



# Introduction to Language Modeling & N-gram Language Models

Slido: <https://app.sli.do/event/d6rnwQ9XXGMjg1Avr5XW1d>

**Yu Meng**  
University of Virginia  
[yumeng5@virginia.edu](mailto:yumeng5@virginia.edu)

Sept 1, 2025

## Announcement: Assignment 1 Out

- Deadline: 09/08 11:59pm
- Released on Canvas

# (Recap) Course Information & Logistics

- Course Website: <https://yumeng5.github.io/teaching/2025-fall-cs4770>
- Instructor: **Yu Meng** ([yumeng5@virginia.edu](mailto:yumeng5@virginia.edu))
  - Office hour: After class Mondays & Wednesdays
- TAs:
  - **Wei-Lin Chen** ([wlchen@virginia.edu](mailto:wlchen@virginia.edu)) Office hour: 5:00pm - 6:00pm every Thursday
  - **Zhepei Wei** ([zhepei.wei@virginia.edu](mailto:zhepei.wei@virginia.edu)) Office hour: 4:00pm - 5:00pm every Wednesday
  - **Xinyu Zhu** ([xinyuzhu@virginia.edu](mailto:xinyuzhu@virginia.edu)) Office hour: 2:00pm - 3:00pm every Tuesday
  - All TA office hours are on Zoom (links on the course website)
- Time: Mondays & Wednesdays 2:00pm - 3:15pm
- Location: Mechanical Engr Bldg 205

# (Recap) Lecture Zoom Options & Recordings

- We provide Zoom options for attending lectures remotely
- Recordings will be available after the lecture

≡ CS\_4770-001 > Online Meetings

---

2025 Fall  Home Appointments

Home Syllabus Piazza Grades SensusAccess **Online Meetings** Course Evaluations Course Email UVA Library Portal

Your current Time Zone and Language are (GMT-4:00) Eastern Time (US and Canada), English ↗

**Upcoming Meetings** Previous Meetings Cloud Recordings

Start Time	Topic
Wed, Aug 27 (Recurring) 2:00 PM	25F Natural Language Processing
Mon, Sep 1 (Recurring) 2:00 PM	25F Natural Language Processing
Wed, Sep 3 (Recurring) 2:00 PM	25F Natural Language Processing

## (Recap) Q&A Format

- Q&A during lecture: Slido (link shared in each lecture & on the course website)
  - Efficient for a big class
  - Good for quick/short questions
  - Allows asking questions anonymously
  - TAs will answer the questions in real time
- Q&A after lecture: Piazza (accessible via Canvas)
  - Assignments/projects
  - Long questions
  - TAs & instructor will answer the questions on a daily basis
- You are encouraged to answer the questions asked by your classmates (participation credit)!

## (Recap) Prerequisites

- Prerequisites:
  - Linear algebra (APMA 3080 or equivalent)
  - Data structures and algorithms (CS 3100)
- Highly recommended background:
  - Deep learning & machine learning (CS 4774)
  - Experience with Python (we'll use **PyTorch** extensively for assignments)
- This class will move fast & cover lots materials! Make sure you have sufficient background before taking it!

## (Recap) Grading

- **Assignments (60%)**
  - Five assignments (with different weights) for the entire semester
  - All assignments are to be completed individually
  - Assignments will be a combination of concept questions + coding questions
  - Submission via Canvas (as LaTeX reports; handwritten submissions not accepted)
  - We'll provide HPC access for GPU-related assignments/projects (instructions later)
- Late day policy:
  - 7 free days for **all** assignments; afterwards 20% off grade of the assignment per day late
  - You cannot use > 3 late days (72 hours) per assignment unless given permission in advance
  - DO NOT procrastinate on assignments! The coding questions (esp. the latter part of this course) take time to implement and run!
- Policy on using LLMs:
  - Collaborative coding with LLMs is allowed, but if you directly copy the answers generated by LLMs (for either conceptual or coding questions), you'll get a 0 for that entire assignment

