

Natural Language Processing



Dan Klein, GSIs: Kevin Lin, Eric Wallace
UC Berkeley

Logistics



Logistics

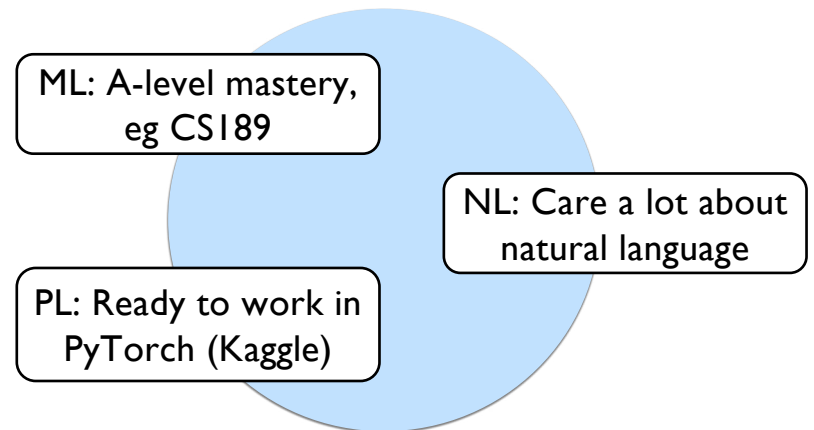
■ Enrollment

- Class is full but we hope to fit waitlist
- We will process enrollments after P1
- Material access doesn't require enrollment

■ Course expectations

- Readings, lectures, 5 projects
- No sections, no exams
- Engagement with the course
- Workload will be high, autonomy required
- Patience: class is under construction!

■ Requirements



■ COVID/Attendance Policies

- Webcasts posted (Berkeley-local)
- Expanded late day policy (14 day)
- Please do not come to lecture sick



Resources and Readings

■ Resources

- Webpage (syllabus, readings, slides, links)
- Ed Discussion (course communication)
- Gradescope (submission and grades)
- Compute via Kaggle notebooks
- Berkeley-internal recordings for main lectures

CS 288	Schedule
Home	Jan 23: Intro (HW1 Released)
Course Info	Jan 25: Language Modeling
Staff	Jan 30: Language Modeling
	Feb 1: Machine Translation

■ Readings (see webpage)

- Individual papers will be linked
- Optional text: Jurafsky & Martin, 3rd (more NL)
- Optional text: Eisenstein (more ML)

ed CS 288 – Ed Discussion

New Thread

COURSES +
CS 288

CATEGORIES
■ General
■ Lectures
■ Assignments

Search

Filter ▾

No threads
Be the first to create a thread!



Projects and Infrastructure

■ Projects

- P1: Language Models
- P2: Machine Translation
- P3: Syntax and Parsing
- P4: Single-task NLP with LLMs
- P5: Many-task NLP with LLMs

■ Infrastructure

- Python / PyTorch
- Compute via Kaggle notebooks
- Grading via Gradescope

The screenshot shows the Kaggle Project 1 interface. On the left is a sidebar with the Kaggle logo and navigation links: Create, Home, Competitions, Datasets, Code (selected), Discussions, Learn, More, Your Work, and Recently Edited (notebook797e85a997). The main area is titled 'Project 1 Draft saved' and has a menu bar with File, Edit, View, Run, Add-ons, and Help. Below the menu bar are icons for adding, deleting, undo, redo, and running cells, along with a 'Run All' button and a 'Markdown' dropdown. The notebook content includes a paragraph about regularization techniques, a section for a write-up describing extensions and modifications, and a list of four requirements for the write-up. Below the text is a light gray box with the instruction: 'Run the cell below in order to train your improved LSTM and evaluate it.' At the bottom are buttons for '+ Code' and '+ Markdown'. A code cell is shown with the following content:

```
[ ]:
## Feel free to copy your original LSTM solution down here to modify for your report if you'd like.
# YOUR CODE [optionally] HERE
##

#lstm_model = LSTMModel()
#lstm_model.train()

print('lstm validation perplexity:', lstm_model.dataset_perplexity(validation_dataset))
save_truncated_distribution(lstm_model, 'lstm_predictions.npy', short=False)
```

What is NLP?



Natural Language Processing



NLP

Goal: Deep Understanding

- Requires knowledge, context, and grounding
- Traditionally out of reach, but huge recent advances



