# DSC291: Advanced Statistical Natural Language Processing

## Overview

**Zhiting Hu**

Lecture 1, March 29, 2022

**UC San Diego**

HALICIOĞLU DATA SCIENCE INSTITUTE

# Logistics

- Class webpage: http://zhiting.ucsd.edu/teaching/dsc291spring2022

DSC291-Spring2022                                                 Logistics   Lectures   Homework   Project

**A**dvanced **S**tatistical **N**atural **L**anguage **P**rocessing

DSC 291 • Spring 2022 • UC San Diego

# Logistics

- Lectures
  - **Time:** Tuesday/Thursday 3:30pm-4:50pm
  - **Location:** HSS 1315

- No discussion session as a DSC 291 class
- Instead: Office hours, Piazza, ad-hoc meetings if needed

# Logistics

Instructor: Zhiting Hu
Email: zhh019@ucsd.edu
Office hours: Thursday 2:30-3:30pm
Location: SDSC E247

TA: Pushkar Bhuse
Email: pbhuse@ucsd.edu
Office hours: TBA
Location: TBA

- Canvas
- Discussion forum: Piazza
- Homework & writeup submission: Gradescope

# Logistics: grading

- 2 Homework assignments (30% of grade)
- Paper presentation (20%)
- Course project (46%)
- Participation (4%)

# Logistics: grading

- 2 Homework assignments (30% of grade)
  - Theory exercises, implementation exercises
  - 3 total late days without penalty
- Paper presentation (20%)
- Course project (46%)
- Participation (4%)

# Logistics: grading

- 2 Homework assignments (30% of grade)
- Paper presentation (20%)
  - Each student will give an oral presentation on a research paper
  - 10 mins = 8 mins presentation + 2 mins QA
  - Discuss both strengths and limitations of the paper
  - Sign up in a google sheet (TBA)
  - Starting TBA
- Course project (46%)
- Participation (4%)

# Logistics: grading

- 2 Homework assignments (30% of grade)
- Paper presentation (20%)
- Course project (46%)
  - 3 or 4-member team to be formed and sign up in a google sheet (TBA)
  - Designed to be as similar as possible to researching and writing a conference-style paper:
    - Due to tight timeline, fine to use synthetic/toy data for proof-of-concept experiments + explanation of theory/intuition of why your approach is likely to work
  - **Proposal** : 2 pages excluding references (10%) -- Due 04/14
    - Overview of project idea, literature review, potential datasets and evaluation, milestones
  - **Midway Report** : 4-5 pages excluding references (20%)
  - **Presentation** : oral presentation, 15-20mins (20%)
  - **Final Report** : 6-8 pages excluding references (50%)

# Logistics: grading

- 2 Homework assignments (30% of grade)
- Paper presentation (20%)
- Course project (46%)
- Participation (4%)
  - Contribution to discussion on Piazza
  - Complete mid-quarter evaluation
  - Any constructive suggestions

# Advanced Statistical Natural Language Processing

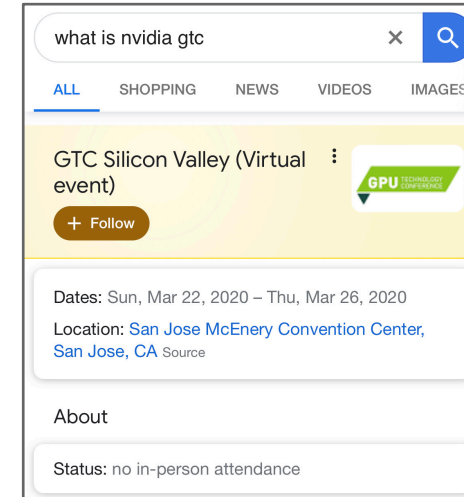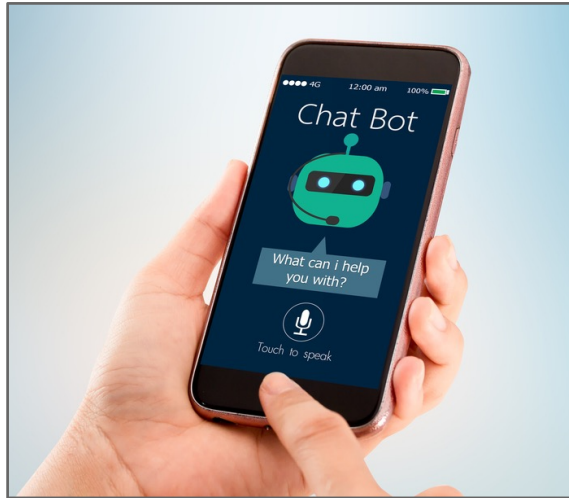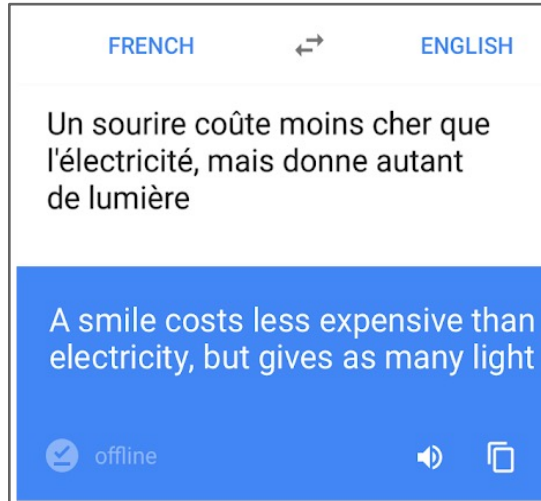# **Advanced <span style="color:blue">Statistical</span> <span style="color:red">Natural Language Processing</span>**

What is NLP?

