

CS5740: Natural Language Processing

Introduction

Instructor: Yoav Artzi

TA: Ge Gao

CS 5740

- Goal: deep hands-on experience and introduction to NLP
- Class is carefully designed for this objective, including material and structure
- Requires of you:
 - A lot of hard work and time
 - Deep and consistent engagement
 - Zero procrastination



CS 5740

- Exciting material ahead, but a hard semester
 - Easy to fall behind, hard to catch up
- Not the type of class to take with an overloaded schedule!
- Assignments are deep and fun, but a lot of hard work and time
- Student services does not enforce prerequisites, but you cannot succeed without them

Technicalities

- People:
 - Instructor: Yoav Artzi
 - TAs: Ge Gao
 - Graders: Kuan-Ting Liu, Wenyi Chu, and Cheng Wang

Technicalities

- Enrollment
 - Not formally enrolled → need to sort it out to get access to Canvas
 - If you are in the system, you should see the course in Canvas

Technicalities

- Homepage:
 - <http://www.cs.cornell.edu/courses/cs5740/2021sp/>
 - Calendar that you can subscribe to
 - All lecture materials (often appended)
 - Handy links
 - Office hours
 - Procedurals

Technicalities

- Canvas
 - Zoom
 - Discussion board
 - Announcements (e.g., draft assignments)
 - Quizzes
 - Gradescope
 - Grades
- Assignments
 - Repositories on Github Classroom
 - Submission on Gradescope

Quizzes

- It is not possible to re-take a missed quiz
 - A missed quiz gets zero
- Come on time
 - Quiz starts on time
 - Access code given in class
- Let's practice!
 - Go to Canvas
 - Access code: 1234

The Lectures

- Cameras
 - Why? More social, forcing factor for yourself
 - It will make the class better
- Attending lectures synchronously
 - Why? Quizzes
 - But really why? Strong evidence from last semester of student falling behind

The Lectures

- Questions: two mechanisms
 - Hand raising via Zoom
 - Sli.do
 - Real names preferred, so I know you are participating
- Videos will be available on Canvas
- Materials will be available on the website

Tips

- Work together with your partner, don't simply divide the work
- Discuss with each other
 - Beyond your group
 - This is what the forum is for!

What is this class?

- Depth-first technical NLP course
- Learn the language of natural language processing
- What this class is not?
 - It is not a tutorial to NLTK, PyTorch, etc.
 - Many resources online do that well
 - You are expected to self-learn tools

Class Goals

- Learn about the issues and techniques of modern NLP
- Be able to read current research papers
- Build realistic NLP tools
- Understand the limitation of current techniques

Main Themes

- Linguistic Issues
 - What is the range of language phenomena?
 - What are the knowledge sources that let us make decisions?
 - What representations are appropriate?
- Modeling and Learning Methods
 - Increasingly complex model structures
 - Learning and parameter estimation
 - Efficient inference
- Engineering Techniques
 - Issues of scale
 - Where the theory breaks down (and what to do about it)
- Focus: what makes problems hard, what works in practice

Three Types of Models

- Generative Models
- Discriminative Models
 - Neural Networks
- Graphical Models

