

CS 4120: Natural Language Processing

Instructor: Prof. Lu Wang
 Northeastern University
 Webpage: www.ccs.neu.edu/home/luwang

1

Time and Location

- **Time:** Tuesdays 11:45 am - 1:25 pm, Thursdays 2:50 pm - 4:30 pm
- **Location:** Snell Library 037

2

Course Webpage

- http://www.ccs.neu.edu/home/luwang/courses/cs4120_sp2020/cs4120_sp2020.html
 - Slides, schedule for lectures and quizzes, exam, assignments (and due dates)
- You can also go to the instructor's web page and find it from there:
 - <http://www.ccs.neu.edu/home/luwang/>

3

Prerequisites

- Programming
 - Being able to write code in some programming languages (Python recommended) proficiently
- Courses
 - Algorithms
 - Probability and statistics
 - Linear algebra (optional but highly recommended)
 - Supervised machine learning (also optional but highly recommended)

4

Prerequisites

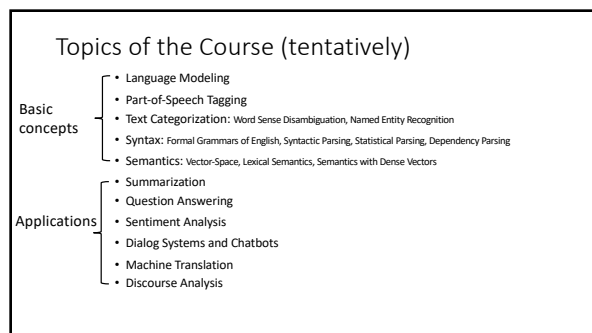
- Quiz 0 (next lecture):
 - 20 simple questions, True or False (relevant to probability, statistics, and linear algebra)
 - The purpose of this quiz is to indicate the expected background of students.
 - 80% of the questions should be easy to answer.
 - **Not counted in your final score!**
- Great notes on probability, statistics, and linear algebra
 - Probability and Statistics for Data Science, by Carlos Fernandez-Granda
 - https://cims.nyu.edu/~cfergranda/pages/stuff/probability_stats_for_DS.pdf
 - No need to be proficient in all aspects!

5

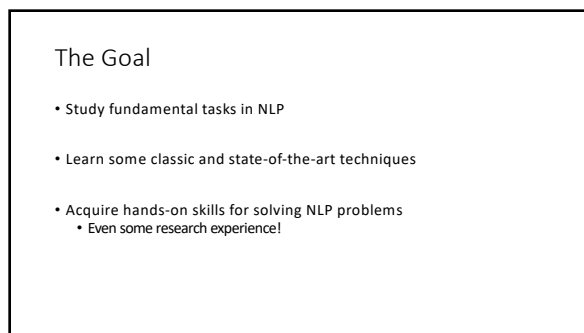
Textbook and References

- Main textbook
 - Dan Jurafsky and James H. Martin, "Speech and Language Processing, 2nd Edition", Prentice Hall, 2009.
 - We will also use some material from 3rd edition (for the available part).
 - <http://web.stanford.edu/~jurafsky/slp3/>
- Other references
 - Chris Manning and Hinrich Schütze, "Foundations of Statistical Natural Language Processing", MIT Press, 1999
 - Machine learning textbooks:
 - Christopher M. Bishop, "Pattern Recognition and Machine Learning", Springer, 2006.
 - Tom Mitchell, "Machine Learning", McGraw Hill, 1997.

