

# CS 124/LINGUIST 180

## From Languages to Information

DAN JURAFSKY

PROFESSOR OF COMPUTER SCIENCE

PROFESSOR OF LINGUISTICS

STANFORD UNIVERSITY

WINTER 2023

INTRODUCTION AND COURSE OVERVIEW

# From Languages to Information

Automatically extracting meaning and structure from:

- Human language (news, social media, etc.)
- Social networks

Interacting with humans via language

- Dialog systems/Chatbots
- Question Answering
- Recommendation Systems

# Commercial World



OpenAI



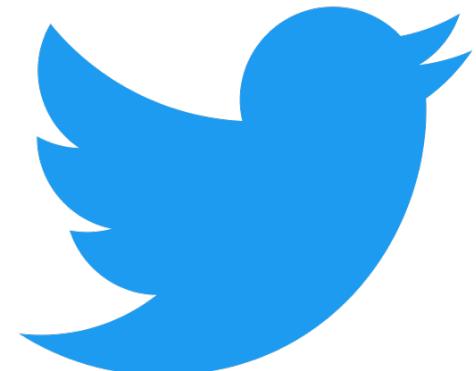
YouTube

*Microsoft*®

∞ Meta



amazon



Social World

Disaster Relief

Helping teachers in the classroom

Improve Police-Community relations via  
Body-Cameras

Training Mental Health Counselors

# 1. Extracting information from language

# Information Retrieval

6,586,013,574 web searches every day (by one estimate)

Text-based information retrieval is likely the most frequently used piece of software in the world

How does it work? Can you build an IR engine?

*Programming Assignment 3: Search!*

# 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 anfom  
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.



*Programming  
Assignment 2: Triage!*

# Extracting Sentiment and Social Meaning

Lots of meaning is in **connotation**

"connotation: an idea or feeling that a word invokes in addition to its literal or primary meaning."

Extracting connotation is generally called  
**sentiment analysis**

*Programming Assignment 2: Sentiment*

A (fictional) application of sentiment analysis for toxicity intervention  
that shows how hard it can be!

# Emotional Spell-Check

# Sentiment in Restaurant Reviews

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

900,000 Yelp reviews online

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 ...

# 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 ...

# Other narratives with this language

A genre using:

Past tense, we/us, negative, people narratives

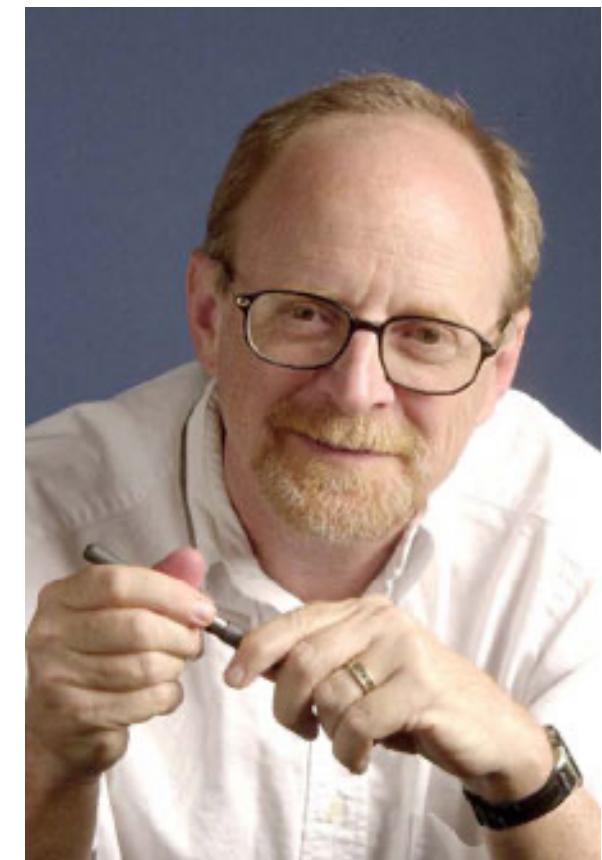
Texts written by **people suffering trauma**

- James Pennebaker lab at UT Austin
- Past tense is used for "distancing"
- Use of “we”: seeking solace in community

**1-star reviews are trauma narratives!**

The lesson of reviews:

**It's all about personal interaction**



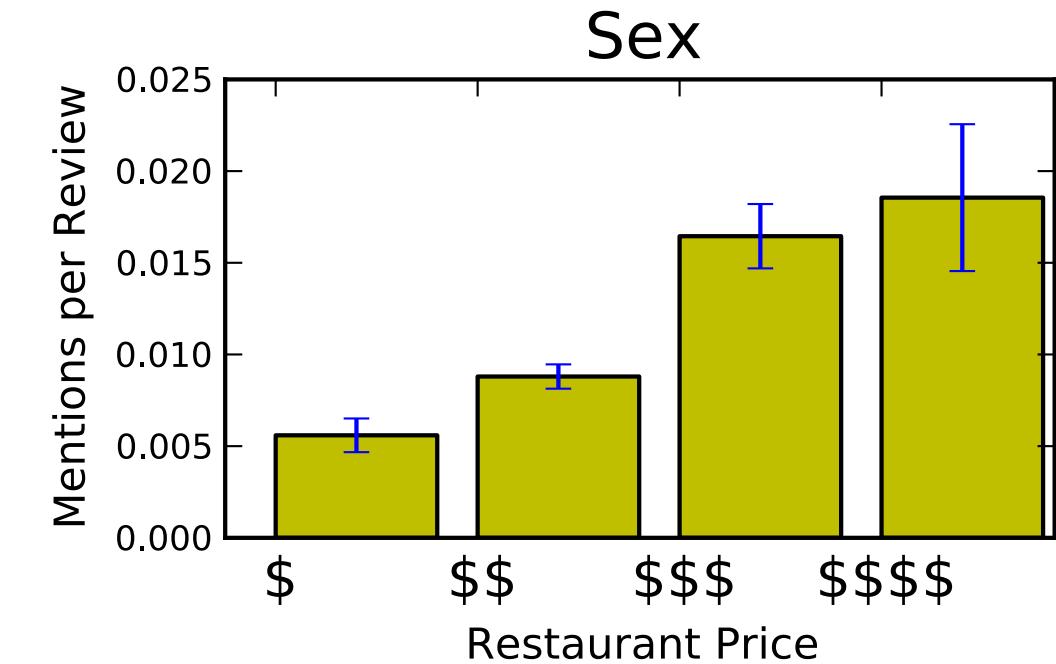
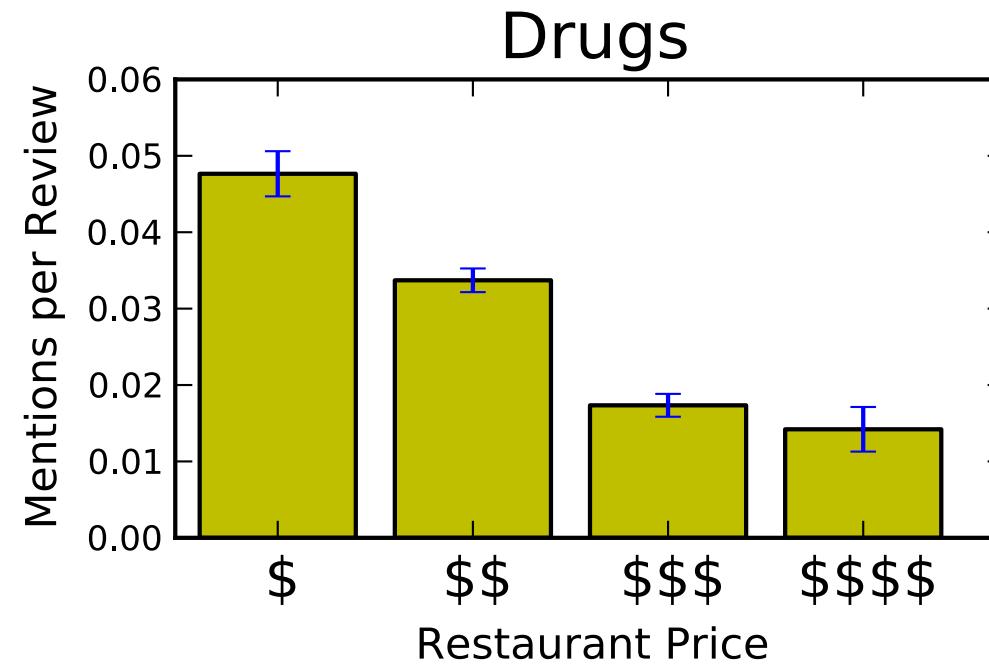
# What about positive reviews? Sex, Drugs, and Dessert

