# Natural Language Processing

Berkeley
N L P

Dan Klein, John DeNero, GSI: David Gaddy
UC Berkeley

---

# Logistics

---
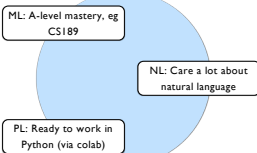
# Logistics

- **Enrollment**
  - Class is currently full
  - Space may open up after P1
  - We'll announce as we go

- **Course expectations**
  - Readings, lectures, ~4 projects
  - No sections, no exams
  - Workload will be high, self-direction
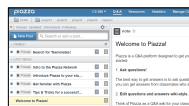  - Patience: class is under construction

- **Requirements**

  ML: A-level mastery, eg CS189

  NL: Care a lot about natural language

  PL: Ready to work in Python (via colab)

---

# Resources and Readings

- **Resources**
  - Webpage (syllabus, readings, slides, links)
  - Piazza (course communication)
  - Gradescope (submission and grades)
  - Compute via Colab notebooks

- **Readings (see webpage)**
  - Individual papers will be linked
  - Optional text: Jurafsky & Martin, 3$^{rd}$ (more NL)
  - Optional text: Eisenstein (more ML)

---

# Projects and Compute

- **Projects**
  - P0: Warm-up
  - P1: Language Models
  - P2: Machine Translation
  - P3: Syntax and Parsing
  - P4: Semantics and Grounding

- **Infrastructure**
  - Python / PyTorch
  - Compute via Colab notebooks
  - Grading via Gradescope

---

# What is NLP?

## Natural Language Processing



**Goal: Deep Understanding**
- Requires context, linguistic structure, meanings…

**Reality: Shallow Matching**
- Requires robustness and scale
- Amazing successes, but fundamental limitations

## NLP History



Neural nets? | Weaver on MT | ALPAC kills MT | Rule-based MT | Search | Neural ASR | Rule-based Semantics | Penn Treebank | Structured ML | Neural MT | Bell Labs ASR | Statistical MT | Neural TTS | Pretraining | Regexps | Grep | CYC

1950   1960   1970   1980   1990   2000   2010   2020

Pre-Compute Era      Symbolic Era      Empirical Era      Scale Era

## Transforming Language

## Speech Systems

- **Automatic Speech Recognition (ASR)**
  - Audio in, text out
  - SOTA: <<1% error for digit strings, 5% conversational speech, still >>20% hard acoustics



Microsoft reaches 'human parity' with new speech recognition system

"Speech Lab"

- **Text to Speech (TTS)**
  - Text in, audio out
  - SOTA: nearly perfect aside from prosody

Google launches more realistic text-to-speech service powered by DeepMind's AI

Speak-N-Spell / Google WaveNet / The Verge

## Machine Translation



- Translate text from one language to another
- Challenges:
  - What's the mapping?  [learning to translate]
  - How to make it efficient?  [fast translation search]
  - Fluency (next class) vs fidelity (later)

Example: Yejin Choi

## Machine Translation



Google Translate 2020

## Spoken Language Translation



Image: Microsoft Skype via Yejin Choi

## Summarization

- Condensing documents
  - Single or multiple docs
  - Extractive or synthetic
  - Aggregative or representative

- Very context-dependent!

- An example of analysis with generation



Image: CNN via Wei Gao

## Understanding Language

## Search, Questions, and Reasoning



## Jeopardy!



Images: Jeopardy Productions

## Question Answering: Watson

## Question Answering: Watson



Jeopardy!
World Champion

US Cities: Its largest airport is named for a World War II hero; its second largest, for a World War II battle.

What is Toronto?????
$36,681
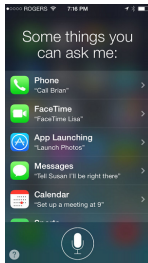
Slide: Yejin Choi

## Language Comprehension?

"The rock was still wet. The animal was glistening, like it was still swimming," recalls Hou Xianguang. Hou discovered the unusual fossil while surveying rocks as a paleontology graduate student in 1984, near the Chinese town of Chengjiang. "My teachers always talked about the Burgess Shale animals. It looked like one of them. My hands began to shake." Hou had indeed found a Naraoia like those from Canada. However, Hou's animal was 15 million years older than its Canadian relatives.

