

# CS371N: Natural Language Processing Lecture 1: Introduction

Greg Durrett  
(he/him)



**TEXAS**  
The University of Texas at Austin





# Administrivia

---

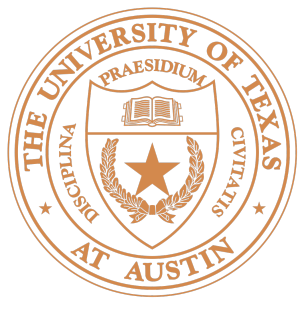
- ▶ Lecture: Tuesdays and Thursdays 9:30am-10:45am in JGB 2.218
  - ▶ Recordings available afterwards on LecturesOnline
- ▶ Course website (including **syllabus**):  
<http://www.cs.utexas.edu/~gdurrett/courses/fa2024/cs371n.shtml>
- ▶ Ed Discussion board: link on Canvas
- ▶ Office hours: see course website and Canvas. Greg's are hybrid, some TA OHs are hybrid too. **Office hours start Thursday after class.**
- ▶ TAs: Juan Diego Rodriguez and Grace Kim.
- ▶ Office hours start today, and I will stay around after this class if you have questions



# Course Requirements

---

- ▶ CS 429
- ▶ Recommended: CS 331, familiarity with probability and linear algebra, programming experience in Python
- ▶ Helpful: Exposure to AI and machine learning (e.g., CS 342/343/363)
- ▶ Assignment 0 is out now (optional):
  - ▶ If this seems like it'll be challenging for you, come and talk to me (this is smaller-scale than the other assignments, which are smaller-scale than the final project)