*addicted to pepper shooters*

*garlic noodles... my drug of choice*

*the fries are like crack*

*orgasmic pastry*  
*sexy food*  
*seductively seared fois gras*



# Computational Biology: Comparing Sequences

**AGGCTATCACCTGACCTCCAGGCCGATGCC**

**TAGCTATCACGACCGCGGGTCGATTGCCCGAC**

-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---

TAG-CTATCAC--GACCGC--GGTCGATTGCCCGAC

# Sequence comparison is key to

- Finding genes
  - Determining function
  - Uncovering evolutionary processes

# This is also how spell checkers work!

# We'll learn: edit distance algorithms (Quiz 1)

# Social Networks

The network formed by your friends or other relations offline or online

- Can we compute properties of these networks?
  - Extract information from them?
- 
- **We'll learn: Network algorithms (Quiz 9)**

# Help improve Police-Community Interaction (week 9)

Problems:

- A flood of viral videos show inappropriate officer use of force
- Black Americans report more negative interactions with police



Could natural language processing help?

- Quantify police-community interactions using body-worn cameras?
- Help develop officer training?
- Reduce the chances of violence?
- I'll talk about work with Prof. Jennifer Eberhardt

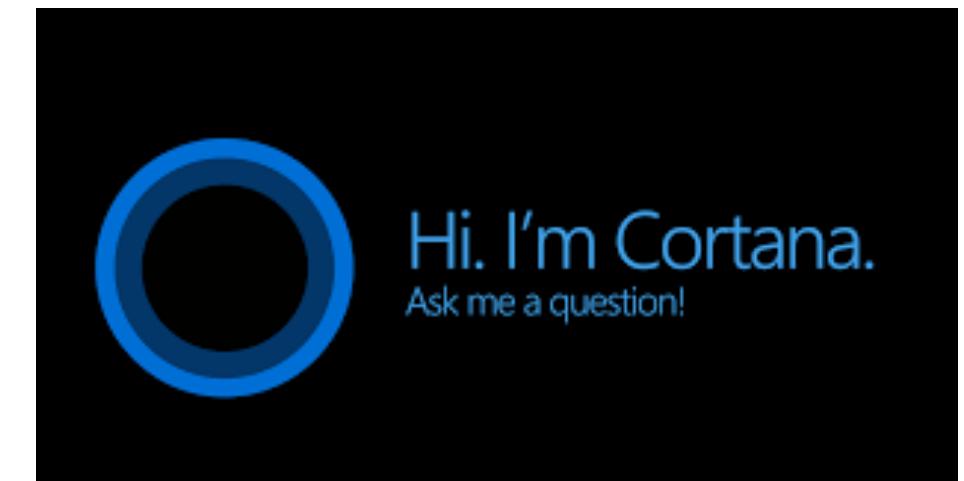
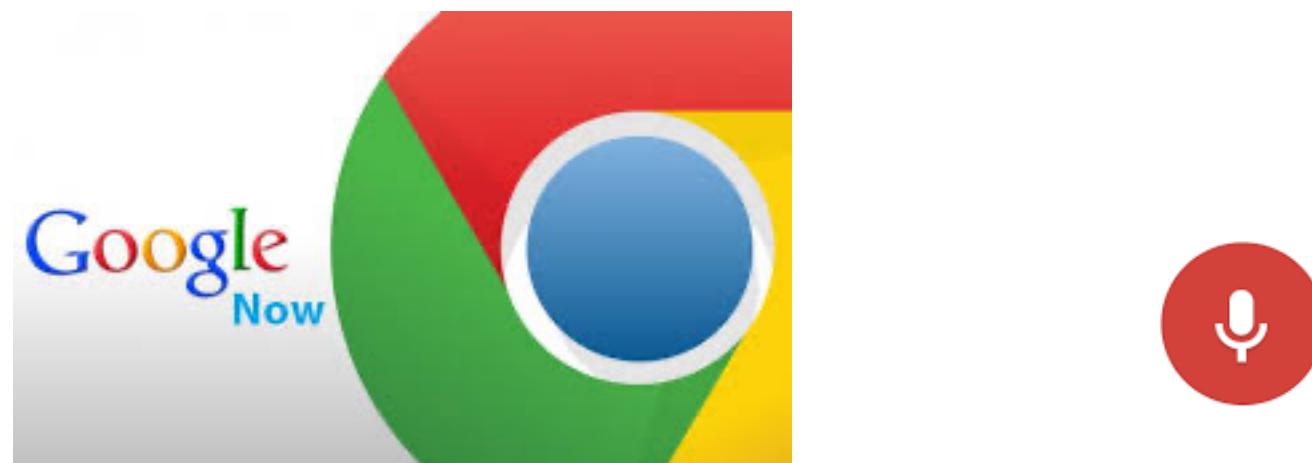