Statistical machine
learning (ML) methods

We'll cover only a subset of
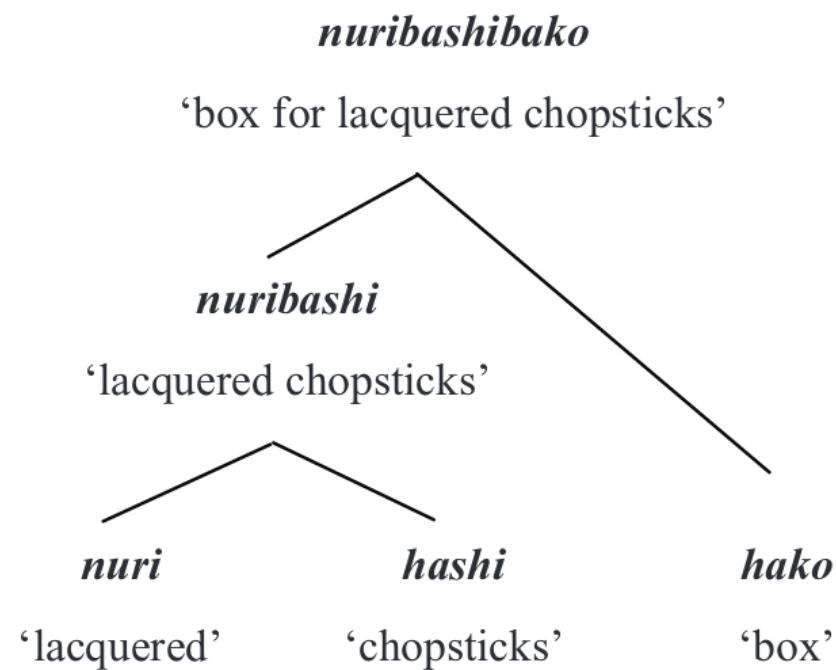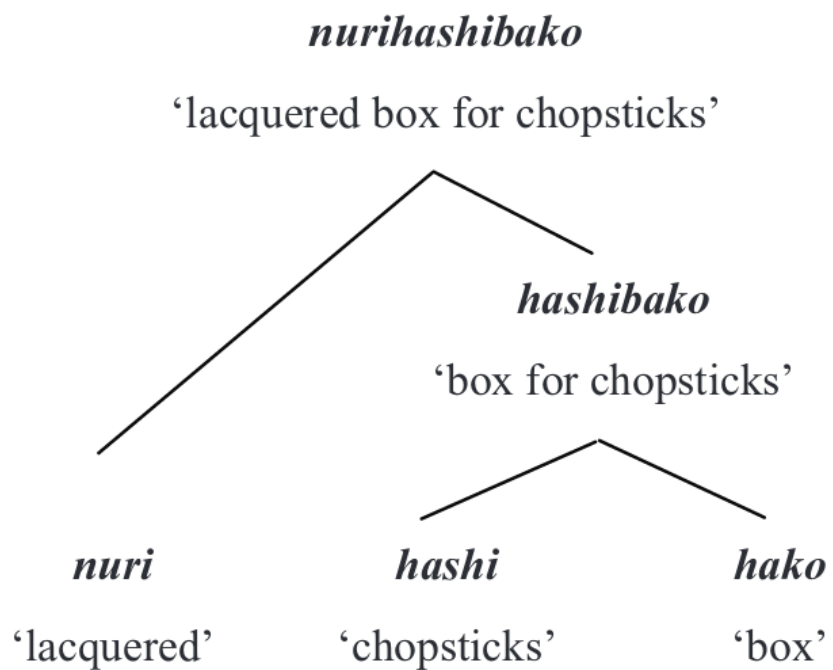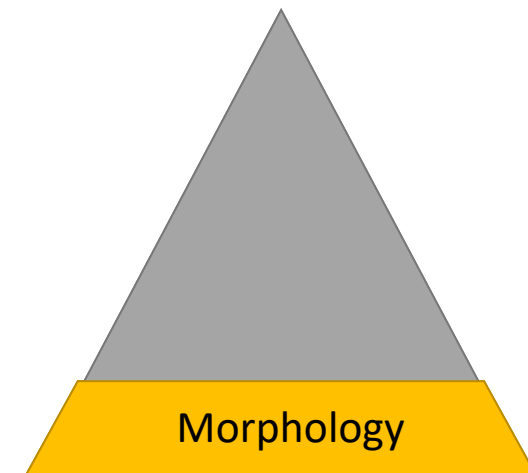advanced, latest methods

# What is NLP

# What is NLP

- NL ∈ { English, German, Chinese, Spanish, Hindi, American Sign Language, . . ., Lushootseed }

- Automation of:
  - analysis or "understanding" (to some degree) what a text means
  - generation of fluent, meaningful, context-appropriate text
  - acquisition of these capabilities from knowledge and data
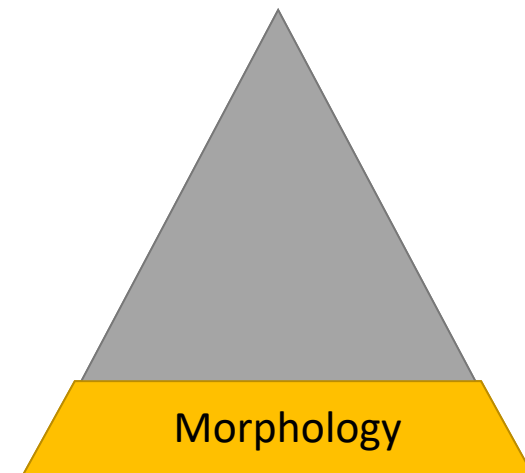
# Language Understanding Pyramid
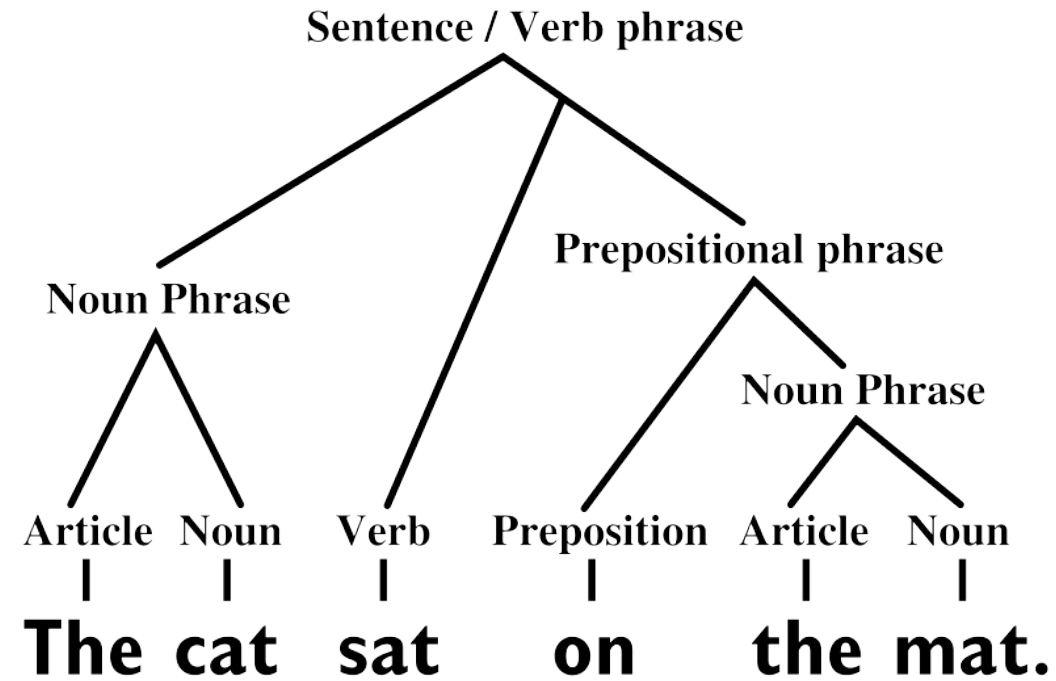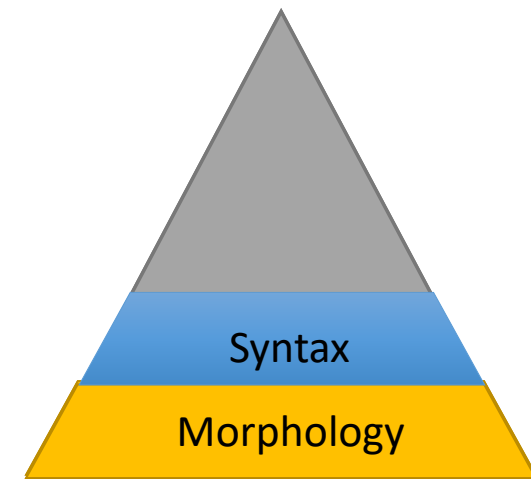
# Morphology