# Format and Accessibility

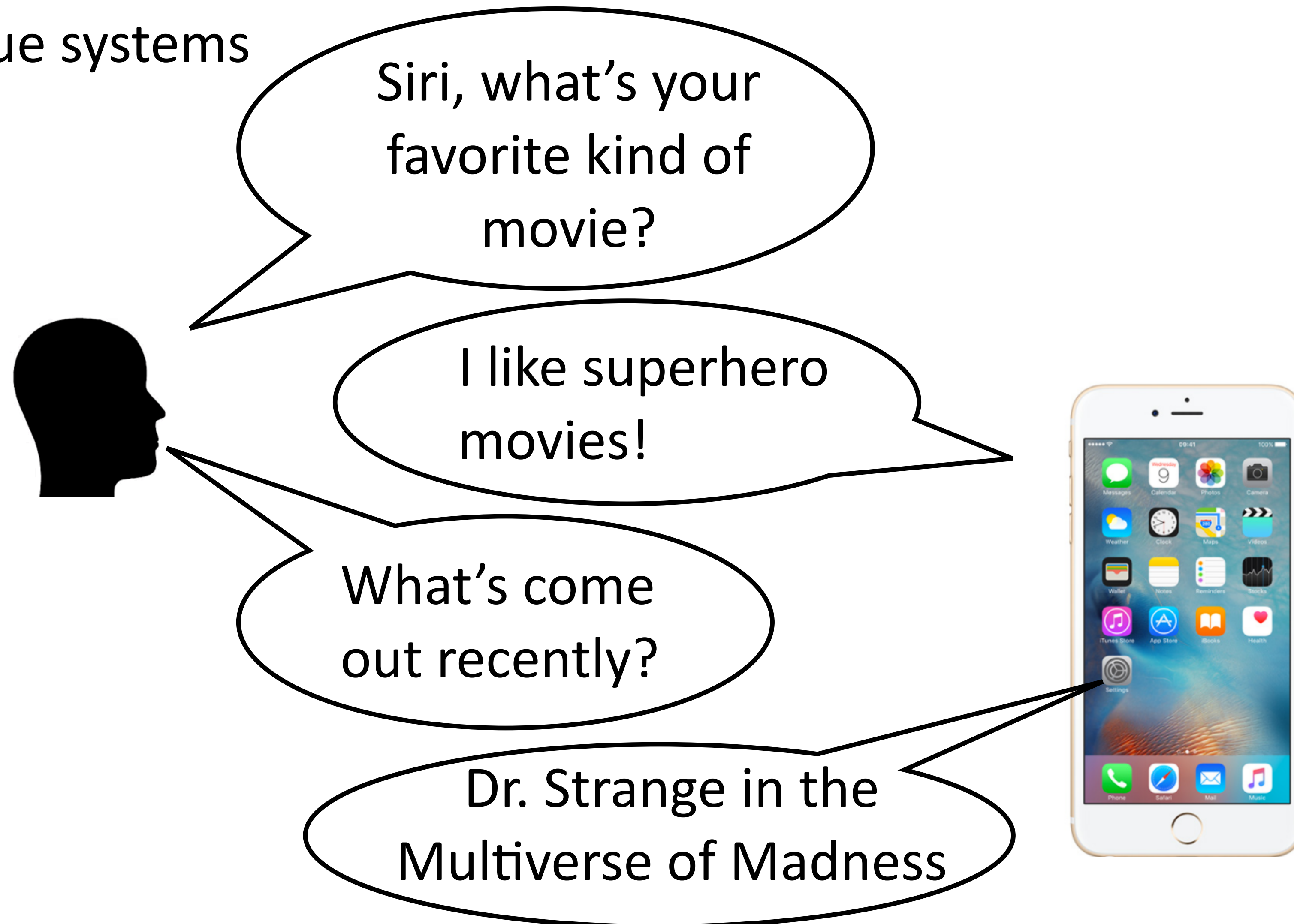
---

- ▶ Lectures will build in time for discussion, in-class exercises, and questions. Additional material is available as videos to watch either before or after lectures
  - ▶ Format: in-person to encourage discussion, but all materials are available asynchronously afterwards
- ▶ Equipment: useful to have a device for lecture to do Instapolls. For homework:
  - ▶ Lab machines available via SSH
  - ▶ A GPU is **not** required to complete the assignments! Having a GPU, GCP credits, or Google Colab access will be helpful for the final project though



# What's the goal of NLP?

- ▶ Be able to solve problems that require deep understanding of text
- ▶ Example: dialogue systems





# Machine Translation

The Political Bureau  
of the CPC Central  
Committee

July 30 hold a meeting

中共中央政治局7月30日召开会议，会议分析研究当前经济形势，部署下半年经济工作。

People's Daily, August 10, 2020

Translate

The Political Bureau of the CPC Central Committee held a meeting on July 30 to analyze and study the current economic situation and plan economic work in the second half of the year.




# Question Answering

When was Abraham Lincoln born?

Name	Birthday	map to Birthday field
Lincoln, Abraham	2/12/1809	→ <b>February 12, 1809</b>
Washington, George	2/22/1732	
Adams, John	10/30/1735	

How many visitors centers are there in Rocky Mountain National Park?



The screenshot shows the Wikipedia article for Rocky Mountain National Park. The article text includes: "Rocky Mountain National Park is an American national park located within the Front Range of the Rocky Mountains. The park is situated on the slopes of the Continental Divide run directly through the center of the park. Features of the park include mountains, alpine lakes and a wide variety of wildlife." The article also states: "The Rocky Mountain National Park Act was signed by President Woodrow Wilson in 1909. The Civilian Conservation Corps built the main automobile road through the park in 1915. In 2018, more than 4.5 million recreation visitors visited the park, ranking it as the third most visited national park in 2015. In 2019, the park has a total of five visitor centers with park headquarters located at the Lloyd Wright School of Architecture at Taliesin West. National Forests include Arapaho National Forest to the north and west, and Arapaho National Forest to the west."

The park has a total of five visitor centers

five



# Generalist Systems?

GD

Translate "中共中央政治局7月30日召开会议，会议分析研究当前经济形势，部署下半年经济工作。" into English.



The translation of "中共中央政治局7月30日召  
当前经济形势，部署下半年经济工作。" into

"The Central Politburo of the Communist Par

GD

What is  $373 * 121$ ?



The product of 373 multiplied by 121 is **45,113.**

GD

When was Abraham Lincoln's birthday?

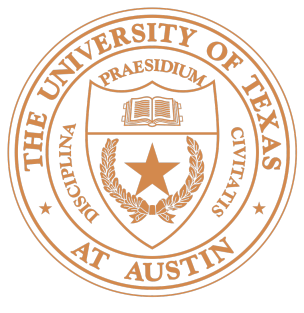


Abraham Lincoln was born on February 12, 1809.