2. Interacting with  
humans via language

# Personal Assistants



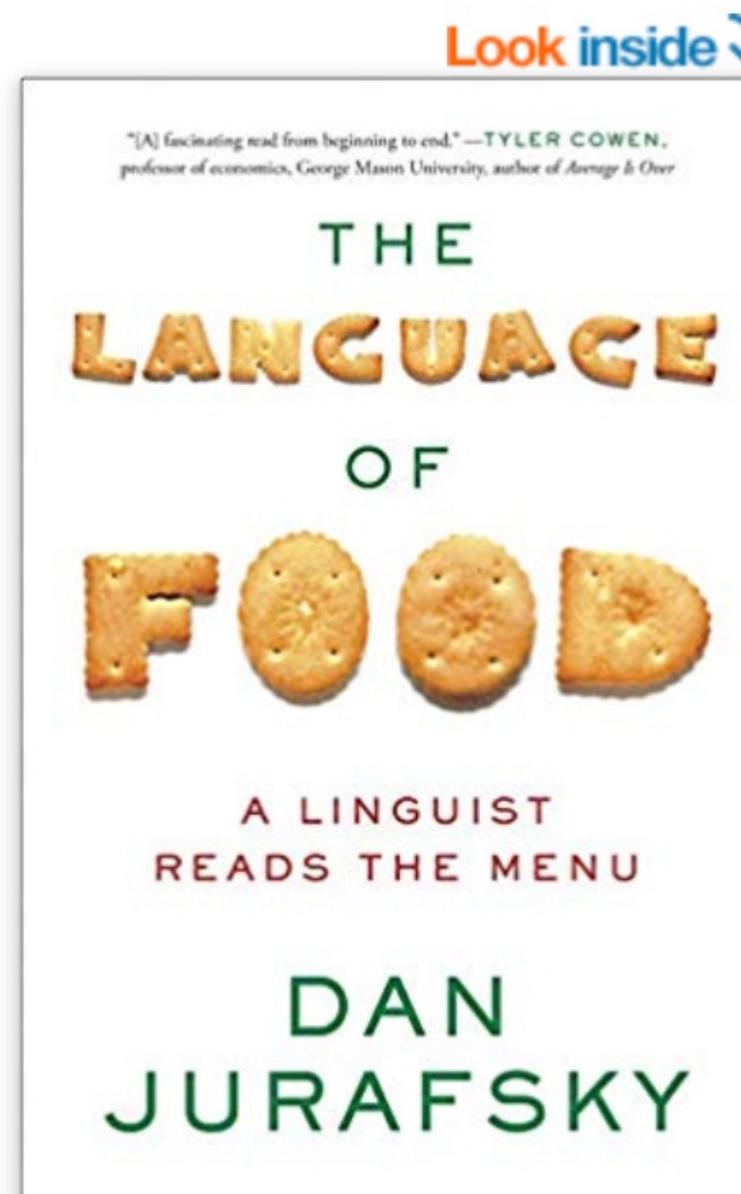
Siri



PA 7 Chatbot!

# Recommendation Engines: The Good

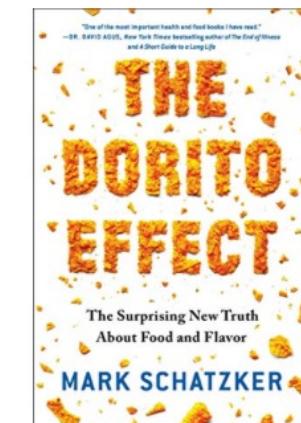
If you bought....



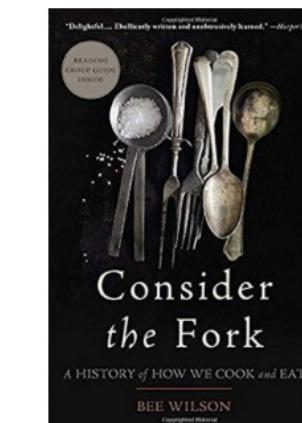
Customers who bought this item also bought



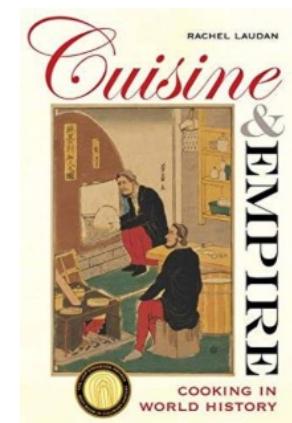
[First Bite: How We Learn to Eat](#)  
by Bee Wilson  
★★★★★ 46  
Paperback  
\$11.37



[The Dorito Effect: The Surprising New Truth About Food and Flavor](#)  
by Mark Schatzker  
★★★★★ 193  
Paperback  
\$9.48



