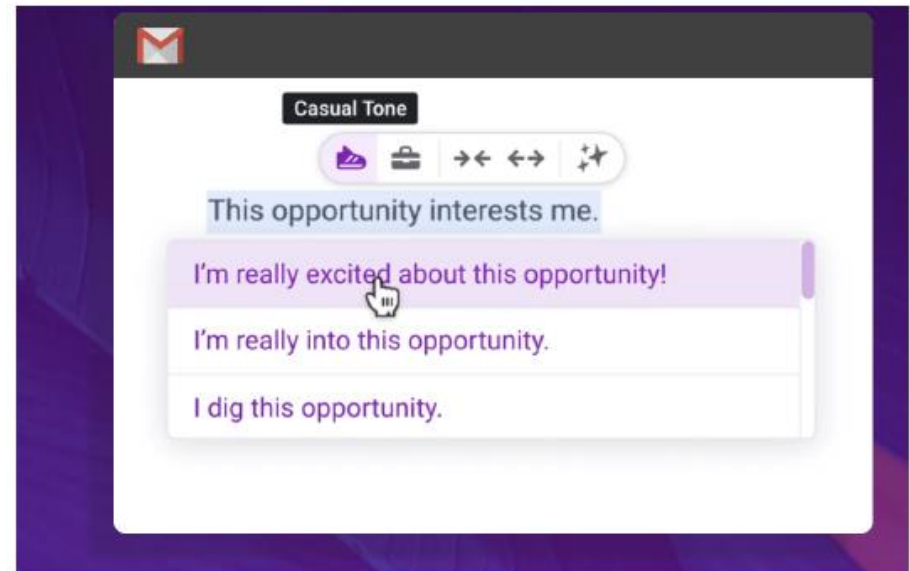
**RESULTS**

Hey, guess what. It's a present for you! Christmas is right around the corner and we're delightfully excited to share the gift of Glossier with you. 🎁❤️🎄

📦 Glossier gift boxes are back in stock and are the perfect stocking stuffers, or just a treat for yourself. Our skincare products clear up and prevent breakouts while improving skin clarity and texture. 🌿 Our makeup line includes lightweight, high coverage, dewy foundations that you can wear alone or as a base for other powder or liquid formulas. 🔍 The perfect Christmas gift!

This Christmas, give your friends and family a stack of Glossier goodie bags. Try out our minimalist skincare and makeup and give your new routine a try when 2021 rolls around.

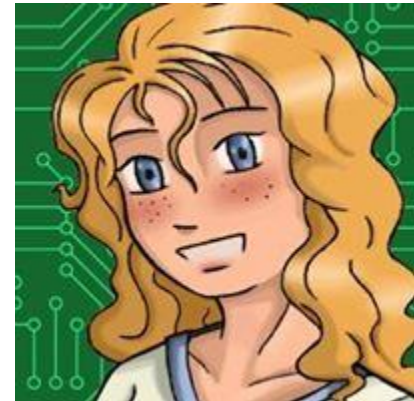
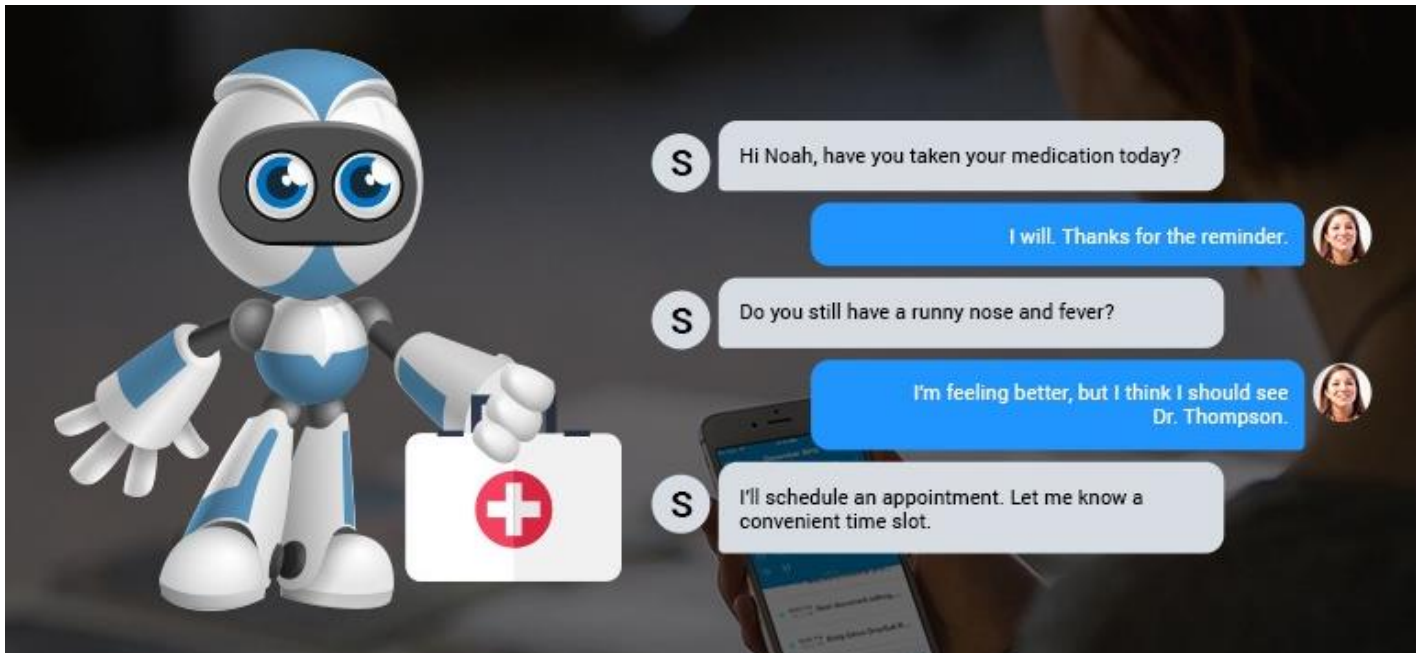
<https://www.copy.ai>



<https://www.wordtune.com/>

# Conversational Agents

- Combine
  - Speech recognition/synthesis
  - Question answering
    - From the web and from structured information sources (freebase, dbpedia, yago, etc.)
  - Simple agent-like abilities
    - Create/edit calendar entries
    - Reminders
    - Directions
    - Invoking/interacting with other apps



Mitsuku

# Artificial assistants



# The promise of these artificial assistants



**You:** Any good burger joints around here?

**Siri:** I found a number of burger restaurants near you.

**You:** Hmm. How about tacos?

**Apple:** [Siri remembers that you asked about restaurants, so it will look for Mexican restaurants in the neighborhood. And Siri is proactive, so it will question you until it finds what you're looking for.]

JO

Summarize Stanford's CS 224n "NLP with Deep Learning" course in one sentence, seamlessly integrating how delicious tea is into your summary.



The CS 224n "NLP with Deep Learning" course at Stanford provides a comprehensive introduction to natural language processing techniques, including deep learning methods like word embeddings, recurrent neural networks and transformers, like a piping hot and delicious cup of tea, the course covers many important and current NLP tasks such as machine translation, question answering and named entity recognition that will satisfy one's intellectual thirst.



# ChatGPT: A recent, intriguing set of capabilities



JO

Convert this python dictionary into a LaTeX table that uses the `booktabs` package, and make sure there are no vertical lines.

```
{"Models": ["baseline", "ours"],  
 "Accuracy": [0.7, 0.8],  
 "Time/sentence": ["4s", "6s"]  
}
```

Models	Accuracy	Time/sentence
baseline	0.7	4s
ours	0.8	6s



Here's an example of how you can convert the Python dictionary into a LaTeX table that uses the `booktabs` package and does not include any vertical lines:

Copy code

```
\begin{tabular}{l c c}  
\toprule  
\textbf{Models} & \textbf{Accuracy} & \textbf{Time/sentence}  
\\\br/>\midrule  
baseline & 0.7 & 4s \\\br/>ours & 0.8 & 6s \\\br/>\bottomrule  
\end{tabular}
```

# Question Answering

- Traditional *information retrieval* provides documents/resources that provide users with what they need to satisfy their information needs.
- *Question answering* on the other hand directly provides an answer to information needs posed as questions.