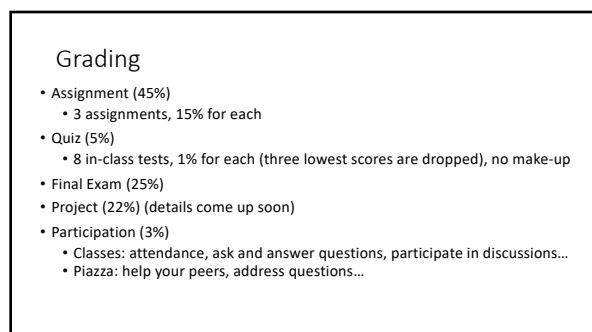
6



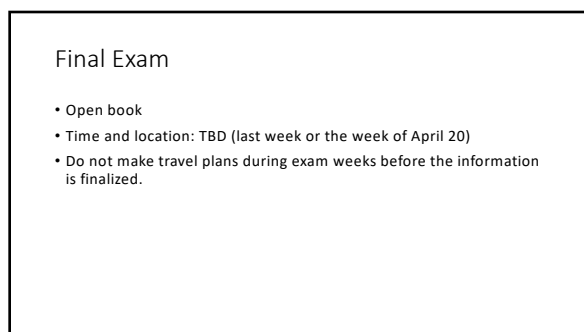
7



8



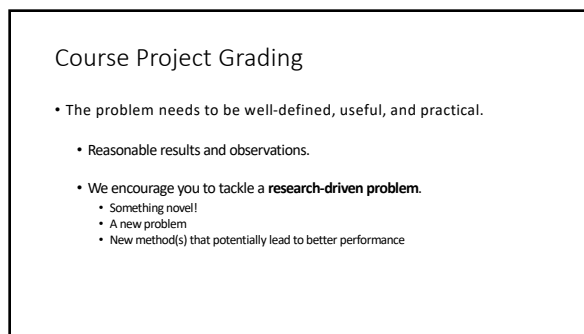
9



10



11



12

Sample Projects

- Neural Semantic Parsing Natural Language into SQL
- Short Passages Reading Comprehension and Question Answering
- Predicting Personality Traits using Tweets
- Android Application for Visual QA
- Novel Summarizer and Keyword Identifier Using Text Rank with Sentence Farn Detection
- Hashtag Similarity based on Tweet Text
- Stance Detection for the Fake News Challenge
- Machine Comprehension Using match-LSTM and Answer-Pointer
- Online Abuse Detection
- Plagiarism Detection Using FP-Growth Algorithm
- An Examination of Influential Framing of Controversial Topics on Twitter
- Paraphrase Identification using Supervised Machine Learning Techniques
- Recognizing Use of Idioms in Natural Language
- Privacy Policy Summarization with a Privacy Score

13

Neural Semantic Parsing Natural Language into SQL

- Enterprise stores data in structured format
- Skilled workforce required to extract knowledge
- What if anyone could ask questions in natural language from structured text?

Pick #	CFL Team	Player	Position	College
27	Hamilton Tiger-Cats	Connor Healy	DB	Wilfrid Laurier
28	Calgary Stampeders	Anthony Forgone	OL	York
29	Ottawa Renegades	L.P. Ladouceur	DT	California
30	Toronto Argonauts	Frank Hoffman	DL	York
...

Question:

SQL:

Result:

An example in WikiSQL. The inputs consist of a table and a question. The outputs consist of a ground truth SQL query and the corresponding result from execution.

14

Short Passages Reading Comprehension and Question Answering

Answer a simple question by reading a short passage.

- Simple question: Factoid QA, < 30 words.
- Short passage: < 300 words.
- The answer is directly given as a **range** of the passage.

- INPUT: A passage, a question.
- OUTPUT: a range of the passage.

15

... John went to school at 9AM and ate a **sandwich** for lunch and came back home at 5PM. ...

What did John have for lunch?

...In 1950, **Alan Turing** published an article titled "Computing Machinery and Intelligence" which proposed what is now called the Turing test as a criterion of intelligence...

Who published "Computing Machinery and Intelligence"?

16

Text Summarization

- Input Article:
 - Prime Minister Bertie Ahern of Ireland called Sunday for a general election on May 24. Mr. Ahern and his centrist party have governed in a coalition government since 1197...Under Irish law, which requires legislative elections every five years, Mr. Ahern had to call elections by midsummer...
- Summary:
 - Prime Min Bertie Ahern of Ireland calls for general election on May 24. He is required by law to call elections by midsummer.