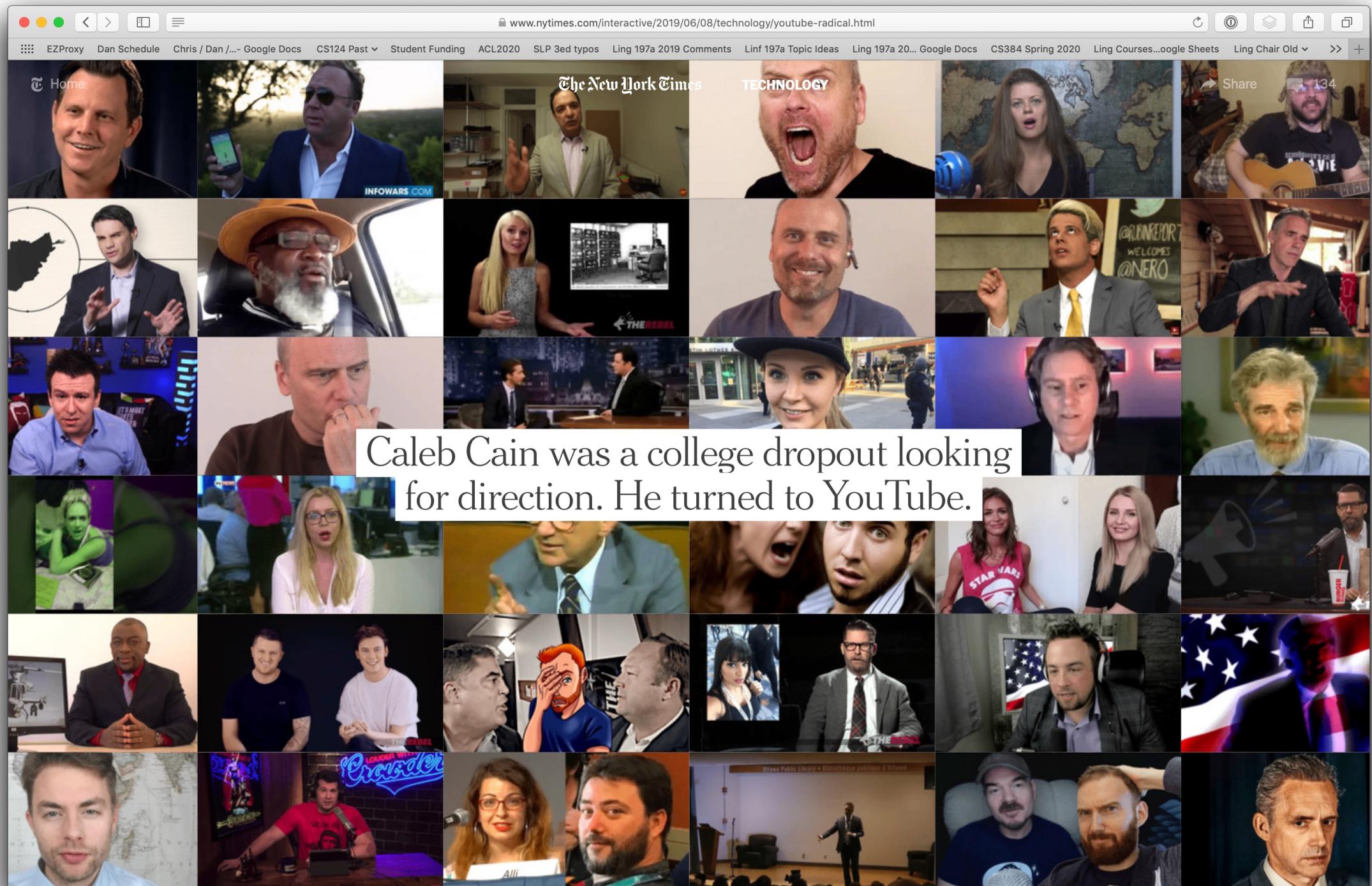
[Consider the Fork: A History of How We Cook and Eat](#)  
by Bee Wilson  
★★★★★ 253  
Paperback  
\$15.65



[Cuisine and Empire: Cooking in World History](#)  
(California Studies in...  
by Rachel Laudan  
★★★★★ 35  
Paperback  
\$16.20

PA 7 and Quiz 8

# The dark side: YouTube Radicalization



# ChatGPT: Large language models

# What makes language interpretation hard?

# Ambiguity

Resolving ambiguity is hard

# Ambiguity

There are at least half a dozen meanings of this sentence:

**The chef made her duck**

Go here and type (and vote for) some definitions

<https://pollev.com/danjurafsky451>



# Ambiguity

## The chef made her duck

The cook cooked waterfowl for a different woman X (person using "she/her" pronouns) to eat

The cook cooked waterfowl belonging to X

The cook cooked waterfowl belonging to the cook

The cook created the (plaster?) waterfowl that X owns

The cook caused X to quickly lower X's head or body

The cook uncovered the true identity of the cook's spy waterfowl

The cook waved their magic wand and turned X into undifferentiated waterfowl

# The chef made her duck

The chef caused X to quickly lower her head or body

**Part of speech:** “duck” can be a Noun or Verb

The chef cooked waterfowl for X (or belonging to X)

**Part of speech:**

“her” is possessive pronoun (“of her”)

“her” is dative pronoun (“for her”)

The chef cooked waterfowl belonging to the chef (vs to X)

**Coreference**

“her” can refer to X or to the Chef

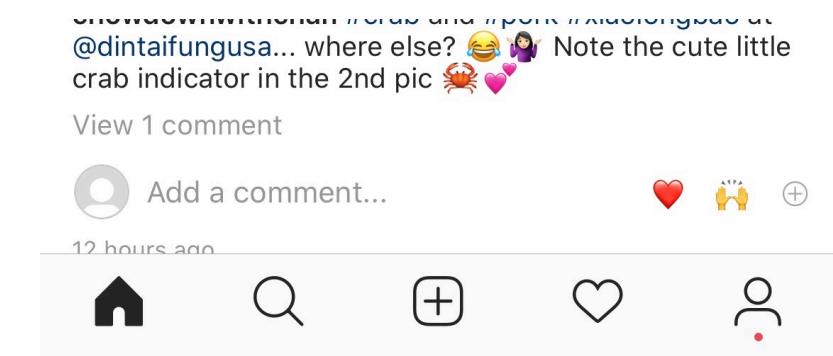
The chef made the (plaster) duck statue X (or the chef) owns

**Word Meaning :** “make” can mean “create” or “cook”

More difficulties:  
Non-standard language,  
emojis, hashtags, names



**chowdownwithchan** #crab and #pork #xiaolongbao at  
@dintaifungusa... where else? 😂🤷‍♀️ Note the cute little  
crab indicator in the 2nd pic 🦀💕



# Models and Tools

Regular Expressions

Edit distance and alignment

Neural word embeddings

Machine Learning classifiers

- Naive Bayes
- Logistic Regression
- Neural Networks

Neural Language Models aka  
"Foundation Models"