## (Recap) Grading

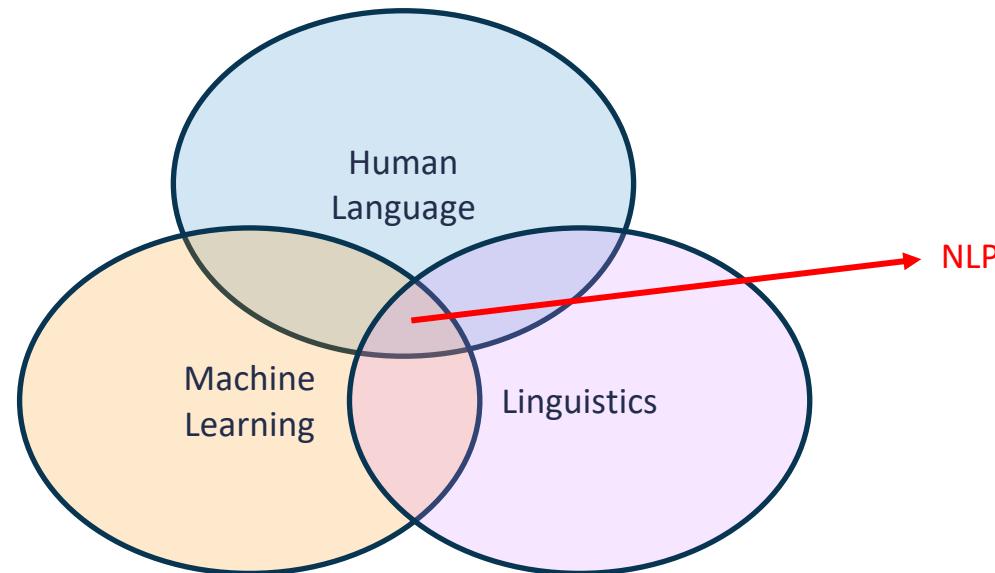
- **Project (35%)**
  - Work in teams of 2–3 students
  - Related to NLP
  - Rule of thumb: demonstrate that you are able to apply the knowledge learned from this course; workload should be more extensive than individual assignments
- Some example project choices:
  - Use word embeddings to analyze sentence semantics (e.g., sentiment analysis)
  - Fine-tune BERT and evaluate its performance for any task you like
  - Benchmark LLMs (either open-weights or proprietary) for challenging tasks
  - Use LLM APIs to create agents for an interesting application (e.g., personal assistants)
  - ...
- Checkpoints (No late dates allowed!)
  - (2%) Project proposal **Deadline:** 09/24
  - (8%) Midterm report **Deadline:** 10/20
  - (25%) Final project presentation **Deadline:** 11/30 + final report **Deadline:** 12/13