Morphology Analysis

Morphology

nurihashibako
'lacquered box for chopsticks'

hashibako
'box for chopsticks'

nuri
'lacquered'

hashi
'chopsticks'

hako
'box'

nuribashibako
'box for lacquered chopsticks'

nuribashi
'lacquered chopsticks'

nuri
'lacquered'

hashi
'chopsticks'

hako
'box'

# Morphology

**Stemming**

adjustable → adjust_
formality → formaliti
formaliti → formal
airliner → airlin_

**Lemmatization**

was → (to) be
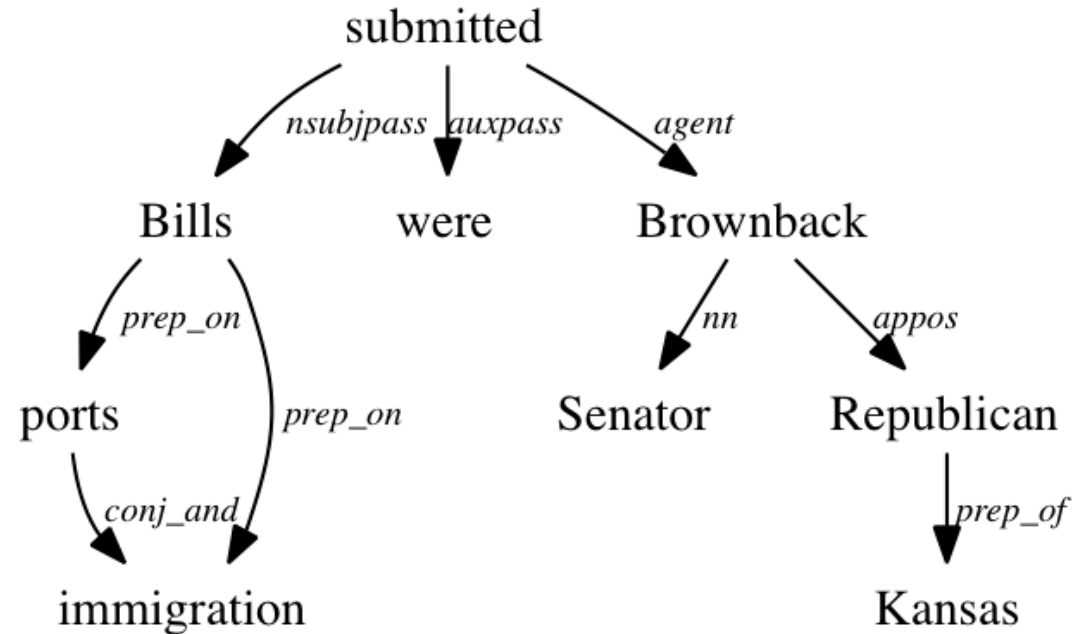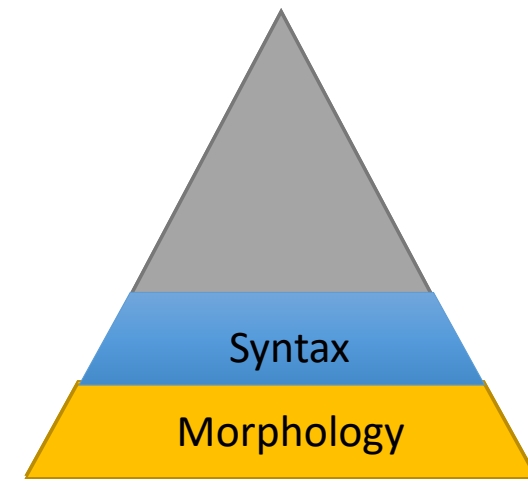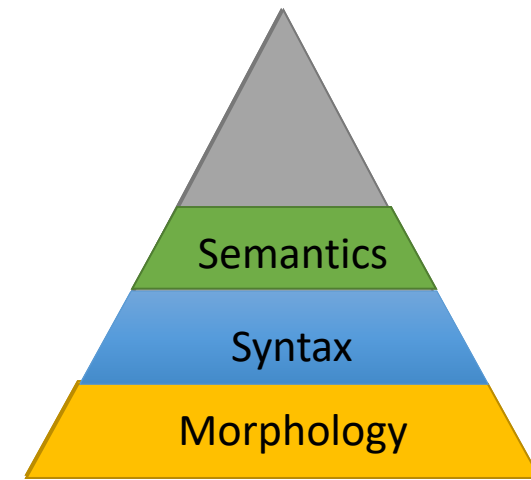better → good
meeting → meeting

Morphology

# Syntax

**Constituent Parsing**

Syntax

Morphology



Sentence / Verb phrase

Noun Phrase

Prepositional phrase

Noun Phrase

Article   Noun   Verb   Preposition   Article   Noun

**The cat   sat   on   the mat.**

17

# Syntax

Dependency Parsing

Syntax

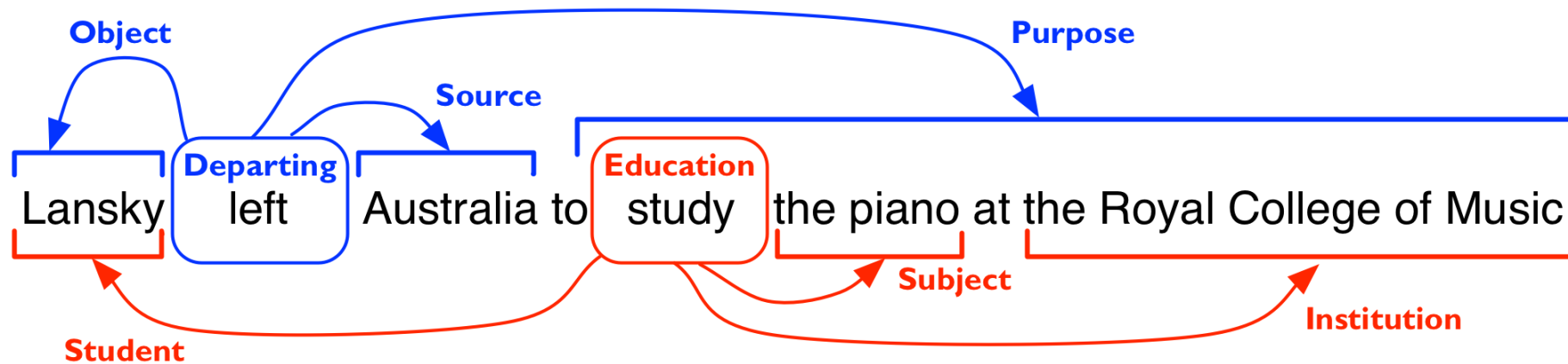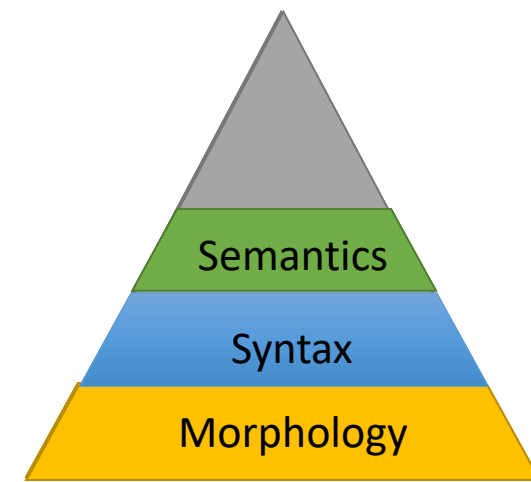Morphology