It can be inferred that Hou Xianguang's "hands began to shake", because he was:

(A) afraid that he might lose the fossil
(B) worried about the implications of his finding
(C) concerned that he might not get credit for his work
(D) uncertain about the authenticity of the fossil
(E) excited about the magnitude of his discovery

## Interactive Language

## Example: Virtual Assistants



Some things you can ask me:

Phone
"Call Brian"

FaceTime
"FaceTime Lisa"

App Launching
"Launch Photos"

Messages
"Tell Susan I'll be right there"

Calendar
"Set up a meeting at 9"

- VAs must do
  - Speech recognition
  - Language analysis
  - Dialog processing
  - Text to speech

Image: Wikipedia

## Conversations with Devices?



What's the best movie to see this weekend

That would probably start an argument. But here's a list of highly-regarded movies:

25 MOVIES
NORTH BY NORTHWEST
Released July 17, 1959    100%
THE TREASURE OF THE SIERRA...
Released January 6, 1948    100%

What can I help you with?

You need to start understanding me Siri

I'll make a note of that.

Yeah you better make a note of that

Noted:

Of that

Slide: Yejin Choi

## Social AIs and Chatbots

Session 1
1  [Human] Hello, Xiaoice.
   [Xiaoice] Hello, your profile pic is interesting. Like I've seen it somewhere.
2  [Human] Really?
   [Xiaoice] Is it a comic character?
3  [Human] How do you know that?
   [Xiaoice] Guess.
4  [Human] Are you human?
   [Xiaoice] Well, you'll know it after chatting with me more.
5  [Human] Aren't you an AI chatbot?
   [Xiaoice] So what? I won't be indifferent to you.
6  [Human] What functions do you have?
   [Xiaoice] I have over 200 skills. Just name one you need.
   -----
Session 6
1  [Human] Have you watched Inuyasha? (Note: A Japanese comic)
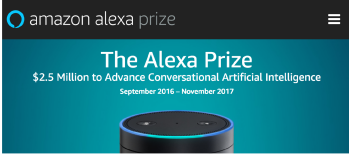   [Xiaoice] Every time Kagome comes back to the modern world, a math test is waiting for her
2  [Human] Bravo, well, do you know how Naraku died?
   [Xiaoice] Naraku is not dead and Shikon Jewel does not disappear.
3  [Human] So amazing!
   [Xiaoice] Why amazed?
4  [Human] Awesome! Surprised that you also like comic
   -----



Microsoft's Xiaolce

Source: Microsoft

## Chatbot Competitions!
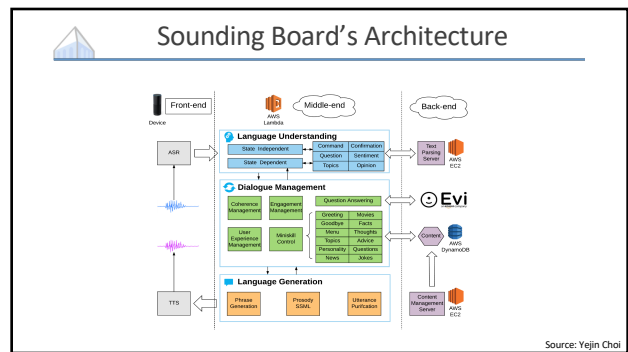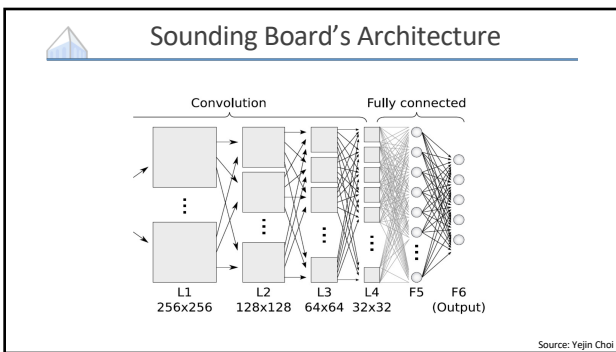


amazon alexa prize

**The Alexa Prize**
$2.5 Million to Advance Conversational Artificial Intelligence
September 2016 – November 2017

- Alexa Prize competition to build chatbots that keep users engaged
  - Winner in 2017: UW's Sounding Board (Fang, Cheng, Holtzman, Ostendorf, Sap, Clark, Choi)
  - Winner in 2018: UC Davis's Gunrock (Zhou Yu et al)
- Compare to the Turing test (eg Loebner Prize) where the goal is to fool people