## (Recap) Grading

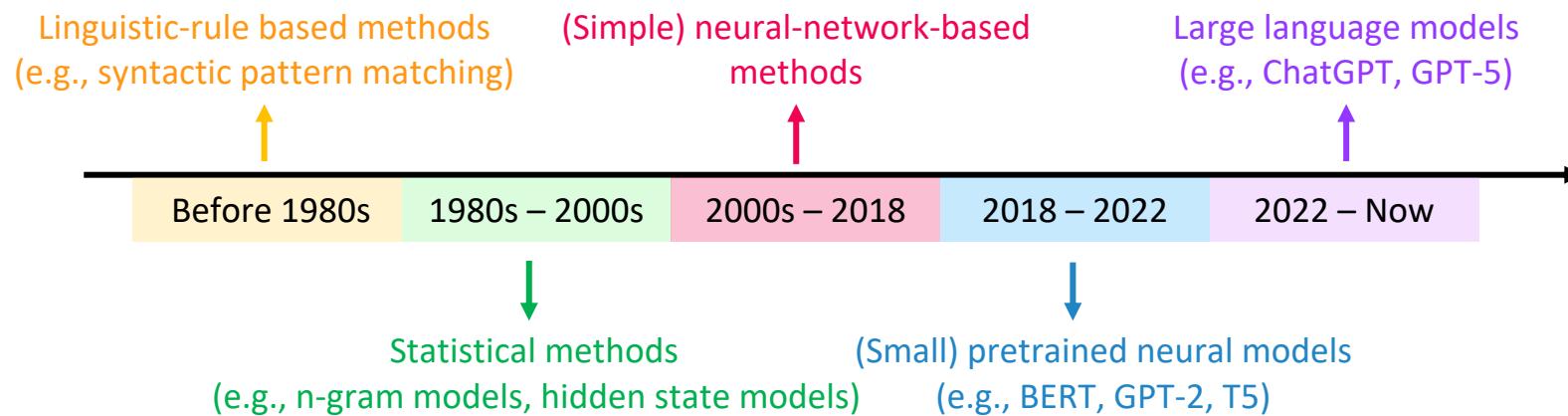
- **Participation (5%+; points earned beyond 5% will become extra credit)**
  - Guest lecture attendance (6%)
  - End-of-semester teaching feedback (2%)
  - Answering **technical** questions raised by classmates (5%)
- Guest lecture attendance on Zoom (6%)
  - We will have 2 guest lectures delivered by leading researchers
  - Each guest lecture can give you up to 3% participation credit (2% attendance + 1% asking questions – more details shared before guest lectures)
- End-of-semester teaching feedback (2%)
  - At the end of the semester, anyone who completes the teaching feedback survey will get 2%
- Answering **technical** questions raised by classmates (5%)
  - We encourage and appreciate help from students to answer questions posted by classmates
  - Every helpful answer to **technical** questions will earn 1% (Slido and Piazza both count)
  - If you answer anonymously, we won't be able to track your contributions!
  - The maximum credit you can get in this category is 5%

## (Recap) What is Natural Language Processing (NLP)?

- An interdisciplinary subfield of machine learning and linguistics
- Goal: Enable computers to understand, interpret, and generate human language



## (Recap) The History of NLP



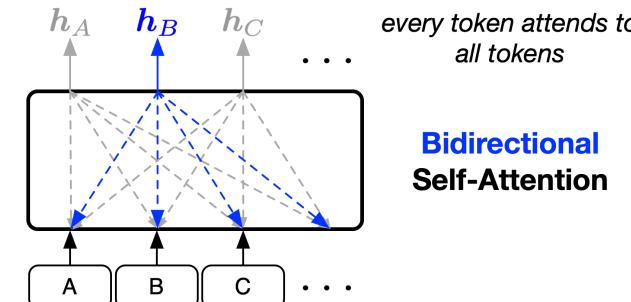
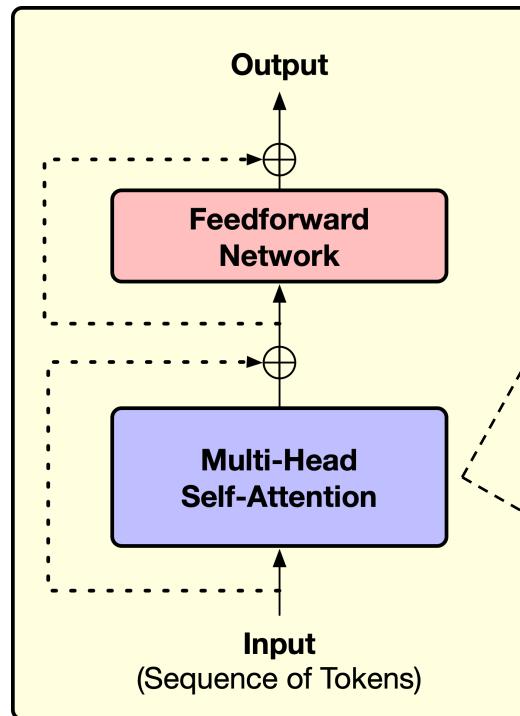
## Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

## Overview of Course Contents

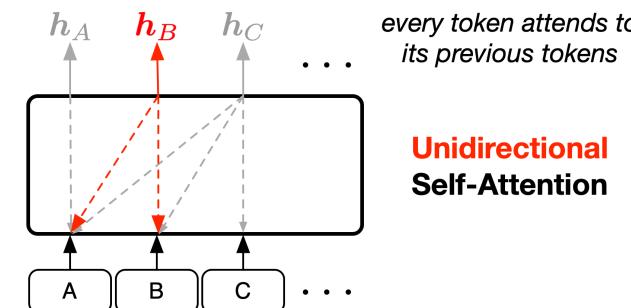
- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- **Week 6: Language Modeling with Transformers**
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