18

# Semantics

Named Entity Recognition

Semantics

Syntax

Morphology

When **Sebastian Thrun** PERSON started at **Google** ORG in **2007** DATE , few people outside of the company took him seriously. "I can tell you very senior CEOs of major **American** NORP car companies would shake my hand and turn away because I wasn't worth talking to," said **Thrun** PERSON , now the co-founder and CEO of online higher education startup Udacity, in an interview with **Recode** ORG **earlier this week** DATE .
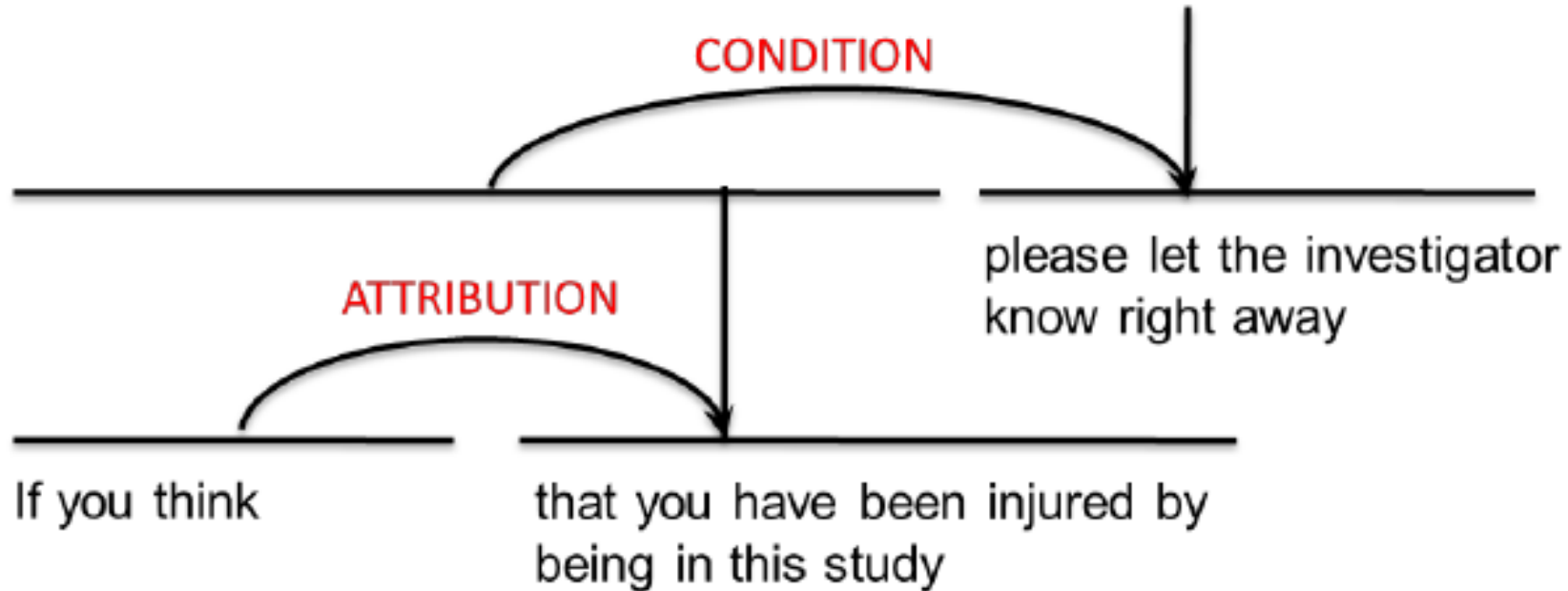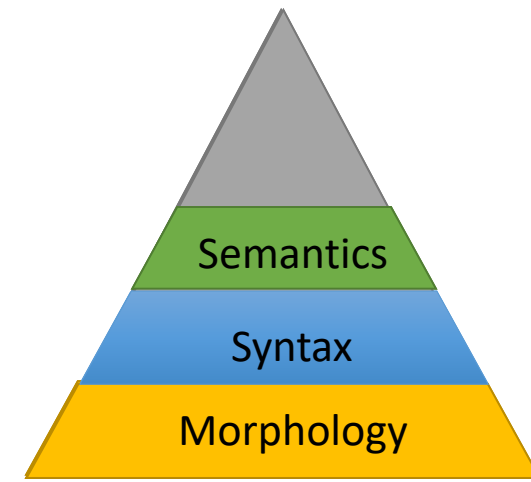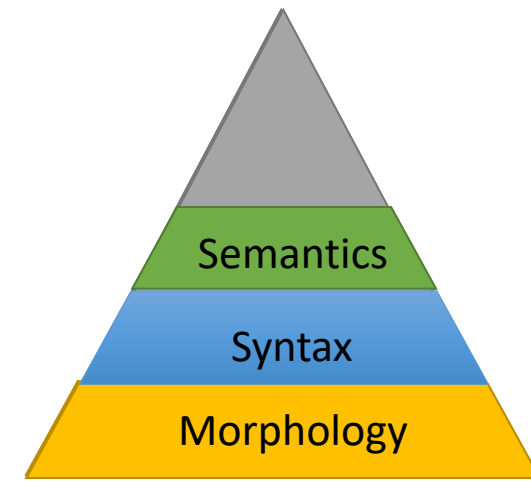
# Semantics

Semantic Roles

Semantics

Syntax

Morphology



Object

Purpose

Source

Departing

Education

Lansky  left  Australia to  study  the piano at the Royal College of Music

Subject

Student

Institution

20

# Semantics

Discourse Parsing

# Semantics

Coreference

Semantics

Syntax

Morphology

"*I* voted for *Nader* because *he* was most aligned with *my* values," *she* said.

# Semantics

Semantics

Syntax

Morphology

**Entity Linking**

Kate Winslet and Leonardo Dicaprio
have definitely created a timeless classic.