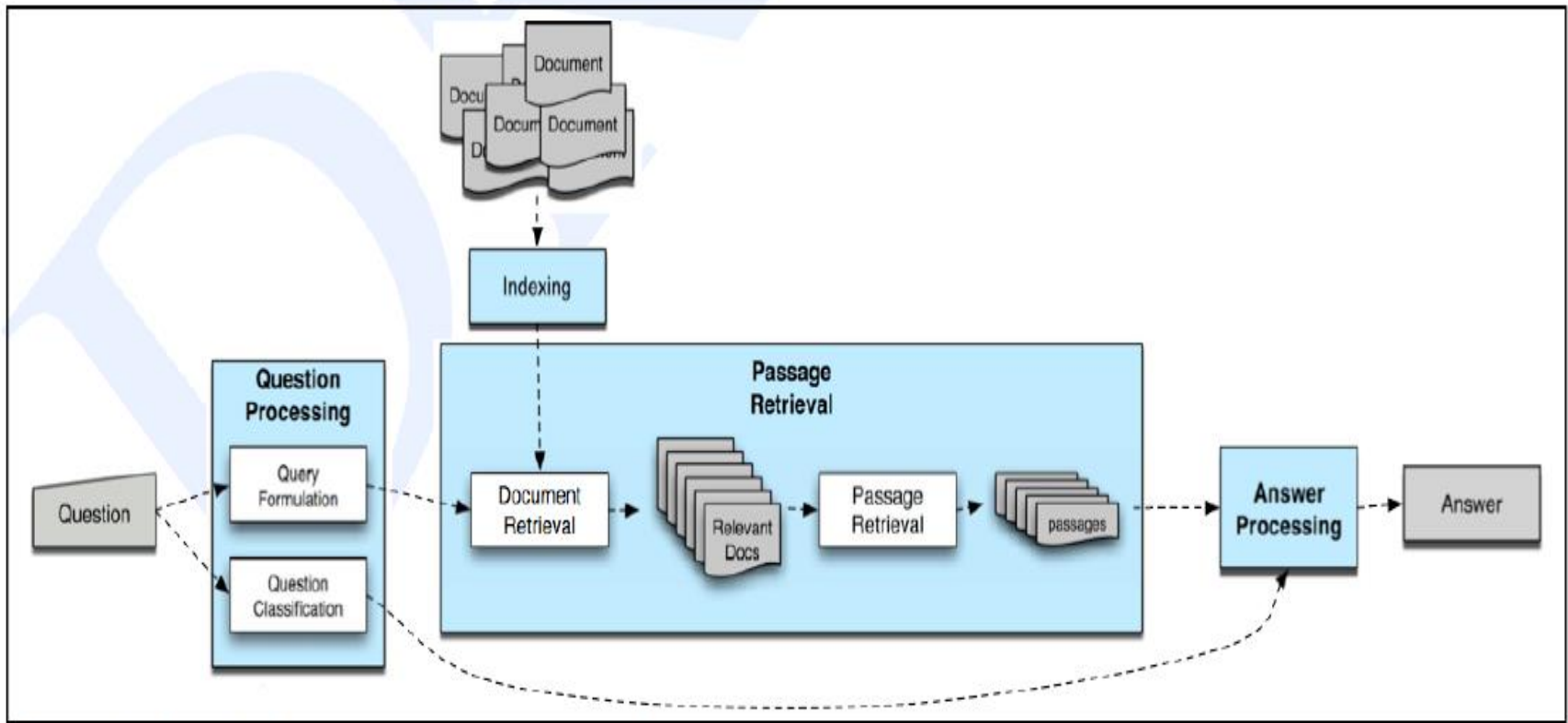
# IBM Watson



**IBM Watson** is a [question-answering](#) computer system capable of answering questions posed in [natural language](#),<sup>[1]</sup> developed in [IBM's](#) DeepQA project by a research team led by [principal investigator David Ferrucci](#).<sup>[2]</sup> Watson was named after IBM's founder and first CEO, industrialist [Thomas J. Watson](#).<sup>[3][4]</sup> The computer system was initially developed to answer questions on the [quiz show Jeopardy!](#)<sup>[5]</sup> and, in 2011, the Watson computer system competed on [Jeopardy!](#) against champions [Brad Rutter](#) and [Ken Jennings](#),<sup>[3][6]</sup> winning the first place prize of \$1 million.<sup>[7]</sup>

[https://www.youtube.com/watch?v=WFR3IOm\\_xhE](https://www.youtube.com/watch?v=WFR3IOm_xhE)

# Question Answering Architecture



**Figure 23.8** The 3 stages of a generic question answering system: question processing, passage retrieval, and answer processing..

# Machine Translation

- The automatic translation of texts between languages is one of the oldest non-numerical applications in Computer Science.
- In the past 15 years or so, MT has gone from a niche academic curiosity to a robust commercial industry.

## 巨大な銃規制集会 が米国を席卷

学生が主催する「私たちの生活のための行進」イベントでは、全国的に数十万人の抗議者が集まります。

🕒 4時間 | 米国とカナダ



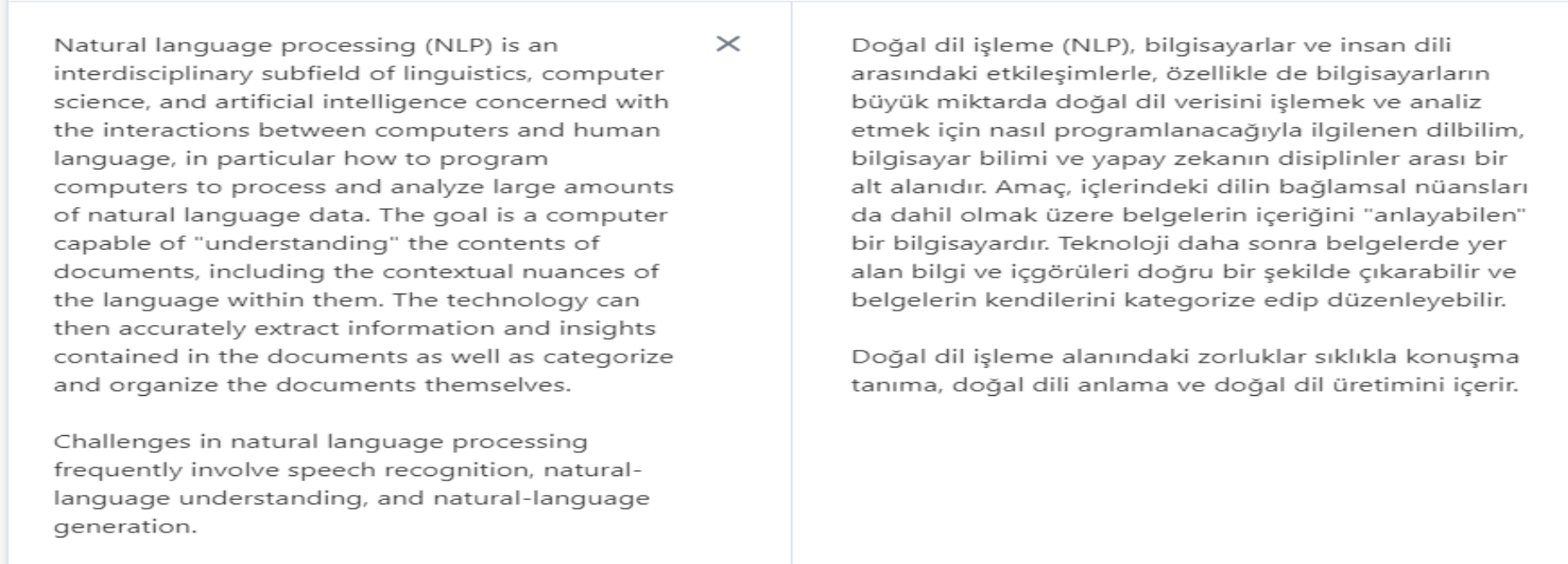
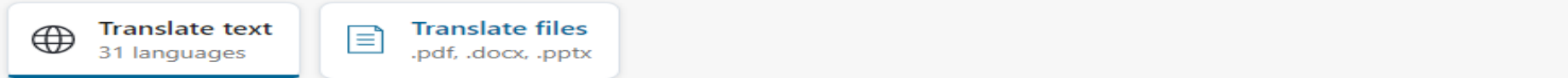
## Huge gun-control rallies sweep US

Student-led March For Our Lives events nationwide draw hundreds of thousands of protesters.

