

# Natural Language Processing

**[anoopsarkar.github.io/nlp-class](https://anoopsarkar.github.io/nlp-class)**

# Programming Languages

C, C++, Java, Python, ...

- unambiguous
- fixed
- designed
- learnable?
- known simple semantics

# Natural Languages

French, English, Korean, Chinese, Tagalog, ...

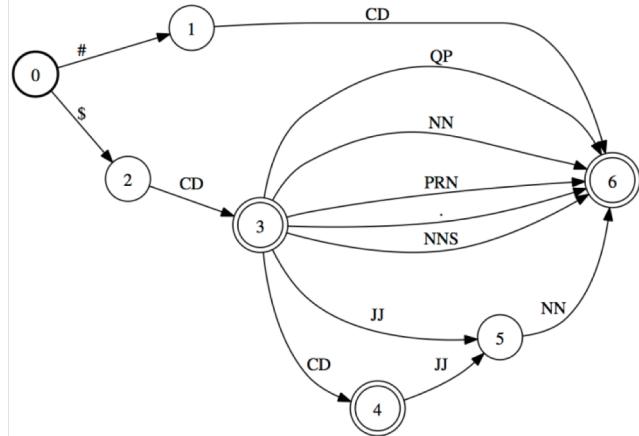
- ambiguous
- evolving
- transmitted
- learnable
- complex semantics

Why is NLP computationally hard?

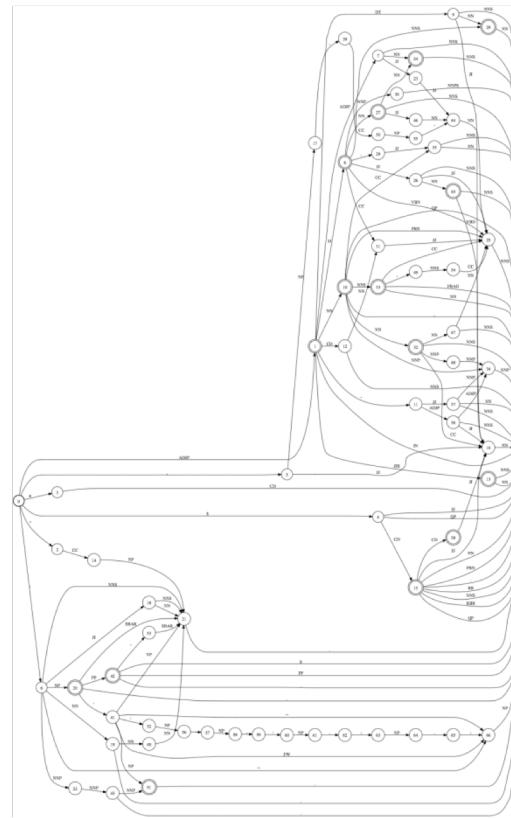
\$ 20 happy meal

\$ CD JJ NN a noun phrase type

10 noun phrase types

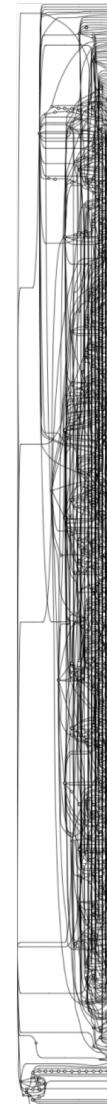


100 noun phrase types



Language  
is complex

1000 noun phrase types



6K noun phrase types

# Language is ambiguous

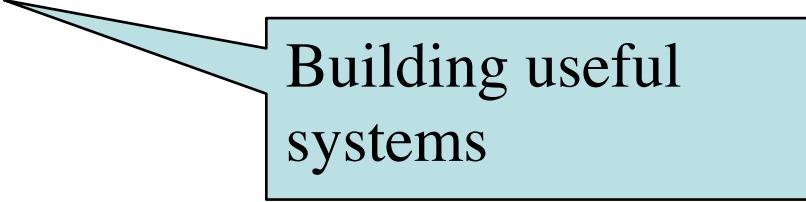
- Lung cancer in women mushrooms
  - Mushrooms is noun or a verb?
- Ban on nude dancing on governor's desk
  - Similar to “if-then-else” ambiguity
- Island Monks Fly in Satellite to Watch Pope Funeral
  - “**fly in**” vs. “**fly [OBJ in Satellite]**” hidden segmentation
- British Left Waffles on Falkland Islands
  - Is it **British/Noun Left/Verb** or **British Left/NP Waffles/Verb**?