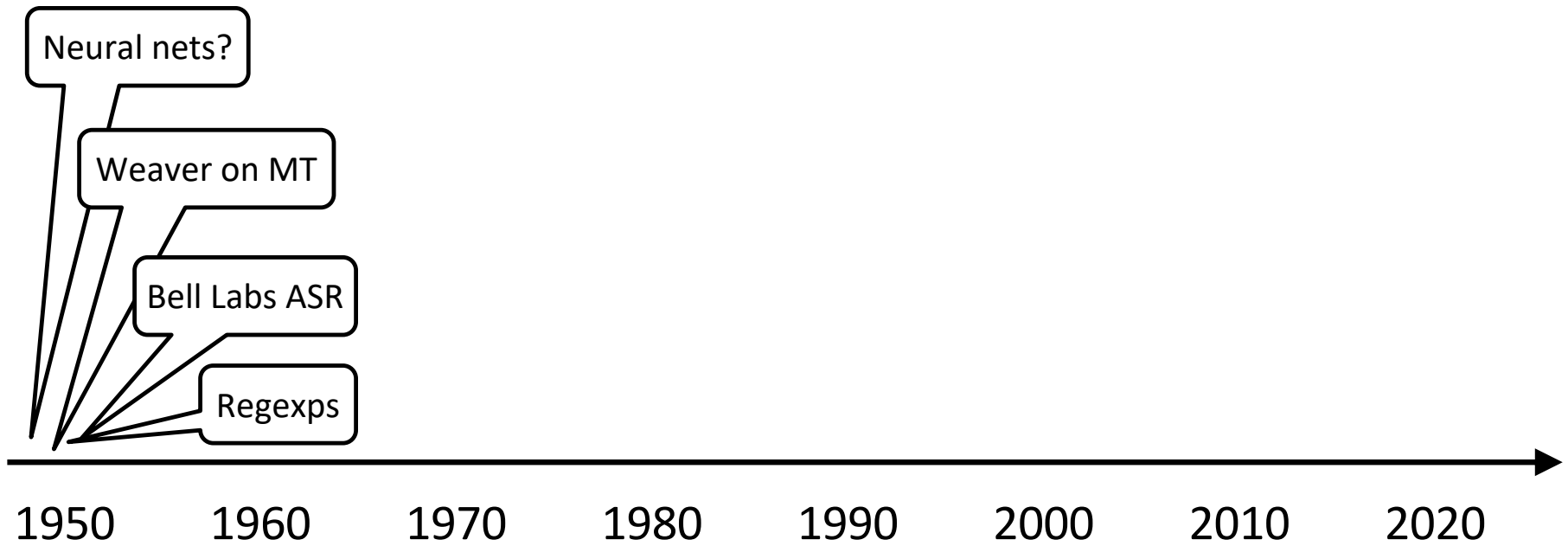
grep

Classic Systems: Shallow Matching

- Focus on robustness, scale, and predictability
- Very useful, but has had fundamental limitations



NLP History



Pre-Compute Era

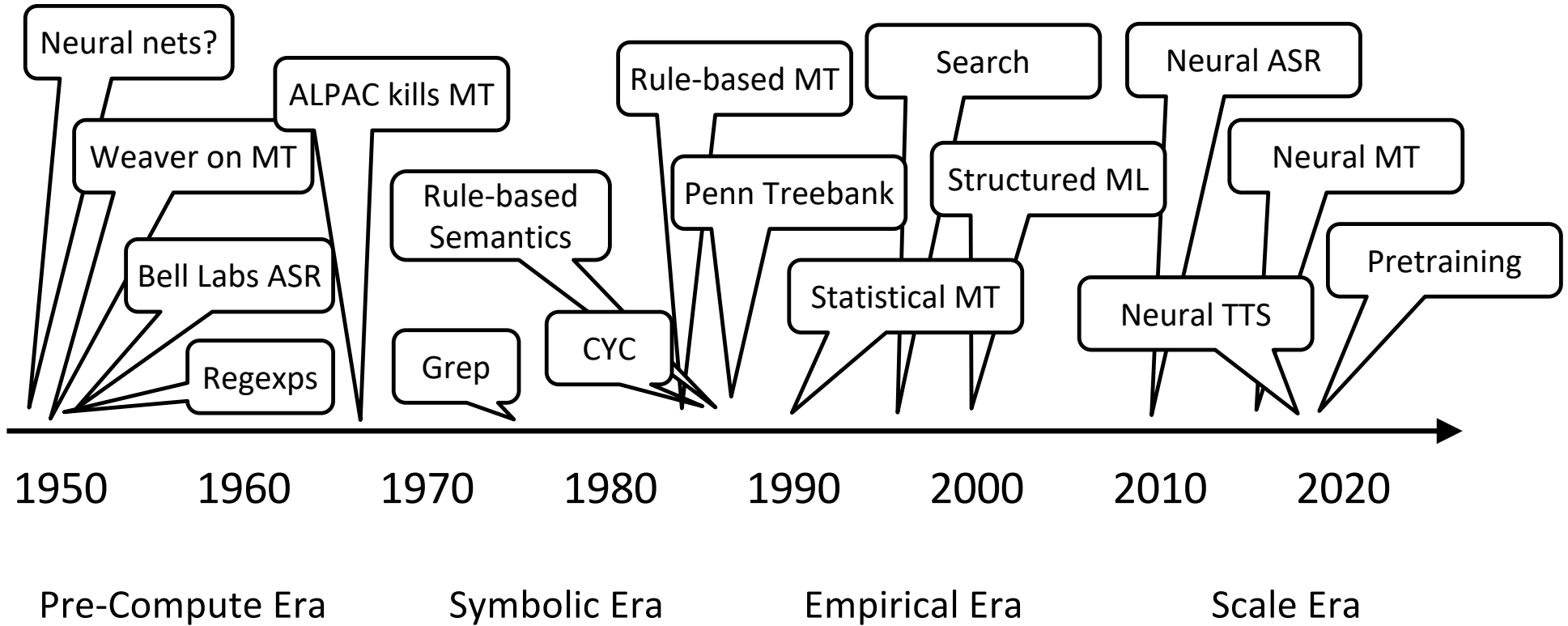


NLP History





NLP History



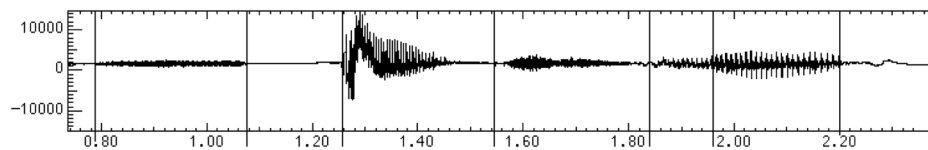
Transforming Language



Speech Systems

■ Automatic Speech Recognition (ASR)

- Audio in, text out
- SOTA: $\ll 1\%$ error for digit strings, 5% conversational speech, still $\gg 20\%$ hard acoustics



“Speech Lab”

MICROSOFT TECH ARTIFICIAL INTELLIGENCE

Microsoft reaches 'human parity' with new speech recognition system

■ Text to Speech (TTS)

- Text in, audio out
- SOTA: nearly perfect aside from prosody



GOOGLE TECH ARTIFICIAL INTELLIGENCE

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



Speak-N-Spell / Google WaveNet / The Verge



Machine Translation

EN CE MOMENT Impôts Kenya Syrie Pakistan Emploi Scandale Prism

Impôt sur le revenu : combien paierez-vous en 2014 ?