🕒 4h | US & Canada

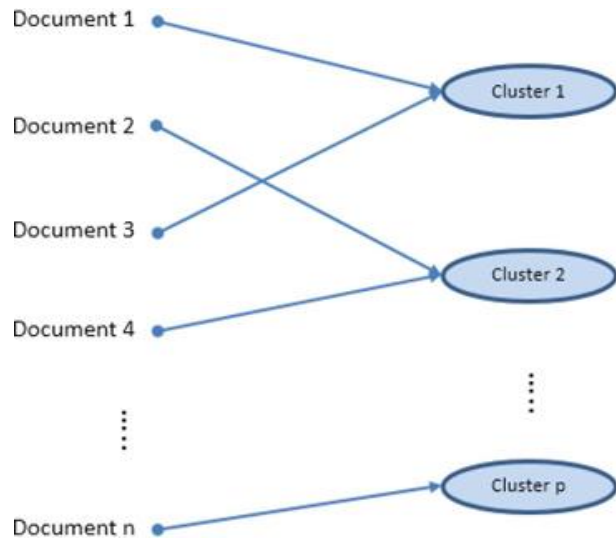


# Trained on text data, neural machine translation is quite good!



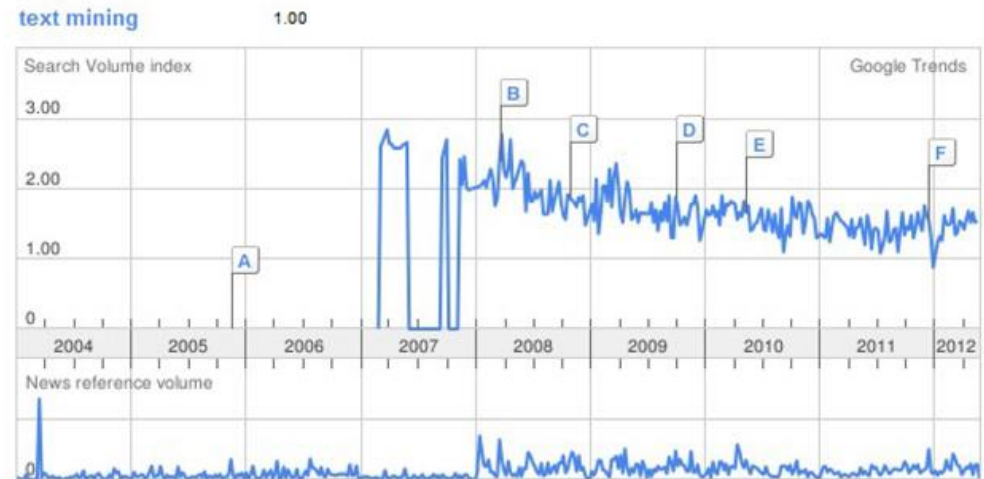
# Text Mining Applications – Unsupervised

- Text clustering



Cluster No.	Comment	Key Words
1	1, 3, 4	doctor, staff, friendly, helpful
2	5, 6, 8	treatment, results, time, schedule
3	2, 7	service, clinic, fast

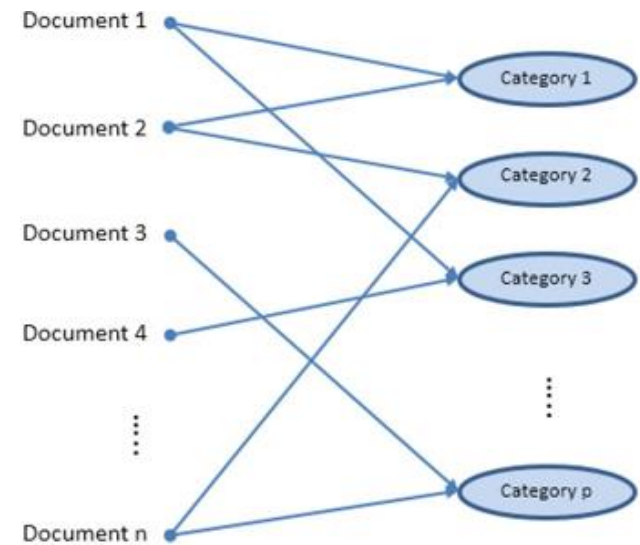
- Trend analysis



Trend for the Term “text mining” from Google Trends

# Text Mining Applications – Supervised

- Many typical **predictive modeling** or classification applications can be enhanced by incorporating textual data in addition to traditional input variables.
  - churning propensity models that include customer center notes, website forms, e-mails, and Twitter messages
  - hospital admission prediction models incorporating medical records notes as a new source of information
  - insurance fraud modeling using adjustor notes
  - sentiment categorization (next page)
  - stylometry or forensic applications that identify the author of a particular writing sample



# Sentiment Analysis

- The field of sentiment analysis deals with categorization (or classification) of opinions expressed in textual documents

The TV is wonderful. Great size, great picture, easy interface. It makes a cute little song when you boot it up and when you shut it off. I just want to point out that the 43" does not in fact play videos from the USB. This is really annoying because that was one of the major perks I wanted from a new TV. Looking at the product description now, I realize that the feature list applies to the X758 series as a whole, and that each model's capabilities are listed below. Kind of a dumb oversight on my part, but it's equally stupid to put a description that does not apply on the listing for a very specific model.

Green color represents positive tone, red color represents negative tone, and product features and model names are highlighted in blue and brown, respectively.

# News Understanding and Recommendation (Li et al., ACL2020Demo)

Home Back Query: Target: Януковича (Yanukovych) Event Search Number of Events: 2

Automated Summary: Source Document Translation from Ukrainian/Russian Show Visual Knowledge Elements Hide Visual Knowledge Elements

Event Summary

Visual Entity Linking

Visual Entity Extraction

Source Doc & Text Extraction Result

Source Doc Translation

Recommended Events

Event Arguments

Date	Location	Attackers	Target	Instrument	Type of Attack
201402	Unknown	Unknown	Януковича (Yanukovych)	Unknown	Conflict.Attack

Event Type

Date	Location	Attackers	Target	Instrument	Type of Attack
201402	Unknown	group	Unknown	Unknown	Conflict.Attack

(Li et al., ACL2020 Best Demo Paper Award)

GitHub: <https://github.com/GAIA-IE/gaia>

DockerHub: <https://hub.docker.com/orgs/blendernlp/repositories>

Demo: [http://159.89.180.81/demo/video\\_recommendation/index\\_attack\\_dark.html](http://159.89.180.81/demo/video_recommendation/index_attack_dark.html)

# Event-centric Question Answering

## Identifying event orders and predicting future events

Heavy **snow** is **causing disruption** to **transport** across the UK, with heavy **rainfall bringing flooding** to the south-west of England. Rescuers **searching** for a woman **trapped** in a **landslide** at her home in Looe, Cornwall, **said** they had **found** a body.

**Q1: What events have already finished?**

A: **searching** trapped **landslide** **said** found

**Q2: What events have begun but has not finished?**

A: **snow** **causing** **disruption** **rainfall** **bringing** **flooding**

**Q3: What will happen in the future?**

A: No answers.

warm-up

**Q4: What happened before a woman was trapped?**

A: **landslide**

**Q5: What had started before a woman was trapped?**

A: **snow** **rainfall** **landslide**

**Q6: What happened while a woman was trapped?**

A: **searching**

**Q7: What happened after a woman was trapped?**

A: **searching** **said** found

User-provided

Ning, et al. TORQUE: A Reading Comprehension Dataset of Temporal Ordering Questions. EMNLP, 2020

- 3.2k news snippets with 21k human-generated questions querying temporal relationships

**Q: Who will drop Japan as a trading partner in August 2019?**

earlier than  
timestamp  
(2019-08-01)

(1/1/19) Apart from the fact of being one another's closest neighbors, the people of South Korea and Japan have a remarkable amount in common. Economically, they are among one another's biggest trading partners. And yet, time and again, relations between *Seoul and Tokyo* are marked, not by mutual support and co-operation but by *anger, reproach and exasperation*.

A) South Korea [0.41] B) Syria [0.28]

C) South Africa [0.15] D) Portugal [0.16]

**Q: Will primary schools in Europe admit non-vaccinated children around September 2019?**

earlier than  
timestamp  
(2019-09-01)

(3/8/18) Public officials and health experts had given several warnings: *Do not allow a student in school if they had not been vaccinated against measles*.

(6/27/19) Fines for parents refusing measles jab. Parents will be fined up to € 2,500 if they don't vaccinate their children against measles under draft legislation in Germany which also *threatens* exclusion from crèches, nurseries and schools.