# Language is Parsed

- Google's Computer Might Betters Translation Tool
  - New York Times March 8, 2010
- Number of Lothian patients made ill by drinking rockets
  - Edinburgh Evening News, March 4, 2010
- Violinist linked to JAL crash blossoms
  - <http://languagelog.ldc.upenn.edu/nll/?p=1693>

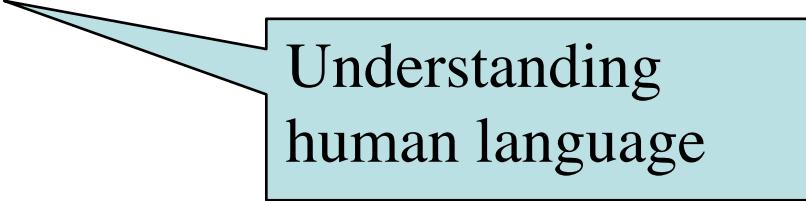
# What is the difference?

NLP = Natural Language Processing



Building useful  
systems

CL = Computational Linguistics



Understanding  
human language

# Different levels of language

- Phonetics: acoustic and perceptual elements
- Phonology: basic sounds (phonemes) and rules for combination
  - e.g. vowel harmony. Anupu is pronunciation of Anoop in Classic Period Mayan
- Morphology: how morphemes combine to form words, relationship of phonemes to meaning
  - e.g. delight-ed vs. de-light-ed



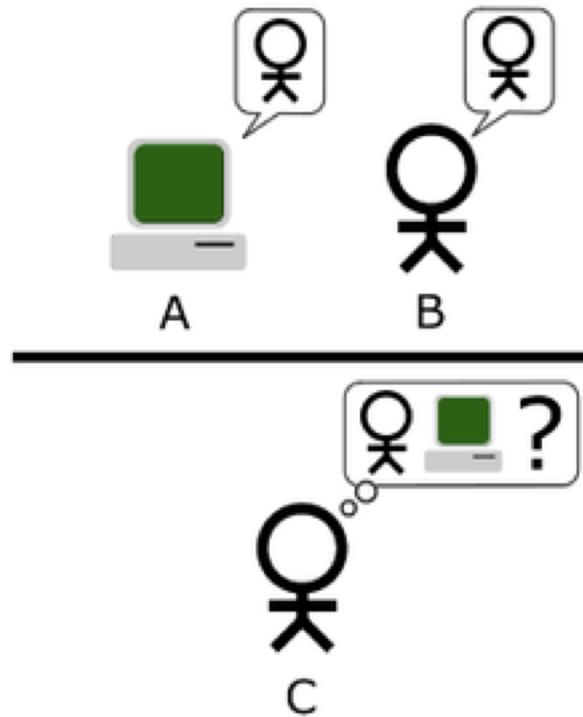
# Different levels of language

- Syntax: sentence formation
  - e.g. The clown who the musician hits watches the ballerina
- Semantics: meaning (from syntax to logical formulas)
  - e.g. Everyone is not here => what does this mean? Nobody / Not everyone is here.
- Pragmatics: meaning that is not part of compositional meaning,
  - e.g. This professor dresses even worse than Anoop!

Imagine an "**Imitation Game**," in which a man and a woman go into separate rooms and guests try to tell them apart by writing a series of questions and reading the typewritten answers sent back. In this game both the man and the woman aim to convince the guests that they are the other.



Alan Turing



We now ask the question, "**What will happen when a machine takes the part of A in this game?**" Will the interrogator decide wrongly as often when the game is played like this as he does when the game is played between a man and a woman? These questions replace our original, "**Can machines think?**"

# Winograd Schema

The town councillors refused to give the angry demonstrators a permit because they feared violence.