17

Detection and Interpretation of English Puns

- Puns are a class of language constructs, in which lexical-semantic ambiguity is a deliberate effect of the communication act.
- For example, "I used to be a banker but I lost *interest*".

18

More Project Samples

- Stanford NLP class
 - <http://web.stanford.edu/class/cs224p>
 - Notice its focus on deep learning
 - Your project can use any machine learning technique(s) on a natural language processing problem, and shouldn't be limited to deep learning only.

19

Course Project Grading

- You are encouraged to talk to the instructor and/or TAs on project topics!
- How to find teammates?
 - Talk to your classmates and see if you share interests!
 - Post on piazza with your background (programming language and skills) and potential project ideas

20

Course Project Grading

- Three reports
 - One-page proposal (3%), due on Jan 28 at 11:59pm.
 - Progress, with code (6%)
 - Final, with code (8%)
- One presentation
 - In class (5%)

21

Audience Award

- **Bonus points!**
 - All teams vote for their favorite project(s) after presentation.
 - The team gets the most votes will be awarded with 1% bonus point!

22

Submission and Late Policy

- Programming language
 - Python (recommended), Java, C/C++
- All submissions are in electronic format.
 - Due on blackboard.

23

Submission and Late Policy

- Assignment or report turned in late will be charged 20 points (out of 100 points) off for each late day (i.e. 24 hours).
- Each student has a budget of **6 days in total** throughout the semester before a late penalty is applied.
- Late days are not applicable to final presentation.
- Each group member is charged with the same number of late days, if any, for their submission.
- You don't have to inform the instructors on the usage of late days. Timestamp of the last submission will be used for automatic grade calculation. You will see the original grades on blackboard.

24

How to find us?

- All materials, including slides, assignments, sample reports, etc, can be found on the course webpage:
 - http://www.ccs.neu.edu/home/luwang/courses/cs4120_sp2020/cs4120_sp2020.html
- Office hours and staffs
 - Prof. Lu Wang: Thursdays, from 1:30pm to 2:30pm, or by appointment, Rm 2211 at 177 Huntington Ave.
 - Note:** to attend OH at 177, you'll need to
 - Put down your name on Piazza by 1pm that Thursday (so that I can enter your name in the guest system!)
 - Bring your photo ID and check in at the front desk when you come in
 - TA Akshay Vasant Dangare, please see piazza
 - All OH starts in the week of Jan 13.
- Piazza
 - <http://piazza.com/northeastern/spring2020/cs4120>, please sign up.
 - All course relevant questions should go here! Also is the best way to reach the instructor and TAs.

25

What is Natural Language Processing?


26

What is Natural Language Processing?

- Allowing machines to communicate with human
- Natural language understanding + natural language generation

27

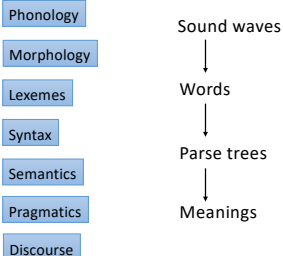
What does it mean to understand a language?



- "Stop"
- "Turn it up"
- "Volume level 6"
- "Repeat that"
- "What can you do?"
- "Play some music"
- "Play music by [artist]"
- "Play dance music on YouTube"
- "Play KEXP radio on TuneIn"
- "Play the latest episode of Radiolab"
- "Pause"
- "Next song"
- "When's my first appointment tomorrow?"
- "Wake me up at 6am tomorrow"
- "Tell me about my day"
- "How long will it take to get to work?"
- "What's the weather today?"

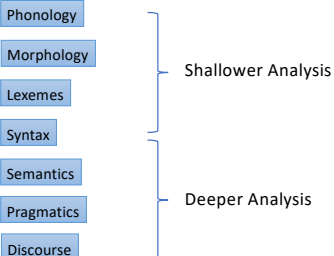
28

What does it mean to understand a language?



29

What does it mean to understand a language?