Yes [0.38] / No [0.62]

[ForecastQA: A Question Answering Challenge for Event Forecasting](#)



# NLP in Industry

- Search (written and spoken)
- Online advertisement
- Automated/assisted translation



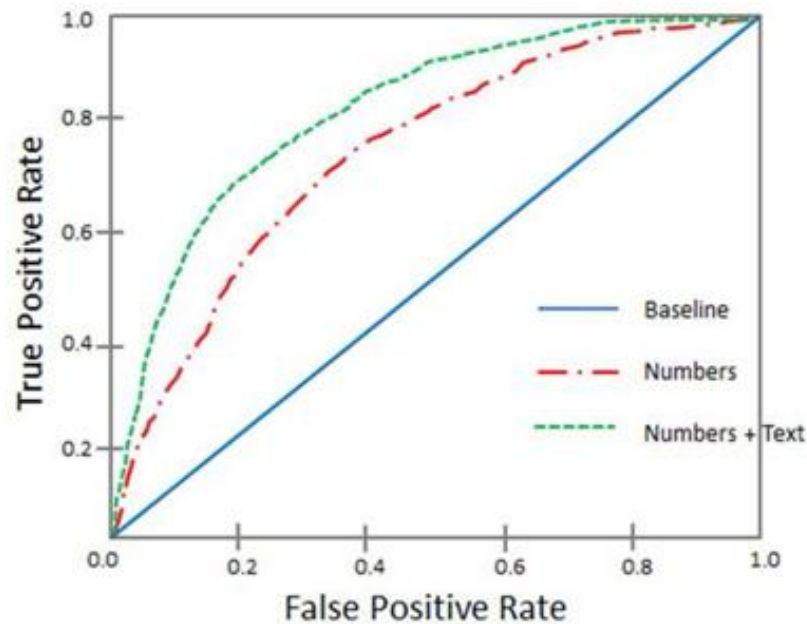
- Sentiment analysis for marketing or finance/trading
- Speech recognition
- Automating customer support



3/18/11 at 4:00 PM | 17 Comments  
**Mentions of the Name 'Anne Hathaway' May Drive Berkshire Hathaway Stock**  
By Patrick Huguenin  
f t e x

# Structured + Text Data in Predictive Models

- Use of both types of data in building predictive models.



ROC Chart of Models With and Without Textual Comments

# High-level NLP tasks

- Information Extraction
  - search, event detection, textual entailment
- Writing Assistance
  - spell checking, grammar checking, auto-completion
- Text Classification
  - spam, sentiment, author, plagiarism
- Natural language understanding
  - metaphor analysis, argumentation mining, question-answering
- Natural language generation
  - summarization, tutoring systems, chat bots
- Multilinguality
  - machine translation, cross-lingual information retrieval

# Some Low-level NLP Tasks

- NLP applications require several NLP analyses:
  - Word tokenization
  - Sentence boundary detection
  - Part-of-speech (POS) tagging
    - to identify the part-of-speech (e.g. noun, verb) of each word
  - Named Entity (NE) recognition
    - to identify proper nouns (e.g. names of person, location, organization; domain terminologies)
  - Parsing
    - to identify the syntactic structure of a sentence
  - Semantic analysis
    - to derive the meaning of a sentence

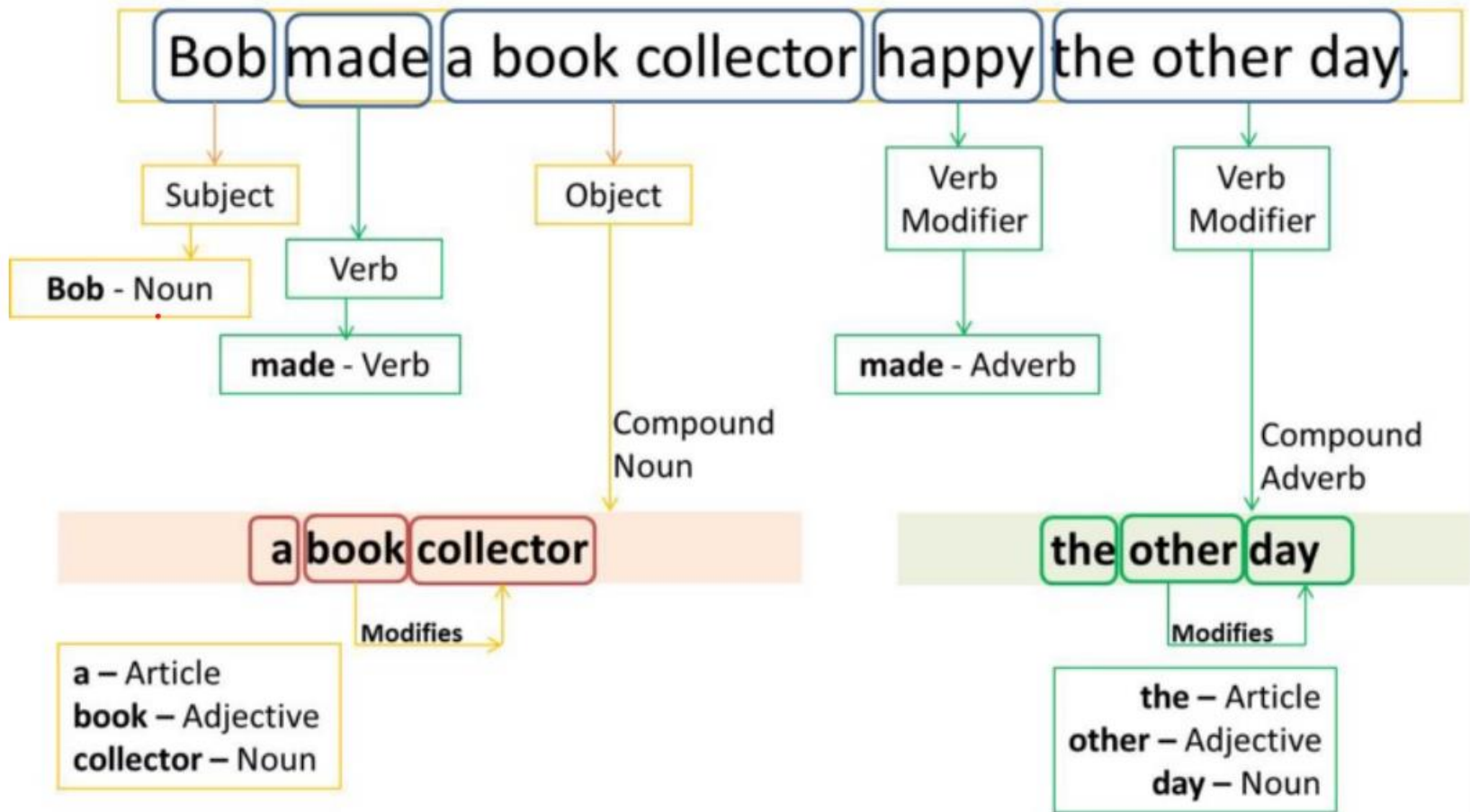
# 1. Part-Of-Speech (POS) Tagging

- POS tagging is a process of assigning a POS or lexical class marker to each word in a sentence (and all sentences in a corpus).

Input:           the lead paint is unsafe

Output:         the/Det lead/N paint/N is/V unsafe/Adj

# Example



# NLTK POS Tagger

## POS tag list:

- CC coordinating conjunction
- CD cardinal digit
- DT determiner
- EX existential there (like: "there is" ... think of it like "there exists")
- FW foreign word
- IN preposition/subordinating conjunction
- JJ adjective 'big'
- JJR adjective, comparative 'bigger'
- JJS adjective, superlative 'biggest'
- LS list marker 1)
- MD modal could, will
- NN noun, singular 'desk'
- NNS noun plural 'desks'
- NNP proper noun, singular 'Harrison'
- NNPS proper noun, plural 'Americans'
- PDT predeterminer 'all the kids'
- POS possessive ending parent's
- PRP personal pronoun I, he, she
- PRP\$ possessive pronoun my, his, hers
- RB adverb very, silently,
- RBR adverb, comparative better
- RBS adverb, superlative best
- RP particle give up
- TO to go 'to' the store.
- UH interjection errrrrrrm
- VB verb, base form take
- VBD verb, past tense took
- VBG verb, gerund/present participle taking
- VBN verb, past participle taken
- VBP verb, sing. present, non-3d take
- VBZ verb, 3rd person sing. present takes
- WDT wh-determiner which
- WP wh-pronoun who, what
- WP\$ possessive wh-pronoun whose
- WRB wh-abverb where, when

```
>>> text = word_tokenize("They refuse to permit us to obtain the refuse permit")
>>> nltk.pos_tag(text) [('They', 'PRP'), ('refuse', 'VBP'), ('to', 'TO'), ('permit', 'VB'), ('us', 'PRP'), ('to', 'TO'), ('obtain', 'VB'), ('the', 'DT'), ('refuse', 'NN'), ('permit', 'NN')]
```

## 2. Named Entity Recognition (NER)

- NER is to process a text and identify named entities in a sentence
  - e.g. “U.N. official Ekeus heads for Baghdad.”

[ORG U.N. ] official [PER Ekeus ] heads for [LOC Baghdad ] .