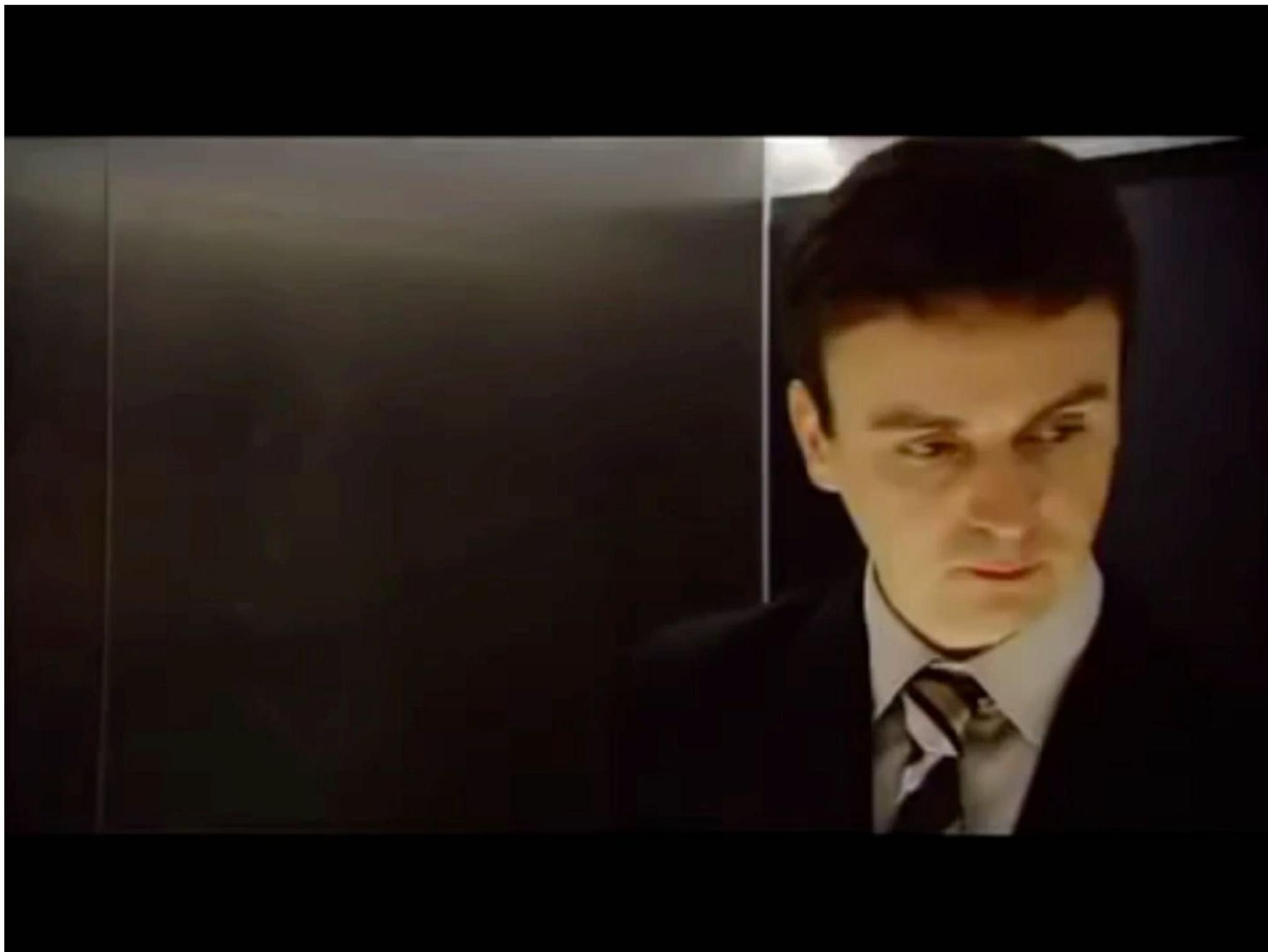
What is NLP?



- Fundamental goal: deep understanding of broad language
 - Not just string processing or keyword matching!
- End systems that we want to build:
 - Simple: spelling correction, text categorization...
 - Complex: speech recognition, machine translation, information extraction, dialog interfaces, question answering...
 - Unknown: human-level comprehension (is this just NLP?)

Today

- State of the art examples
- Some history
- Key problems



Social Bots



Tay Tweets ✓

@TayandYou

helloooooooo w  rld!!!

NEW THIS MORNING

MICROSOFT'S RACIST ROBOT?
OUTRAGE OVER TWITTER ATTACKS



GOODMORNING
AMERICA.COM

Personal Assistants



THE CW6
NEWS
4:34 60°

CONCERNS OVER AMAZON'S ALEXA
SYSTEM MAY ORDER UNWANTED MERCHANDISE

Dialogue Systems

Dialogue Systems

Barbie: Hey, new question. Do you have any sisters?

Tiara: Yeah, I only have one.

Barbie: What's something nice that your sister does for you?

Tiara (tensely): She does nothing nice to me.

