

CS 7650: Natural Language Processing

Alan Ritter

Administrivia

- ▶ Course website:
<https://aritter.github.io/CS-7650-sp25/>
- ▶ Piazza and Gradescope: links on the course website
 - ▶ We will do our best to answer questions within 24 hours (or Monday for questions asked over the weekend).
 - ▶ **Please Sign Yourself up for Piazza**
- ▶ TA Office hours:
 - ▶ See spreadsheet
<https://tinyurl.com/7650-TAs>

Instructor



[Alan Ritter](#)

alan.ritter@cc.gatech.edu

Teaching Assistants

[Duong Le](#)

dminh6@gatech.edu

[Jeongrok Yu](#)

jyu673@gatech.edu

[Junmo Kang](#)

junmo.kang@gatech.edu

[Mohit Raghavendra](#)

mraghavendra6@gatech.edu

[Ray Yang](#)

ryang396@gatech.edu

Prerequisites

- ▶ Probability/Statistics
- ▶ Linear Algebra
- ▶ Multivariable Calculus
- ▶ Programming / Python experience
- ▶ A prior Machine Learning course will be very helpful

There will be a lot of math and programming!

Coursework

- ▶ 4 Programming Projects (fairly substantial implementation effort)
 - ▶ Logistic Regression
 - ▶ Text classification
 - ▶ Sequence Labeling
 - ▶ Neural chatbot
- ▶ 3 written assignments + midterm exam
 - ▶ Mostly math problems related to ML / NLP
- ▶ Final project
- ▶ Problem Set 0 (background review) is out now and **due Thursday**.

Final Project

- ▶ Final project (20%)
 - ▶ Groups of 3-4.
 - ▶ 1 is possible, but will require more work. Grading requirements are the same for individual projects.
 - ▶ 4 page report.
 - ▶ We will discuss more and have a “project kickoff” later in the semester.

Problem Set 0 (Background Review)

- ▶ Due this Thursday.
- ▶ Background review on probability, linear algebra, calculus.
- ▶ **Waitlisted students:** please submit PS0 by Thursday if you plan to enroll in the course.
 - ▶ We can't predict whether or not you will get in, as this depends on other students dropping the class...
- ▶ Submit on Gradescope
 - ▶ If you don't have access to Gradescope, send the course staff a private message on Piazza, and we will give you the access code.

Project 0 is also out (please look!)



▼ Logistic Regression

CS 4650 "Natural Language Processing" Project 0

Georgia Tech, Fall 2023 (Instructor: Alan Ritter)

In this assignment, we will walk you through the process of implementing logistic regression from scratch. You will also apply your implemented logistic regression model to a small dataset and predict whether a student will be admitted to a university. This dataset will allow you to visualize the data and debug more easily. You may find [this documentation](#) very helpful (although it uses Octave rather than Python).

This assignment also serves as a programming preparation test. We will use [NumPy](#) -- a popular Python package for scientific computing and implementing machine learning algorithms. It provides very good support for matrix and vector operations. You need to feel comfortable working with matrices, vectors, and tensors in order to complete all the programming projects in CS 4650.

To start, first make a copy of this notebook to your local drive, so you can edit it.

▼ 4. Gradient Computation [5 points]

Implement the gradient computations for logistic regression.

```
▶ def gradient_update(theta, x, y):
    """ The gradient update for logistic regression"""
    #####
    # Compute the gradient update #
    #####
```

Free Textbooks!



- ▶ 2 excellent textbooks for NLP
 - ▶ There will be assigned readings from both
 - ▶ Both freely available online

Natural Language Processing

Speech and Language Processing (3rd ed. draft)

[Dan Jurafsky](#) and [James H. Martin](#)

Jacob Eisenstein

Not free: GPUs



- ▶ Modern NLP methods require non-trivial computation
 - ▶ Training neural networks with many parameters can take a long time (it is a very good idea to start working on the assignments early!)
 - ▶ This is a big part of modern NLP methods. It is important to get experience training these networks.
 - ▶ You will need to use GPUs to complete the programming assignments.
 - ▶ Google Colab: has free GPUs, but with some big limitations that will make the assignments very difficult to complete.
 - ▶ The programming projects are designed with Colab in mind
 - ▶ Colab Pro subscription (\$10/month). This is highly recommended once we start working with PyTorch.

Outline of the Course

- Machine Learning Review (Naive Bayes, Log. Reg. SVMs, Neural Nets)
- Sequence Models (HMMs, CRFs)
- Word Embeddings
- Neural Networks in NLP (NBOW, RNNs, CNNs, Transformers)
- Pre-trained Transformers (e.g. BERT, BART, T5, GPT)
- Machine Translation
- Dialogue
- Question Answering

Course Goals

- ▶ Cover fundamental machine learning techniques used in NLP
 - ▶ Deeper understanding of algorithms beyond “how to use ML/NLP libraries”.
- ▶ Know about modern NLP methods: what is the state-of-the-art in 2024?
- ▶ Make you a “producer” rather than a “consumer” of NLP tools