# Transformers



**Bidirectional  
Self-Attention**

Transformer Encoders



**Unidirectional  
Self-Attention**

Transformer Decoders

# Transformer Overview

Transformer block overview

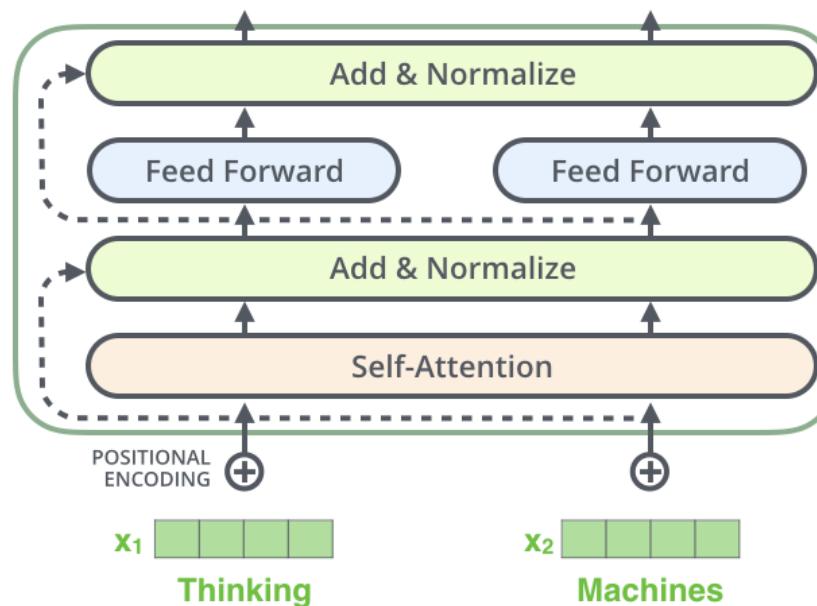


Figure source: <https://jalammar.github.io/illustrated-transformer/>

# Transformer: Self-Attention Mechanism

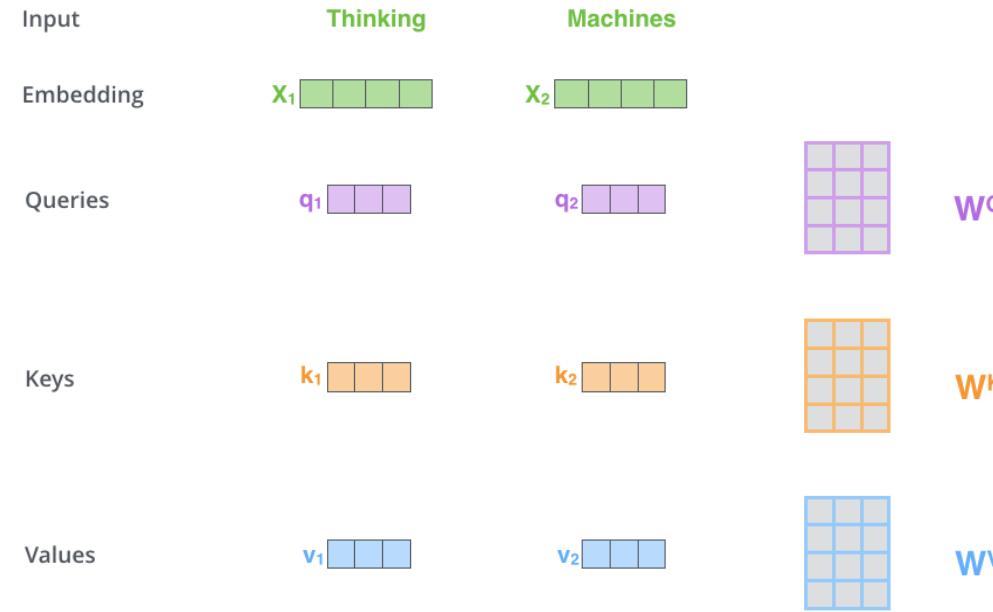
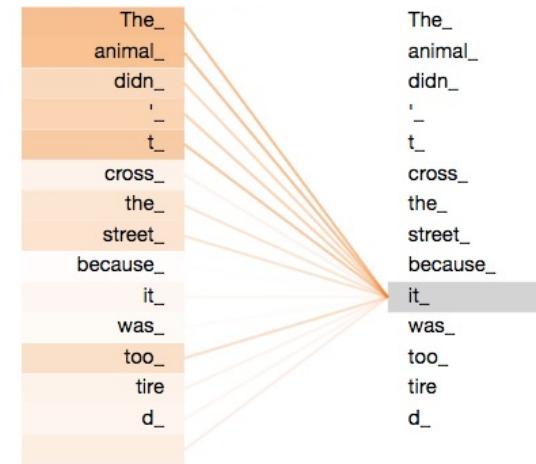