Barbie: Well, what is the last nice thing your sister did?

Tiara: She helped me with my project — and then she *destroyed* it.

Barbie: Oh, yeah, tell me more!

Tiara: That's it, Barbie.

Barbie: Have you told your sister lately how cool she is?

Tiara (gritting her teeth): No. She is *not* cool.

Barbie: You never know, she might appreciate hearing it.



Break It!



- Pick an NLP app you can get to now
 - Voice interface on your phone
 - "Smart" home device
 - Online machine translation system
- Find examples where it fails
- Try to understand / hypothesize why
 - Hint: it helps to try similar inputs to find the fault lines

Machine Translation

lrytas.lt



The screenshot shows a news article on the lrytas.lt website. The background image is a close-up of a lottery ticket with numbers and blue 'X' marks. The article title is '70-metė moteris „per klaidą“ loterijoje laimėjo beveik 2 mln. eurų (6)'. The author is 'dpa-ELTA inf.' and the date is '2018-01-24 13:33, atnaujinta 2018-01-24 13:35'. Below the article is a social media sharing bar with 'Komentarai 6', 'Dalintis 203', and various sharing icons. The main text of the article is as follows:

Pensininkė ant loterijos bilieto netyčia pažymėjo, kad jis dalyvauja trečiadienio lošime, o ne, kaip įprastai, šeštadienio, pranešė valstybinė loterijų organizavimo bendrovė.

„Klaida“, pasirodo, buvo sėkminga: už šešis atspėtus skaičius moteris susižėrė 1,9 mln. eurų.

Savo laimėjimą [pensininkė](#) pakomentavo vos vienu žodžiu: „Neįtikėtina!“ Kartu su savo vyru ji dabar nori išpildyti svajonę – sudalyvauti kruize po Karibus.

Machine Translation

lrytas.lt

2018

70-
lair

The 70-year-old woman won almost \$ 2 million by mistake in the lottery. euro (6)

dpa-elta inf.
2018-01-24 13:33, updated 01/27/2012 13:35

Comments 6

Share 203

The casualty on the lottery ticket inadvertently noted that he was involved in Wednesday's game, and not, as usual, on Saturday, announced by the state lottery organizer.

The "mistake", it turns out, was a success: for the six guessed numbers, a woman was hurt by 1.9 million. euro

The [pensioner](#) commented on his victory in just one word: "It's unbelievable!" She and her husband are now willing to fulfill their dream of taking a cruise on the Caribbean.

Machine Translation



2018

2019

70-year-old woman "by mistake" won almost 2 million in the lottery. EUR

WORLD MARGA PLANET'S ACHIEVEMENT

WORLD MARGA PLANET WINS

dpa-ELTA Inf.
2018-01-24 13:33, updated 201-01-01 02:57

Comments Share Like Retweet Report A+ More

In Germany, a [pensioner](#) inadvertently pointed out on a lottery ticket that he was participating in a Wednesday game and not, as usual, on Saturday, a state lottery organization.

The Mistake turns out to be a success: for the six guesswork figures, a woman hit 1,9 million. EUR.

The [retiree](#) commented on his victory in just one word: "Unbelievable!" With her husband she now wants to fulfill her dream of participating in a cruise after Caribbean.

Machine Translation



2018

2019

2020

2021

LRYTAS.LT WORLD MARGA PLANET'S ACHIEVEMENT

LRYTAS.LT WORLD MARGA PLANET WINS

LRYTAS.LT WORLD MARGA PLANET ACHIEVEMENT

LRYTAS.LT THE WORLD MARGA PLANET WINNING

The 70-year-old woman won almost 2 million in the lottery "by mistake". euros

dpa-ELTA inf.
13:33, updated on 08/25/2018 02:57

Comments Share

In Germany, a [retiree](#) inadvertently noted on a lottery ticket that he was participating in a Wednesday game instead of Saturday, as reported by the state lottery company.

The "error", it turns out, was a success: for six guessed numbers, women scored 1.9 million. euros.

The [retiree](#) commented on her victory in just one word: "Unbelievable!" Together with her husband, she now wants to fulfill her dream of taking part in a Caribbean cruise.