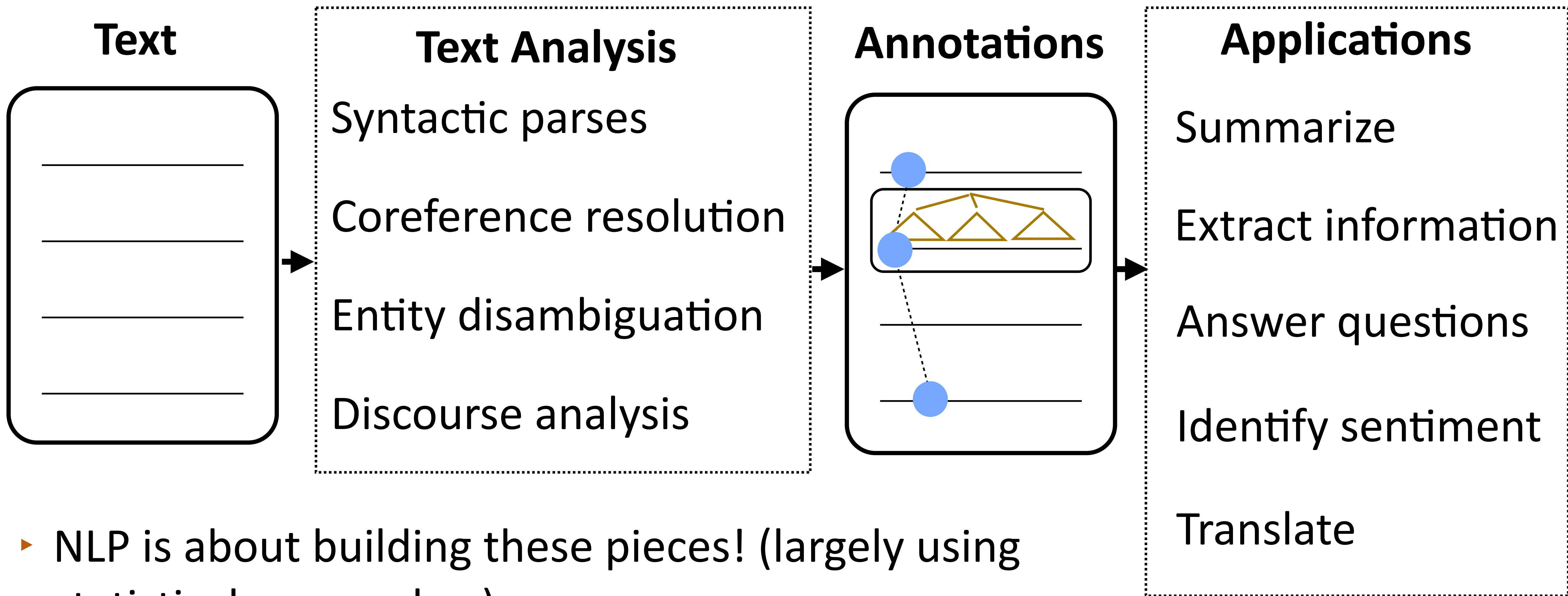
**45,133**  
is correct

Still useful to think  
about capabilities along  
different tasks/domains.





# Classical NLP Analysis Pipeline



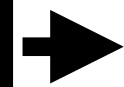
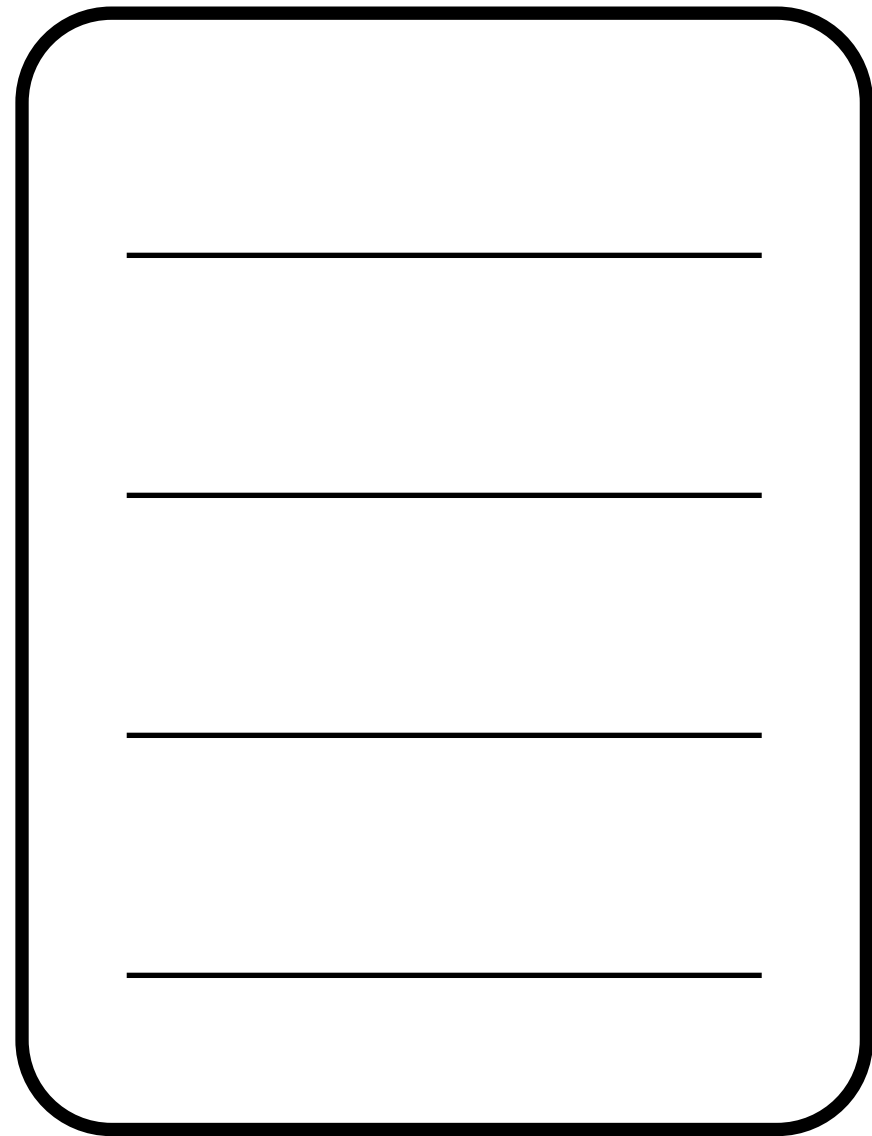
▶ NLP is about building these pieces! (largely using statistical approaches)

▶ Lots of this is done end-to-end with neural nets. But analysis is still useful...