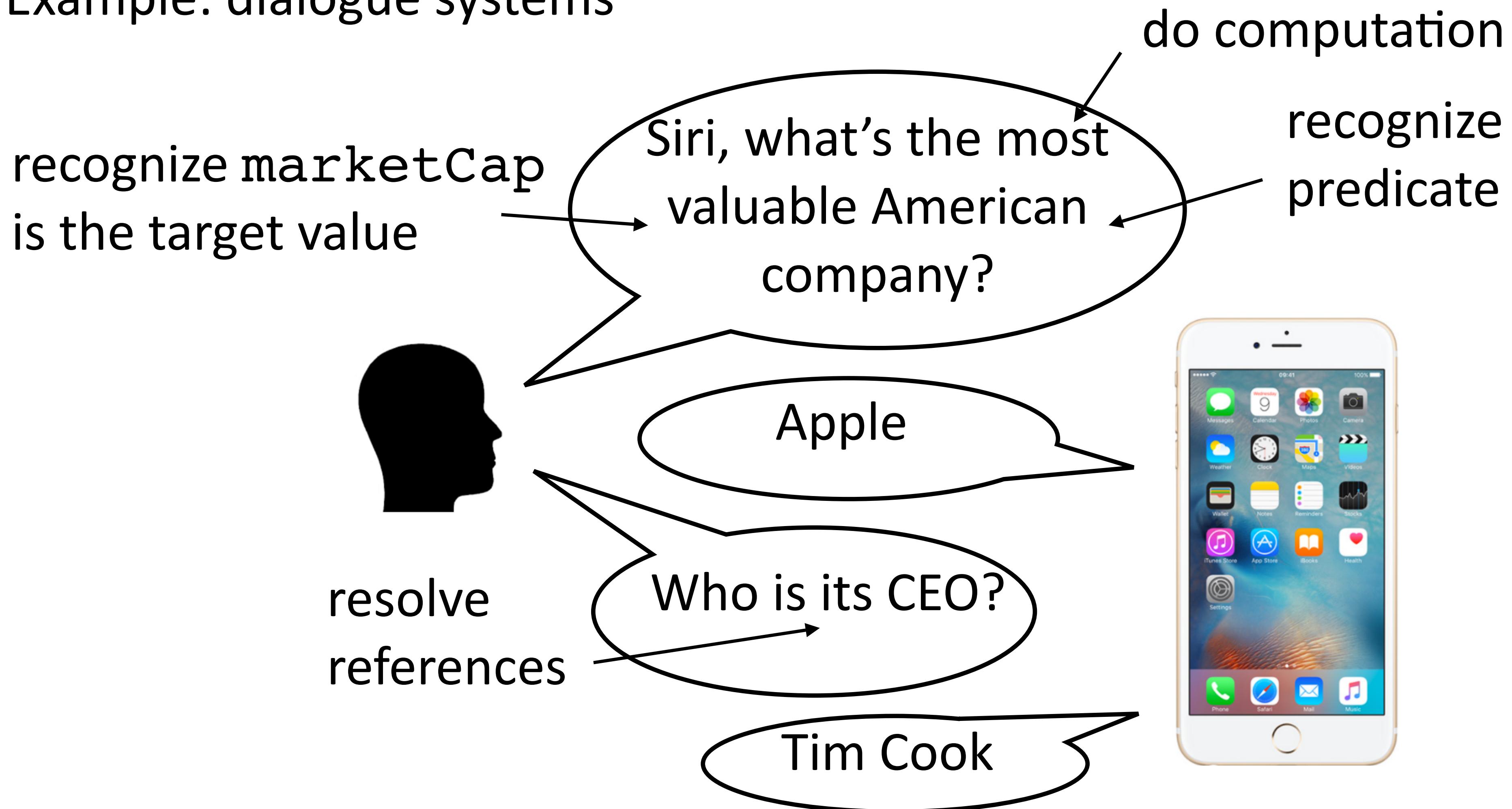
Programming Assignments

- ▶ 3 Programming Assignments
 - ▶ Implementation-oriented
 - ▶ ~2 weeks per assignment

These projects require understanding of the concepts, ability to write performant code, and ability to think about how to debug complex systems. **They are challenging, so start early!**

What's the goal of NLP?

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



Automatic Summarization

POLITICS

Google Critic Ousted From Think Tank Funded by the Tech Giant

WASHINGTON — In the hours after European antitrust regulators levied a record [\\$2.7 billion fine](#) against Google in late June, an influential Washington think tank learned what can happen when a tech giant that shapes public policy debates with its enormous wealth is criticized.

•••

But not long after one of New America's scholars [posted a statement](#) on the think tank's website praising the European Union's penalty against Google, Mr. Schmidt, who had been chairman of New America until 2016, communicated his displeasure with the statement to the group's president,

Anne-Marie Slaughter, according to the scholar.

•••

Ms. Slaughter told Mr. Lynn that “the time has come for Open Markets and New America to part ways,” according to an email from Ms. Slaughter to Mr. Lynn. The email suggested that the entire Open Markets team — nearly 10 full-time employees and unpaid fellows — would be [exiled](#) from New America.

compress
text

provide missing
context

One of New America's writers posted a statement critical of Google. Eric Schmidt, [Google's CEO](#), was displeased.

The writer and his team were dismissed.

paraphrase to provide clarity

Machine Translation



THE WALL STREET JOURNAL.

美国众议院议长选举大戏落幕，
共和党议员重点转向支出及中国
问题

7 小时前



CHINESE (SIMPLIFIED) - DETECTED

CHINESE (SIMPLIFIED)

HINDI

FRE ▾



ENGLISH

SPANISH

ARABIC



美国众议院议长选举大戏落幕，共和党议员重点转向支出及中国问题



U.S. House speaker race ends as Republican lawmakers focus on spending, China



Language Models



Translate this sentence to English: 中国和俄罗斯已在推进一种非正式的政治和经济联盟，共同对抗西方。现在，两国正在通过越来越具挑衅性的联合军事演习来加大双边军事合作。



The sentence translates to:

"China and Russia have been advancing an informal political and economic alliance to jointly oppose the West. Now, the two countries are increasing bilateral military cooperation through increasingly provocative joint military exercises."



Summarize this article in two sentences:



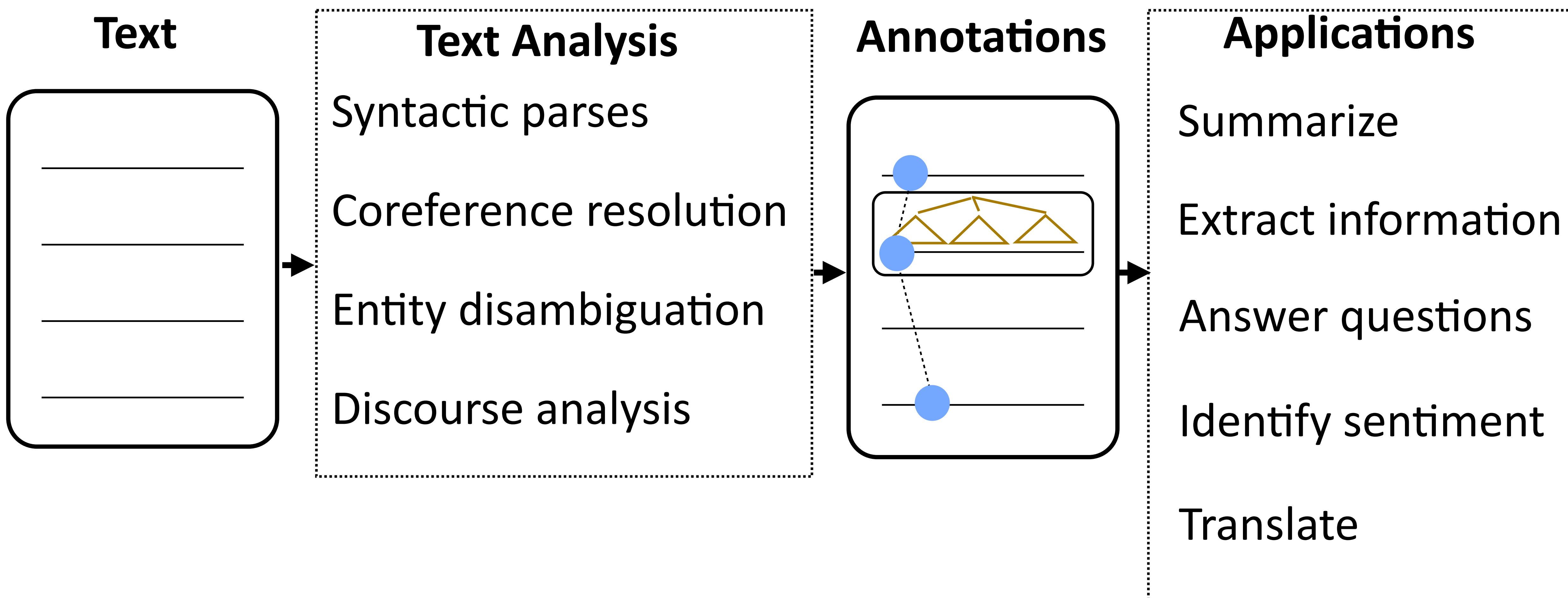
China and Russia have deepened their military cooperation, conducting joint exercises near Alaska and the South China Sea as part of a broader alignment against the United States. This collaboration underscores their shared intent to challenge U.S. dominance, with China also pursuing trade alliances and influence-building to counter American power.