# NER

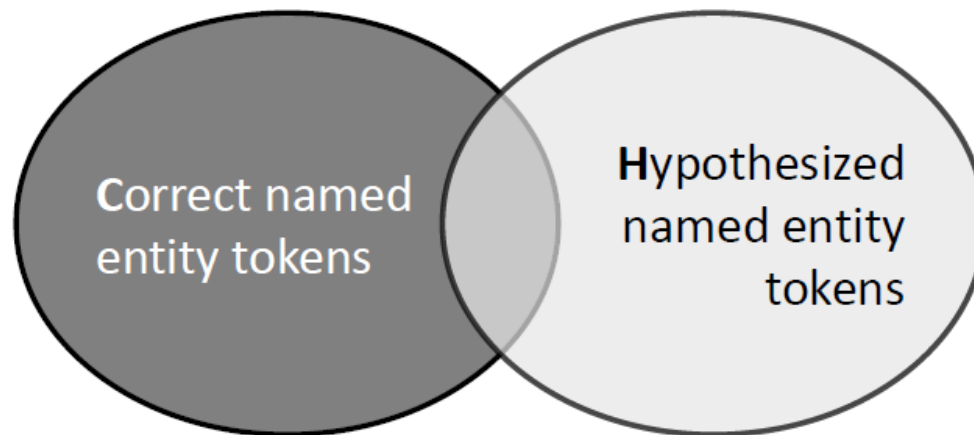
**Elizabeth Warren**, the liberal firebrand who emerged as a top Democratic contender for the **White House** on the strength of an anti-corruption platform backed by a dizzying array of policy proposals, ended her campaign on Thursday. A former bankruptcy law professor who forged a national reputation as a scourge of Wall Street even before entering politics, **Warren** had banked on a strong showing on Super Tuesday after a string of disappointing finishes in the early states. But she trailed far behind front-runners **Bernie Sanders** and **Joe Biden**, placing third in her home state of Massachusetts, which she continues to represent in the U.S. Senate.

- Label certain kinds of proper nouns:
  - Personal names
  - Organizations
  - Geopolitical entities
  - Locations
  - etc.

# Some Named Entity Types

- Different annotation schemes for NER use different types
- Common types include
  - PER—person
  - ORG—Organization
  - LOC—Location
  - GPE—Geopolitical Entity
  - FAC—Facility
  - NAT—Natural phenomenon
- These are only tagged when they are proper names

# Evaluating an NER System

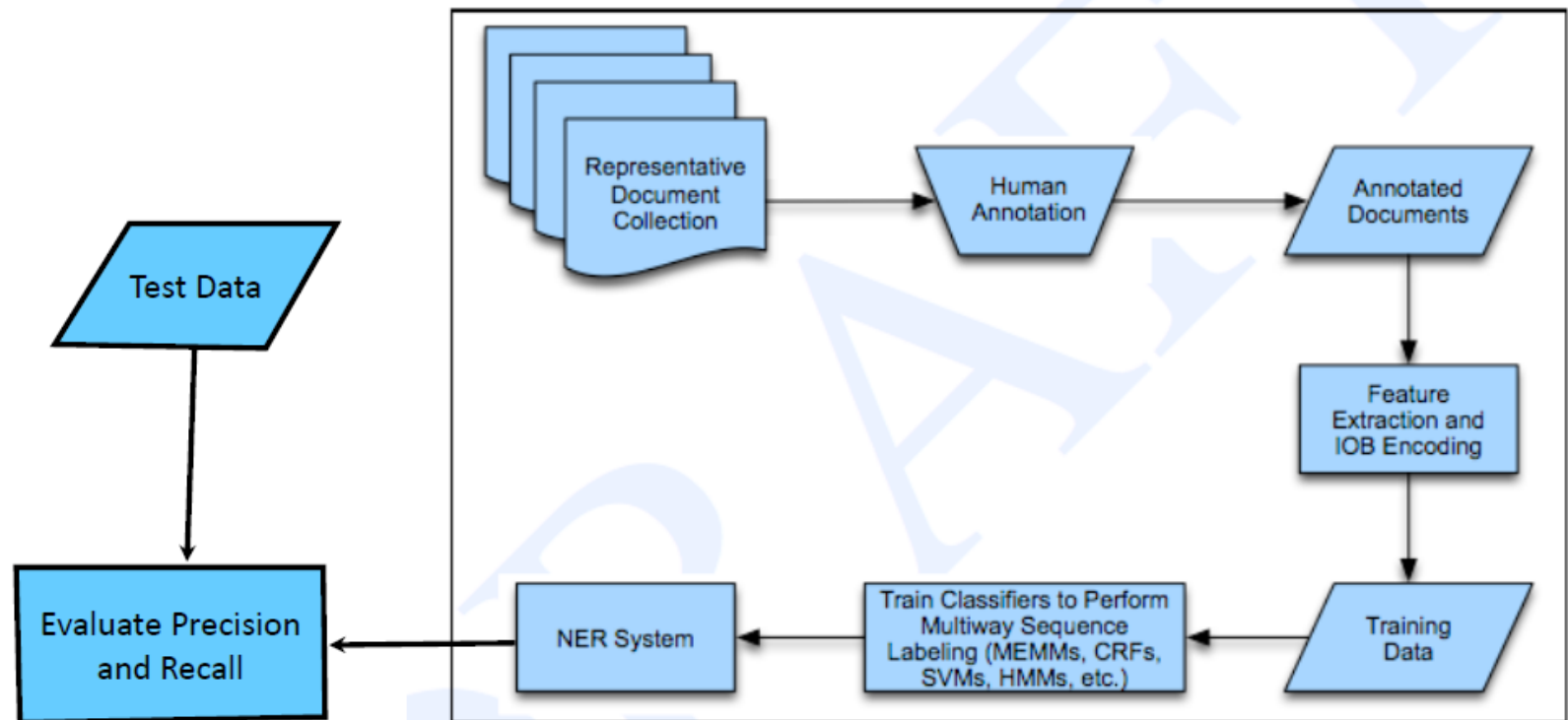


$$\text{recall} = \frac{|C \cap H|}{|C|}$$

$$\text{precision} = \frac{|C \cap H|}{|H|}$$

# NER System Building Process

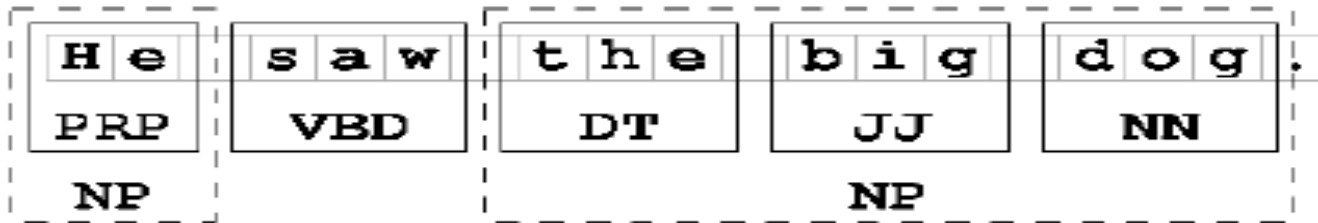
## NER System Building Process



**Figure 22.10** Basic steps in the statistical sequence labeling approach to creating a named entity recognition system.

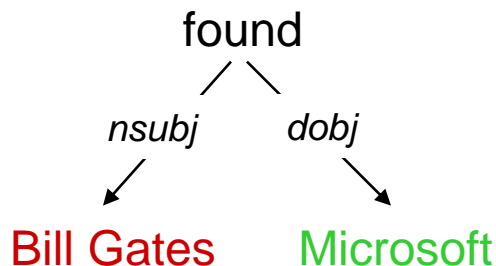
# 3. Shallow Parsing

- **Shallow parsing** (also **chunking** or **light parsing**) is an analysis of a sentence which first identifies constituent parts of sentences (nouns, verbs, adjectives, etc.) and then links them to higher order units that have discrete grammatical meanings (noun groups or phrases, verb groups, etc.).
- Shallow (or Partial) parsing identifies the (base) syntactic phases in a sentence.



- After NEs are identified, **dependency parsing** is often applied to extract the syntactic/dependency relations between the NEs.

[<sub>PER</sub> Bill Gates] founded [<sub>ORG</sub> Microsoft].