Network algorithms

- PageRank

Recommendation  
algorithms

- Collaborative filtering

Linguistic tools

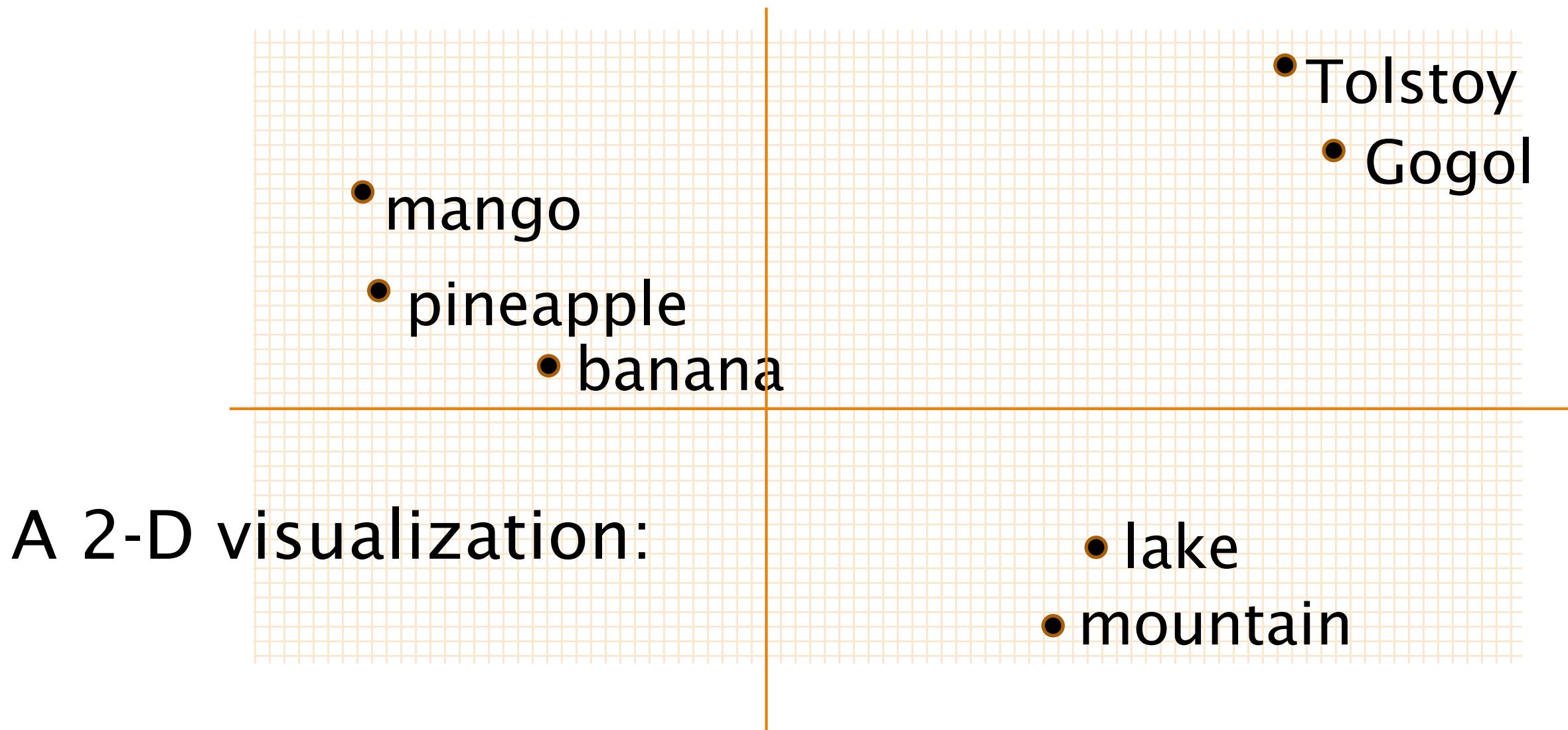
- Sentiment lexicons
- Emotion lexicons

The GUS chatbot architecture  
(Siri, Alexa)

Neural chatbots (ChatGPT)

# Core of modern NLP: Neural "word embeddings"

A word's meaning is a point in (say) 300-dimensional space



# Problem: Embeddings reflect cultural bias!

Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. 2016. "Man is to computer programmer as woman is to homemaker? Debiasing word embeddings." In *NeurIPS 2016*, pp. 4349-4357.

Ask “Paris : France :: Tokyo : x”

- x = Japan

Ask “father : doctor :: mother : x”

- x = nurse

Ask “man : computer programmer :: woman : x”

- x = homemaker

What can we do about this problem? Week 5!

# Logistics: Instructor

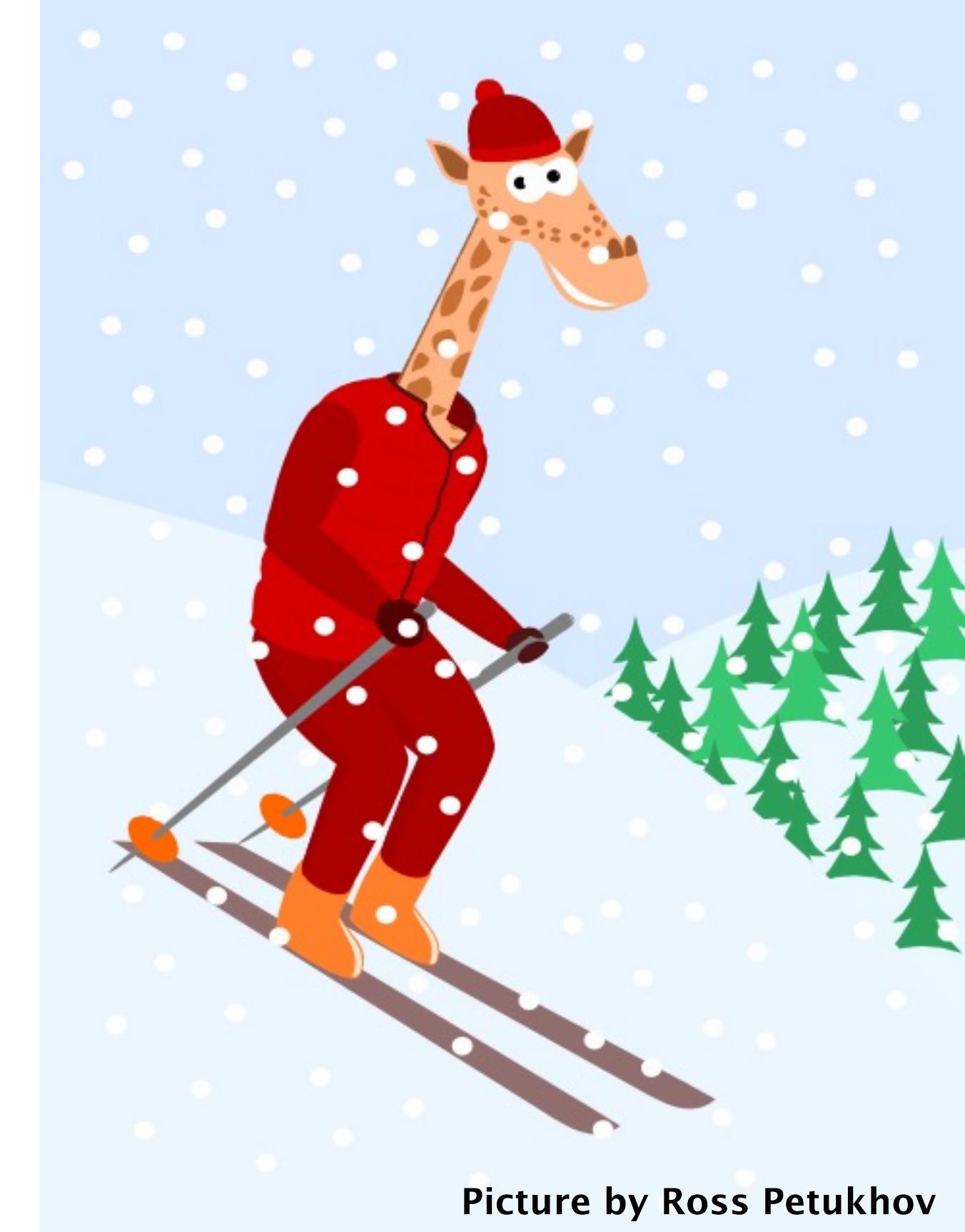
Instructor: Dan Jurafsky (he/him)