DÉCOUVREZ TOUS LES SERVICES ABONNÉS
S'abonner au Monde à partir de 1 €



APPEL À TÉMOIGNAGES
Adhérent(e) d'Europe Ecology-Les Verts, partagez-vous le constat sévère de Noël Mamère sur EELV ?
Partagez votre expérience

En continu

- 19:53 Budget : des dépenses incompressibles
- 19:36 Entendu à la "fashion week" de Paris
- 19:19 Commande géante à Airbus
- 19:04 Plainte contre "Valeurs actuelles"
- 19:01 Venezuela : 17 personnes arrêtées
- 18:59 Vidberg : le budget nouveau est annoncé
- 18:50 La "mission noble" de la NSA
- 18:38 Roms : Joutes entre Bruxelles et Paris

Sélectionnez votre revenu et votre situation familiale pour découvrir si vous bénéficiez de la pause fiscale.

- Comment le budget pour 2014 est-il réparti ? VISUEL INTERACTIF
- Un budget 2014 soumis aux critiques
- Budget : ces dépenses qu'aucun gouvernement ne peut diminuer
- Budget 2014 : le détail des économies VISUEL INTERACTIF

AT THIS MOMENT Taxes Kenya Syrie Pakistan Use Prism scandal

Income tax: how much do you pay in 2014?

DÉCOUVREZ TOUS LES SERVICES ABONNÉS
S'abonner au Monde à partir de 1 €



CALL FOR EVIDENCE
Member (s) of Europe Ecology-Greens, do you share the finding of severe Christmas Mamère EELV?
Share your experience

Continuous

- Expenses
- "Week" in Paris
- Actual Values"
- able arrested
- udget came
- of the NSA
- ween Brussels €

Select your income

- How is the budget allocated?
- A 2014 budget criticism



- 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

CINÉMA • MÉDIAS

Partage 

Disney décide de changer le nom du légendaire studio de cinéma Fox

Le nouveau propriétaire va faire disparaître la mention « Fox » à la 20th Century pour que le nom du studio ne soit pas associé à l'image ultra-conservatrice de la chaîne de télévision Fox News.

Le Monde avec AFP • Publié le 18 janvier 2020 à 06h46

 Lecture 1 min.

Révolution dans le monde d'Hollywood. Le célèbre studio de cinéma 20th Century Fox va voir une partie de son nom modifiée et devenir 20th Century Studios, décision prise par son nouveau propriétaire, Disney.

Selon plusieurs médias américains, Disney aurait décidé de modifier le nom du studio pour qu'il ne soit plus associé à la grande chaîne Fox mais aussi et surtout à Fox News, la chaîne d'information en continu.

CINEMA • MEDIA

Sharing 

Disney decides to change the name of the legendary Fox film studio

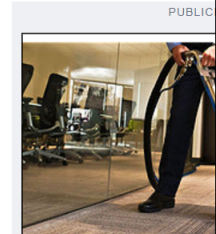
The new owner will remove the mention "Fox" in the 20th Century so that the name of the studio is not associated with the ultra-conservative image of the television channel Fox News.

Le Monde avec AFP • Posted on January 18, 2020 at 06:46

 Reading time 1 min.

Revolution in the Hollywood world. The famous film studio 20th Century Fox will see part of its name changed and become 20th Century Studios, decision taken by its new owner, Disney.

According to several American media, Disney has decided to change the name of the studio so that it is no longer associated with the big chain Fox but also and especially with Fox News, the news channel.





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

The image shows a CNN news article titled "Official: Egypt balloon explosion probe can take 2 weeks". The article is dated February 27, 2013, and is updated at 16:14 GMT (00:14 HKT). The article is by Adam Makary, Saad Abedine, and Mariano Castillo. The article is categorized under "Africa" and "U.S.". The article is written in English, but there is a link to read a version in Arabic.

STORY HIGHLIGHTS

- **NEW:** No foul play is suspected, official says
- The Tuesday accident was the world's deadliest hot air balloon accident in at least 20 years
- Officials: Passengers in the balloon included 19 foreign tourists
- Luxor province bans all hot air balloon flights until further notice

Read a version of this story in Arabic.

Cairo (CNN) -- An official investigation into the cause of a balloon accident that killed 19 people in Egypt could take two weeks, the governor of Luxor province said Wednesday.

The Tuesday accident was the world's deadliest hot air balloon accident in at least 20 years.

Preliminary investigations confirmed no foul play was involved when gas canisters aboard the balloon exploded, causing it to plummet about 1,000 feet (300 meters) to the ground, Gov. Ezzat Saad said.

CNN iReport: After tragedy, vacationers recall glorious balloon rides in Egypt

Passengers in the balloon included 19 foreign tourists: nine from Hong Kong, four from Japan, three from Britain, two from France and one from Hungary, officials said.

How safe is hot air ballooning?

An Egyptian pilot and another Egyptian were also aboard, Luxor province spokesman Badawi al-Masri said.

Balloon rides offering panoramic aerial views of the Nile River and the ancient temples of Karnak and Hatshepsut are a popular tourist attraction in Luxor, about a nine-hour drive southeast of Cairo.

"You can see Valley of the Kings in the background bordered by farmland," Pauline Liang of Vancouver, Canada, told CNN's iReport last year. "Below were banana farms, and behind us was the city of Luxor. There was a great contrast between desert landscape, lush farmland and urban development."