## Kate Winslet

WIKIPEDIA
The Free Encyclopedia

From Wikipedia, the free encyclopedia

**Kate Elizabeth Winslet** CBE (born 5 October 1975) is an English actress. She is particularly known for her work in period dramas, and is often drawn to portraying angst-ridden women. Winslet is the recipient of various accolades, including three British Academy Film Awards, and is among the few performers to have won Academy, Emmy, and Grammy Awards.

Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia
Wikipedia store

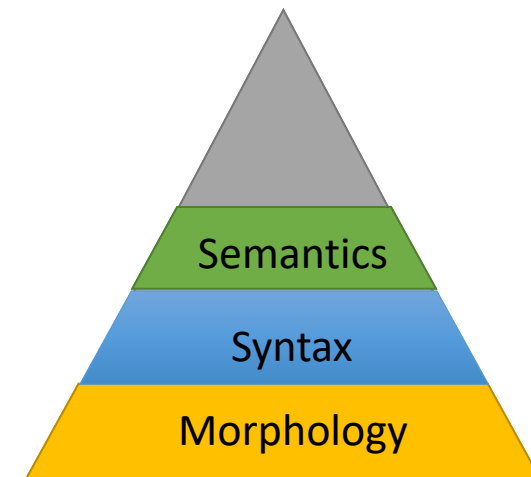Interaction

Help

IMDb ≡ Menu IMDb TV All ▾ Search IMDb

### Leonardo DiCaprio
Actor | Producer | Writer

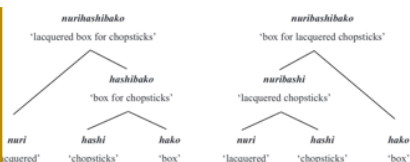1:08 | Interview

Few actors in the world have had a career quite as diverse as Leona has gone from relatively humble beginnings, as a supporting cast m Growing Pains (1985) and low budget horror movies, such as Critter teenage heartthrob in the 1990s, as the hunky lead actor in movies See full bio »
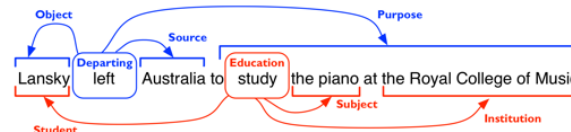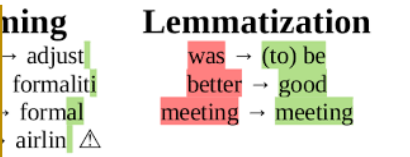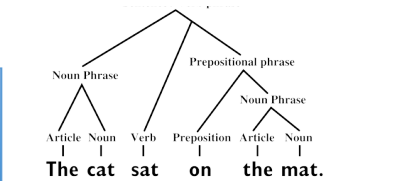
# Language Understanding Pyramid



Morphology

Lemmatize

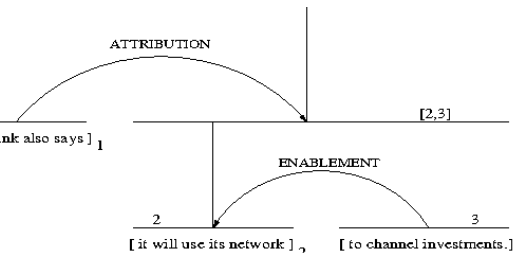Constituent Parse

Dependency Parse

Semantic Parsing

NER

Discourse Parsing

Coreference

Entity Linking

Semantics

Syntax

Morphology

# Pragmatics

# Why NLP is Hard

- Ambiguity
  - A string may have many possible interpretations in different contexts, and resolving ambiguity correctly may rely on knowing a lot about the world.

*We saw the woman with the telescope wrapped in paper.*

- Who has the telescope?
- Who or what is wrapped in paper?
- An event of perception, or an assault?

# Why NLP is Hard

- Ambiguity
  - A string may have many possible interpretations in different contexts, and resolving ambiguity correctly may rely on knowing a lot about the world.
  - Richness: any meaning may be expressed many ways, and there are immeasurably many meanings.
  - Linguistic diversity across languages, dialects, genres, styles, …
- Appropriateness of a representation depends on the application
- Typically, representation of language is a theorized construct, not directly observable, or it is encoded numerically (vectors, matrices, tensors) and inscrutable
- There are many sources of variation and noise in linguistic input

# Advanced Statistical Natural Language Processing

What is NLP?

Statistical machine
learning (ML) methods

We'll cover only a subset of
advanced, latest methods

# Machine Learning

- Computational methods that enable machines to learn concepts and improve performance from **experiences**.

# Experiences of all kinds


Data examples

Type-2 diabetes is 90% more common than type-1

Rules/Constraints


Knowledge graphs


Rewards


Auxiliary agents


Adversaries


Teachers

...

And all combinations thereof

# Experiences of all kinds



Type-2

missing

Data exampl...

Auxiliary agents

...ations thereof

should be conceived
as a kind of intimate reverie

Adversaries

Master classes

# Experiences of all kinds



Data examples

Type-2 diabetes is 90% more common than type-1

Rules/Constraints

Knowledge graphs

SCORE: 107

Rewards

Auxiliary agents

Adversaries

should be conceived as a kind of intimate reverie

Master classes

...

And all combinations thereof

# Experiences: (massive) data examples


Image classification


Machine translation



Language modeling
(BERT, GPT-2, **GPT-3**, …)

45TB of text data: CommonCrawl, WebText,
Wikipedia, corpus of books, …

# Experiences: (massive) data examples

**TECH** \ **ARTIFICIAL INTELLIGENCE**

## OpenAI's text-generating system GPT-3 is now spewing out 4.5 billion words a day

*Robot-generated writing looks set to be the next big thing*

By James Vincent | Mar 29, 2021, 8:24am EDT

**Loud and clear**

Speech-recognition word-error rate, selected benchmarks, %

*Log scale*

Switchboard

Switchboard cellular

Meeting speech

Broadcast speech

IBM, Switchboard

Microsoft, Switchboard — 5.9%

The Switchboard corpus is a collection of recorded telephone conversations widely used to train and test speech-recognition systems