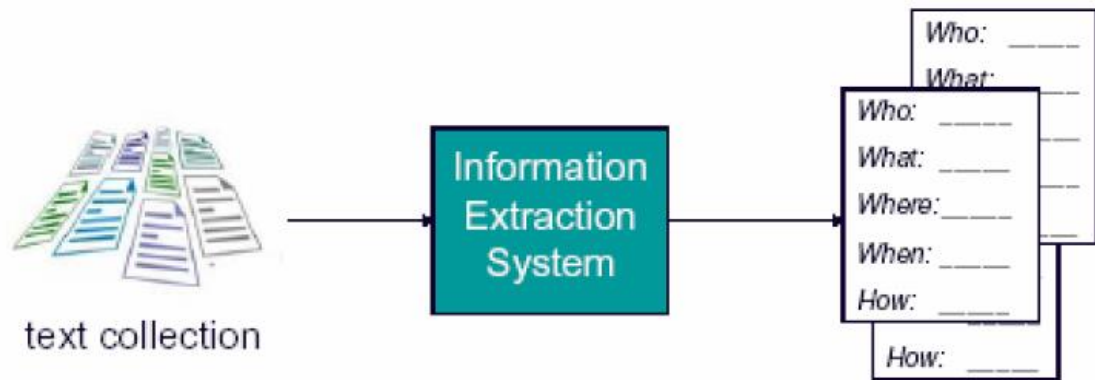


## Dependency Relations

nsubj(Bill Gates, found)  
dobj(found, Microsoft)

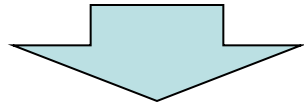
# 4. Information Extraction (IE)

- Identify specific pieces of information (data) in an unstructured or semi-structured text
- Transform unstructured information in a corpus of texts or web pages into a structured database (or templates)
- Applied to various types of text, e.g.
  - Newspaper articles
  - Scientific articles
  - Web pages
  - etc.

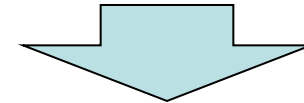


Bridgestone Sports Co. said Friday it had set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be supplied to Japan.

The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with production of 20,000 iron and “metal wood” clubs a month.



*template filling*



### **TIE-UP-1**

**Relationship:** TIE-UP

**Entities:** “Bridgestone Sport Co.”

“a local concern”

“a Japanese trading house”

**Joint Venture Company:**

“Bridgestone Sports Taiwan Co.”

**Activity:** **ACTIVITY-1**

**Amount:** NT\$2000000000

### **ACTIVITY-1**

**Activity:** PRODUCTION

**Company:**

“Bridgestone Sports Taiwan Co.”

**Product:**

“iron and ‘metal wood’ clubs”

**Start Date:**

DURING: January 1990

## 5. Reference resolution

**Elizabeth Warren**, the liberal firebrand who emerged as a top Democratic contender for the **White House** on the strength of an anti-corruption platform backed by a dizzying array of policy proposals, ended her campaign on Thursday. A former bankruptcy law professor who forged a national reputation as a scourge of Wall Street even before entering politics, **Warren** had banked on a strong showing on Super Tuesday after a string of disappointing finishes in the early states. But she trailed far behind front-runners **Bernie Sanders** and **Joe Biden**, placing third in her home state of Massachusetts, which she continues to represent in the U.S. Senate.



**Reference**, in NLP, is a linguistic process where one word in a sentence or discourse may refer to another word or entity. The task of resolving such references is known as **Reference Resolution**.



# 7. Relation detection/extraction

WASHINGTON/SELMA, Ala. (Reuters) - Democratic U.S. presidential front-runner Bernie Sanders raised \$46.5 million in February, his campaign said on Sunday, and will launch new television ad buys in nine states with primaries later this month after this week's Super Tuesday contests. Joe Biden's campaign reported raising \$5 million the day of the South Carolina primary. His February haul was \$18 million, spokesman Michael Gwin said. Meanwhile, rival Elizabeth Warren, who struggled to a fifth-place finish in South Carolina, raised more than \$29 million in February, her campaign manager Roger Lau said in a memo to supporters on Sunday.

Candidate	member_of
Bernie Sanders	Democratic Party
Joe Biden	Democratic Party
Elizabeth Warren	Democratic Party

- Relationship detection/extraction is the task of extracting semantic relationships from a text.
- Extracted relationships usually occur between two or more entities of
  - a certain type (e.g. Person, Organisation, Location) and fall into
  - a number of semantic categories (e.g. married to, employed by, lives in).

# Examples of Relations

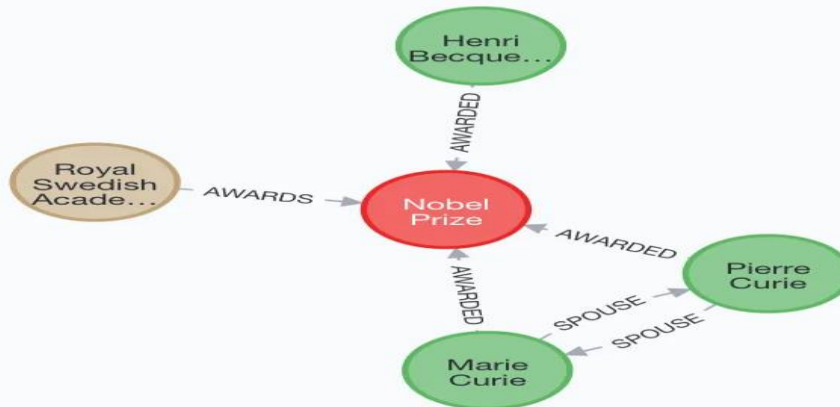
Relations	Examples	Types
Affiliations		
Personal	<i>married to, mother of</i>	PER → PER
Organizational	<i>spokesman for, president of</i>	PER → ORG
Artifactual	<i>owns, invented, produces</i>	(PER   ORG) → ART
Geospatial		
Proximity	<i>near, on outskirts</i>	LOC → LOC
Directional	<i>southeast of</i>	LOC → LOC
Part-Of		
Organizational	<i>a unit of, parent of</i>	ORG → ORG
Political	<i>annexed, acquired</i>	GPE → GPE

**Figure 22.11** Semantic relations with examples and the named entity types they involve.

# Ex: Named Entity Recognition and Relationship Extraction

- Building Named Entity Recognition and Relationship Extraction Components with HuggingFace Transformers (<https://odsc.medium.com/building-named-entity-recognition-and-relationship-extraction-components-with-huggingface-77d233e27e65>)

In December 1903 DATE the Royal Swedish Academy of Sciences ORG awarded Marie PERSON and Pierre Curie PERSON , along with Henri Becquerel PERSON , the Nobel Prize in Physics WORK\_OF\_ART .



Relations detected by a Relation Extraction (RE) component

## But NLP very is hard..

- Understanding natural languages is hard ... because of inherent *ambiguity*
- Engineering NLP systems is also hard ... because of:
  - Huge amount of data resources needed (e.g. grammar, dictionary, documents to extract statistics from)
  - Computational complexity (intractable) of analyzing a sentence

# Ambiguity (1)

“Get the cat with the gloves.”



# Ambiguity (2)

Find at least 5 meanings of this sentence:

*“I made her duck”*

1. I cooked waterfowl for her benefit (to eat)
2. I cooked waterfowl belonging to her
3. I created the (plaster?) duck she owns
4. I caused her to quickly lower her head or body
5. I waved my magic wand and turned her into undifferentiated waterfowl

# Ambiguity (3)


## Some ambiguous headlines

- Juvenile Court to Try Shooting Defendant
- Teacher Strikes Idle Kids
- Kids Make Nutritious Snacks
- Bush Wins on Budget, but More Lies Ahead
- Hospitals are Sued by 7 Foot Doctors

# Ambiguity is Pervasive

- **Phonetics**

- I mate or duck
- I'm eight or duck
- Eye maid; her duck
- Aye mate, her duck
- I maid her duck
- I'm aid her duck
- I mate her duck
- I'm ate her duck
- I'm ate or duck
- I mate or duck



Sound like  
*"I made her duck"*

- **Lexical category** (part-of-speech)
  - “duck” as a noun or a verb
- **Lexical Semantics** (word meaning)
  - “duck” as an animal or a plaster duck statue
- Compound nouns
  - e.g. “dog food”, “Intelligent design scores ...”
- **Syntactic ambiguity**

“I saw a man on the hill with a telescope”

- [But semantics can sometimes help disambiguate]

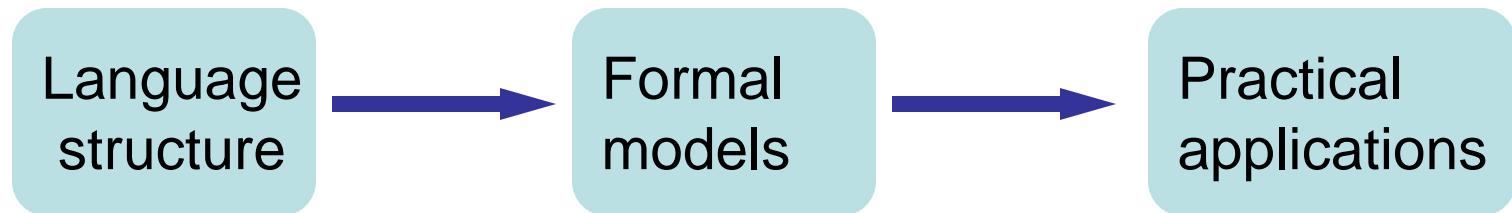
“I saw a man on the hill with a hat”

# The Bottom Line

- Complete NL Understanding (thus general intelligence) is impossible.
- But we can make incremental progress.
- Also we have made successes in **limited domains**.