Who feared violence?



Answer 0: the town councillors

Answer 1: the angry demonstrators

---

The town councillors refused to give the angry demonstrators a permit because they advocated violence.

Who advocated violence?

Answer 0: the town councillors



Answer 1: the angry demonstrators

# Some examples of NLP achievements

# IBM Watson plays Jeopardy



# Speech to Speech Translation



KIT Lecture Translator



NICT Speech Translator



Skype Translator

# Some examples of NLP tasks

# Information Extraction

Así lo explicó hoy el presidente del Gobierno español, José María Aznar, en la conferencia de prensa con la que  
concluyó la XIII Cumbre Hispano-francesa, celebrada en Santander, con asistencia del presidente francés,  
Jacques Chirac ; del primer ministro, Lionel Jospin, y trece miembros de ambos gabinetes.

```
graph TD; Gobierno[ORG] --- Gobierno; Aznar[PER] --- Aznar; Cumbre[MISC] --- Cumbre; Santander[LOC] --- Santander; Chirac[PER] --- Chirac; Jospin[PER] --- Jospin; Miembros[PER] --- Miembros;
```

# Information Extraction



Association of N-glycosylation of apolipoprotein B-100 with plasma cholesterol levels in Watanabe heritable hyperlipidemic rabbits.



Terry Pratchett lives in England.

LocalizationID

1022

PSID

10126

- 1) Select "valid" if the passage contains strong evidence of an experimentally determined localization.

[PubMed Entrez](#)

PMID 9811664

[PubMed Central](#)

PMCID 107680

The cytoplasmic membrane proteins ExbB and ExbD support TonB-dependent active transport of iron siderophores and vitamin B12 across the essentially unenergized outer membrane of *Escherichia coli*.

Valid

Invalid

Maybe

Reviewer

Comments

- 2) If the passage is valid then select whether the protein, organism, and location names are also valid. (If you want to defer your decision then select neither valid nor invalid)

Protein:

ExbB

Valid

Invalid

Organism:

Escherichia coli

Valid

Invalid

Location:

cytoplasmic membrane

Valid

Invalid

# SQuASH: SFU QA Summarization System

**Input:** 25 news articles, Complex question    **Output:** 250-word summary

Q. Describe developments in the movement for the independence of Quebec from Canada.

A. Canadian Prime Minister Jean Chretien has ruled out the possibility that the French-speaking province of Quebec could declare independence even if a majority of Quebecers vote for secession in the next referendum. Premier Lucien Bouchard and the separatist Parti Quebecois have retained control of Canada's largest province, possibly setting the stage for a bitter confrontation with the rest of Canada, and possibly another referendum on independence. Canada's Supreme Court managed to rule on the question of Quebec secession without infuriating leaders on either side of a bitter ideological divide. . . .

# Headline Generation

Headline A: US launches air raids in Somalia

Headline B: Somalia says dozens killed in US attack

Headline C: Many dead after US strike in Somalia

Headline D: US Launches New Attacks in Somalia

Headline E: US strikes terrorist targets in Somalia

Cluster of headlines for an event on Google News

# Headline Generation

Headline Candidate	Score
Bush to sign of	-22.614
Bush to sign bill on	-26.652
Bush to sign of the	-26.835
the House of The Internet gambling	-29.946
The bill of the Internet gambling	-29.982
Bush to end of the Internet gambling	-32.576
Bush to sign bill on the Internet gambling	-35.746
Bush to sign bill on the Internet gambling law	-39.710
Bush to end of the Internet gambling on The Senate bill	-46.988
Bush to sign bill on the Internet gambling site of The law	-50.912

Table 5.9: Top Headlines for “Law on Internet Gambling” news story

TEAM	WIN	LOSS	PTS	FG_PCT	RB	AS ...
Heat	11	12	103	49	47	27
Hawks	7	15	95	43	33	20