- ▶ LLMs are very general models that can follow natural language instructions and do many NLP tasks.

- ▶ Still useful to think about capabilities in specific tasks / domains.

Traditional NLP Analysis Pipeline



- ▶ All of these components were modeled with statistical approaches trained with machine learning

How do we represent language?

Text

Labels

the movie was good +

Beyoncé had one of the best videos of all time subjective

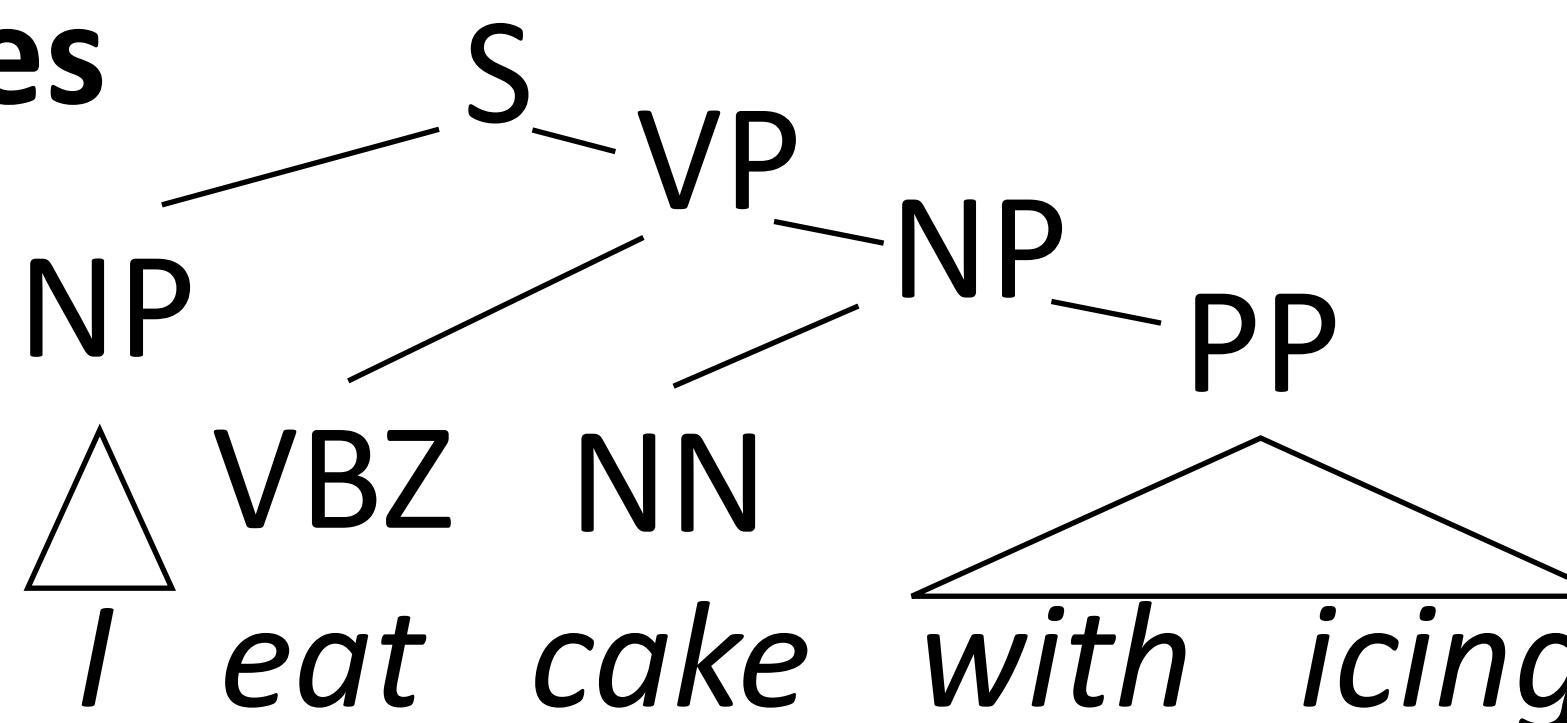
Sequences/tags

PERSON

Tom Cruise stars in the new Mission Impossible film

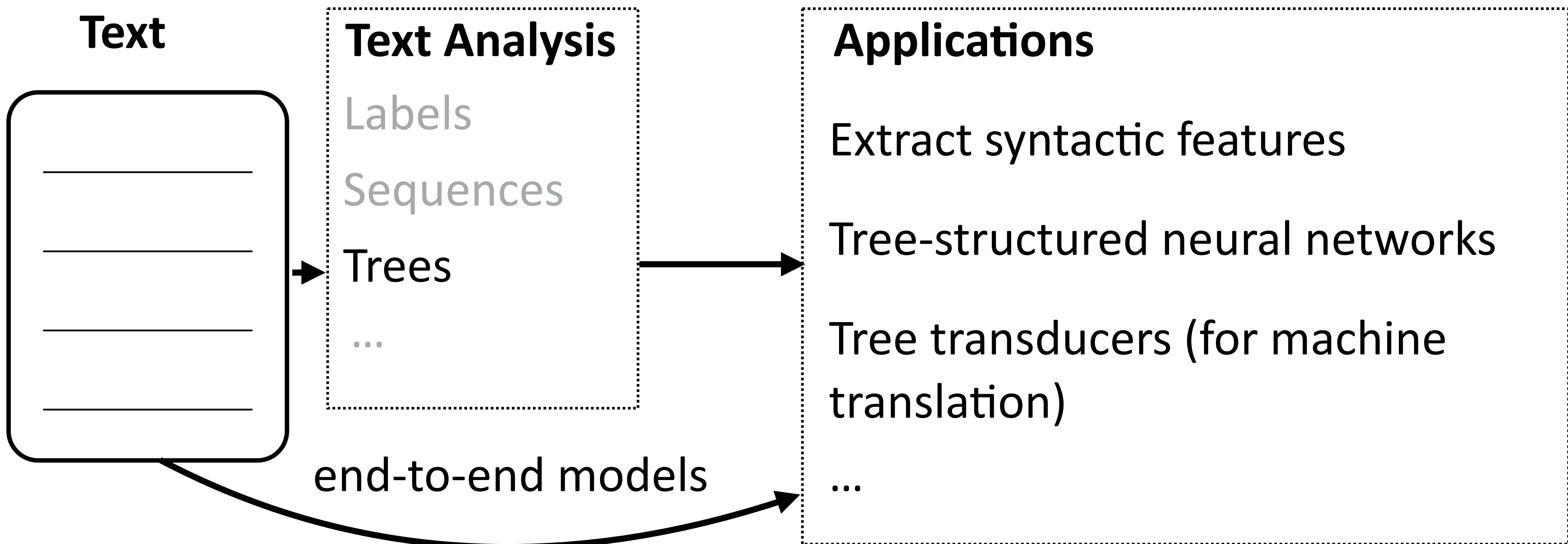
MOVIE

Trees



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

How do we use these representations?



- ▶ Main question: What representations do we need for language? What do we want to know about it?