30

Syntax, Semantics, Pragmatics

- Syntax concerns the proper ordering of words and its affect on meaning.
 - The dog bit the boy.
 - The boy bit the dog.
 - Bit boy dog the the.
- Semantics concerns the (literal) meaning of words, phrases, and sentences.
 - "plant" as a photosynthetic organism
 - "plant" as a manufacturing facility
 - "plant" as the act of sowing
- Pragmatics concerns the overall communicative and social context and its effect on interpretation.
 - Honest or dishonest?
 - Context 1: Kyle and Ellen would like to see a movie. Kyle has \$20 in his pocket. Tickets cost \$8 each.
 - Context 2: Kyle and Ellen would like to see a movie. Kyle has \$20 in his pocket. Tickets cost \$10 each.

31

Syntax, Semantics, Pragmatics

- Syntax concerns the proper ordering of words and its affect on meaning.
 - The dog bit the boy.
 - The boy bit the dog.
 - Bit boy dog the the.
- Semantics concerns the (literal) meaning of words, phrases, and sentences.
 - "plant" as a photosynthetic organism
 - "plant" as a manufacturing facility
 - "plant" as the act of sowing
- Pragmatics concerns the overall communicative and social context and its effect on interpretation.
 - Honest or dishonest?
 - Context 1: Kyle and Ellen would like to see a movie. Kyle has \$20 in his pocket. Tickets cost \$8 each.
 - Context 2: Kyle and Ellen would like to see a movie. Kyle has \$20 in his pocket. Tickets cost \$10 each.
 - Kyle: "I have \$8."

32

Where NLP is used?

33

Commercial World



34

Social World

- Disaster Relief
- Chatbots for Mental Health
- Detecting abusive language in online posts


35

Text Classification: Disaster Response

- Haiti Earthquake 2010
- Classifying SMS messages

Mwen thomassin 32 nan pyron
 mwen ta renmen jwen yon ti dlo
 gras a dieu bo lakay mwen anform
 se sel dlo nou bezwen

I am in Thomassin number 32, in the area named Pyron. I would like to have some water. Thank God we are fine, but we desperately need water.



36

Extracting Social Meaning from Language

- Uncertainty (students in tutoring)
- Annoyance (callers to dialogue systems)
- Anger (police-community interaction)
- Deception
- Emotion
- Intoxication
- Flirtation, Romantic interest

37

Sentiment in Restaurant Reviews

A very bad (one-star) review:

The bartender... absolutely horrible... we waited 10 min before we even got her attention... and then we had to wait 45 - FORTY FIVE! - minutes for our entrees... stalk the waitress to get the cheque... she didn't make eye contact or even break her stride to wait for a response ...

Dan Jurafsky, Victor Chahuneau, Bryan R. Routledge, and Noah A. Smith, 2014, Narrative framing of consumer sentiment in online restaurant reviews. First Monday 19:4

38

What is the language of bad reviews?

- Negative sentiment language
horrible awful terrible bad disgusting
- Past narratives about people
waited, didn't, was he, she, his, her, manager, customer, waitress, waiter
- Frequent mentions of *we* and *us*
... we were ignored until we flagged down a waiter to get our waitress ...

39

Personal Assistants

Siri: What can I help you with?

amazon alexa

Google Now

Facebook M: A personal assistant inside Messenger

Hi, I'm Cortana. How can I help you?

40

Question Answering: IBM's Watson

41

Recommendation Engines

If you bought...

Look inside ↓

Customers who bought this item also bought

THE LANGUAGE OF FOOD
A LINGUIST READS THE MENU
DAN JURAFSKY

First Bite: How We Learn to Eat
By Ben Wilson
★ ★ ★ ★ ☆ 46
\$11.37 -prime

THE DORITO EFFECT: The Surprising New Truth About Food and Flavor
By Mark Schutler
★ ★ ★ ★ ☆ 183
\$9.48 -prime

Consider the Fork: A History of How We Cook and Eat
By Ben Wilson
★ ★ ★ ★ ☆ 253
\$15.65 -prime

Cuisine and Empire: Cooking in World History
By Rachel Laudan
★ ★ ★ ★ ☆ 25
\$16.20 -prime

42

Why NLP is challenging?