# How do we represent language?

## Text



## Labels

*the movie was good* +

*Beyoncé had one of the best videos of all time* **subjective**

## Sequences

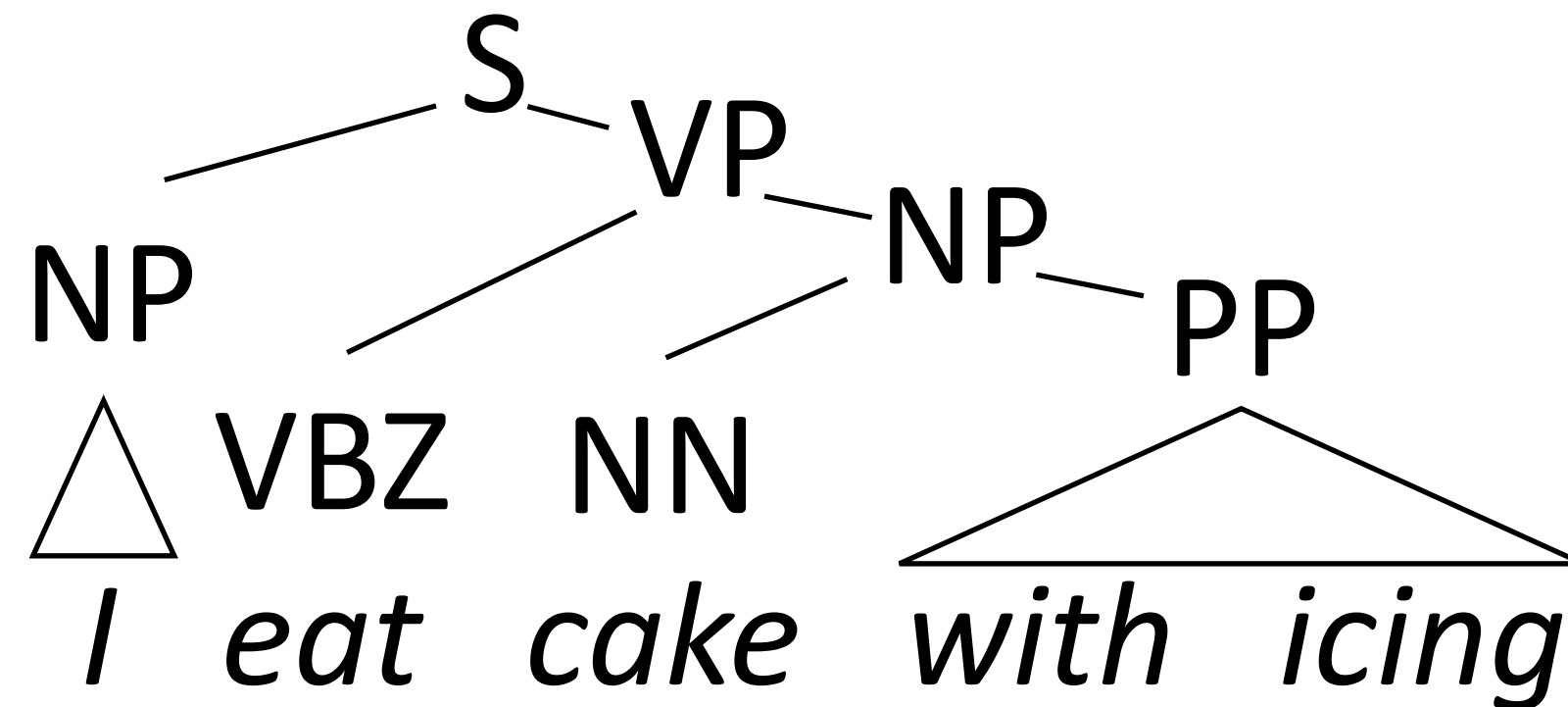
/tags

**PERSON**

*Tom Cruise stars in the new Mission Impossible film*

**WORK\_OF\_ART**

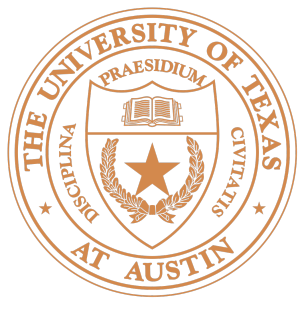
## Trees



$\lambda x. \text{flight}(x) \wedge \text{dest}(x)=\text{Miami}$   
*flights to Miami*

- ▶ Question: What ambiguities do these representations need to help us resolve?

Why is language hard?  
(and how can we handle that?)