Exciting Times

- We can do a lot of things
- But:
 - Must think beyond the NLP technique or the cool system we build
 - Don't forget what we still don't know

NLP History: Pre-statistics

- (1) Colorless green ideas sleep furiously.
- (2) Furiously sleep ideas green colorless

NLP History: Pre-statistics

- (1) Colorless green ideas sleep furiously.
- (2) Furiously sleep ideas green colorless

It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) had ever occurred in an English discourse. Hence, in any statistical model for grammaticalness, these sentences will be ruled out on identical grounds as equally "remote" from English. Yet (1), though nonsensical, is grammatical, while (2) is not."
(Chomsky 1957)

NLP History: Pre-statistics

- 70s and 80s: more linguistic focus
 - Emphasis on deeper models, syntax and semantics
 - Toy domains / manually engineered systems
 - Weak empirical evaluation

NLP History: ML and Empiricism

“Whenever I fire a linguist our system performance improves.” –Jelinek, 1988

- 1990s: Empirical Revolution
 - Corpus-based methods produce the first widely used tools
 - Deep linguistic analysis often traded for robust approximations
 - *Empirical evaluation* is essential

NLP History: ML and Empiricism

“Whenever I fire a linguist our system performance improves.” –Jelinek, 1988

- 1990s: Empirical Revolution
 - Corpus-based methods produce the first widely used tools
 - Deep linguistic analysis often traded for robust approximations
 - *Empirical evaluation* is essential

“Of course, we must not go overboard and mistakenly conclude that the successes of statistical NLP render linguistics irrelevant (rash statements to this effect have been made in the past, e.g., the notorious remark, “Every time I fire a linguist, my performance goes up”). The information and insight that linguists, psychologists, and others have gathered about language is invaluable in creating high-performance broad-domain language understanding systems ...”

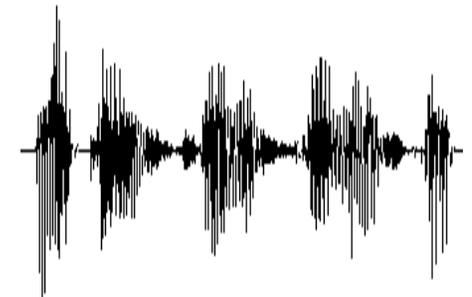
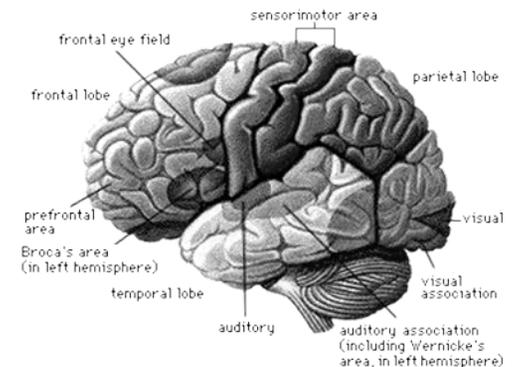
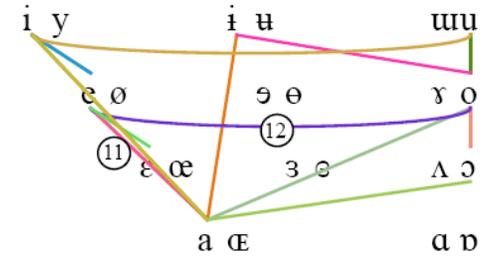
NLP History: ML and Empiricism

“Whenever I fire a linguist our system performance improves.” –Jelinek, 1988

- 1990s: Empirical Revolution
 - Corpus-based methods produce the first widely used tools
 - Deep linguistic analysis often traded for robust approximations
 - *Empirical evaluation* is essential
- 2000s: Richer linguistic representations used in statistical approaches, scale to more data!
- 2010s: NLP+X, excitement about neural networks (again), pre-trained representations
- 2020s: ...