43

Ambiguity is Ubiquitous

- Speech Recognition
 - "recognize speech" vs. "wreck a nice beach"
 - "youth in Asia" vs. "euthanasia"

44

Ambiguity is Ubiquitous

- Speech Recognition
 - "recognize speech" vs. "wreck a nice beach"
 - "youth in Asia" vs. "euthanasia"
- Syntactic Analysis
 - "I ate spaghetti **with** chopsticks" vs. "I ate spaghetti **with** meatballs."

45

Ambiguity is Ubiquitous

- Speech Recognition
 - "recognize speech" vs. "wreck a nice beach"
 - "youth in Asia" vs. "euthanasia"
- Syntactic Analysis
 - "I ate spaghetti **with** chopsticks" vs. "I ate spaghetti **with** meatballs."
- Semantic Analysis
 - "The dog is in the **pen**." vs. "The ink is in the **pen**."
 - "I put the **plant** in the window" vs. "Ford put the **plant** in Mexico"

46

Ambiguity is Ubiquitous

- Speech Recognition
 - "recognize speech" vs. "wreck a nice beach"
 - "youth in Asia" vs. "euthanasia"
- Syntactic Analysis
 - "I ate spaghetti **with** chopsticks" vs. "I ate spaghetti **with** meatballs."
- Semantic Analysis
 - "The dog is in the **pen**." vs. "The ink is in the **pen**."
 - "I put the **plant** in the window" vs. "Ford put the **plant** in Mexico"
- Pragmatic Analysis
 - From "The Pink Panther Strikes Again":
 - **Clouseau**: Does your dog bite?
 - **Hotel Clerk**: No.
 - **Clouseau**: [Bowing down to pet the dog] Nice doggie.
 - [Dog barks and bites Clouseau in the hand]
 - **Clouseau**: I thought you said your dog did not bite!
 - **Hotel Clerk**: That is not my dog.

47

Ambiguity

Find at least 6 meanings of this sentence:

I made her duck

48

Ambiguity

Find at least 6 meanings of this sentence:

I made her duck

- I cooked waterfowl for her benefit (to eat)
- I cooked waterfowl belonging to her
- I created the (plaster?) waterfowl she owns
- I caused her to quickly lower her head or body
- I recognized the true identity of her spy waterfowl
- I waved my magic wand and turned her into undifferentiated waterfowl

49

Ambiguity

I caused her to quickly lower her head or body

Part of speech: "duck" can be a Noun or Verb

I cooked waterfowl belonging to her.

Part of speech:

"her" is possessive pronoun ("of her")

"her" is dative pronoun ("for her")

I made the (plaster) duck statue she owns

Word Meaning: "make" can mean "create" or "cook"

50

Ambiguity is Explosive

- Ambiguities compound to generate enormous numbers of possible interpretations.
- In English, a sentence ending in n prepositional phrases has *over 2ⁿ* syntactic interpretations
 - "I saw the man with the telescope": 2 parses
 - "I saw the man on the hill with the telescope": 5 parses
 - "I saw the man on the hill in Texas with the telescope": 14 parses
 - "I saw the man on the hill in Texas with the telescope at noon": 42 parses
 - "I saw the man on the hill in Texas with the telescope at noon on Monday": 132 parses

51

Humor and Ambiguity

- Many jokes rely on the ambiguity of language:
 - Policeman to little boy: "We are looking for a thief with a bicycle." Little boy: "Wouldn't you be better using your eyes?"
 - Why is the teacher wearing sun-glasses. Because the class is so bright.
 - Groucho Marx: One morning I shot an elephant in my pajamas. How he got into my pajamas, I'll never know.
 - She criticized my apartment, so I knocked her flat.
 - Noah took all of the animals on the ark in pairs. Except the worms, they came in apples.

52

Why is Language Ambiguous?

53

Why is Language Ambiguous?