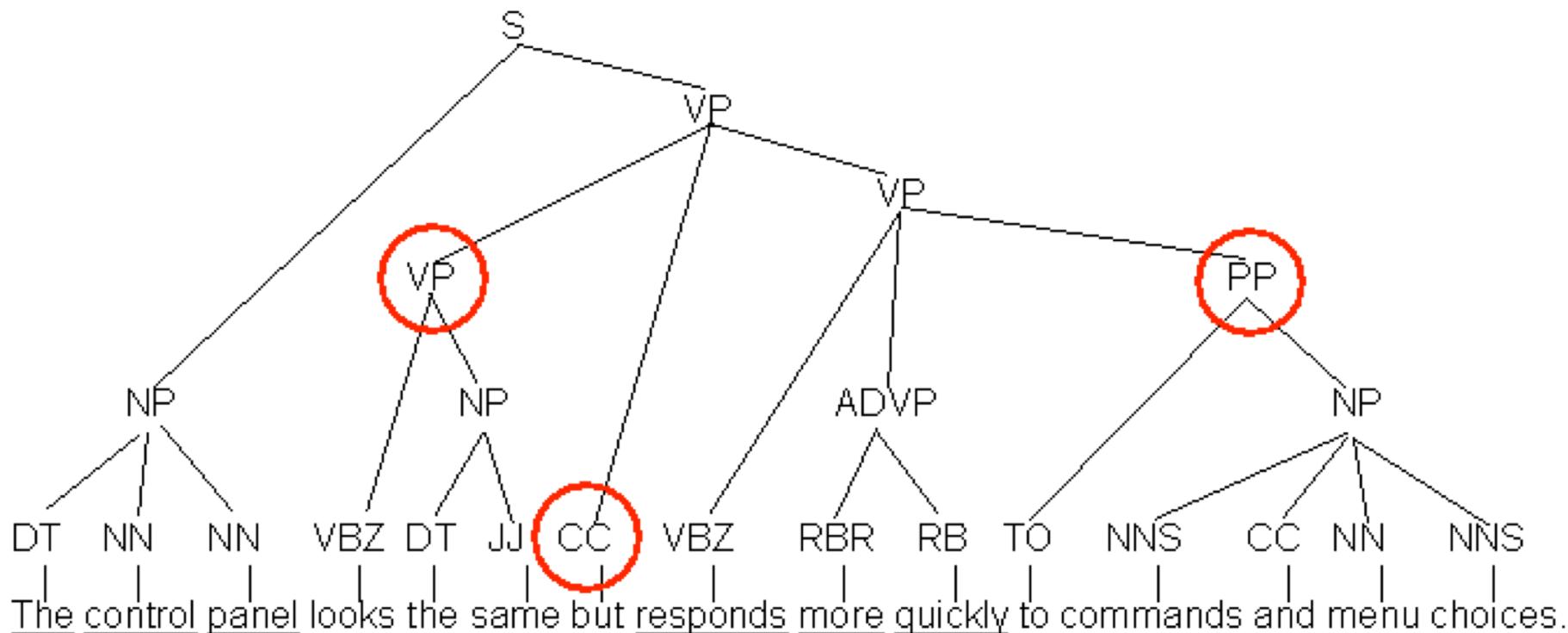
PLAYER	AS	RB	PT	FG	FGA	CITY ...
Tyler Johnson	5	2	27	8	16	Miami
Dwight Howard	11	17	23	9	11	Atlanta
Paul Millsap	2	9	21	8	8	
Goran Dragic	4	2	21	8	8	



# Natural Language Generation

The Atlanta Hawks defeated the Miami Heat, 103 - 95, at Philips Arena on Wednesday. Atlanta was in desperate need of a win and they were able to take care of a shorthanded Miami team here. Defense was key for the Hawks, as they held the Heat to 42 percent shooting and forced them to commit 16 turnovers. Atlanta also dominated in the paint, winning the rebounding battle, 47 - 34, and outscoring them in the paint 58 - 26. The Hawks shot 49 percent from the field and assisted on 27 of their 43 made baskets. This was a near wire-to-wire win for the Hawks, as Miami held just one lead in the first five minutes. Miami ( 7 - 15 ) are as beat-up as anyone right now and it's taking a toll on the heavily used starters. Hassan Whiteside really struggled in this game, as he amassed eight points, 12 rebounds and one blocks on 4 - of - 12 shooting ...