Tuesday's crash prompted the governor to ban all hot air balloon flights until further notice.

Twitter

- CNN: official investigation into yesterday air balloon accident in Luxor could take 2 weeks
- Governor bans all hot air balloon flights until further notice
- Foul play not suspected in fatal balloon accident
- Official: Egypt balloon explosion probe can take 2 weeks
- Egypt balloon explosion
- An official investigation into the cause of a balloon accident that killed 19 people in Egypt could take two w...
- Egypt: Balloon probe could take weeks

Click to play

Tourists killed in hot air balloon blast

2009 balloon crash survivor speaks

How to stay safe in a hot air balloon

Image: CNN via Wei Gao

Understanding Language



Search, Questions, and Reasoning

Google

Web Images Groups News Froogle Local more »

any US states' capitals are also their largest cities? Search

Google

Web

Tip: Have a question?

Web definitions

Cheyenne wordnet

Scholarly articles

The Culture Capital

MSN Encarta

Cheyenne (city) ... In 1869 the

encarta.msn.com/encyclopedia_761556413/Cheyenne_(city_Wyoming)

Your search - **How many US states' capitals are also their largest cities?** - did not match any documents.

Suggestions:

- Make sure all words are spelled correctly.
- Try different keywords.
- Try more general keywords.
- Try fewer keywords.

Google Home - - Business Solutions

How many US states' capitals are also their largest cities?

About 39,100,000 results (0.84 seconds)

17 States

The 17 **States** Whose Capital is also its **Largest City**. For the majority of **US States**, the **state capital** is NOT the **largest city** in the state. In some cases, **capitals** were established differently from the **largest city**.

quizriot.com › geography › states › capital cities › largest cities › the 17 states whose capital is also the largest city

The 17 States Whose Capital is also its Largest City

How many US states' capitals are their second largest cities?

About 23,500,000 results (1.01 seconds)

Did you mean: How many US states' capitals are **there** second largest cities?

www.factmonster.com › us › states › state-capitals-and-largest-cities

State Capitals and Largest Cities - Fact Monster

State **Capitals** and **Largest Cities**. The following table lists the **capital** and **largest city** of every state in the **United States**. Featured Videos. ADVERTISEMENT.

Missing: second | Must include: second



Jeopardy!



Images: Jeopardy Productions



The screenshot shows a web browser window displaying the Wiktionary page for the phrase "a camel is a horse designed by a committee". The page layout includes a sidebar on the left with navigation links, a main content area with the title and a discussion forum, and a search bar at the top right. The title "a camel is a horse designed by a committee" is highlighted in a blue box. The discussion forum section shows a post by Ruben P. Mendez on April 16, 2004, asking for the origin of the maxim. Below the post, there are two links: "Re: A camel is a horse designed by committee" by SR dated 16/April/04, and "Re: A camel is a horse designed by committee" by Henry dated 18/April/04.

Wiktionary
a multi-lingual free encyclopedia

Log in / create account

Entry Discussion Read Edit History Search

a camel is a horse designed by a committee

Contents [hide]

1 English

1.1 Alternative forms

1.2 Proverb

The Phrase Finder

e > Discussion Forum

Google Custom Search Search

A camel is a horse designed by committee

Posted by Ruben P. Mendez on April 16, 2004

Does anyone know the origin of this maxim? I heard it way back at the United Nations, which is chockfull of committees. It may have originated there, but I'd like an authoritative explanation. Thanks

- Re: A camel is a horse designed by committee SR 16/April/04
 - Re: A camel is a horse designed by committee Henry 18/April/04

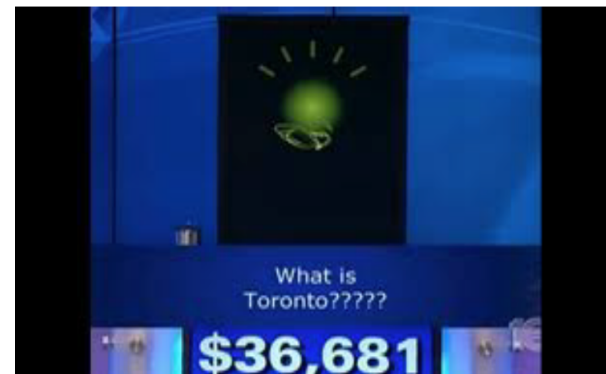


Question Answering: Watson



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.





Language Comprehension?

Opera refers to a dramatic art form, originating in Europe, in which the emotional content is conveyed to the audience as much through music, both vocal and instrumental, as it is through the lyrics. By contrast, in musical theater an actor's dramatic performance is primary, and the music plays a lesser role. The drama in opera is presented using the primary elements of theater such as scenery, costumes, and acting. However, the words of the opera, or libretto, are sung rather than spoken. The singers are accompanied by a musical ensemble ranging from a small instrumental ensemble to a full symphonic orchestra.

1. It is pointed out in the reading that opera -----.

- ☐ A) has developed under the influence of musical theater
- ☐ B) is a drama sung with the accompaniment of an orchestra
- ☐ C) is not a high-budget production
- ☐ D) is often performed in Europe
- ☐ E) is the most complex of all the performing arts

[From GrammarBank]



Language Comprehension: LLMs

Opera refers to a dramatic art form, originating in Europe, in which the emotional content is conveyed to the audience as much through music, both vocal and instrumental, as it is through the lyrics. By contrast, in musical theater an actor's dramatic performance is primary, and the music plays a lesser role. The drama in opera is presented using the primary elements of theater such as scenery, costumes, and acting. However, the words of the opera, or libretto, are sung rather than spoken. The singers are accompanied by a musical ensemble ranging from a small instrumental ensemble to a full symphonic orchestra.