- Having a unique linguistic expression for every possible conceptualization that could be conveyed would make language overly complex and linguistic expressions unnecessarily long.
- Allowing resolvable ambiguity permits shorter linguistic expressions, i.e. data compression.
- Language relies on people's ability to use their knowledge and inference abilities to properly resolve ambiguities.
- Infrequently, disambiguation fails, i.e. the compression is lossy.

54

More difficulties: Non-standard language

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either

And neologisms:

- unfriend
- retweet
- bromance

55

Some NLP Tasks

56

Syntactic Tasks

57

Word Segmentation

- Breaking a string of characters into a sequence of words.
- In some written languages (e.g. Chinese) words are not separated by spaces.
- Even in English, characters other than white-space can be used to separate words [e.g., ;, -, : ()]
- Examples from English URLs:
 - jumptheshark.com ⇒ jump the shark .com
 - twitter.com/realdonaldtrump ⇒ real donald trump .com
 - myspace.com/pluckerswingbar
 - ⇒ myspace .com pluckers wing bar
 - ⇒ myspace .com plucker swing bar

58

Morphological Analysis

- **Morphology** is the field of linguistics that studies the internal structure of words. (Wikipedia)
- A **morpheme** is the smallest linguistic unit that has semantic meaning (Wikipedia)
 - e.g. "carry", "pre", "ed", "ly", "s"
- Morphological analysis is the task of segmenting a word into its morphemes:
 - carried ⇒ carry + ed (past tense)
 - independently ⇒ in + (depend + ent) + ly
 - Googlers ⇒ (Google + er) + s (plural)
 - unlockable ⇒ un + (lock + able) ?
 - ⇒ (un + lock) + able ?

59

Part Of Speech (POS) Tagging

- Annotate each word in a sentence with a part-of-speech.

I ate the spaghetti with meatballs.
 Pro V Det N Prep N
 John saw the saw and decided to take it to the table.
 PN V Det N Con V Part V Pro Prep Det N

- Useful for subsequent syntactic parsing and word sense disambiguation.

60

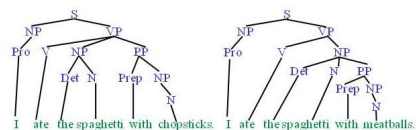
Phrase Chunking

- Find all non-recursive noun phrases (NPs) and verb phrases (VPs) in a sentence.
 - [NP I] [VP ate] [NP the spaghetti] [PP with] [NP meatballs].
 - [NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to] [NP only # 1.8 billion] [PP in] [NP September]

61

Syntactic Parsing

- Produce the correct syntactic parse tree for a sentence.



62

Semantic Tasks

63

Word Sense Disambiguation (WSD)

- Words in natural language usually have a fair number of different possible meanings.
 - Ellen has a strong **interest** in computational linguistics.
 - Ellen pays a large amount of **interest** on her credit card.
- For many tasks (question answering, translation), the proper sense of each ambiguous word in a sentence must be determined.

64

Semantic Role Labeling (SRL)

- For each clause, determine the semantic role played by each noun phrase that is an argument to the verb.
 - agent patient source destination instrument
 - John drove Mary from Austin to Dallas in his Toyota Prius.
 - The hammer broke the window.
- Also referred to a "case role analysis," "thematic analysis," and "shallow semantic parsing"

65

Semantic Parsing

- A **semantic parser** maps a natural-language sentence to a complete, detailed semantic representation (**logical form**).
- For many applications, the desired output is immediately executable by another program.
- Example: Mapping an English database query to Prolog:


```
How many cities are there in the US?
answer(A, count(B, (city(B), loc(B, C),
                    const(C, countryid(USA))),
A))
```

66

Textual Entailment

- Determine whether one natural language sentence entails (implies) another under an ordinary interpretation.
- E.g., "A soccer game with multiple males playing. -> Some men are playing a sport."

67

Pragmatics/Discourse Tasks

68

Anaphora Resolution/Co-Reference

- Determine which phrases in a document refer to the same underlying entity.
 - John put the carrot on the plate and ate it.
- Bush started the war in Iraq. But the president needed the consent of Congress.
- Some cases require difficult reasoning.
 - Today was Jack's birthday. Penny and Janet went to the store. They were going to get presents. Janet decided to get a kite. "Don't do that," said Penny. "Jack has a kite. He will make you take it back."