- ▶ Boils down to: what ambiguities do we need to 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 _____ violence

they advocated

they feared

- ▶ This is so complicated that it's an AI challenge problem! (AI-complete)
- ▶ Referential/semantic ambiguity

Language is Ambiguous!

- ▶ Ambiguous News Headlines:
 - ▶ Teacher Strikes Idle Kids
 - ▶ Hospitals 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
- ▶ Syntactic/semantic ambiguity: 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
He makes truly boyfriend
It fact actually handsome

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

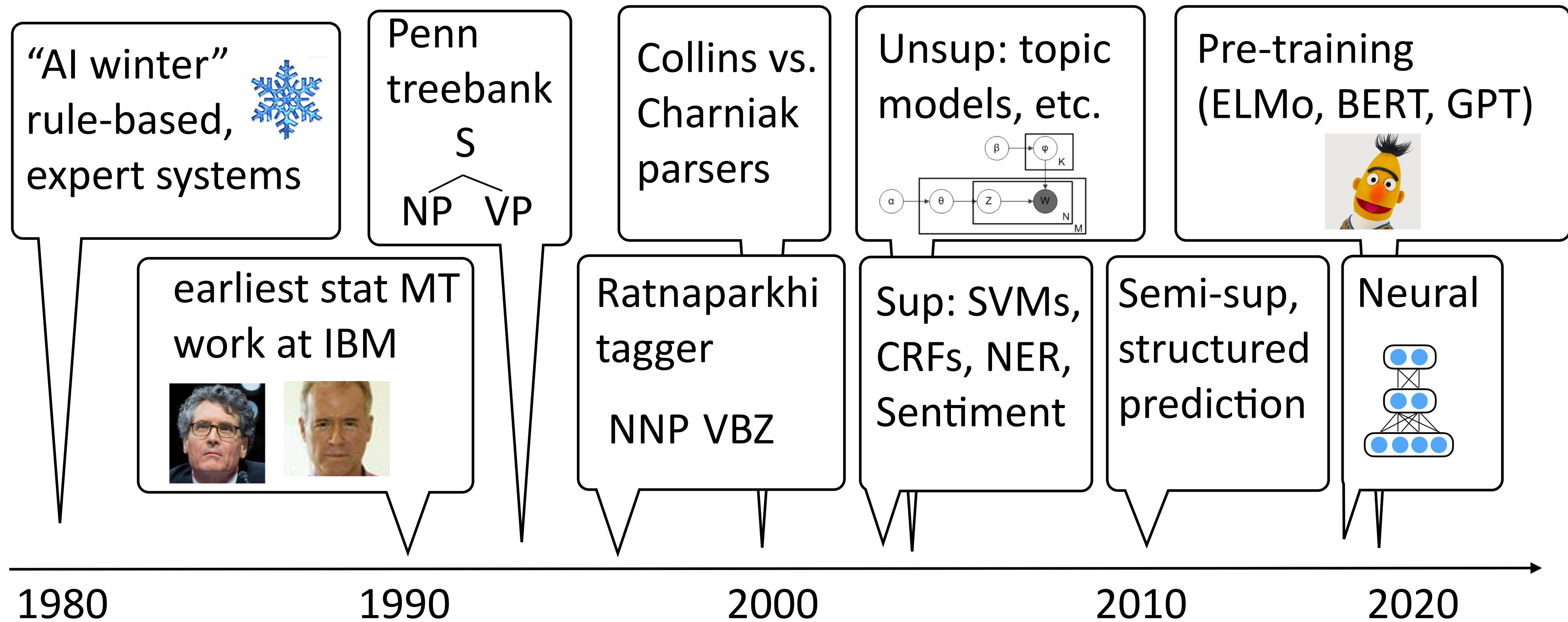
What do we need to understand language?

- ▶ Lots 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] [.]

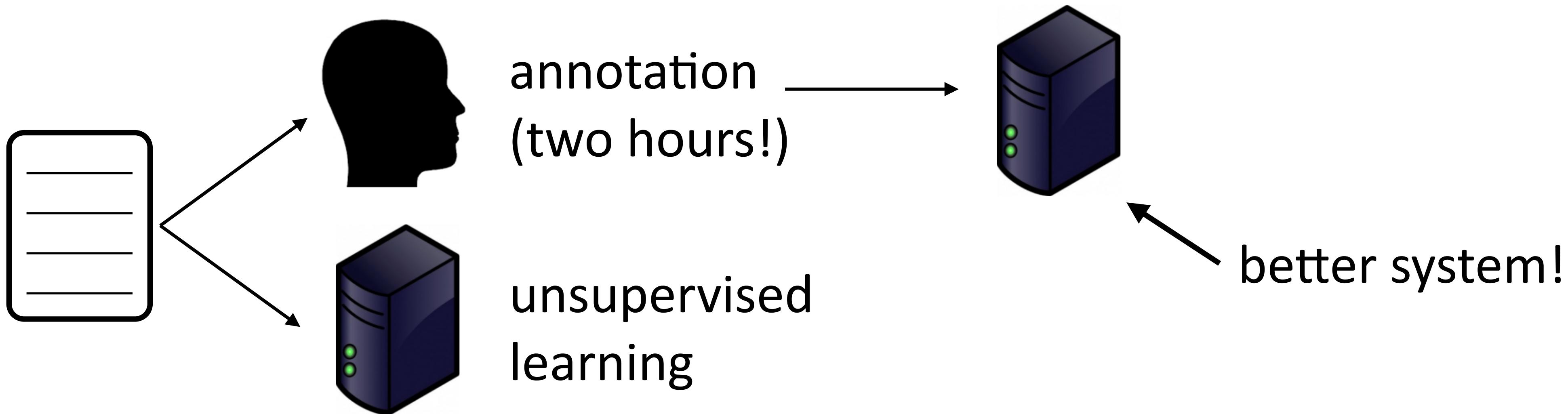
What techniques do we use?
(to combine data, knowledge, linguistics, etc.)