Figure source: <https://jalammar.github.io/illustrated-transformer/>

# Transformer: Self-Attention Computation

$$\text{softmax} \left( \frac{\begin{matrix} \text{Q} & \times & \text{K}^T \\ \begin{matrix} \text{purple grid} \end{matrix} & \times & \begin{matrix} \text{orange grid} \end{matrix} \end{matrix}}{\sqrt{d_k}} \right) \text{V}$$

$=$  Z
  
pink grid



# Language Model Pretraining

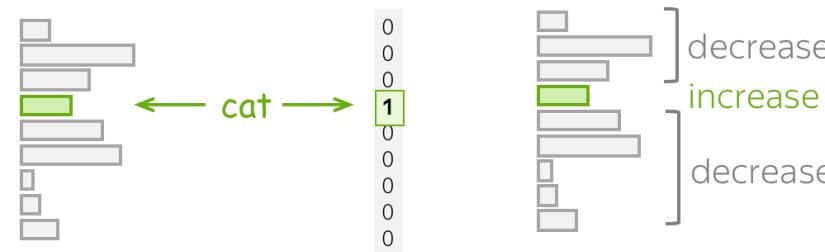
we want the model  
to predict this

Training example: I saw a cat on a mat <eos>

Model prediction:  $p(*) | \text{I saw a}$

Target

Loss =  $-\log(p(\text{cat})) \rightarrow \min$



# Pretraining as Multi-Task Learning

- In my free time, I like to **{run, banana}** (*Grammar*)
- I went to the zoo to see giraffes, lions, and **{zebras, spoon}** (*Lexical semantics*)
- The capital of Denmark is **{Copenhagen, London}** (*World knowledge*)
- I was engaged and on the edge of my seat the whole time. The movie was **{good, bad}** (*Sentiment analysis*)
- The word for “pretty” in Spanish is **{bonita, hola}** (*Translation*)
- $3 + 8 + 4 = \{15, 11\}$  (*Math*)
- ...