It is pointed out in the reading that opera

- A) has developed under the influence of musical theater
- B) is a drama sung with the accompaniment of an orchestra
- C) is not a high-budget production
- D) is often performed in Europe
- E) is the most complex of all the performing arts

B) is a drama sung with the accompaniment of an orchestra

[From GPT-3 davinci]

Interactive Language



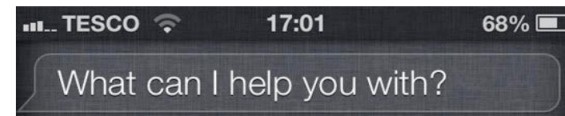
Example: Virtual Assistants



- Virtual assistants must do
 - Speech recognition
 - Language analysis
 - Dialog processing
 - Text to speech
 - ... and back-end actions!



Conversations with Devices?





ELIZA

```
ELIZA: You're not really talking about me -- are you ?
YOU:   What?
ELIZA: What do you think ?
YOU:   I mean, I feel happy to be writing right now, to be typing this as we spe
ak.  Lately, I've had writer's block.
ELIZA: You say you mean ?
YOU:   That's what I said.
ELIZA: Can you elaborate on that ?
YOU:   I only really feel happy when I'm writing.
ELIZA: How long have you been writing ?
YOU:   Always, but I've been writing seriously since I was eighteen.
ELIZA: What incident are you thinking of ?
YOU:   Um.
ELIZA: What does that suggest to you ?
YOU:   I don't think it's enough, to only be happy when you write.
ELIZA: Does that trouble you ?
YOU:   ...Um, because you can't write all the time.  And then there's the real w
orld, and you still have to deal with real people.
ELIZA: What are your feelings now ?
YOU:   I've been thinking a lot of that poem by Yeats:  "The Man and the Echo."
        Do you know it?
ELIZA: Do you say you've been thinking a lot of that poem by yeats: "the man and
        the echo for some special reason ?
YOU:   Yes
```




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

.....

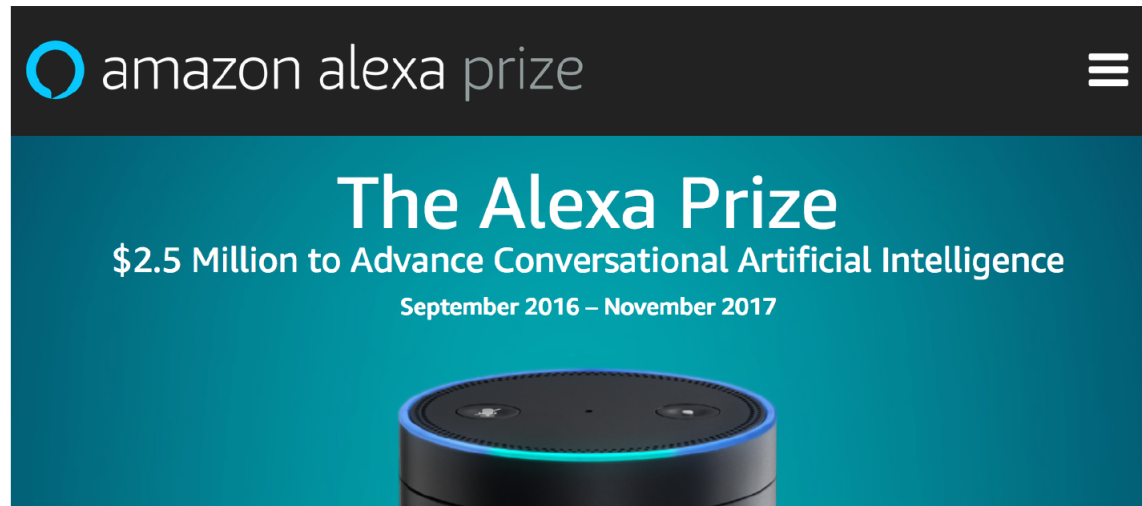


Xiaoice, developed
by Microsoft

Source: Microsoft



Chatbot Competitions!



- Alexa Prize competition to build chatbots that keep users engaged
 - First winner in 2017: UW's Sounding Board (Fang, Cheng, Holtzman, Ostendorf, Sap, Clark, Choi)
- 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 <u>five</u> .	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 of the movie is your favorite?
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