## SoundingBoard Example

| Turn | User | Sounding Board |
|---|---|---|
| 1 | Let's chat | Hi! This is an Alexa Prize socialbot. How's your day? |
| 2 | I'm five. | Glad to hear it! We could talk about robots, batman, or superman. Say "next", to chat about other things. |
| 3 | Superman. | It looks like you wanna chat about superman. I was high up in the cloud when I realized: If people don't recognize Clark Kent as Superman because of his glasses, does his eye doctor know his identity? |
| 4 | I guess so. | Did you know that Henry Cavill almost missed the call for the role of Superman cause he was playing "World of Warcraft" |
| 5 | Really, I didn't know that. | Weird, right? Speaking of superman, Did you see: Superman? It was released in 1997. It's a comedy. |
| 6 | Yes, it was hilarious. | I'm glad you feel this is hilarious. Which part about this movie? |
| 7 | The part when he met lewis leah. | Interesting. meccartin, and raffi co-directed this film. The movie has a 6.3 out of 10 on IMDB. which seems pretty good! do you like the movie's director? |
| 8 | Yeah, he's cool. | Cool! By the way, do you wanna know something interesting about raffi? |
| 9 | Okay, go ahead. | Super! My friend told me that Raffi is a worldwide children's rights activist and has fought against exploitative advertising to children. |

Negotiation
Thought
Movie
Fact

Source: Mari Ostendorf

## Sounding Board's Architecture



Convolution          Fully connected

L1          L2          L3        L4    F5    F6
256x256   128x128    64x64    32x32        (Output)

Source: Yejin Choi

## Sounding Board's Architecture



Source: Yejin Choi

## Related Areas

## What is Nearby NLP?

- **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 Processing**
  - Mapping audio signals to text
  - Traditionally separate from NLP, converging

## Example: NLP Meets CL

/kentrum/ (la)

u → o / some context
m → / some context

/t͡ʃentro/ (yl)

/sentro/ (ib)    /t͡ʃentro/ (it)

/sentro/ (es)    /semtru/ (pt)

| Gloss | Latin | Italian | Spanish | Portuguese |
|-------|-------|---------|---------|------------|
| Word/verb | verbum | verbo | verbo | verbu |
| Center | centrum | centro | centro | centro |

- Example: Language change, reconstructing ancient forms, phylogenies
  … just one example of the kinds of linguistic models we can build

---

## Why is Language Hard?

---

## Problem: Ambiguity

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

- Why are these funny?

---

## What Do We Need to Understand Language?

---

## We Need Representation: Linguistic Structure

```
N     N    V   N              NP
N     V   ADJ  N                  PP
Teacher Strikes Idle Kids              NP

                      N  P
body/      body/    Ban on Nude Dancing on Governor's Desk
position   weapon    N
Iraqi Head Seeks Arms        PP       PP
                              NP
```

▸ Syntactic and semantic ambiguities: parsing needed to resolve these, but need context to figure out which parse is correct

Slide: Greg Durrett

---

## Example: Syntactic Analysis

Hurricane Emily howled toward Mexico 's Caribbean coast on Sunday
packing 135 mph winds and torrential rain and causing panic in Cancun,
where frightened tourists squeezed into musty shelters .

Accuracy: 95+

## We Need Data



## We Need Lots of Data: MT

| SOURCE | Cela constituerait une solution transitoire qui permettrait de conduire à terme à une charte à valeur contraignante. |
|---|---|
| HUMAN | That would be an interim solution which would make it possible to work towards a binding charter in the long term . |
| 1x DATA | [this] [constituerait] [assistance] [transitoire] [who] [permettrait] [licences] [to] [terme] [to] [a] [charter] [to] [value] [contraignante] [.] |
| 10x DATA | [it] [would] [a solution] [transitional] [which] [would] [of] [lead] [to] [term] [to a] [charter] [to] [value] [binding] [.] |
| 100x DATA | [this] [would be] [a transitional solution] [which would] [lead to] [a charter] [legally binding] [.] |
| 1000x DATA | [that would be] [a transitional solution] [which would] [eventually lead to] [a binding charter] [.] |

## We Need Models: Data Alone Isn't Enough!

**CLASSIC SOUPS**                    Sm.    Lg.

| | | |
|---|---|---|
| 57. | House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot) | 1.50  2.75 |
| 58. | Chicken Rice Soup | 1.85  3.25 |
| 59. | Chicken Noodle Soup | 1.85  3.25 |
| 60. | Cantonese Wonton Soup | 1.50  2.75 |
| 61. | Tomato Clear Egg Drop Soup | 1.65  2.95 |
| 62. | Regular Wonton Soup | 1.10  2.10 |
| 63. | Hot & Sour Soup | 1.10  2.10 |
| 64. | Egg Drop Soup | 1.10  2.10 |
| 65. | Egg Drop Wonton Mix | 1.10  2.10 |
| 66. | Tofu Vegetable Soup | NA  3.50 |
| 67. | Chicken Corn Cream Soup | NA  3.50 |
| 68. | Crab Meat Corn Cream Soup | NA  3.50 |
| 69. | Seafood Soup | NA  3.50 |

*Example from Adam Lopez*

## We Need World Knowledge

‣ World knowledge: have access to information beyond the training data



DOJ greenlights Disney - Fox merger

Department of Justice

metaphor; "approves"

‣ What is a green light? How do we understand what "green lighting" does?