Professor in CS and Linguistics

My office hours:

- Tuesday after class 4:30-5:50
- Margaret Jacks Hall 117

\*How to pronounce my name



Picture by Ross Petukhov

# Course Staff



Dan Jurafsky  
Professor



Deveshi Buch  
Head TA



Amelie Byun  
Course Manager



John Cho  
Course Coordinator

# Course Assistants



Niki Agrawal



Jason Ah Chuen



Naomi Eigbe



Amelia Hardy



David Lim



Alexis Lowber



Dwight Moore



John Nguyen



Tolu Oyeniyi



Samantha Silverstein



Anooshree Sengupta



Alisa Wang



Michelle Xu

# Evidence Based Pedagogy!

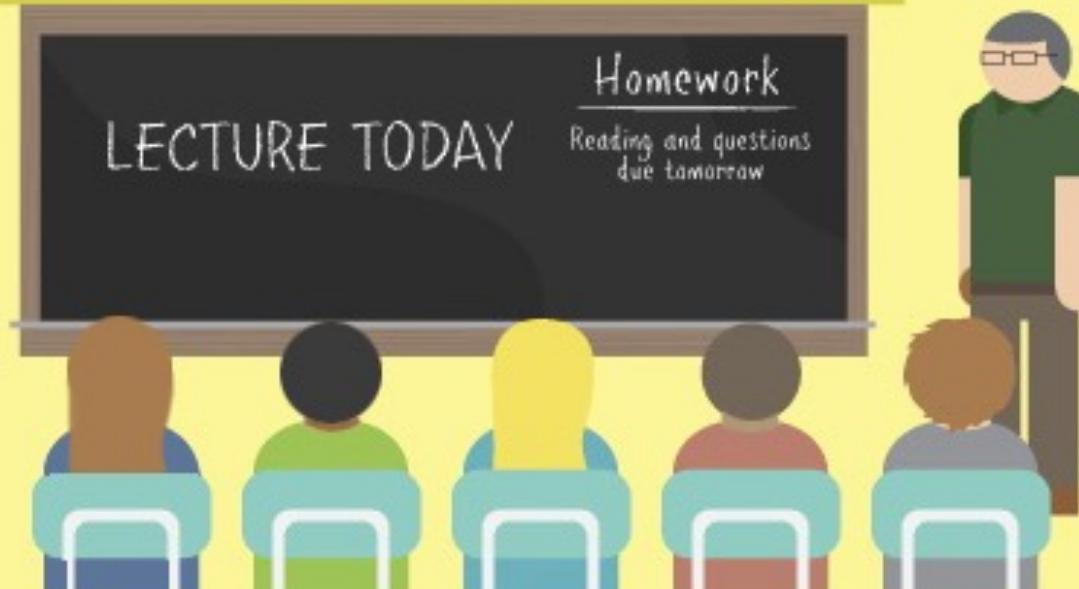
# WHAT IS THE FLIPPED CLASSROOM?

The flipped classroom inverts traditional teaching methods, delivering instruction online outside of class and moving “homework” into the classroom.

## THE INVERSION

### The Traditional Classroom

Teacher's Role: Sage on the Stage



### The Flipped Classroom

Teacher's Role: Guide on the Side



# Why the flipped classroom (1)

**Mastery learning:** Learn until you master

Benjamin Bloom, 1968



# Bloom's mastery learning

Personalized, **goal-driven practice**, driven by **feedback**

1. Watch (and re-watch) lectures at your own pace and learn when it's best for you
2. Videos have embedded miniquizzes. If you get it wrong, it gives you feedback about why you misunderstood.
3. You have **infinite** chances at each weekly Tuesday Quiz, so you can go back to the lecture and retake them.
4. With programming assignments you can see your performance on the training and dev set to see what you might be doing wrong on the test set!

# Why the videos have embedded quizzes: “summative” vs “formative” assessment

## **Summative assessment**

- Final exams/midterms: goal is grading

## **Formative assessment**

- Along the way: goal is for **you** to find out what you don’t know so you can learn

# Why the flipped classroom (2)

Attention span: everyone spaces out during long lectures

- Middendorf and Kalish, 1995, Johnstone and Percival 1976, Burns 1985

“the class started 1:00. The student sitting in front of me took copious notes until 1:20. Then he just nodded off... motionless, with eyes shut for about a minute and a half, pen still poised. Then he awoke and continued his rapid note-taking as if he hadn’t missed a beat.”