A brief history of (modern) NLP



How Much Training Data do we Need?

- ▶ All of these techniques are data-driven! Some data is naturally occurring, but may need to label
- ▶ Supervised techniques work well on very little data



- ▶ Even neural nets can do pretty well!

“Learning a Part-of-Speech Tagger from Two Hours of Annotation”
Garrette and Baldridge (2013)

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

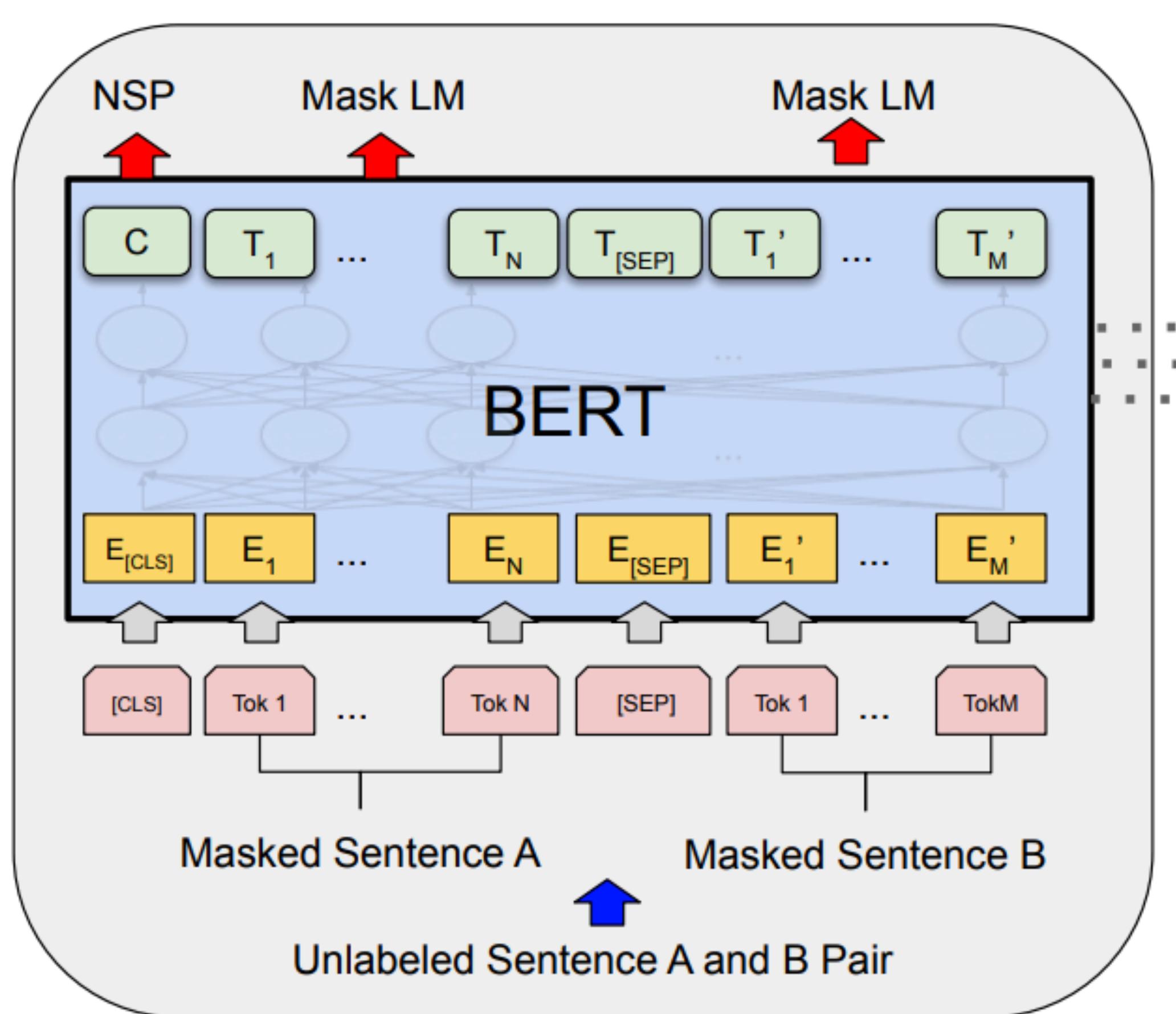
$P(w | \text{the acting was horrible, I think the movie was}) = 0.1 \text{ bad}$

0.001 good

- ▶ Model understands some sentiment?

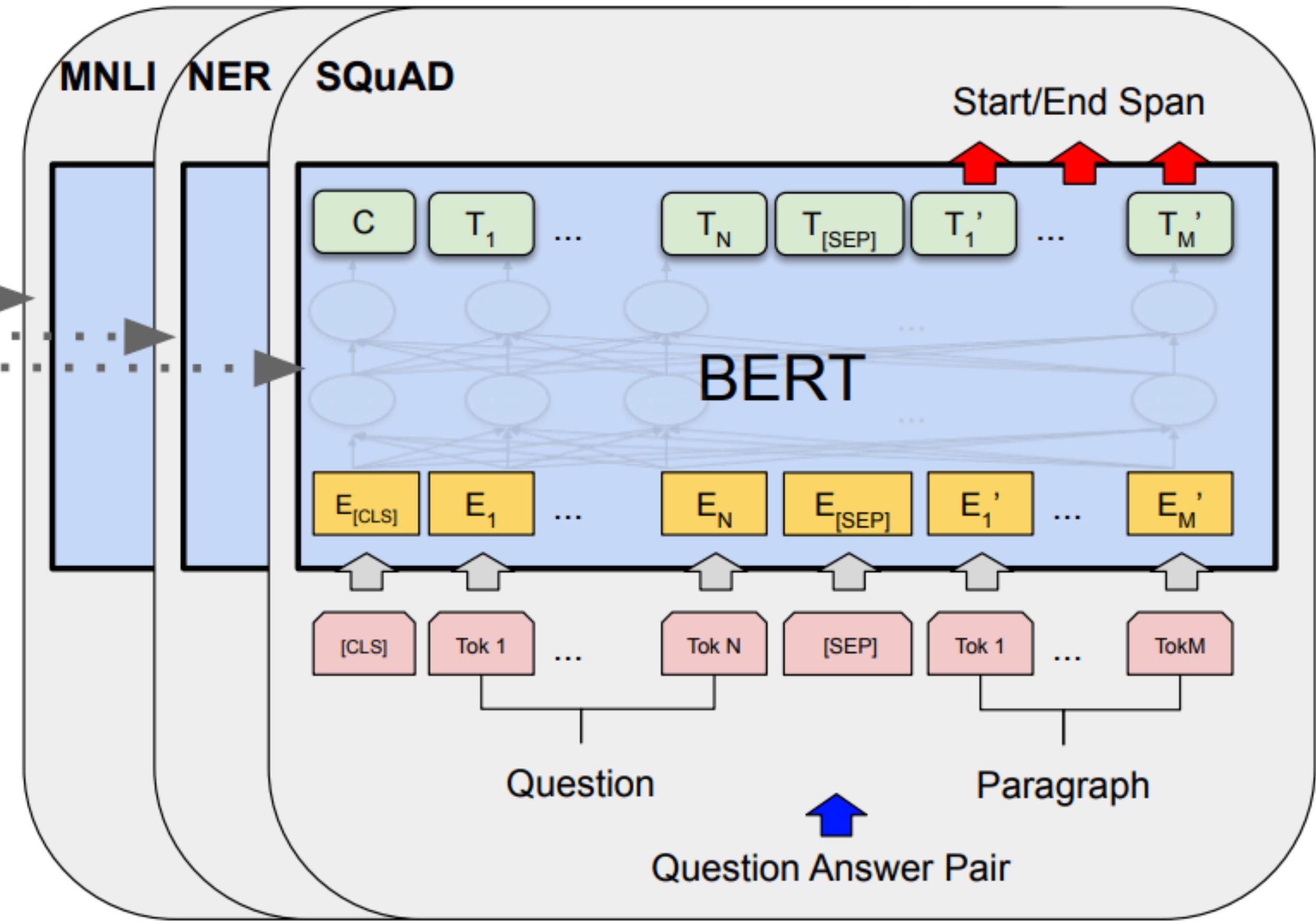
- ▶ Train a neural network to do language modeling on massive unlabeled text, fine-tune it to do {tagging, sentiment, question answering, ...}