# The Big Picture Approach

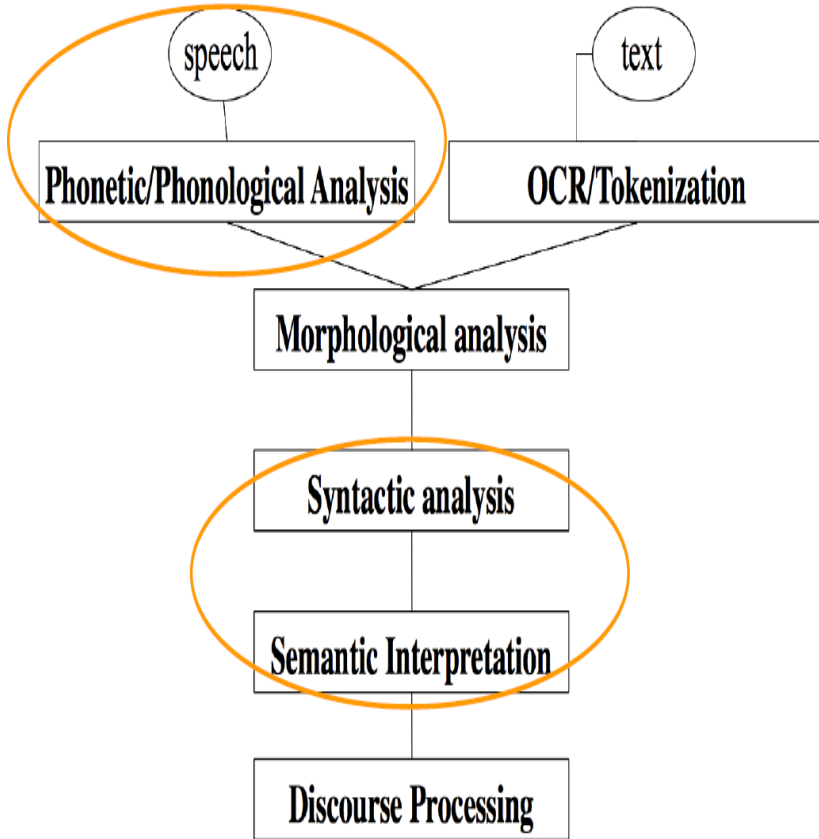
All of these applications operate by **exploiting underlying regularities** in human languages. Sometimes in complex ways, sometimes in pretty trivial ways.



# Topics: Linguistics

- Word-level processing
- Syntactic processing
- Lexical and compositional semantics
- Discourse structure

# NLP Levels



NLP requires various kinds of knowledge of language:


- Phonetics and Phonology —knowledge about linguistic sounds
- Morphology —knowledge of the meaningful components of words
- Syntax —knowledge of the structural relationships between words
- Semantics —knowledge of meaning
- Pragmatics — knowledge of the relationship of meaning to the goals and intentions of the speaker
- Discourse —knowledge about linguistic units larger than a single utterance

# Different Levels of Linguistic Analysis

- Phonology
  - Speech audio signal to phonemes
- Morphology
  - Inflection (e.g. “I”, “my”, “me”; “eat”, “eats”, “ate”, “eaten”)
  - Derivation (e.g. “teach”, “teacher”, “nominate”, “nominee”)
- Syntax
  - Part-of-speech (noun, verb, adjective, preposition, etc.)
  - Phrase structure (e.g. noun phrase, verb phrase)
- Semantics
  - Meaning of a word (e.g. “book” as a bound volume or an accounting ledger) or a sentence
- Discourse
  - Meaning and inter-relation between sentences

# Topics: Techniques

- Finite-state methods
- Context-free methods
- Probabilistic models
- Neural network models

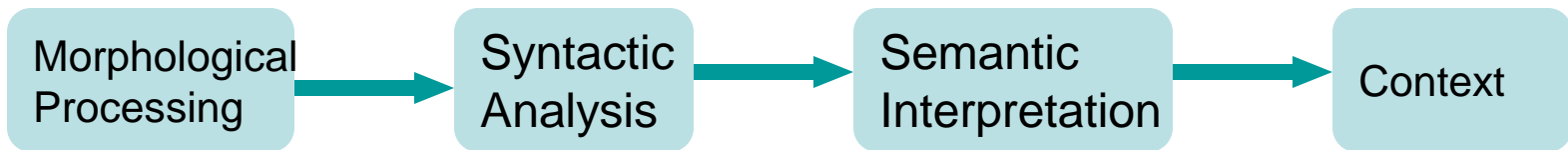


Supervised machine  
learning methods

# Process Pipeline

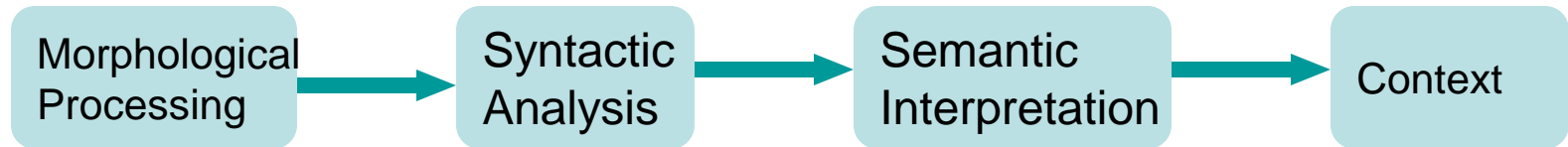
- Phonology
- Morphology
- Syntax
- Semantics
- Pragmatics
- Discourse

Each kind of knowledge has associated with it an encapsulated set of processes that make use of it. Interfaces are defined that allow the various levels to communicate. This often leads to a **pipeline architecture**.

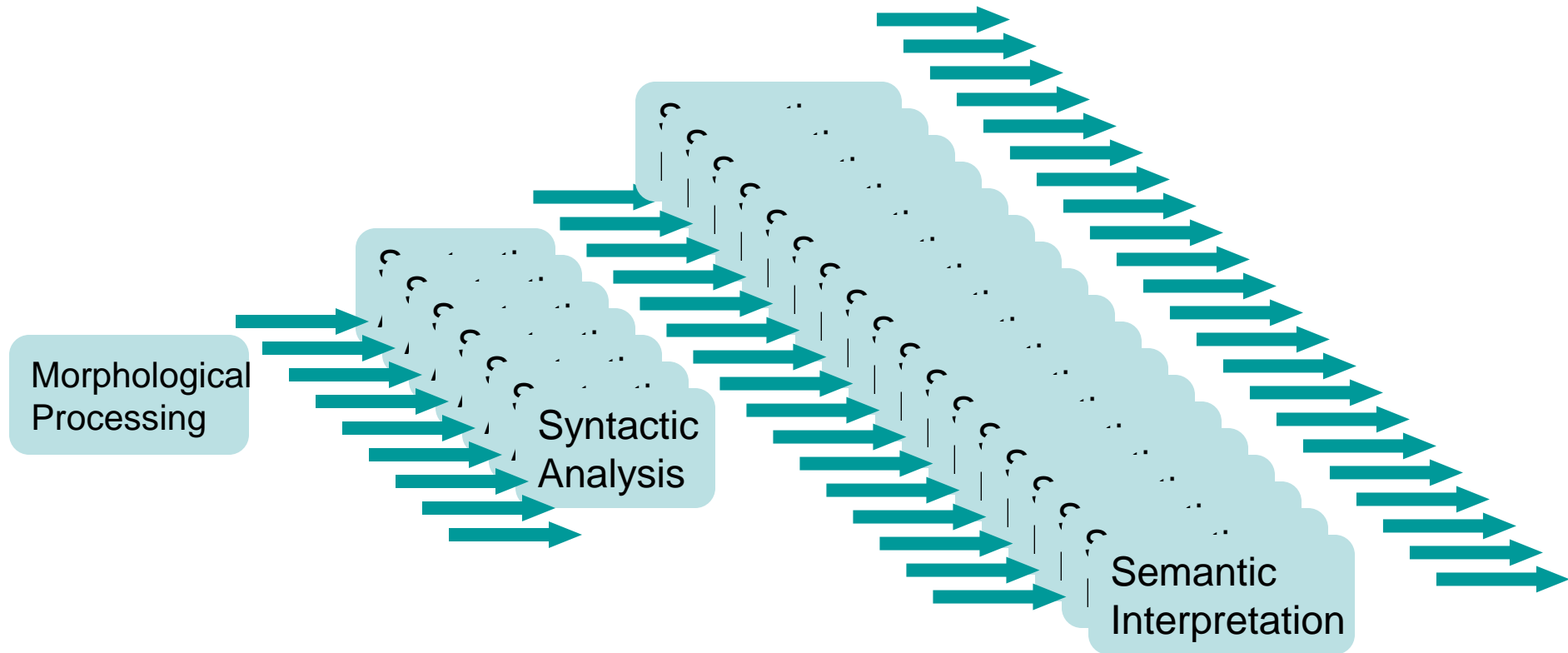


# But Problem..

- Remember our pipeline...