69

More Application-driven Tasks

70

Information Extraction (IE)

- Identify phrases in language that refer to specific types of entities and relations in text.
- Named entity recognition is task of identifying names of people, places, organizations, etc. in text.
 - people organizations places
 - Michael Dell is the CEO of Dell Computer Corporation and lives in Austin Texas.
- Relation extraction identifies specific relations between entities.
 - Michael Dell is the CEO of Dell Computer Corporation and lives in Austin Texas.
 - Michael Dell is the CEO of Dell Computer Corporation and lives in Austin Texas.

71

Question Answering

- Directly answer natural language questions based on information presented in a corpora of textual documents (e.g. the web).
 - Who is the president of United States?
 - Donald Trump
 - What is the popular of Massachusetts?
 - 6.8 million

72

Text Summarization

- Produce a short summary of one or many longer document(s).
- **Article:** An international team of scientists studied diet and mortality in 135,335 people between 35 and 70 years old in 18 countries, following them for an average of more than seven years. Diet information depended on self-reports, and the scientists controlled for factors including age, sex, smoking, physical activity and body mass index. The study is in The Lancet. Compared with people who ate the lowest 20 percent of carbohydrates, those who ate the highest 20 percent had a 28 percent increased risk of death. But high carbohydrate intake was not associated with cardiovascular death. ...
- **Summary:** Researchers found that people who ate higher amounts of carbohydrates had a higher risk of dying than those who ate more fats.

73

Spoken Dialogue Systems -- Chatbots

- Q: Is it going to rain today?
- A: It will be mostly sunny. No rain is expected.



74

Machine Translation

- Translate a sentence from one natural language to another.
- 我喜欢汉堡 → I like burgers.

75

Ambiguity Resolution is Required for Translation

- Syntactic and semantic ambiguities must be properly resolved for correct translation:
 - "John plays the guitar." → "John 弹 吉他"
 - "John plays soccer." → "John 踢 足球"

76

Ambiguity Resolution is Required for Translation

- Syntactic and semantic ambiguities must be properly resolved for correct translation:
 - "John plays the guitar." → "John 弹 吉他"
 - "John plays soccer." → "John 踢 足球"
- An apocryphal story is that an early MT system gave the following results when translating from English to Russian and then back to English:
 - "The spirit is willing but the flesh is weak." → "The liquor is good but the meat is spoiled."
 - "Out of sight, out of mind." → "Invisible idiot."

77

Resolving Ambiguity

- Choosing the correct interpretation of linguistic utterances requires (commonsense) knowledge of:
 - **Syntax**
 - An agent is typically the subject of the verb
 - **Semantics**
 - Michael and Ellen are names of people
 - August is the name of a month (and of a person)
 - Toyota is a car company and Prius is a brand of car
 - **Pragmatics**
 - Some social norm, communicative goals
 - Asking a question, expecting an answer
 - **World knowledge**
 - Credit cards require users to pay financial interest
 - Agents must be animate and a hammer is not animate

78

State-of-the-Arts

- Learning from large amounts of text data (cf. rule-based methods)
 - Supervised learning or unsupervised learning
- Statistical machine learning-based methods
 - The probabilistic knowledge acquired allows robust processing that handles linguistic regularities as well as exceptions.
- Now with neural network-based methods mostly

79

Related Fields

- Artificial Intelligence
- Machine Learning
- Linguistics
- Cognitive science
- Logic
- Data science
- Political science
- Education
- Economics
- ...many more

80

Relevant Scientific Conferences and Journals

- Association for Computational Linguistics (ACL)
- North American Association for Computational Linguistics (NAACL)
- Empirical Methods in Natural Language Processing (EMNLP)
- International Conference on Computational Linguistics (COLING)
- Conference on Computational Natural Language Learning (CoNLL)
- Transactions of the Association for Computational Linguistics (TACL)
- Journal of Computational Linguistics (CL)

81