# Language is Ambiguous!

- ▶ Hector Levesque (2011): “Winograd schema challenge” (named after Terry Winograd, the creator of SHRDLU)

The city council refused the demonstrators a permit because they advocated violence

The city council refused the demonstrators a permit because they feared violence

The city council refused the demonstrators a permit because they \_\_\_\_\_ violence

- ▶ >5 datasets in the last few years examining this problem and commonsense reasoning
- ▶ Referential ambiguity



# Language is Ambiguous!

---

Teacher Strikes Idle Kids

Ban on Nude Dancing on Governor's Desk

Iraqi Head Seeks Arms

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



# Language is Really Ambiguous!

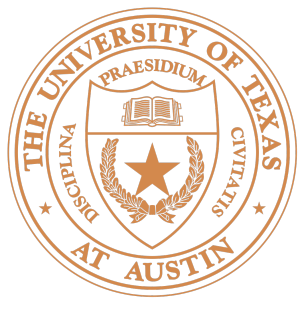
---

- ▶ There aren't just one or two possibilities which are resolved pragmatically

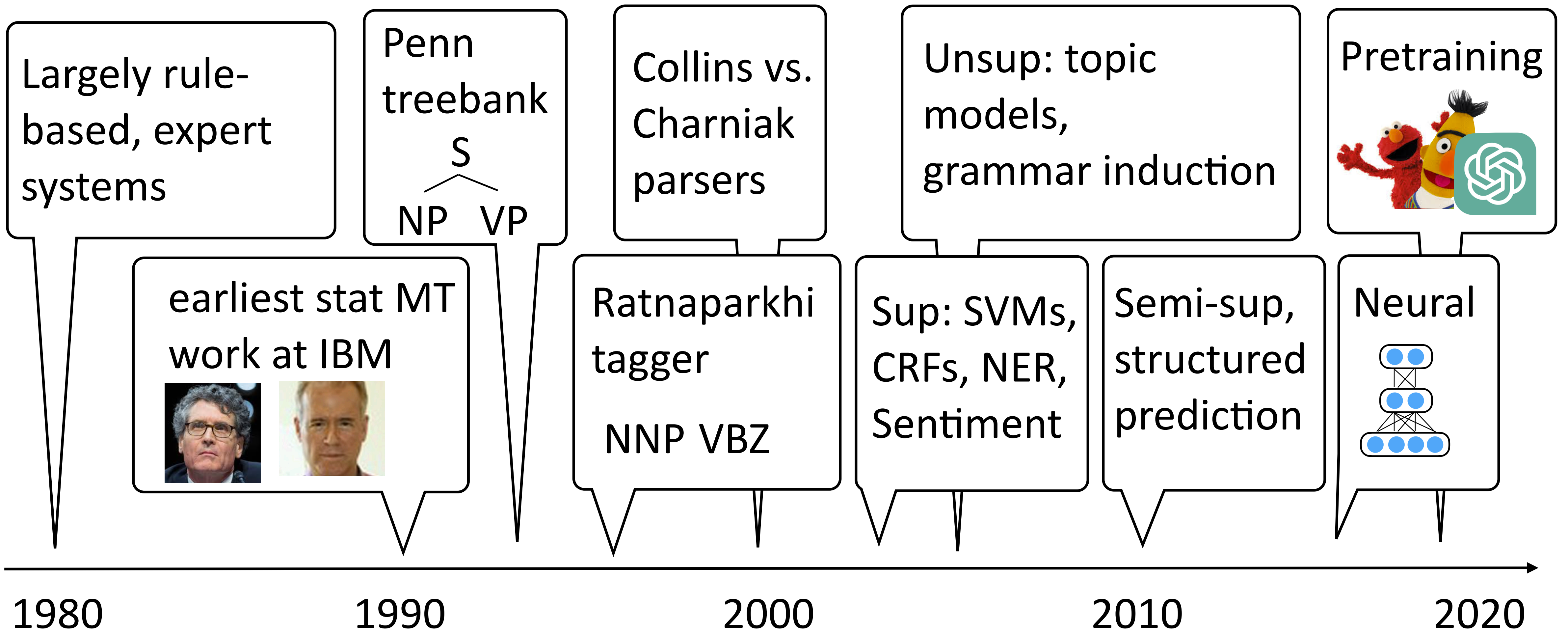
*il fait vraiment beau* → It is really nice out  
It's really nice  
The weather is beautiful  
It is really beautiful outside  
He makes truly beautiful  
It fact actually handsome

- ▶ Combinatorially many possibilities, many you won't even register as ambiguities, but systems still have to resolve them

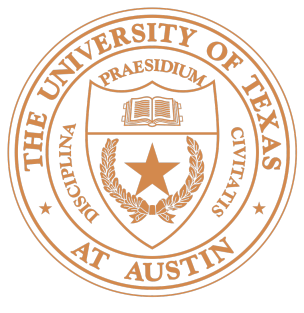
What techniques do we use?  
(to combine data, knowledge, linguistics, etc.)