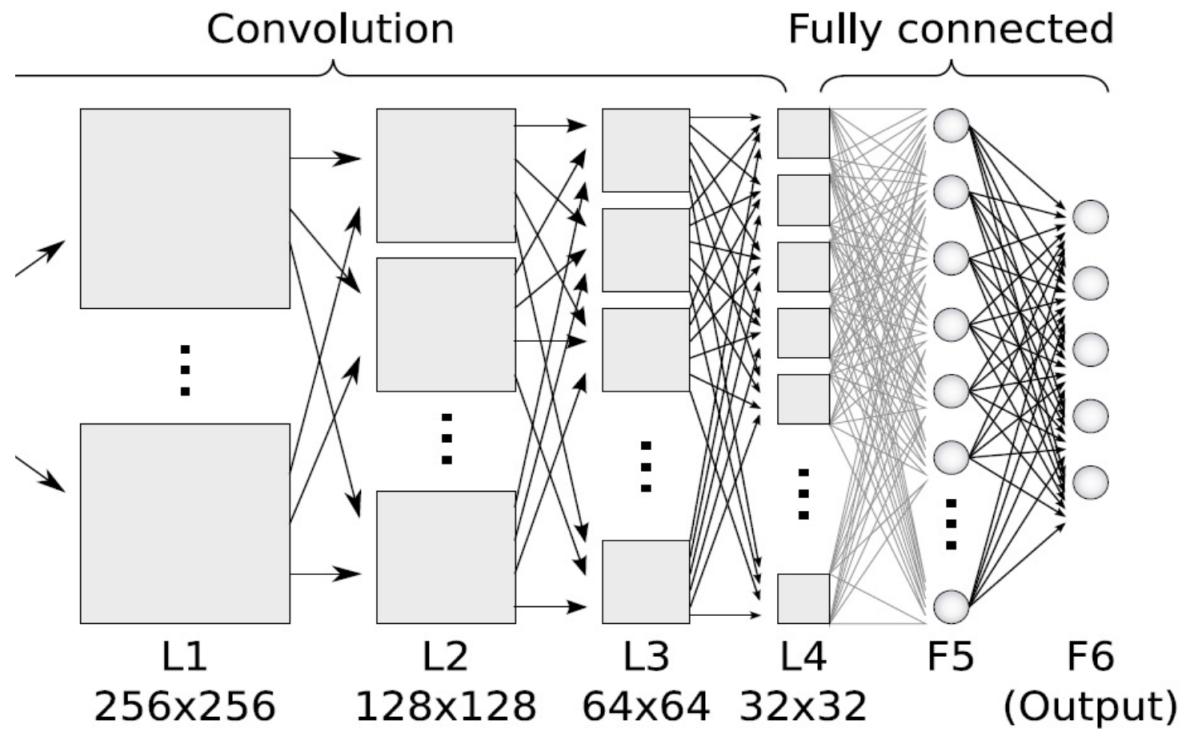
Movie

Fact

Source: Mari Ostendorf



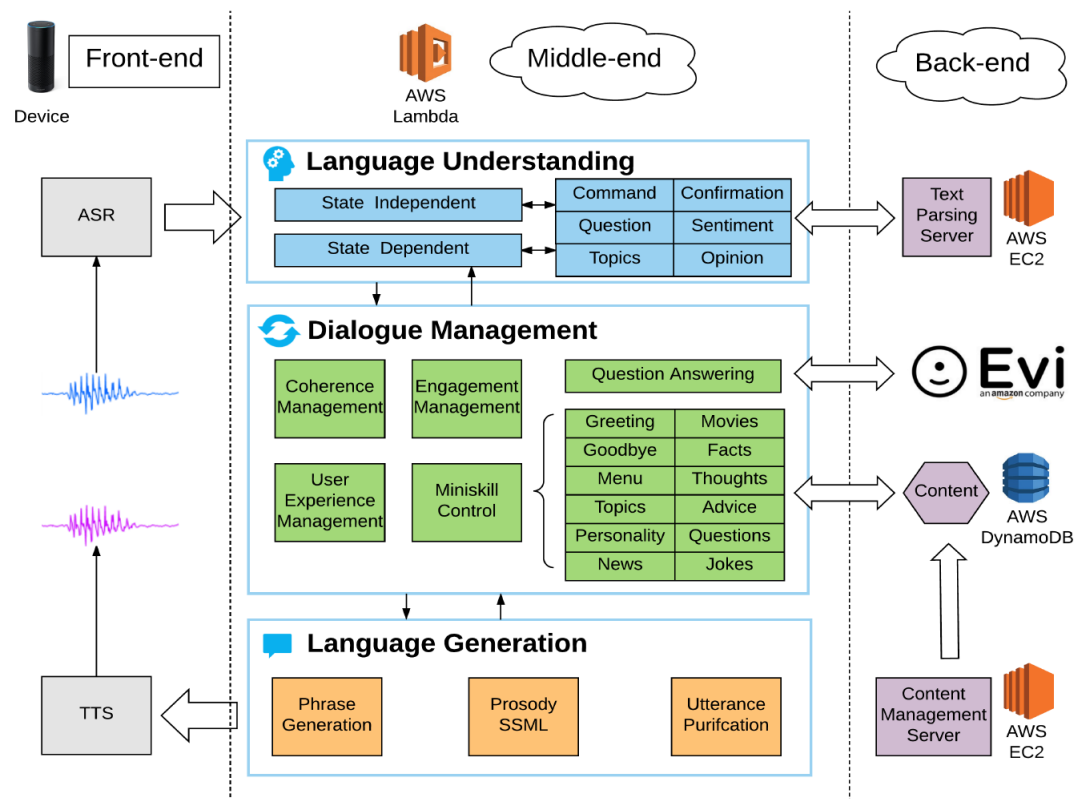
Sounding Board's Architecture?



Source: Yejin Choi



Sounding Board's Architecture



Source: Yejin Choi

Why is Language Hard?



Ambiguities



Live TV



Stevie Wonder announces he'll be having kidney surgery during London concert

By Amir Vera, CNN

Updated 11:16 PM EDT, Sat July 06, 2019



(CNN) — [Stevie Wonder](#) will be taking a break from music.

The legendary singer-songwriter announced during [a concert in London](#) Saturday that he will be undergoing kidney surgery.



Problem: Ambiguity

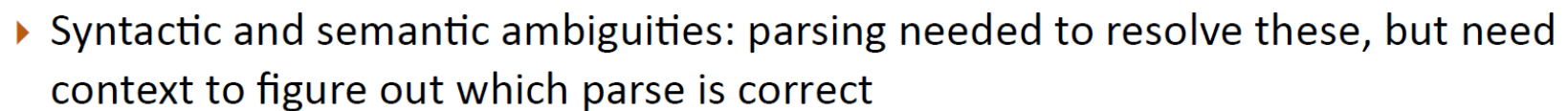
- Headlines:
 - Teacher Strikes Idle Kids
 - Hospitals Are Sued by 7 Foot Doctors
 - Ban on Nude Dancing on Governor's Desk
 - Stolen Painting Found by Tree
 - Kids Make Nutritious Snacks
 - Iraqi Head Seeks Arms
 - Local HS Dropouts Cut in Half

- Why are these funny?