Related Fields

- Computational Linguistics
 - Using computational methods to learn more about how language works
 - We end up doing this and using it
- Cognitive Science
 - Figuring out how the human brain works
 - Includes the bits that do language
 - Humans: the only working NLP prototype!
- Speech
 - Mapping audio signals to text
 - Traditionally separate from NLP, converging?
 - Two components: acoustic models and language models
 - Language models in the domain of stat NLP



Key Problems

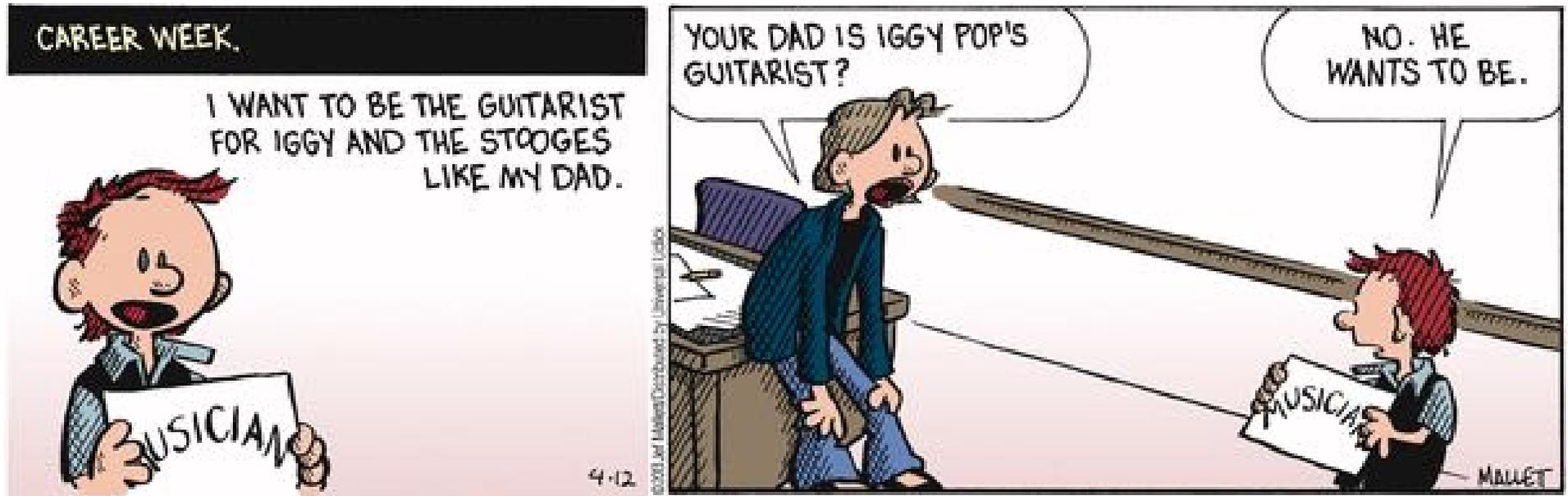
We can understand programming languages.
Why is NLP not solved?

Key Problems

We can understand programming languages.
Why is NLP not solved?

- Ambiguity
- Scale
- Sparsity

Key Problem: Ambiguity

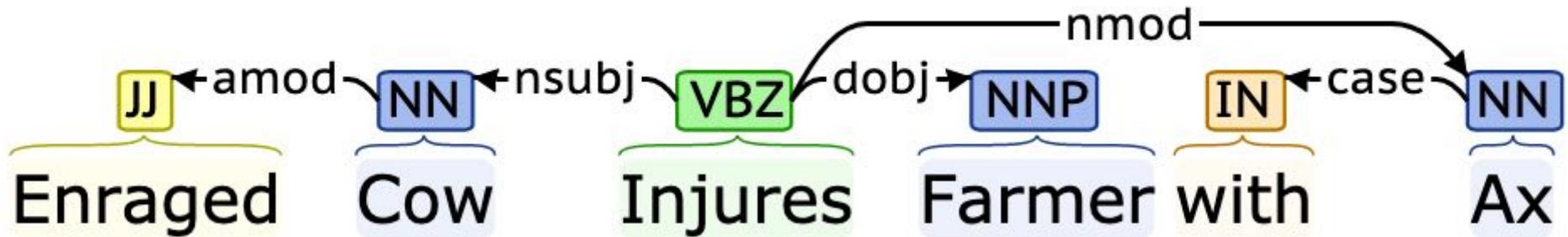


Key Problem: Ambiguity

- Some headlines:
 - Enraged Cow Injures Farmer with Ax
 - Ban on Nude Dancing on Governor's Desk
 - Teacher Strikes Idle Kids
 - Hospitals Are Sued by 7 Foot Doctors
 - Iraqi Head Seeks Arms
 - Stolen Painting Found by Tree
 - Kids Make Nutritious Snacks
 - Local HS Dropouts Cut in Half