BERT



Pre-training

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



Fine-Tuning

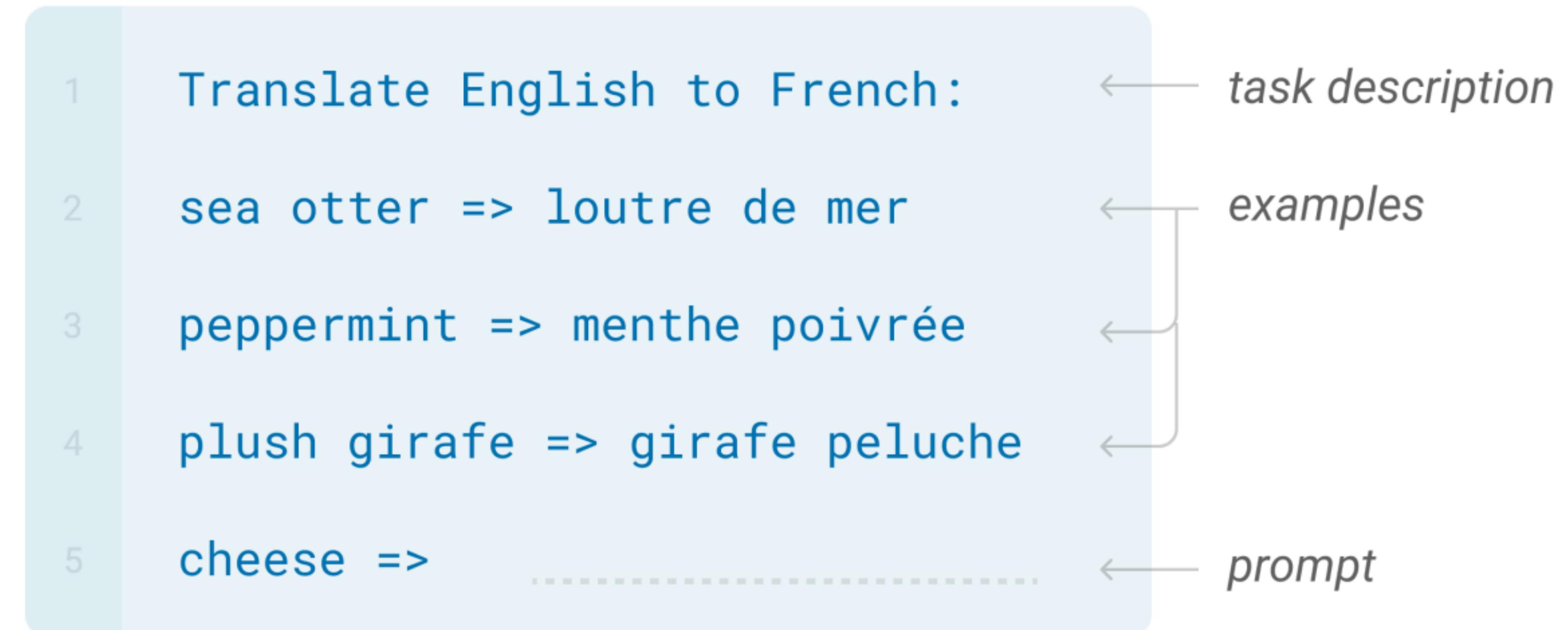
Devlin et al. (2019)

GPT and In-Context Learning

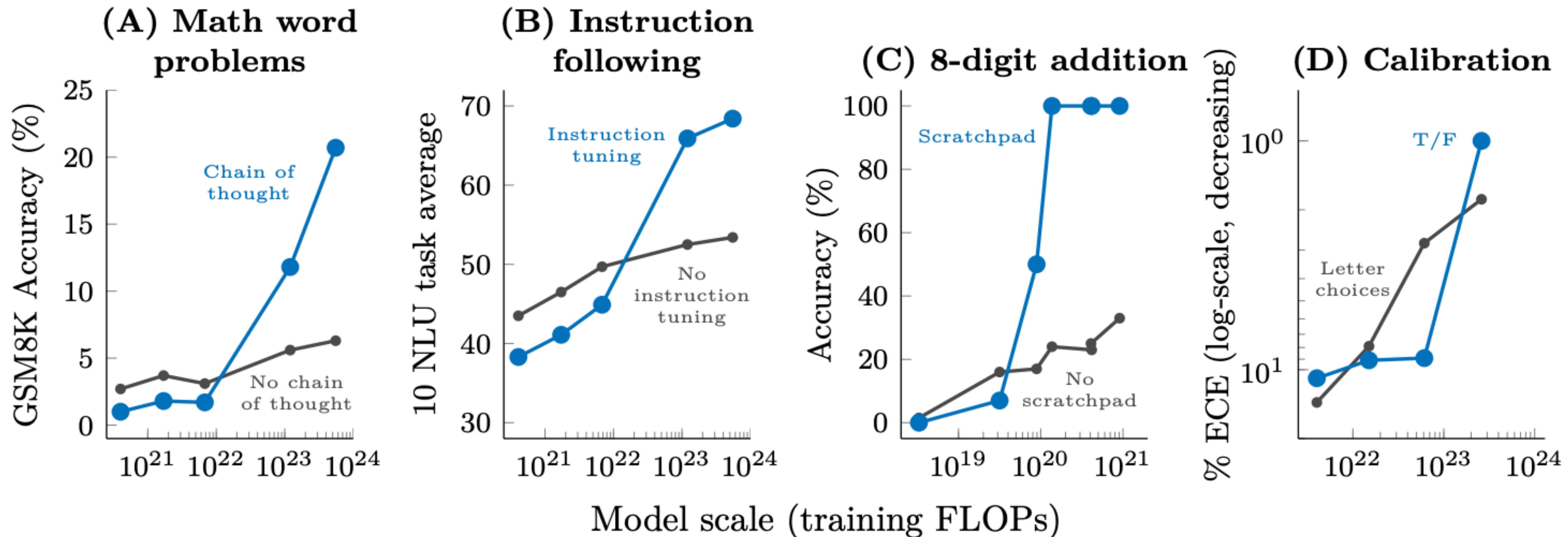
- ▶ Even more “extreme” setting: no gradient updates to model, instead large language models “learn” from examples in their context
- ▶ Many papers studying why this works. We will read some!