100

10

1

1993  96  98  2000  02  04  06  08  10  12  14  16

Sources: Microsoft; research papers

**Speak easy**

Human scorers' rating* of Google Translate and human translation

Translation method | Phrase-based† | Neural-network† | Human

3    4    5    Perfect translation=6

English → Spanish
English → French
English → Chinese

Spanish → English
French → English
Chinese → English

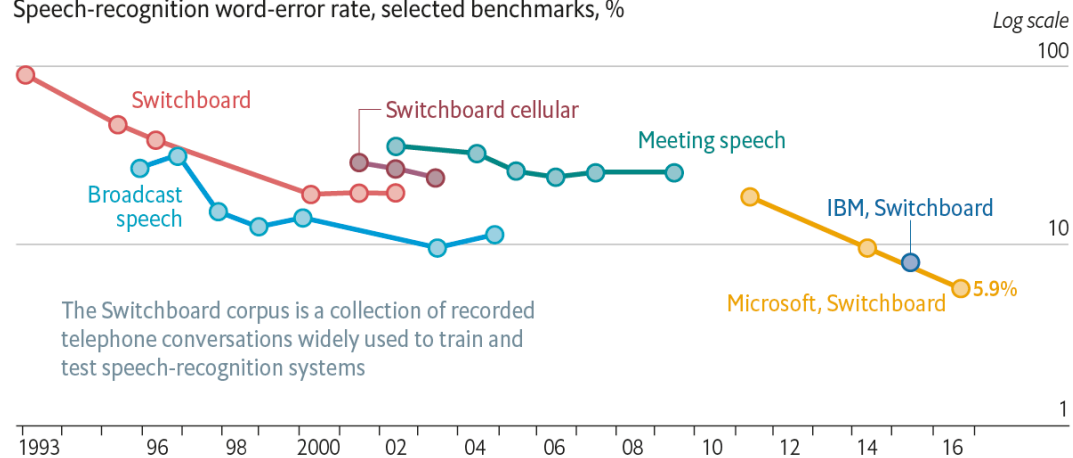**Input sentence** — Pour l'ancienne secrétaire d'Etat, il s'agit de faire oublier un mois de cafouillages et de convaincre l'auditoire que M. Trump n'a pas l'étoffe d'un président

**Phrase-based†**
For the former secretary of state, this is to forget a month of bungling and convince the audience that Mr Trump has not the makings of a president

**Neural-network†**
For the former secretary of state, it is a question of forgetting a month of muddles and convincing the audience that Mr Trump does not have the stuff of a president

**Human**
The former secretary of state has to put behind her a month of setbacks and convince the audience that Mr Trump does not have what it takes to be a president

Source: Google          *0=completely nonsense translation, 6=perfect translation  †Machine translation

[The Economist]

34

# Problems with few data (labels)

- Privacy, security issues

Assistive diagnosis

Normal findings

Abnormal findings

``*The heart size and mediastinal contours appear within normal limits. There is blunting of the right lateral costophrenic sulcus which could be secondary to a small effusion versus scarring ...*''

# Problems with few data (labels)

- Expense to collect/annotate

Controllable content generation

*Controlling sentiment*

Pos — *The film is full of imagination!*

↓

Neg — *The film is strictly routine!*

*Controlling writing style*

Plain — *LeBron James contributed 26 points, 8 rebounds, 7 assists.*

↓

Elaborate — *LeBron James rounded out the box score with an all around impressive performance, scoring 26 points, grabbing 8 rebounds and dishing out 7 assists.*

Applications: personalized chatbot, live sports commentary production

# Problems with few data (labels)

- Difficult / expertise-demanding to annotate

Adversarial attack

"entailment"  "neutral"  "contradiction"

Entailment classifier

The Old One always comforted Ca'daan, except today.

Your gift is appreciated by each and every student …

At the other end of Pennsylvania Avenue, people …

The person saint-pierre-et-saint-paul is ..

premises

hypothesis (attack)

Applications: test model robustness

# Problems with few data (labels)

- Difficult / expertise-demanding to annotate

Prompt generation: automatically generating prompts to steer pretrained LMs



Pretrained LM
(e.g., GPT3)

Generate a story about cat: once upon a time,    ...

prompt                              input        continuation

# Problems with few data (labels)

- Specific domain    Low-resource languages

    ~7K languages in the world

# Problems with few data (labels)

- Specific domain    Low-resource languages



Written languages
(3.5K)

All languages
(7K)

Languages with
NER Annotation
(30?)

42

# Problems with few data (labels)

- Specific domain    Low-resource languages



Written languages (3.5K)

All languages (7K)