# A brief history of (modern) NLP







# Pretraining

- ▶ Language modeling: predict the next word in a text  $P(w_i | w_1, \dots, w_{i-1})$

$P(w | \text{I want to go to}) = 0.01 \text{ Hawai'i}$

0.005 LA

0.0001 class



: use this model for other purposes

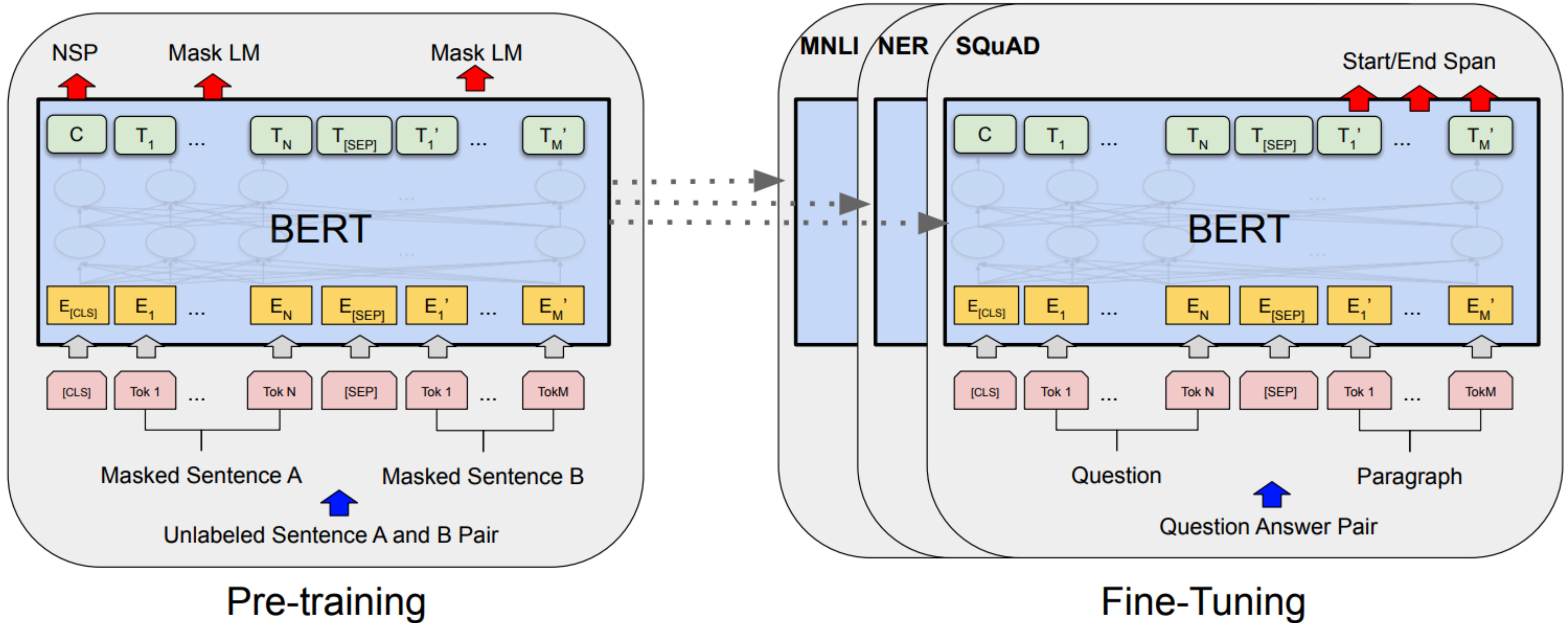
$P(w | \text{the acting was horrible, I think the movie was}) = 0.1 \text{ bad}$

0.001 good

- ▶ Model understands some sentiment?
- ▶ Train a neural network to do language modeling on massive unlabeled text, fine-tune it to do {tagging, sentiment, question answering, ...}



# BERT



- ▶ Key parts which we will study: (1) Transformer architecture; (2) what data is used (both for pre-training and fine-tuning)