Few-shot

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

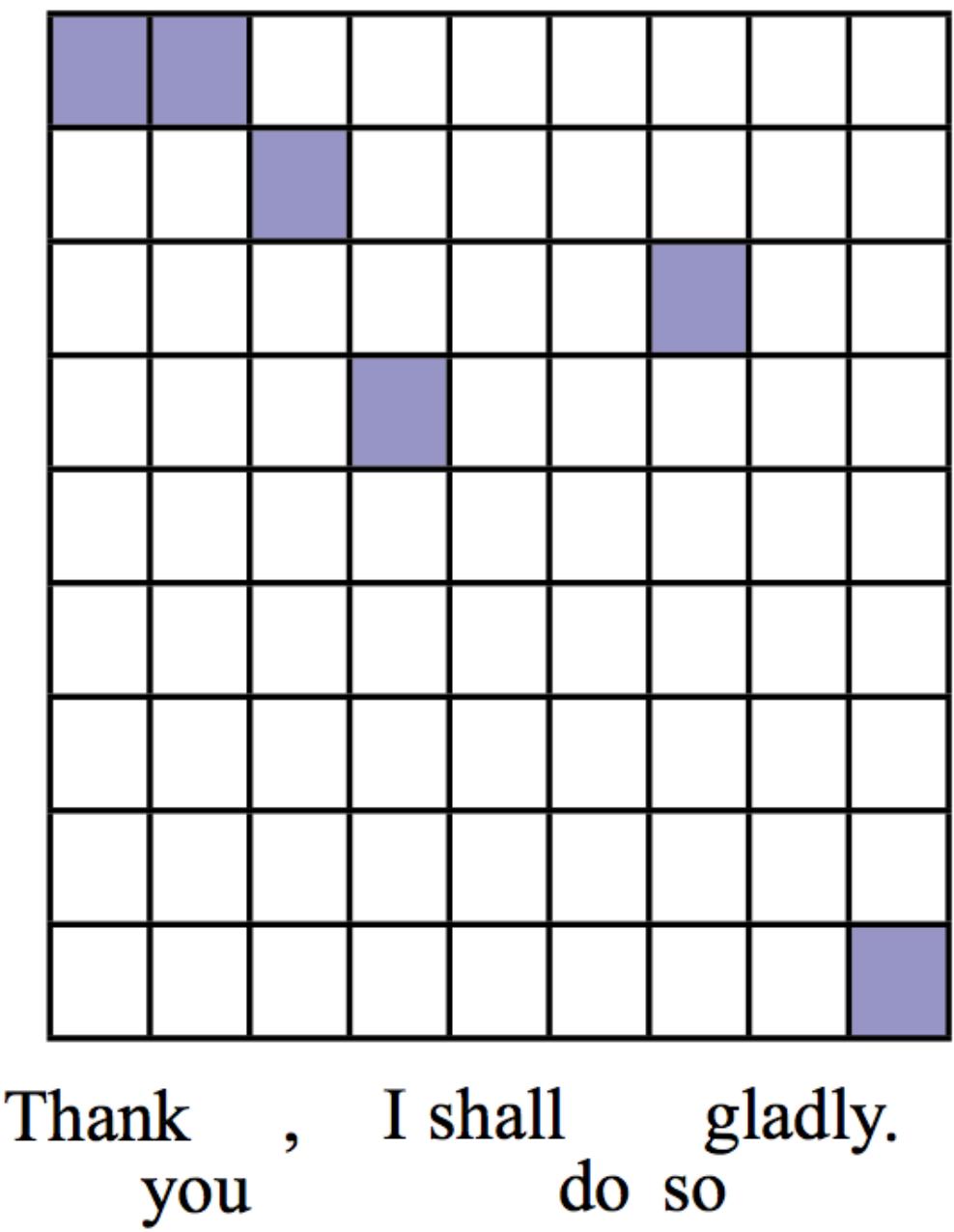


Scaling Laws

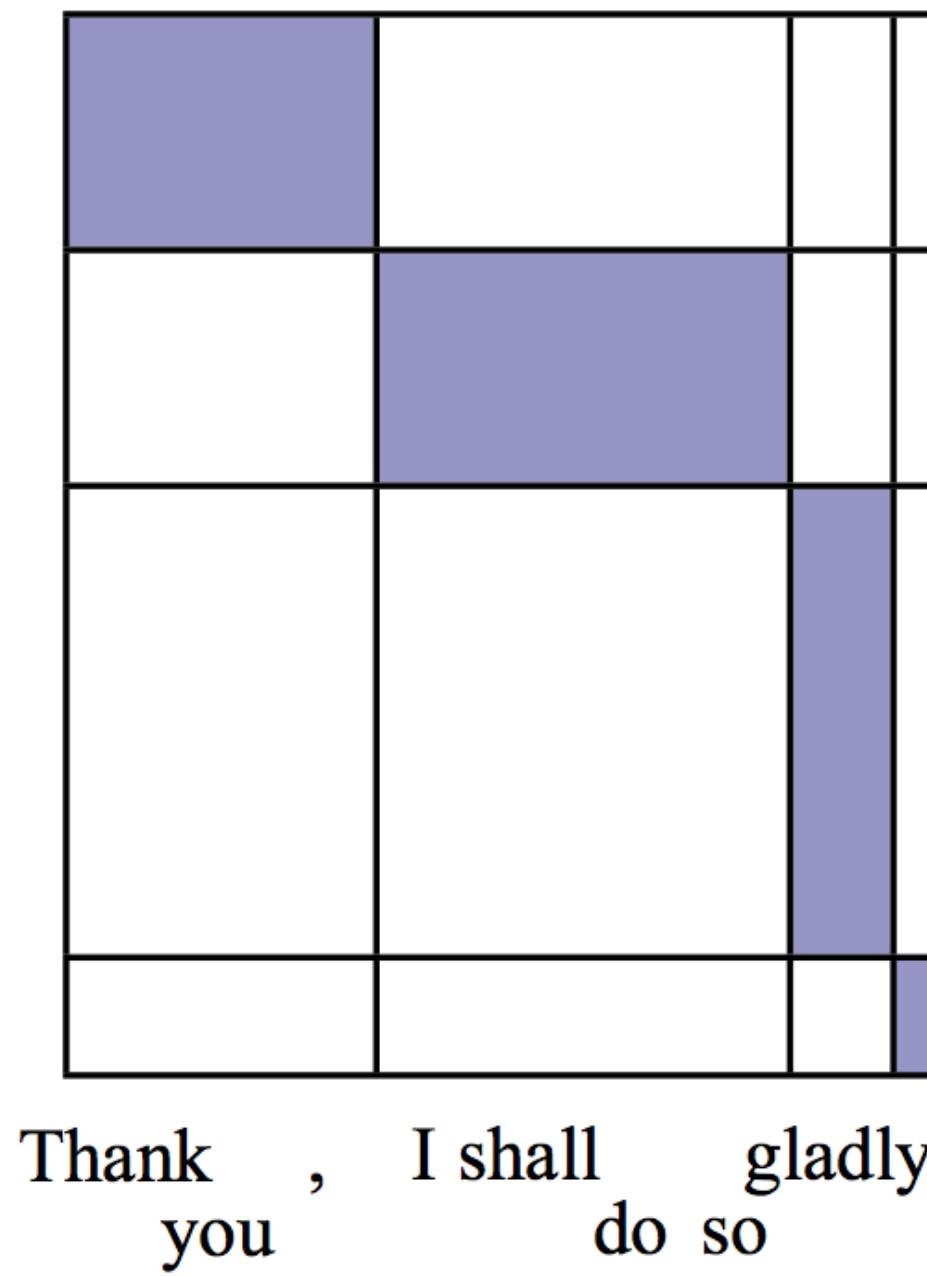


- ▶ Many of the ideas that are big in 2023 only make sense and only work because the models are so big!

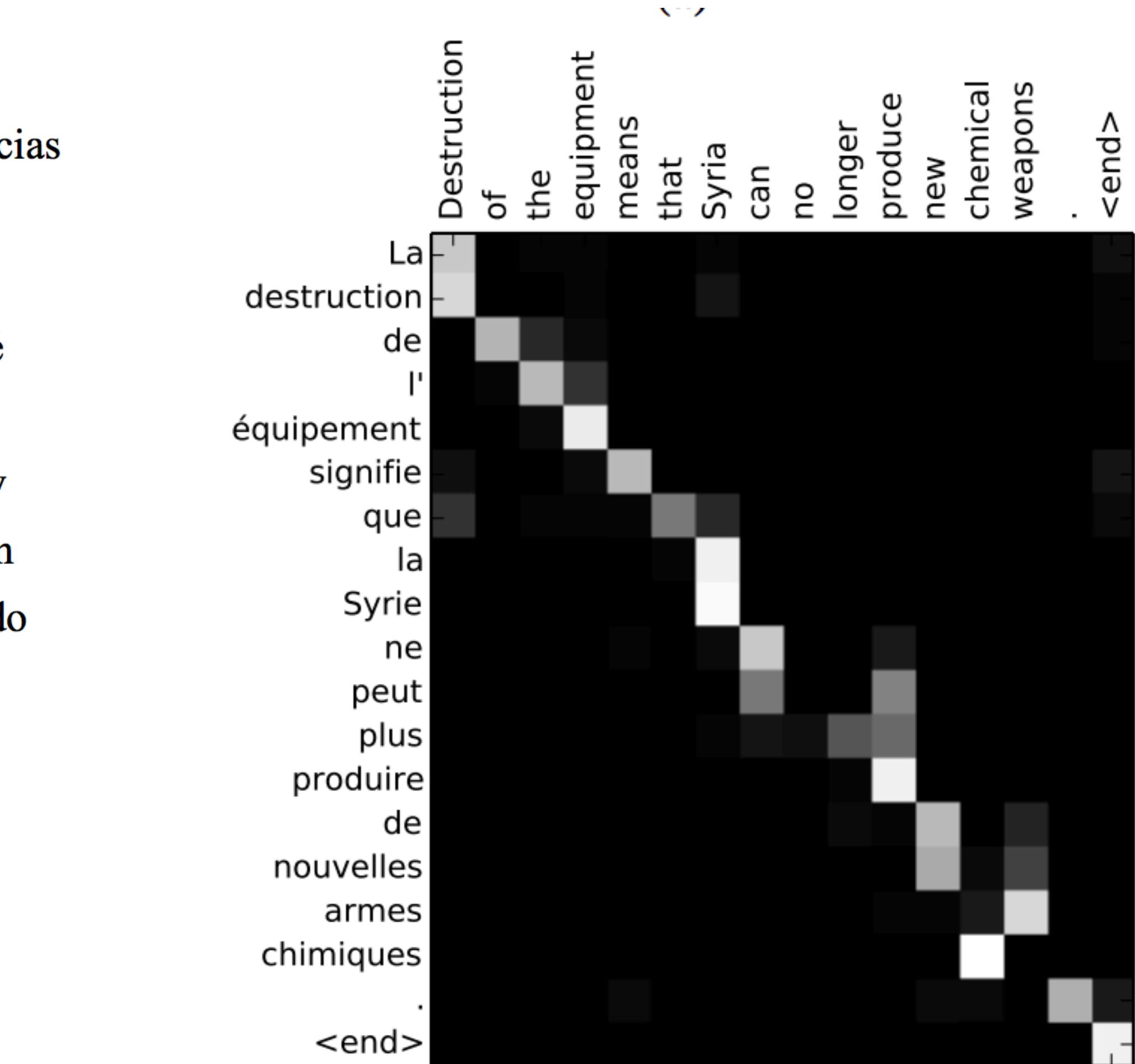
Less Manual Structure?



(a) example word alignment



(b) example phrase alignment



Where are we?

- ▶ NLP consists of: analyzing and building representations for text, solving problems involving text
- ▶ These problems are hard because language is ambiguous, requires drawing on data, knowledge, and linguistics to solve
- ▶ Knowing which techniques to use requires understanding dataset size, problem complexity, and a lot of tricks!
- ▶ NLP encompasses all of these things