Can we translate the annotation to other languages?
Requires parallel data for training
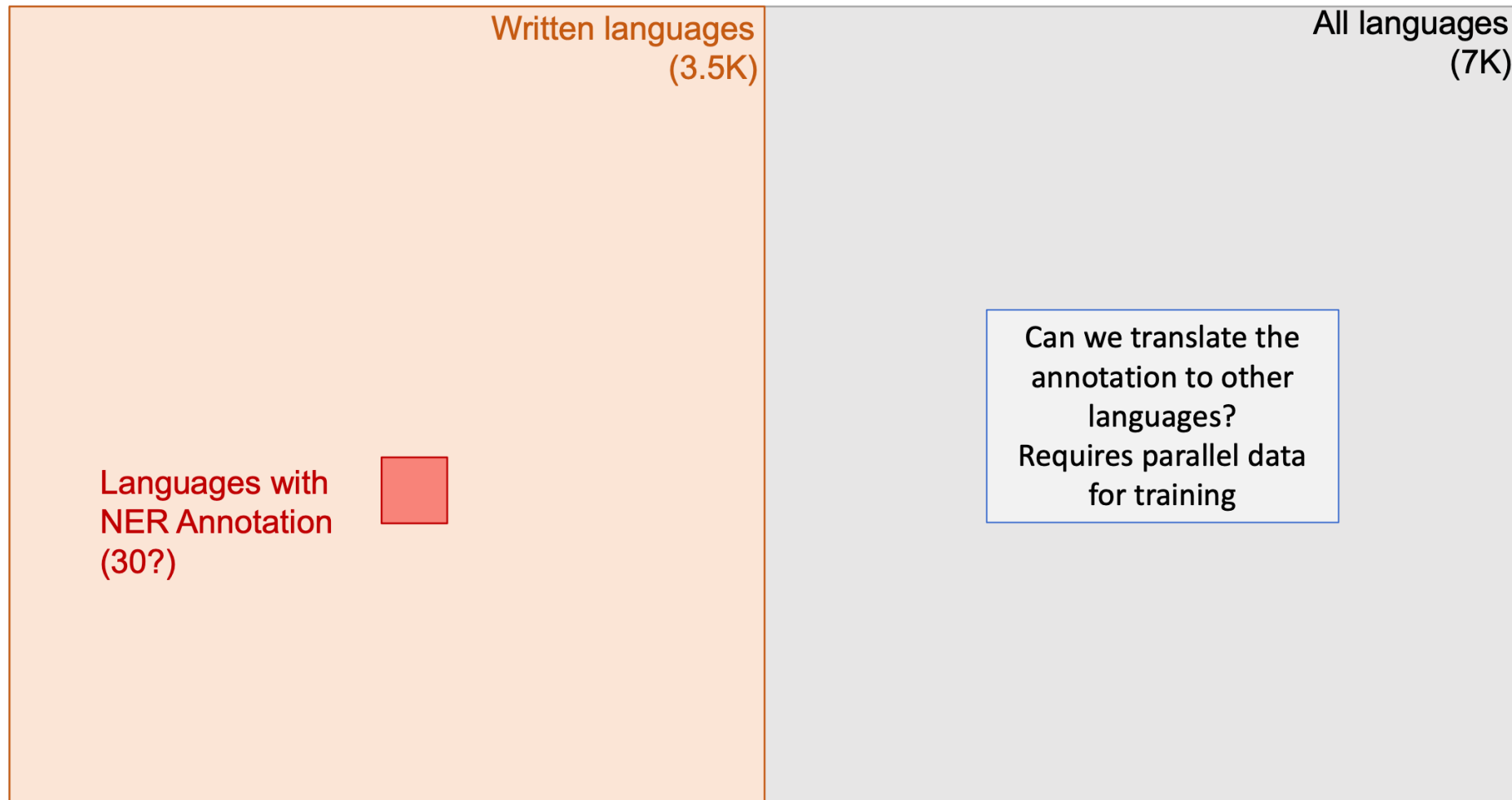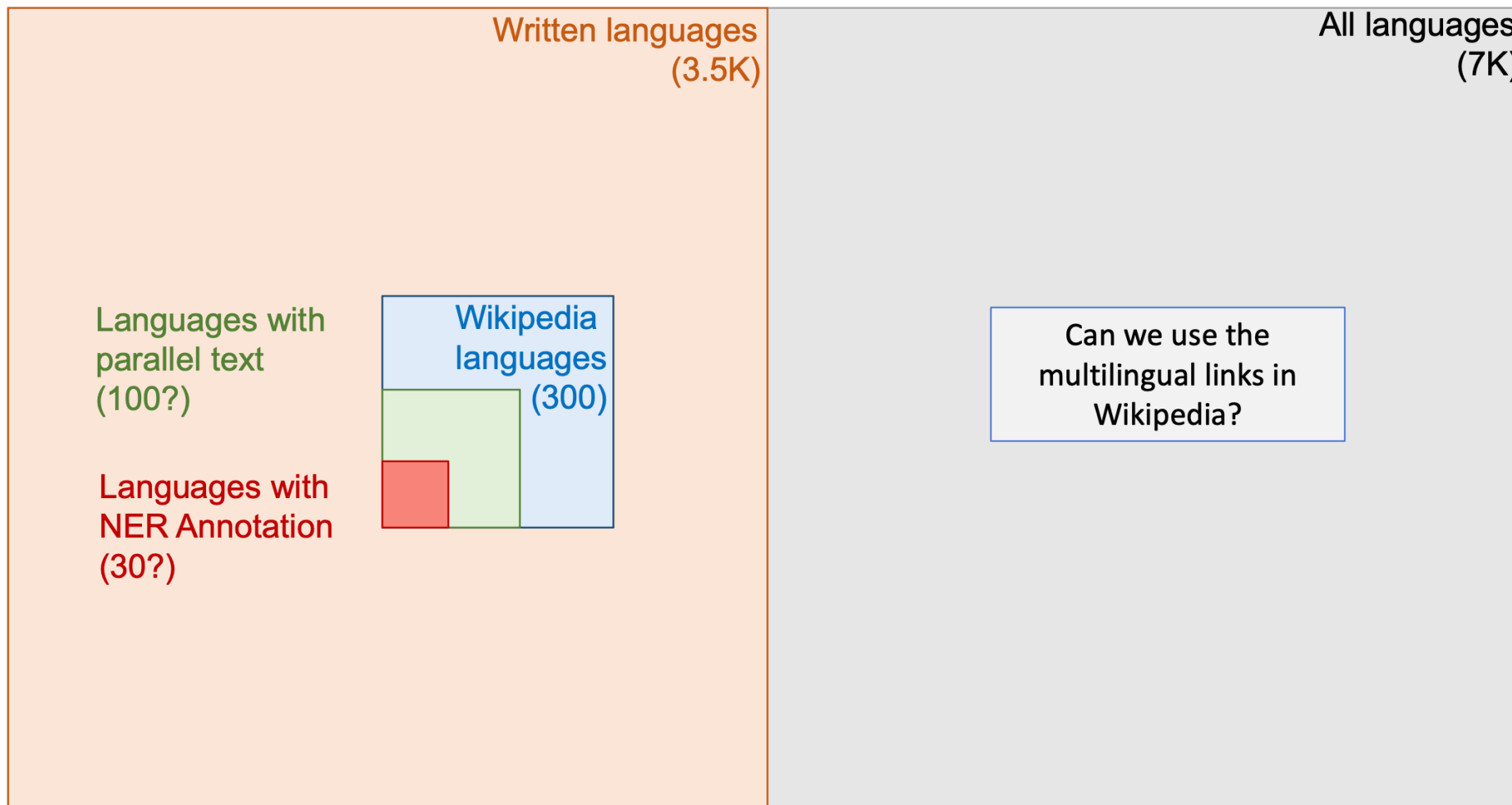
Languages with NER Annotation (30?)
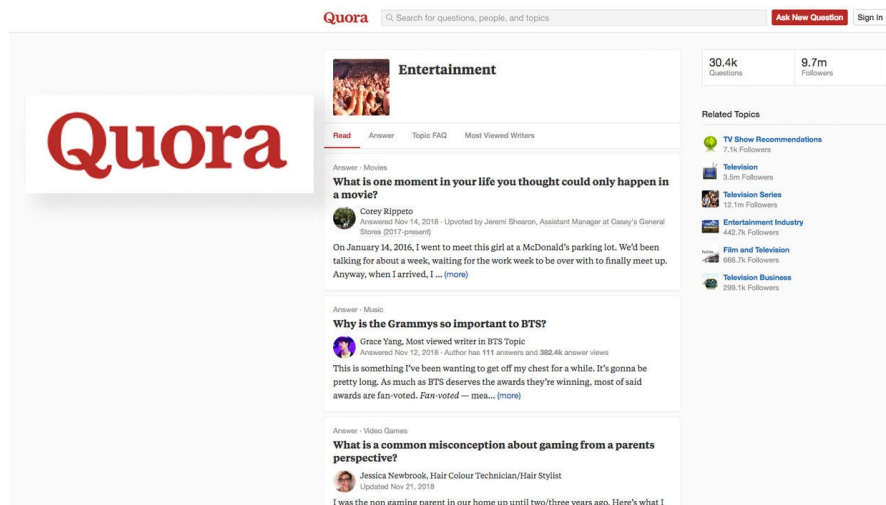
# Problems with few data (labels)
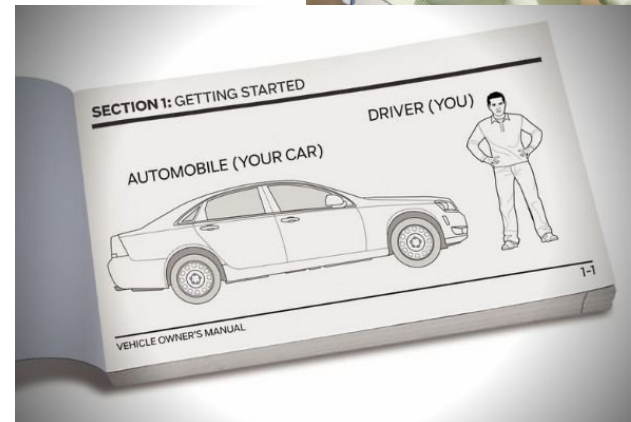
- Specific domain   Low-resource languages

# Problems with few data (labels)

- Specific domain

Question answering

QA based on car manual?