# Sentence Compression



# Paraphrasing

- open borders imply increasing racial fragmentation in *european countries* .
- open borders imply increasing racial fragmentation in *the countries of europe* .
- open borders imply increasing racial fragmentation in *european states* .
- open borders imply increasing racial fragmentation in *europe* .
- open borders imply increasing racial fragmentation in *european nations* .
- open borders imply increasing racial fragmentation in *the european countries* .

Why is paraphrasing useful?

# Sentiment detection

Annotate tweets using labels from [http://en.wikipedia.org/wiki/List\\_of\\_emoticons](http://en.wikipedia.org/wiki/List_of_emoticons)

## 10 Happiest Tweets

- @WRiTExMiND no doubt! <--guess who I got tht from? Bwahaha anyway doe I like surprising people it's kinda my thing so ur welcome! And hi :)
- @skvillain yeh wiz is dope, got his own lil wave poppin! I'm fuccin wid big sean too he signed to kanye label g.o.o.d music
- And @pumahbeatz opened for @MarshaAmbrosius & blazed! So proud of him! Go bro! & Marsha was absolutely amazing! Awesome night all around. =)
- Awesome! RT @robscoms: Great 24 hours with nephews. Watched Tron, homemade mac & cheese for dinner, Wii, pancakes & Despicable Me this am!
- Good Morning 2 U Too RT @mzmonique718: Morningggg twitt birds!...up and getting ready for church...have a good day and LETS GO GIANTS!
- Goodmorning #cleveland, have a blessed day stay focused and be productive and thank god for life
- AMEN!!!>>>RT @DrSanlare: Daddy looks soooo good!!! God is amazing! To GOD be the glory and victory #TeamJesus Glad I serve an awesome God
- AGREED!! RT @ILoveElizCruz: Amen to dat... We're some awesome people! RT @itsVonnell\_Mars: @ILoveElizCruz gotta love my sign lol
- #word thanks! :) RT @Steph0e: @IBtunes HAppy Birthday love!!! =) still a fan of ya movement... yay you get another year to be dope!!! YES!!
- Happy bday isaannRT @isan\_coy: Selamatt ulang tahun yaaa RT @Phitz\_bow: Selamat siangg RT @isan\_coy: Slamat pagiiii

# Sentiment detection

Annotate tweets using labels from [http://en.wikipedia.org/wiki/List\\_of\\_emoticons](http://en.wikipedia.org/wiki/List_of_emoticons)

## 10 Saddest Tweets

- Migraine, sore throat, cough & stomach pains. Why me God?
- Ik moet werken omg !! Ik lig nog in bed en ben zo moe .. Moet alleen opstaan en tis koud buiten :(
- I Feel Horrible ' My Voice Is Gone Nd I'm Coughing Every 5 Minutes ' I Hate Feeling Like This :-/
- SMFH !!! Stomach Hurting ; Aggy ; Upset ; Tired ;; Madd Mixxy Shyt Yo !
- Worrying about my dad got me feeling sick I hate this!! I wish I could solve all these problems but I am only 1 person & can do so much..
- Malam2 menggil+ga bs napas+sakit kepala....badan remuk redam \*I miss my husband's hug....#nangismanja#
- Waking up with a sore throat = no bueno. Hoping someone didn't get me ill and it's just from sleeping. D:
- Aaaa ini tenggorokan gak enak, idung gatel bgt bawaannya pengen bersin terus. Calon2 mau sakit nih -\_\_\_-
- I'm scared of being alone, I can't see to breathe when I am lost in this dream, I need you to hold me?
- Why the hell is suzie so afraid of evelyn! Smfh no bitch is gonna hav me scared I dnt see it being possible its not!

## Word Segmentation (in Chinese)

北京大学学生体育馆

- 北京 (Beijing) 大学生 (university students) 体育馆 (gym)  
The gym for university students in Beijing.
- 北京大学 (Peking University) 生 (give birth to) 体育馆 (gym)  
Peking University gave birth to the gym?