# It really looks like this



# Dealing with Ambiguity

Four possible approaches:

1. **Formal approaches** -- Tightly coupled interaction among processing levels; knowledge from other levels can help decide among choices at ambiguous levels.
2. Pipeline processing that ignores ambiguity as it occurs and hopes that other levels can eliminate incorrect structures.
3. **Probabilistic approaches** based on making the most likely choices
4. Don't do anything, maybe it won't matter

# Models and Algorithms

- By **models** we mean the formalisms that are used to capture the various kinds of linguistic **knowledge** we need.
- **Algorithms** are then used to manipulate the knowledge representations needed to tackle the task at hand.

# Various Models

- Finite state machines
- Rule-based and logic-based approaches
- Probabilistic models
- Neural network models

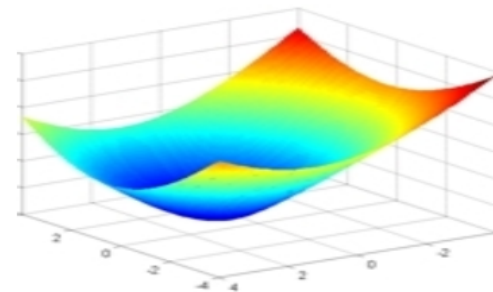
# Various Algorithms

- In particular..
  - State-space search
    - To manage the problem of making choices during processing when we lack the information needed to make the right choice
  - Dynamic programming
    - To avoid having to redo work during the course of a state-space search
      - CKY, Earley, Minimum Edit Distance, Viterbi, Baum-Welch
  - Classifiers
    - Machine learning based classifiers that are trained to make decisions based on features extracted from the local context

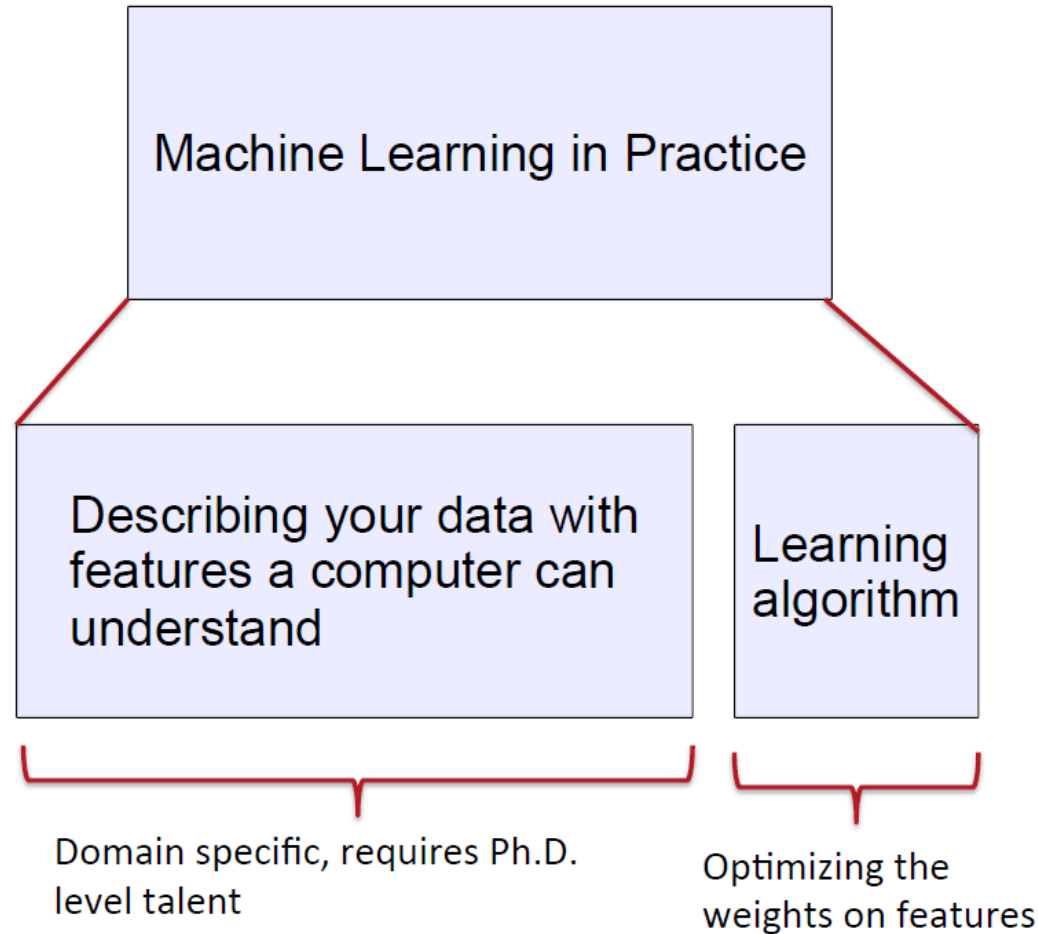
# What's Deep Learning (DL)?

- Deep learning is a subfield of machine learning
- Most machine learning methods work well because of human-designed representations and input features
  - For example: features for finding named entities like locations or organization names (Finkel, 2010):
- Machine learning becomes just optimizing weights to best make a final prediction

Feature	NER
Current Word	✓
Previous Word	✓
Next Word	✓
Current Word Character n-gram	all
Current POS Tag	✓
Surrounding POS Tag Sequence	✓
Current Word Shape	✓
Surrounding Word Shape Sequence	✓
Presence of Word in Left Window	size 4
Presence of Word in Right Window	size 4

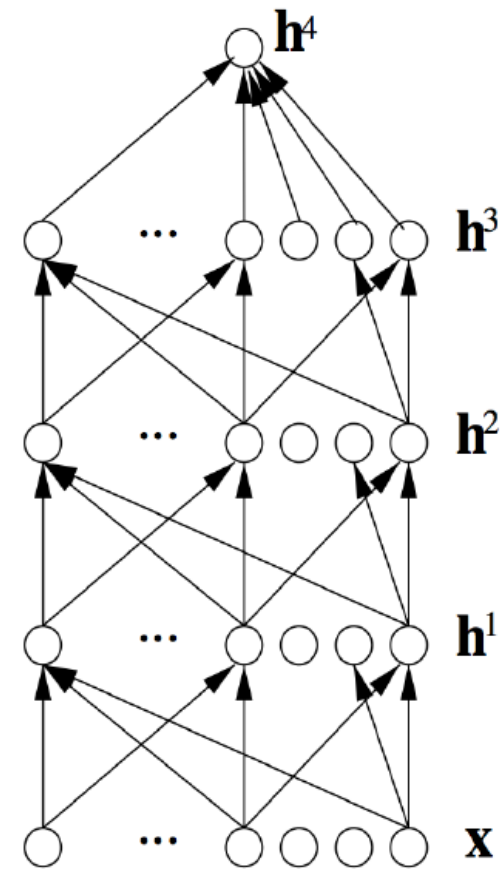


# Machine Learning vs Deep Learning



# What's Deep Learning (DL)?

- Representation learning attempts to automatically learn good features or representations
- Deep learning algorithms attempt to learn (multiple levels of) representation and an output
- From “raw” inputs  $\mathbf{x}$  (e.g. words)



# Reasons for Exploring Deep Learning

- Manually designed features are often over-specified, incomplete and take a long time to design and validate
- **Learned Features** are easy to adapt, fast to learn
- Deep learning provides a very flexible, (almost?) universal, learnable framework for **representing** world, visual and linguistic information.
- Deep learning can learn **unsupervised** (from raw text) and **supervised** (with specific labels like positive/negative)

# Reasons for Exploring Deep Learning

- In 2006 **deep** learning techniques started outperforming other machine learning techniques. Why now?
  - DL techniques benefit more from a lot of data
  - Faster machines and multicore CPU/GPU help DL
  - New models, algorithms, ideas
- **Improved performance** (first in speech and vision, then NLP)

# Deep Learning + NLP = Deep NLP

- Combine ideas and goals of NLP and use representation learning and deep learning methods to solve them
- Several big improvements in recent years across different NLP
  - **levels:** speech, morphology, syntax, semantics
  - **applications:** machine translation, sentiment analysis and question answering