Student remembered only the first 15-20 minutes

# Why the flipped classroom (3)

**Active learning:** Be in charge of your learning

- Most important: programming assignments
- Active learning (“constructivism”), learning by doing

**Collaborative learning:** Learn from each other

- Use class time for group problem-solving
- “Small group active learning”
- You must do PA6 in groups

# cs124: Flipped classroom

## **1. Prerecorded video lectures on video:**

- About ~90 minutes/week of video lectures
- Some people watch it speeded up

## **2. Live sessions:**

- 2 required lectures
- 5 required in-class group works (“active learning”)
  - Group Work #1 next Tuesday is **required live**
  - Group works #2, #3, #4, #5 are recommended to be done live, but if you miss class that day, make sure you still do it at home; the material will be on the midterms

# Logistics More Specifically

Online Video Lectures w/embedded non-graded questions (watch before class)

20 pages of reading a week (up to you when to read)

Weekly online Quizzes (Tue of following week)

7 Python homeworks (mostly due Fri of following week)

Two midterms during class time (but we will have alternate times on both days)

- Feb 16
- Mar 16

# Learning Goals

At the end of this course, you will be able to:

# Learning goals

Write efficient regular expressions to solve any kind of text-based extraction task

# Learning goals

Apply the edit distance algorithm to all sorts of text sequence problems

# Learning goals

Build a supervised classifier to solve problems like sentiment classification

# Learning goals

Build a neural network and train it using stochastic gradient descent

# Learning goals

Build a search engine

# Learning goals

Build a recommendation engine

# Learning goals

Build a computational model of word meaning  
(using lexicons and neural word embeddings)

# Learning goals

Build a chatbot

# Learning goals

Understand and implement PageRank and other social network functions

# Learning goals

Understand neural language models and be able to reason both about what they can do, and also their social implications

# This class is the undergrad intro to:

Win 2023: cs224N Natural Language Processing w/Deep Learning (Manning)

Win 2023: cs246 Mining Massive Data Sets (Leskovec)

Spr 2023: cs222U Natural Language Understanding (Potts)

Win 2023: cs224C: NLP for Computational Social Science (Yang)

Spr 2023: cs346 Ethical and Social Issues in NLP (Jurafsky)

Spr 2023: cs 224??: Human-Centered NLP

Aut 2023: cs224W Machine Learning with Graphs

Aut 2023: cs221 Artificial Intelligence

? cs224S Spoken Language Processing

? cs276 Information Retrieval and Web Search

# Should I take 124 or 224X?

## **CS124 is designed for sophomores or juniors**

- It's gentle (I explain everything) and broad (covering many topics, not just NLP but also recommendation engines, IR, social networks)
- Mastery learning, quizzes and programming homeworks
- No research project, but a fun chatbot final homework

## **CS224N and 224U are deeper, more focused, grad courses**

- They assume you are familiar with machine learning and will jump right into optimization and do advanced stuff
- Learning via research: novel research projects as a large component

CS224?? (Human Center NLP) and CS346 (Social and Ethical Issues in NLP) require that you have already had 224N or 224U

CS224C is a different topic: Computational Social Science

(You should of course take all of them!!)

# Syllabus

[cs124.stanford.edu](http://cs124.stanford.edu)

# Where do I find all the programming assignments and quizzes and readings?

Everything is on the webpage `cs124.stanford.edu`

Except the videos which are on canvas!

In other words:

- Lectures slides: webpage
- Group work instructions: webpage
- Tutorial information: webpage
- Programming assignments: webpage (points to git where they live)
- Weekly quizzes: webpage (points to gradescope where they live)
- Practice midterms: webpage
- Midterms: gradescope
- Videos: canvas

# Coming up this week: Thursday

**Optional tutorial** on jupyter notebooks and PA0, getting ready for PA1

Come to class **with your laptops** and we'll go through PA0 together!

This tutorial will be led by amazing head TA Deveshi Buch!!!  
But I and many other CAs will be there!

# Action Items Before Thursdays class!

- 1) Read the syllabus webpage at [cs124.stanford.edu](http://cs124.stanford.edu)
- 2) Look at PAO (you can find it from the webpage)
- 3) Watch Canvas Videos on "PAO Mac Setup" (or  
"PAO Windows Setup")