# Problems with few data (labels)

- Privacy, security issues
- Expensive to collect/annotate
- Difficult / expertise-demanding to annotate
- Specific domain

# Machine learning solutions given few data (labels)

- How can we make more efficient use of the data?
  - Clean but small-size
  - Noisy
  - Out-of-domain

- Can we incorporate other types of experiences in learning?

Data examples

Type-2 diabetes is 90% more common than type-1

Rules/Constraints

Knowledge graphs

Rewards

Auxiliary agents

Adversaries

Master classes

... *And all combinations thereof*

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver
- Model architecture

$$\min_\theta \mathcal{L}(\theta, \mathcal{E})$$

Optimization solver    Loss    Model architecture    Experience

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver
- Model architecture

$$\min_\theta \mathcal{L}(\theta, \mathcal{E})$$

Optimization solver  Loss  Model architecture  Experience

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver
- Model architecture

Model of certain architecture whose parameters are the subject to be learned, $p_\theta(\boldsymbol{x}, \boldsymbol{y})$ or $p_\theta(\boldsymbol{y}|\boldsymbol{x})$

- ○ Neural networks
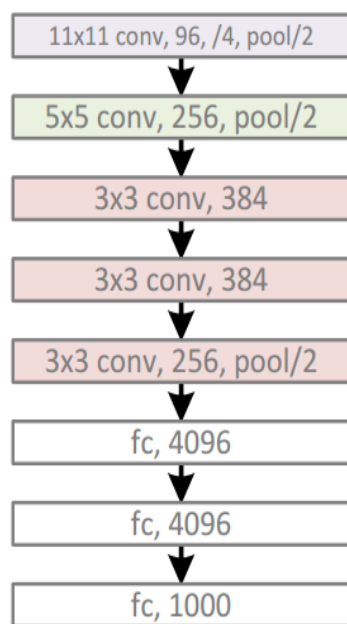- ○ Graphical models
- ○ Compositional architectures

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver
- **Model architecture**

Model of certain architecture whose parameters are the subject to be learned, $p_\theta(x, y)$ or $p_\theta(y|x)$

- ○ Neural networks
- ○ Graphical models
- ○ Compositional architectures



Convolutional networks
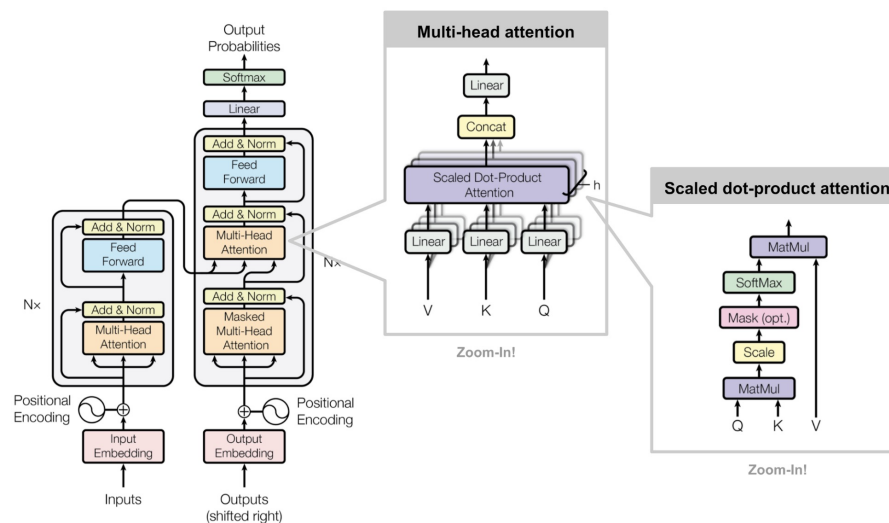


Transformers

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver

- Model architecture

Model of certain architecture whose parameters are the subject to be learned, $p_\theta(x, y)$ or $p_\theta(y|x)$

- Neural networks
- Graphical models
- Compositional architectures



Naive Bayes → SEQUENCE → HMMs → GENERAL GRAPHS → Generative directed models

CONDITIONAL

Logistic Regression → SEQUENCE → Linear-chain CRFs → GENERAL GRAPHS → General CRFs

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver
- Model architecture

Assuming you know basic procedures:
- ○ (Stochastic) gradient descent
- ○ Backpropagation
- ○ Lagrange multiplier
- ○ …

$$\min_\theta \mathcal{L}(\theta, \mathcal{E})$$

Optimization solver

Loss

Model architecture

Experience

# Components of a ML solution (roughly)

- Loss
- Experience
- Optimization solver
- Model architecture

Core of most learning algorithms

$$\min_\theta \mathcal{L}(\theta, \mathcal{E})$$

Optimization solver

Loss

Model architecture

Experience

# Machine learning solutions

- (1) How can we make more efficient use of the data?
  - Clean but small-size, Noisy, Out-of-domain

- (2) Can we incorporate other types of experiences in learning?



*Data examples*     *Rules/Constraints*     *Knowledge graphs*     *Rewards*     *Auxiliary agents*



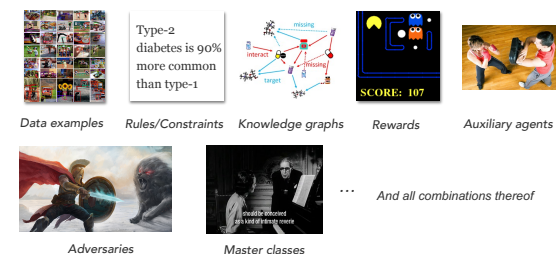*Adversaries*          *Master classes*          ...     *And all combinations thereof*

# Machine learning solutions

- (1) How can we make more efficient use of the data?
  - Clean but small-size, Noisy, Out-of-domain, …
- Algorithms

  - **Supervised learning**: MLE, maximum entropy principle

  - **Unsupervised learning**: EM, variational inference, VAEs

  - **Self-supervised learning**: successful instances, e.g., BERT, GPT-3, contrastive learning, applications to downstream tasks

  - **Distant/weakly supervised learning**: successful instances

  - **Data manipulation:** augmentation, re-weighting, curriculum learning, …

  - **Meta-learning**

# Machine learning solutions

- (2) Can we incorporate other types of experiences in learning?

  - Learning from auxiliary models, e.g., adversarial models:
    - Generative adversarial learning (GANs and variants), co-training, …

  - Learning from structured knowledge
    - Posterior regularization, constraint-driven learning, …

  - Learning from rewards
    - Reinforcement learning: model-free vs model-based, policy-based vs value-based, on-policy vs off-policy, extrinsic reward vs intrinsic reward, …

  - Learning in dynamic environment
    - Online learning, lifelong/continual learning, …



Data examples    Rules/Constraints    Knowledge graphs    Rewards    Auxiliary agents

Adversaries    Master classes    … *And all combinations thereof*

# Questions?