What Do We Need to Understand Language?

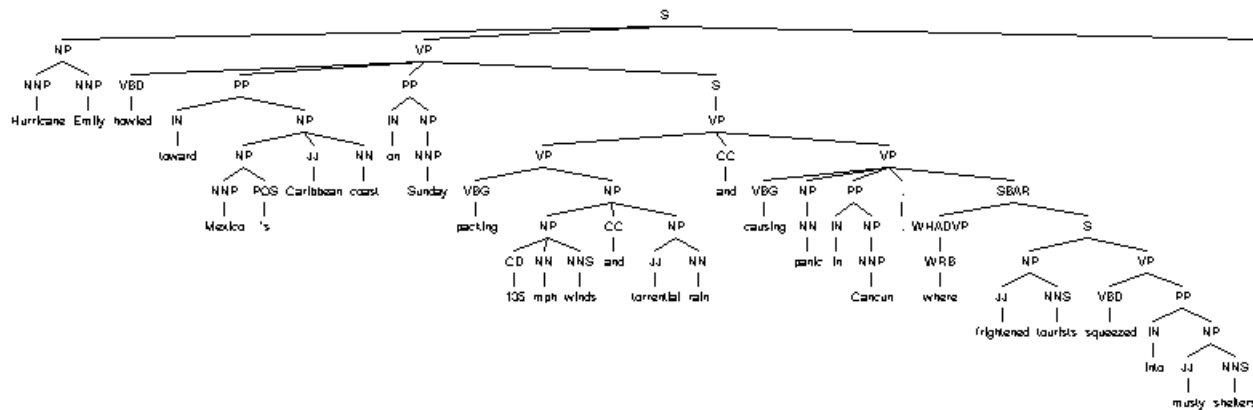


body/ position	body/ weapon
Iraqi Head Seeks Arms	





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+



Latent Representations

Personal Pronouns (PRP)

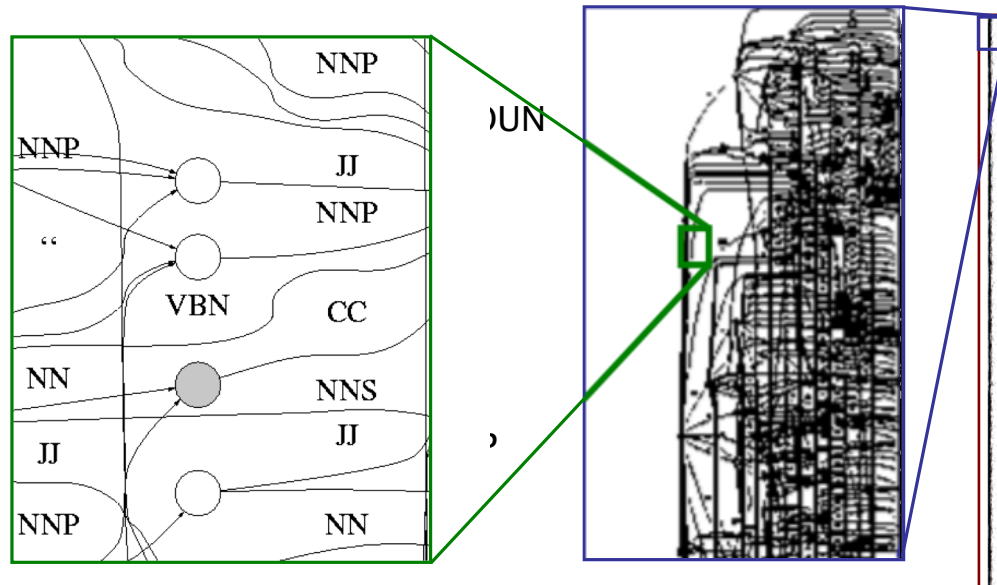
PRP-1	it	them	him
PRP-2	it	he	they
PRP-3	It	He	I