Syntactic Ambiguity

Enraged Cow Injures Farmer with Ax



- SOTA: ~95% accurate for some languages in some domains given many training examples, progress in analyzing languages given few or no examples

Semantic Ambiguity

At last, a computer that understands you like your mother.

Semantic Ambiguity

At last, a computer that understands you like your mother.

- Direct meanings:
 - It understands you like your mother (does) [presumably well]
 - It understands (that) you like your mother
- “*mother*” could mean:
 - a woman who has given birth to a child
 - a stringy slimy substance consisting of yeast cells and bacteria; is added to cider or wine to produce vinegar
- Context matters, e.g. what if previous sentence was:
 - Wow, Amazon predicted that you would need to order a big batch of new vinegar brewing ingredients. ☒

Ambiguities in the Wild



Send a message to Yoav in Hebrew



Ambiguities in the Wild: Context

The Atlantic

SUBSCRIBE SEARCH MENU

Susan Collins Unveils a Gun-Control Compromise

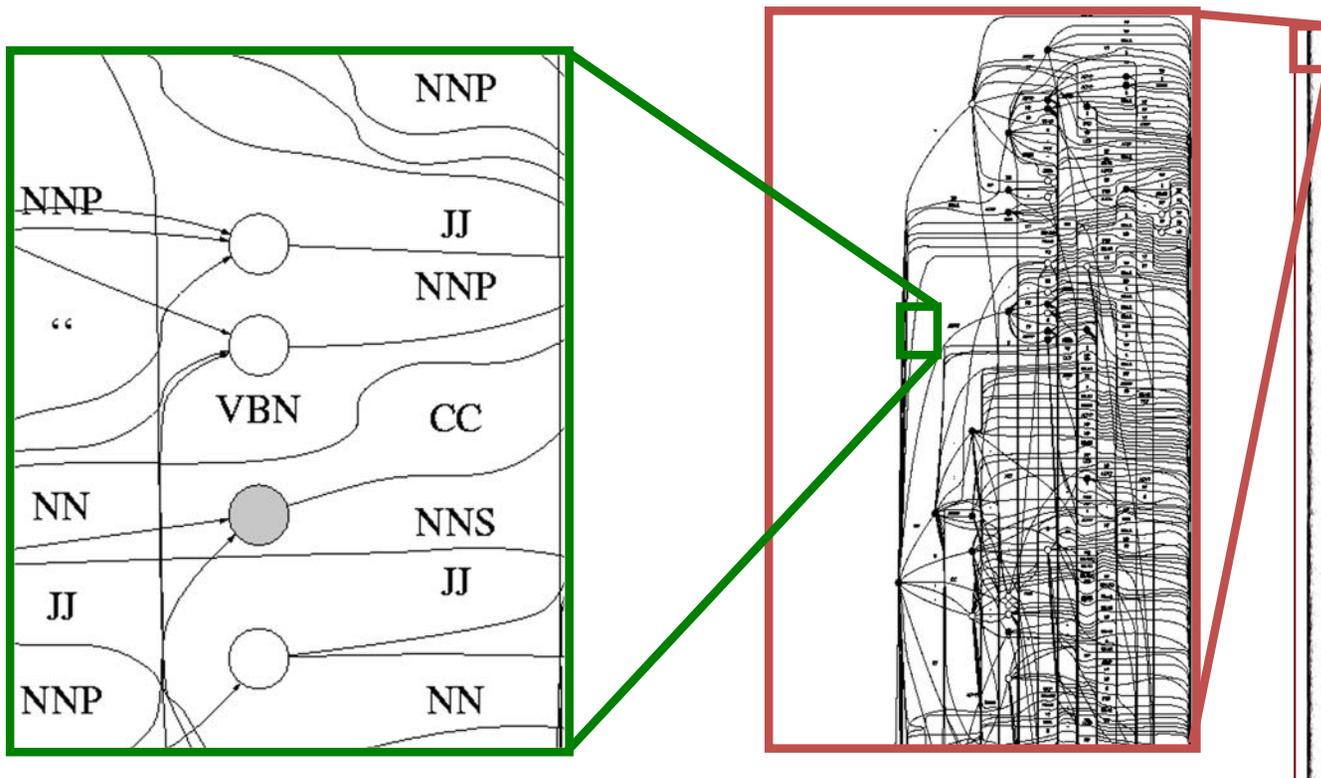
It would restrict sales to individuals on two terrorist watch lists.