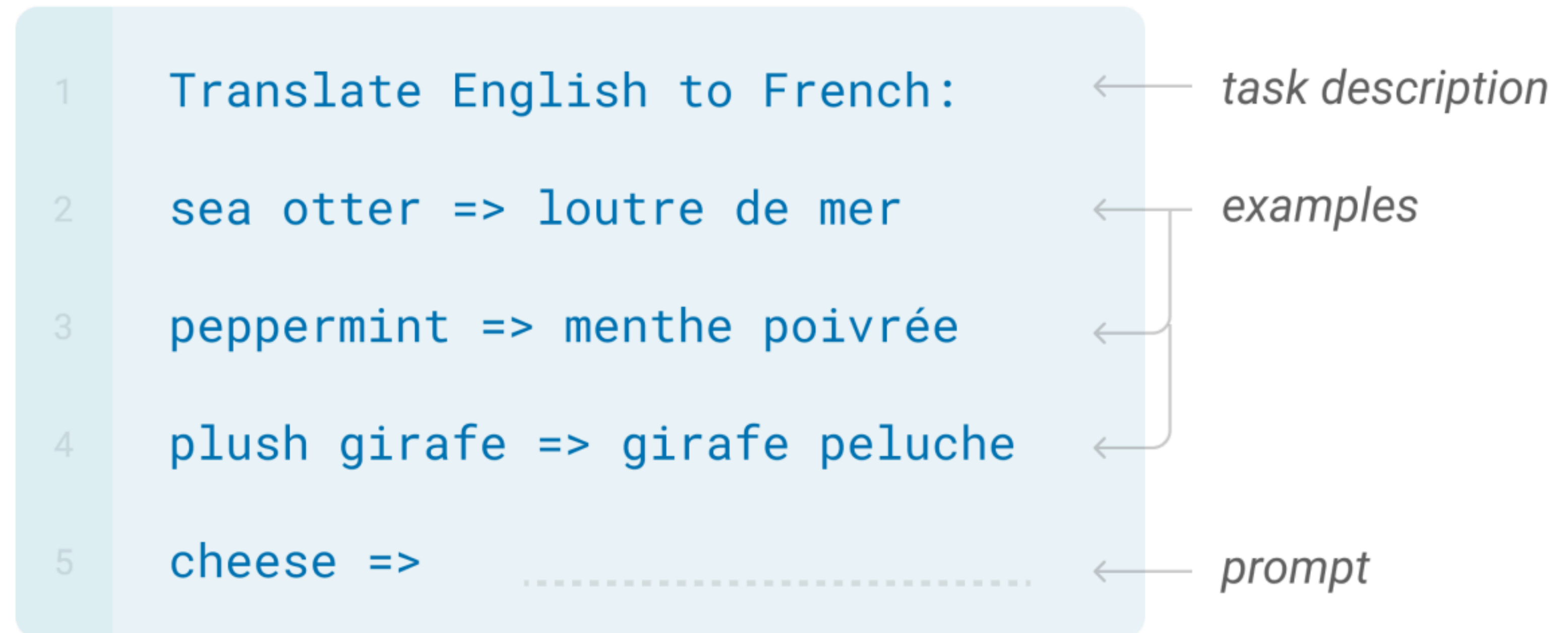


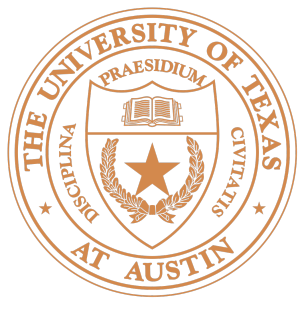
# GPT and In-Context Learning

- ▶ Even more “extreme” setting: no gradient updates to model, instead large language models “learn” from examples in their context
- ▶ Many papers studying why this works. We will read some!

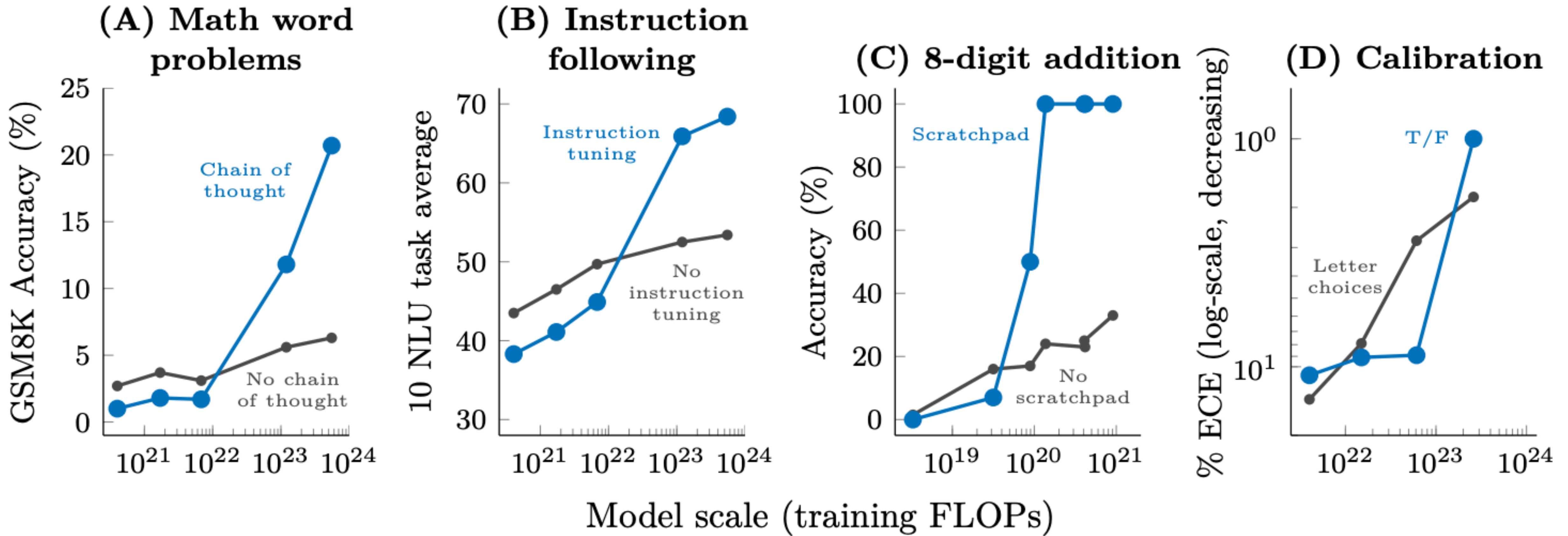
## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.