# Statistical Machine Translation

**SMT** uses parallel corpora to automatically learn a translation

SOURCE: 目前，某些 西方 国家 已经 宣布 终止 对 津巴布韦 的 经济援助 .

H1: at present , some western nations have already announced their  
termination of economic aid to zimbabwe .

H2: at present , certain western countries have already suspended their economic  
aids to zimbabwe .

H3: so far , some western countries have declared ending economic aid to zimbabwe .

H4: some western countries have already halted economic aid to zinbarbwe at present .

SYSTEM: at present , some western countries have announced the\* end\* of the\*  
financial\* assistance\* to zimbabwe .

Open Source Machine Translation! [www.statmt.org](http://www.statmt.org)

# Visualization of Information

[Clear all constraints](#)

[Facets](#)   [Timeline](#)   [Map](#)

**2013 CE**

April 15 – Two bombs explode at the Boston Marathon, in Boston, Massachusetts in the United States, killing three and injuring 183.

**2013 CE**

March 27 – Canada becomes the first country to withdraw from the United Nations Convention to Combat Desertification.

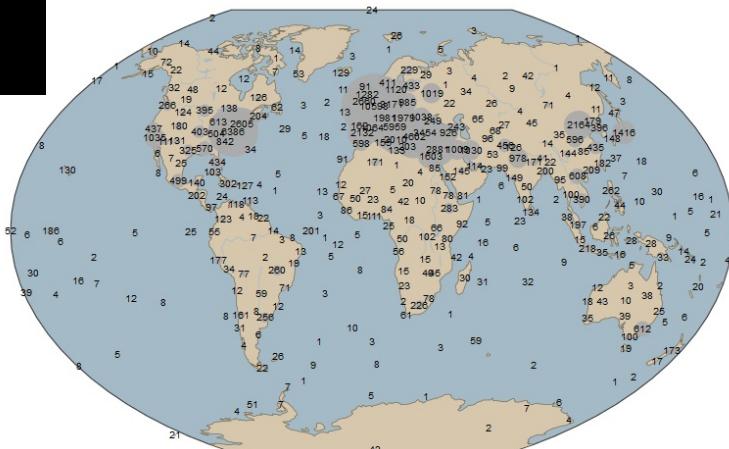
**Faceted Browsing**

Over the last century, there have been 1,491 incidents. The incident, along with a coincidental flyby of an asteroid, prompts international concern regarding the vulnerability of the planet to meteor strikes.

Role	Person	Current country	Location
<a href="#">Clear selection</a>	<a href="#">Clear selection</a>	<a href="#">Clear selection</a>	<a href="#">Clear selection</a>
<a href="#">agent [8981]</a>	Alexander [595]	United States [11993]	England [2813]
<a href="#">theme [6381]</a>	Henry [360]	United Kingdom [8297]	Rome [2811]
<a href="#">patient [4015]</a>	Antiochus [356]	Italy [6997]	United States [2375]
<a href="#">topic [2653]</a>	Charles [349]	Greece [4438]	U.S. [1945]
<a href="#">thing set [2393]</a>	Hannibal [319]	France [4302]	France [1931]
<a href="#">corpse [2321]</a>	John [318]	Turkey [3482]	Athens [1653]
<a href="#">entity changing [2077]</a>	Ptolemy [245]	China [2897]	China [1530]
<a href="#">theme(-creation) [1769]</a>	Demetrius [231]	Germany [2856]	Italy [1370]
<a href="#">killer [1534]</a>	Demosthenes [206]	Spain [2313]	London [1329]
<a href="#">entity defeated [1331]</a>	Antiochus III [203]	Egypt [1704]	Japan [1150]
<a href="#">entity victorious [1208]</a>	Constantine [202]	Japan [1471]	Egypt [1111]
<a href="#">creator [1196]</a>	Philip [198]	Russia [1301]	Germany [1068]
<a href="#">thing gotten [1009]</a>	Franks [196]	India [1174]	Constantinople [1003]
<a href="#">taker [957]</a>	Pericles [184]	The Netherlands [1168]	Spain [1001]
<a href="#">leader [857]</a>	James [170]	Iran [1150]	India [814]
<a href="#">thing taken [852]</a>	Edward [164]	Tunisia [1102]	United Kingdom [810]
<a href="#">task [761]</a>	Pyrrhus [157]	Israel [1097]	Sicily [797]
<a href="#">entity succeeding [753]</a>	Richard [152]	Syria [943]	Greece [773]