‣ Need commonsense knowledge

Slide: Greg Durrett

## Data and Knowledge

▪ Classic knowledge representation worries: How will a machine ever know that...
  ▪ Ice is frozen water?
  ▪ Beige looks like this:
  ▪ Chairs are solid?

▪ Answers:
  ▪ 1980: write it all down
  ▪ 2000: get by without it
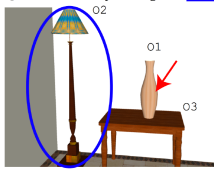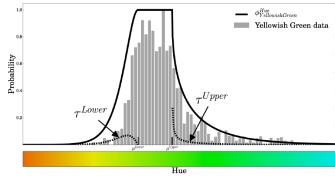  ▪ 2020: learn it from data

## Learning Latent Syntax

### Personal Pronouns (PRP)

| | | | |
|---|---|---|---|
| PRP-1 | it | them | him |
| PRP-2 | it | he | they |
| PRP-3 | It | He | I |

### Proper Nouns (NNP)

| | | | |
|---|---|---|---|
| NNP-14 | Oct. | Nov. | Sept. |
| NNP-12 | John | Robert | James |
| NNP-2 | J. | E. | L. |
| NNP-1 | Bush | Noriega | Peters |
| NNP-15 | New | San | Wall |
| NNP-3 | York | Francisco | Street |

## We Need Grounding

Grounding: linking linguistic concepts to non-linguistic ones

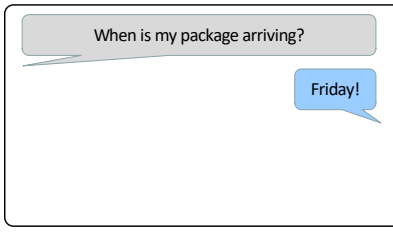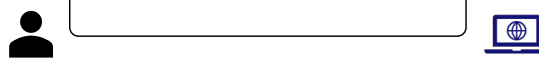**Question:** What object is **right of** O2 ?

Golland et al. (2010)    McMahan and Stone (2015)

Slide: Greg Durrett

## Example: Grounded Dialog
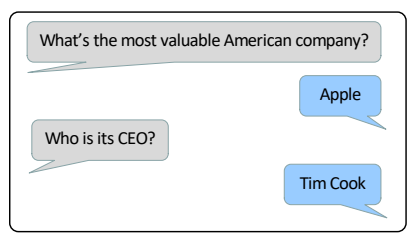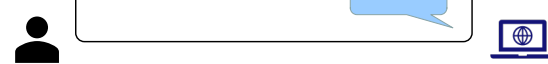
When is my package arriving?

Friday!

## Example: Grounded Dialog

What's the most valuable American company?

Apple

Who is its CEO?

Tim Cook

## Why is Language Hard?

- We Need:
  - Representations
  - Models
  - Data
  - Machine Learning
  - Scale
  - Efficient Algorithms
  - Grounding

- … and often we need all these things at the same time

## What is this Class?

## What is this Class?

- Three aspects to the course:
  - Linguistic Issues
    - What are the range of language phenomena?
    - What are the knowledge sources that let us disambiguate?
    - What representations are appropriate?
    - How do you know what to model and what not to model?
  - Modeling Methods
    - Increasingly sophisticated model structures
    - Learning and parameter estimation
    - Efficient inference: dynamic programming, search, sampling
  - Engineering Methods
    - Issues of scale
    - Where the theory breaks down (and what to do about it)
- We'll focus on what makes the problems hard, and what works in practice…

## Class Requirements and Goals

- Class requirements
  - Uses a variety of skills / knowledge:
    - Probability and statistics, graphical models (parts of cs281a)
    - Basic linguistics background (ling100)
    - Strong coding skills (Python, ML libraries)
  - Most people are probably missing one of the above
  - You will often have to work on your own to fill the gaps

- Class goals
  - Learn the issues and techniques of modern NLP
  - Build realistic NLP tools
  - Be able to read current research papers in the field
  - See where the holes in the field still are!

- This semester: new projects, new topics, lots under construction!