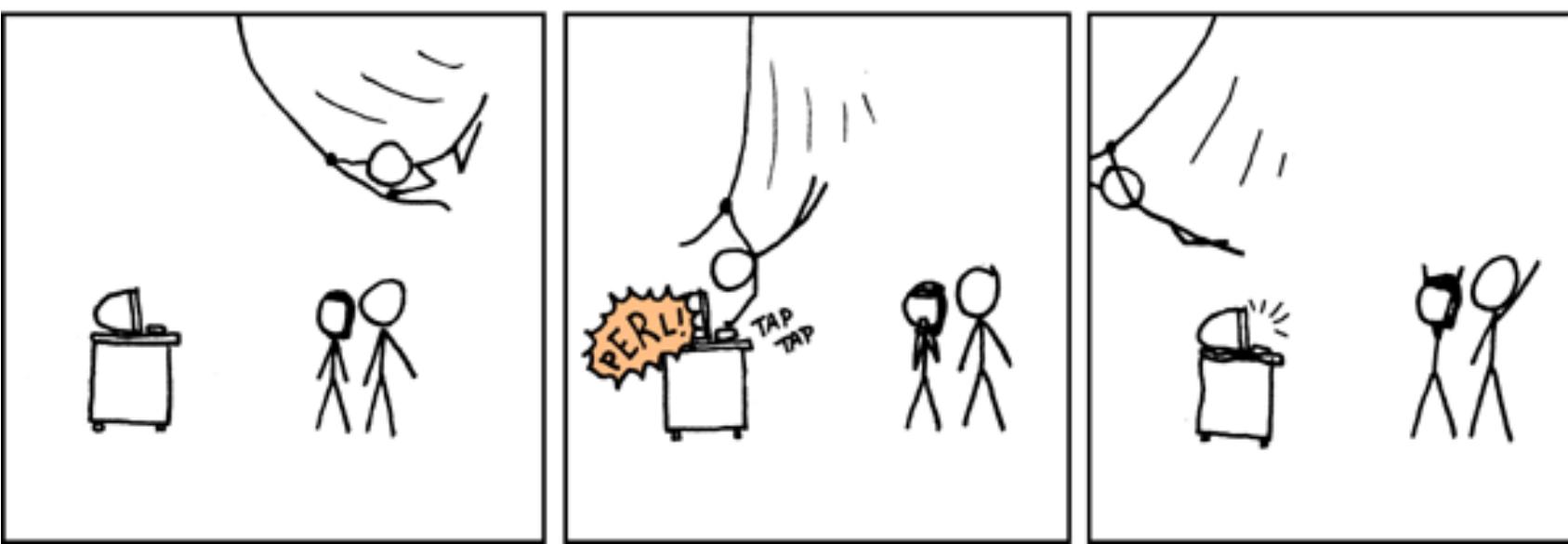
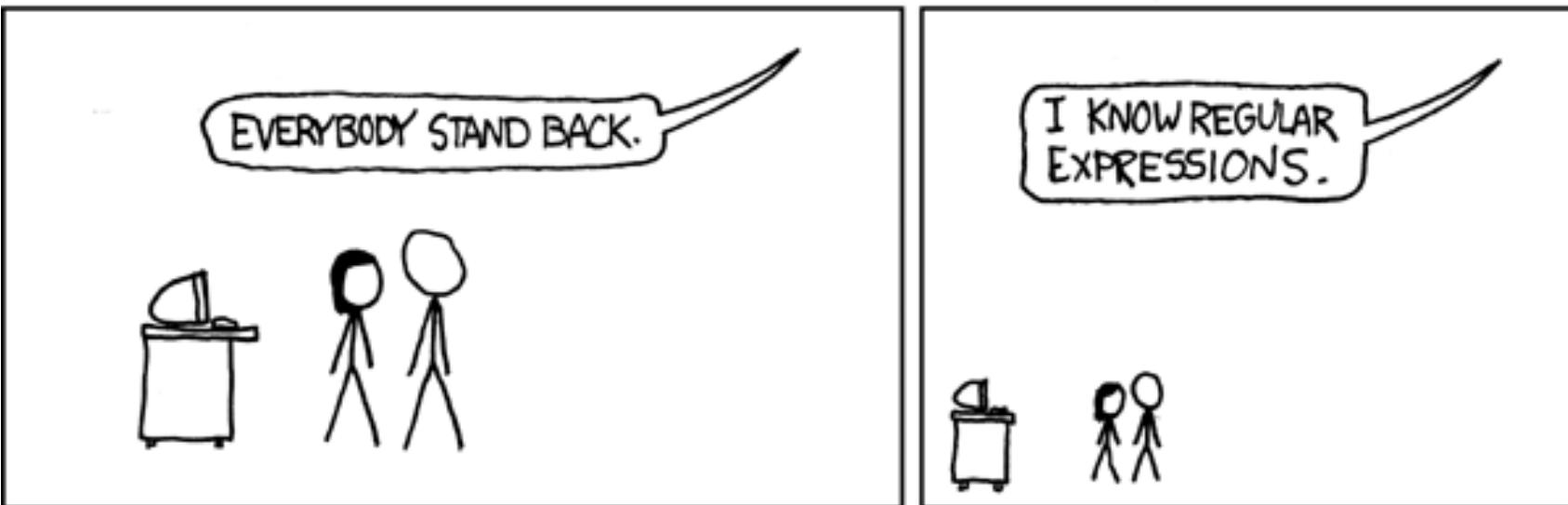
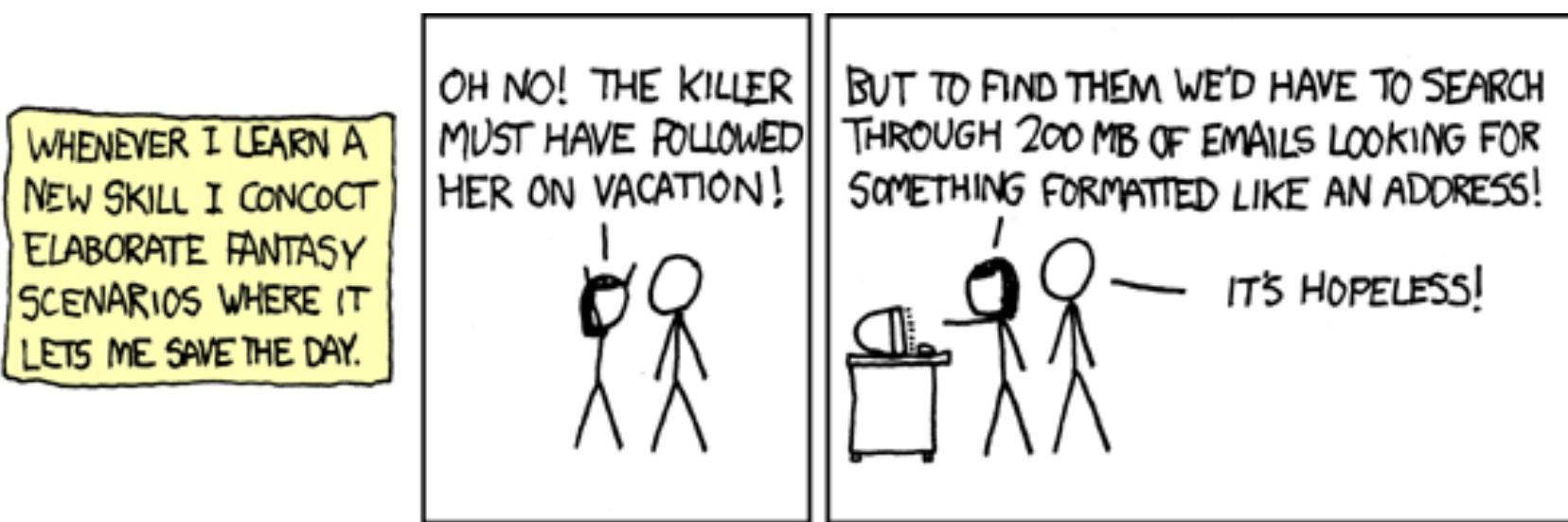
# Coming up next week (Tuesday)

"Unix for poets":

grep

sort

Key UNIX tools for  
dealing with text files and  
regular expressions.



# Action Items Before Tuesday's class!

1) Watch the "week 1" videos on Canvas before class (even earlier, since the quiz is also due Tuesday)

3) Download this file to your laptop

[http://cs124.stanford.edu/nyt\\_200811.txt](http://cs124.stanford.edu/nyt_200811.txt)

4) If you don't know UNIX yet (haven't had cs107):

- For people using a Windows 10 machine, if you don't have Ubuntu on your machine:
  - Watch the first 9 minutes of Bryan's lovely pa0 video about how to download and install Ubuntu:
  - <https://canvas.stanford.edu/courses/144170/modules/items/981067>
- Watch Chris Gregg's excellent UNIX videos here: Logging in, first 7 File System, and first 8 useful commands

<https://web.stanford.edu/class/archive/cs/cs107/cs107.1186/unixref/>

# PA1: Spam Lord!

Write regular expressions to spread evil\* throughout the galaxy!

By extracting email addresses and phone numbers from the web!

jur a fs ky at st anford dot e d u

Goes live Friday 5pm!

\*Just kidding; don't be evil

YOU KNOW HOW SOMETIMES PEOPLE  
PUT A SPACE IN THEIR EMAIL ADDRESS  
TO MAKE IT HARDER TO HARVEST?

YEAH?

THEY HAVE A TOOL THAT  
CAN DELETE THE SPACE!

OH MY GOD.



LESS-DRAMATIC REVELATIONS  
FROM THE CIA HACKING DUMP