[Clear selection](#)   [Toggle](#)   [Drag](#)   [Pan](#)

## Map



## Timeline



# Identifying confusable drug names

G. Kondrak and B. Dorr

**Table 4** Top 8 names that are most similar to *Toradol* according to the BI-SIM similarity measure, and the corresponding recall values

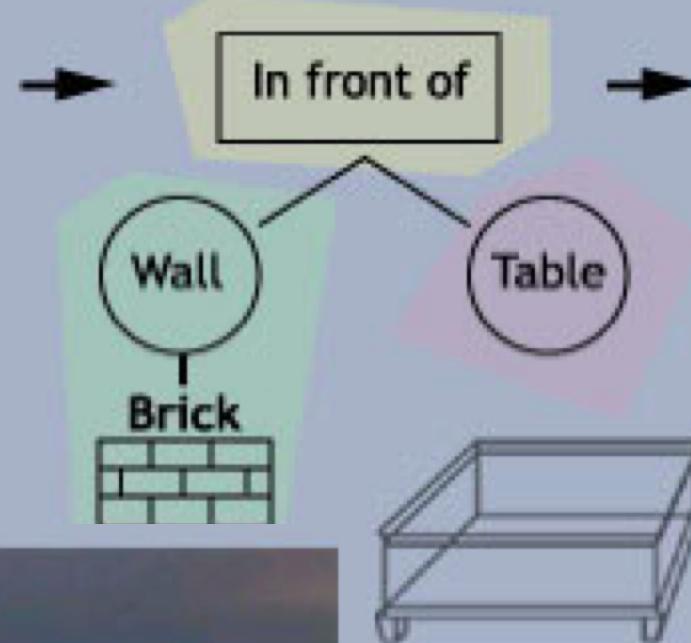
	Name	Score	+/-	Recall
1.	<i>Tramadol</i>	0.6875	+	0.25
2.	<i>Tobradex</i>	0.6250	-	0.25
3.	<i>Torecan</i>	0.5714	+	0.50
4.	<i>Stadol</i>	0.5714	-	0.50
5.	<i>Torsemide</i>	0.5000	-	0.50
6.	<i>Theraflu</i>	0.5000	-	0.50
7.	<i>Tegretol</i>	0.5000	+	0.75
8.	<i>Taxol</i>	0.5000	-	0.75

# Holy Grail: Understanding Language

- Can we *generate* language from our knowledge of language?
- Can we convert a natural language utterance into a *model* (or some other fancy logic thing)
- Can we map it into a *database*?
- Can we map it into a *mental picture* (or a *real* one?)
- Demo: WordsEye (from Richard Sproat's group at AT&T)

# Text to semantic model to image

The vase is on the Richard Sproat coffee table. The table is in front of the brick wall. The Van Gogh picture is on the wall. The Matisse sofa is next to the table. Mary is sitting on the sofa. She is playing the violin. She is wearing a straw hat.





The Devil is  
in the details

# Topics in NLP research

- Dialogue and Interactive Systems
- Discourse and Pragmatics
- Generation and Summarization
- Information Extraction and Question Answering
- Information Retrieval
- Language Resources and Evaluation
- Language and Vision
- Linguistic and Psycholinguistic Aspects of CL
- Machine Learning for NLP
- Machine Translation
- NLP for Web, Social Media and Social Sciences
- NLP-enabled Technology
- Phonology, Morphology and Word Segmentation
- Semantics
- Sentiment Analysis and Opinion Mining
- Spoken Language Processing
- Tagging, Chunking, Syntax and Parsing
- Text Categorization and Topic Models