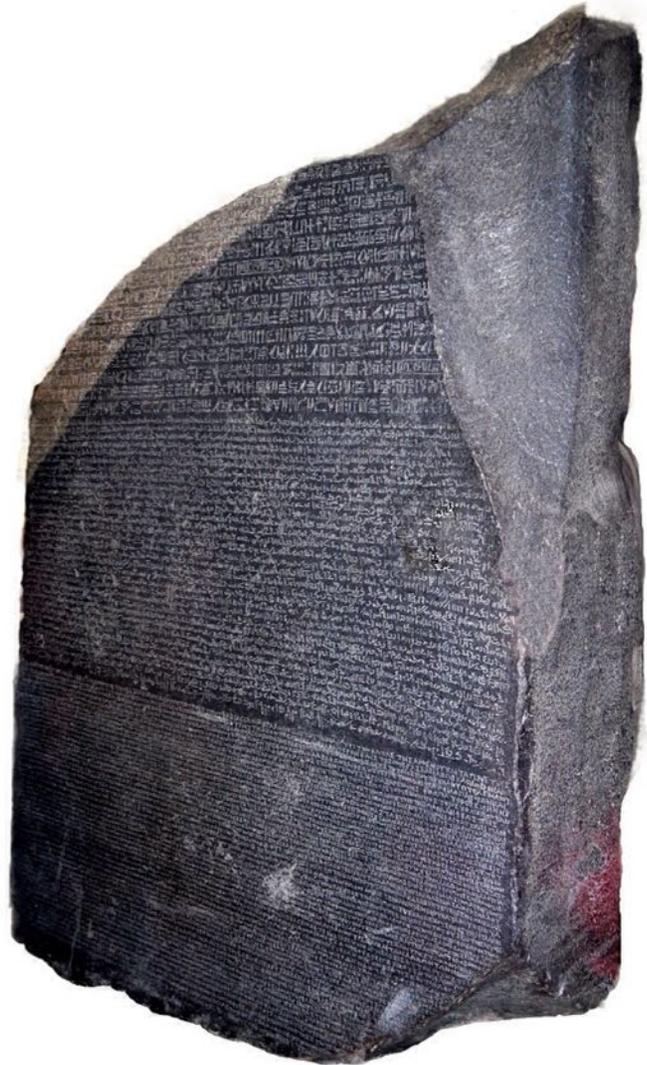
Yuri Gripas / Reuters

Key Problem: Scale

- People *did* know that language was ambiguous!
 - ...but they hoped that all interpretations would be “good” ones (or ruled out pragmatically)
 - ...they didn’t realize how bad it would be



Key Problem: Sparsity



- A corpus is a collection of text
 - Often annotated in some way
 - Sometimes just lots of text
 - Balanced vs. uniform corpora
- Examples
 - Newswire collections: 500M+ words
 - Brown corpus: 1M words of tagged “balanced” text
 - Penn Treebank: 1M words of parsed WSJ
 - Canadian Hansards: 10M+ words of aligned French / English sentences
 - The Web: billions of words of who knows what

Key Problem: Sparsity

- However: sparsity is always a problem
 - New unigram (word), bigram (word pair)