# Scaling Laws



- ▶ Many of the methods that work in LLMs today only make sense and only work because the models are so big!



# Where are we?

---

- ▶ We have very powerful neural models that can fit lots of datasets
- ▶ Data: we need data that is not just correctly labeled, but reflects what we actually want to be able to do
- ▶ Users: systems are not useful unless they do something we want
- ▶ Language/outreach: who are we building this for? What languages/dialects do they speak?



# Social Impact

- ▶ NLP systems are increasingly used in the world



...and increasingly we have to reckon with their impact



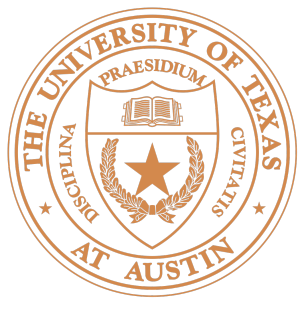
- ▶ This lecture: let's warm up by thinking about these issues a bit



# Social Impact

---

- ▶ Rate your awareness of the social impact of NLP, AI, and machine learning from 1 to 5, where 1 is little awareness and 5 is strong awareness (5 = you feel like you could write a blog post about a current issue).
- ▶ Describe one scenario where you think deployment of an NLP system might pose ethical challenges *due to the application* itself (i.e., using NLP to do “bad stuff”)
- ▶ Describe one scenario where you think deployment of an NLP system might pose ethical challenges due to *unintended* consequences (e.g., unfairness, indirectly causing bad things to happen, etc.).

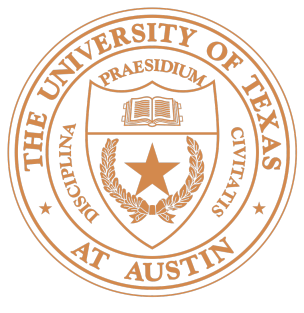


# Outline of the Course

---

- ▶ Classification: linear and neural, word representations (3.5 weeks)
- ▶ Language modeling, Transformers, pre-training (2.5 weeks)
- ▶ Tagging, parsing, and linguistic structure (2 weeks, ending in midterm)
- ▶ Modern pre-trained models, ChatGPT, etc. (2.5 weeks)
- ▶ Applications, modern topics, and ethics (2.5 weeks)
- ▶ Goals:
  - ▶ Cover fundamental techniques used in NLP
  - ▶ Understand how to look at language data and approach linguistic phenomena
  - ▶ Cover modern NLP problems encountered in the literature: what are the active research topics in 2023?





# Coursework

---

- ▶ Five assignments, worth 40% of grade
  - ▶ Mix of writing and implementation;
  - ▶ Assignment 0 is out now, optional diagnostic
  - ▶ ~2 weeks per assignment except for A4
  - ▶ 5 “slip days” throughout the semester to turn in assignments 24 hours late
  - ▶ Submission on Gradescope

These assignments require understanding the concepts, writing performant code, and thinking about how to debug complex systems. **They are challenging; start early!**

Office hours: please come! However, **the course staff are not here to debug your code!** We **will** help you understand the concepts and come up with debugging strategies!



# Coursework

---

- ▶ Midterm (25% of grade), take-home
  - ▶ Similar to written homework problems
- ▶ Final project (25% of grade)
  - ▶ Groups of 1 or 2
  - ▶ Standard project: understanding dataset biases
  - ▶ Independent projects are possible: these must be proposed earlier (to get you thinking early) and will be held to a high standard!
- ▶ Social Impact Responses, UT Instapoll (10% of the grade)
  - ▶ These will be done online and can be done during or after class



# Academic Honesty

---

- ▶ You may work in groups, but your final writeup and code **must be your own**
- ▶ Don't share code with others!



# Conduct



**YOU  
BELONG  
HERE**

**A climate conducive to learning and creating knowledge is the right of every person in our community.** Bias, harassment and discrimination of any sort have no place here.



The University of Texas at Austin  
College of Natural Sciences

*The College of Natural Sciences is steadfastly committed to enriching and transformative educational and research experiences for every member of our community. Find more resources to support a diverse, equitable and welcoming community within Texas Science and share your experiences at [cns.utexas.edu/diversity](https://cns.utexas.edu/diversity)*



# Survey

---

- ▶ See Instapoll (you can answer later as well)