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 Data





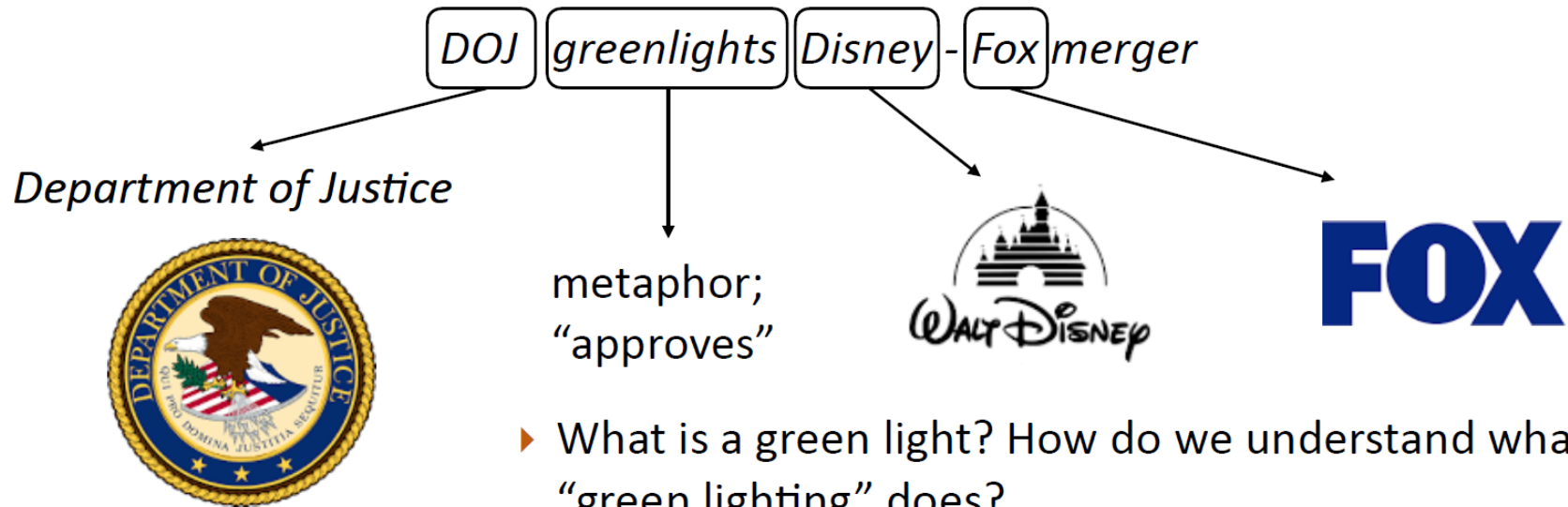
We Need A Lot of Data

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 Knowledge

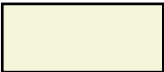
- ▶ World knowledge: have access to information beyond the training data



- ▶ What is a green light? How do we understand what "green lighting" does?
- ▶ Need commonsense knowledge



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 web data



Knowledge from Large LMs?

AI2 Allen Institute for AI

AllenNLP

🔍 Annotate a sentence ^

Semantic Role Labeling

Named Entity

Recognition

Constituency Parsing

Dependency Parsing

Open Information

Extraction

Sentiment Analysis

📄 Annotate a passage ^

Coreference Resolution

❓ Answer a question ^

Reading

Comprehension

🔧 Semantic parsing ^

Language Modeling

This demonstration uses the public 345M parameter [OpenAI GPT-2](#) language model to generate sentences.

Enter some initial text and the model will generate the most likely next words. You can click on one of those words to choose it and continue or just keep typing. Click the left arrow at the bottom to undo your last choice.

Sentence:

At high temperatures, ice becomes

Predictions:

32.0% [The](#)

10.3% [Ice](#)

5.2% [↩](#)

3.7% [Water](#)

2.7% [When](#)

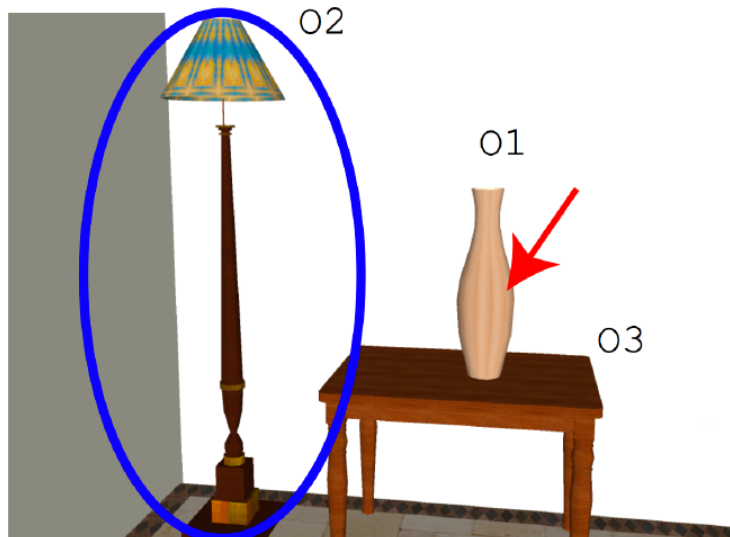
[← Undo](#)



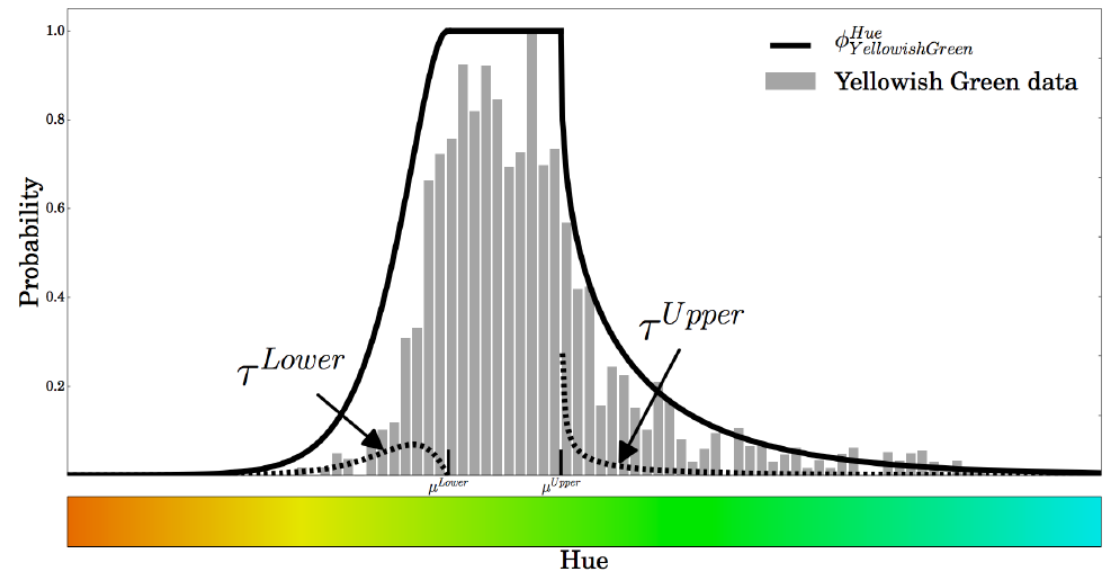
We Need Grounding

Grounding: linking linguistic concepts to non-linguistic ones

Question: What object is right of O2 ?



Golland et al. (2010)



McMahan and Stone (2015)



Example: Grounded Dialog

When is my package arriving?

Friday!





Why is Language Hard?

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

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