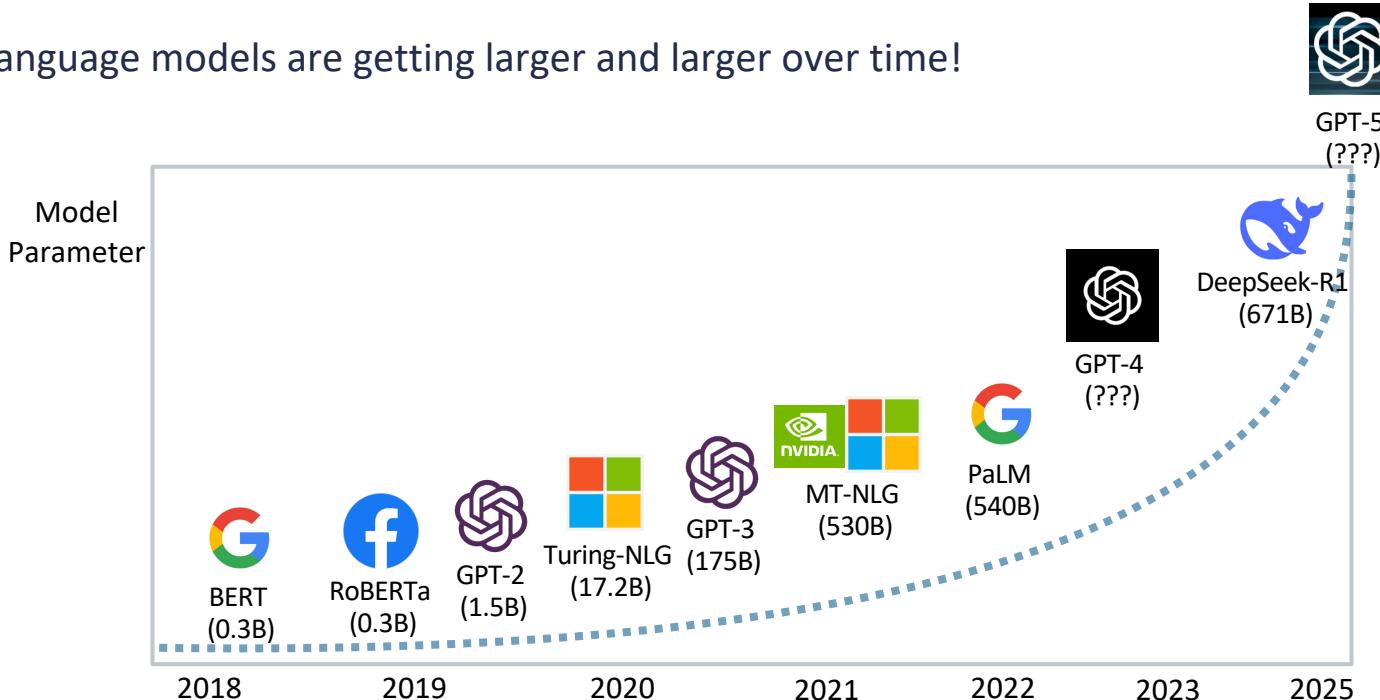
Examples from: [https://docs.google.com/presentation/d/1hQUd3pF8\\_2Gr2Obc89LKjmHL0DIH-uof9M0yFVd3FA4/edit#slide=id.g28e2e9aa709\\_0\\_1](https://docs.google.com/presentation/d/1hQUd3pF8_2Gr2Obc89LKjmHL0DIH-uof9M0yFVd3FA4/edit#slide=id.g28e2e9aa709_0_1)

## Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- **Week 9: Large Language Models (LLMs) & In-context Learning**
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

# Large Language Models (LLMs)

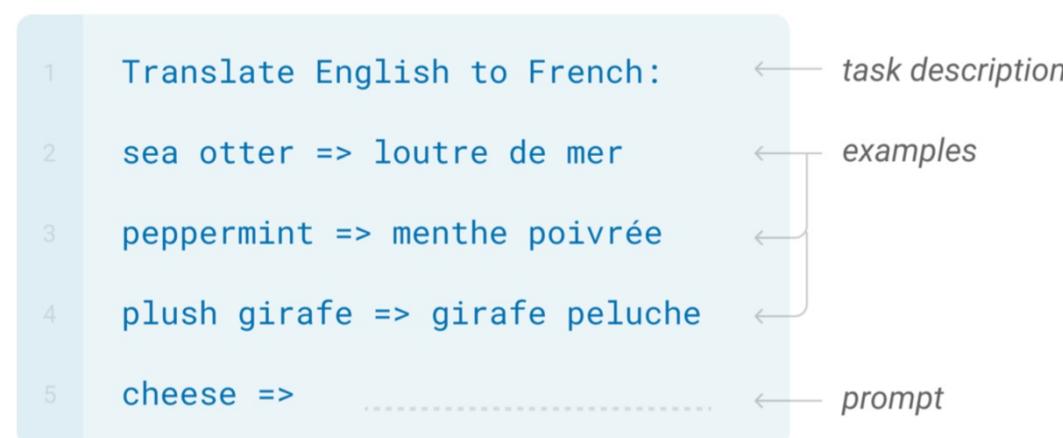
Language models are getting larger and larger over time!



# In-Context Learning

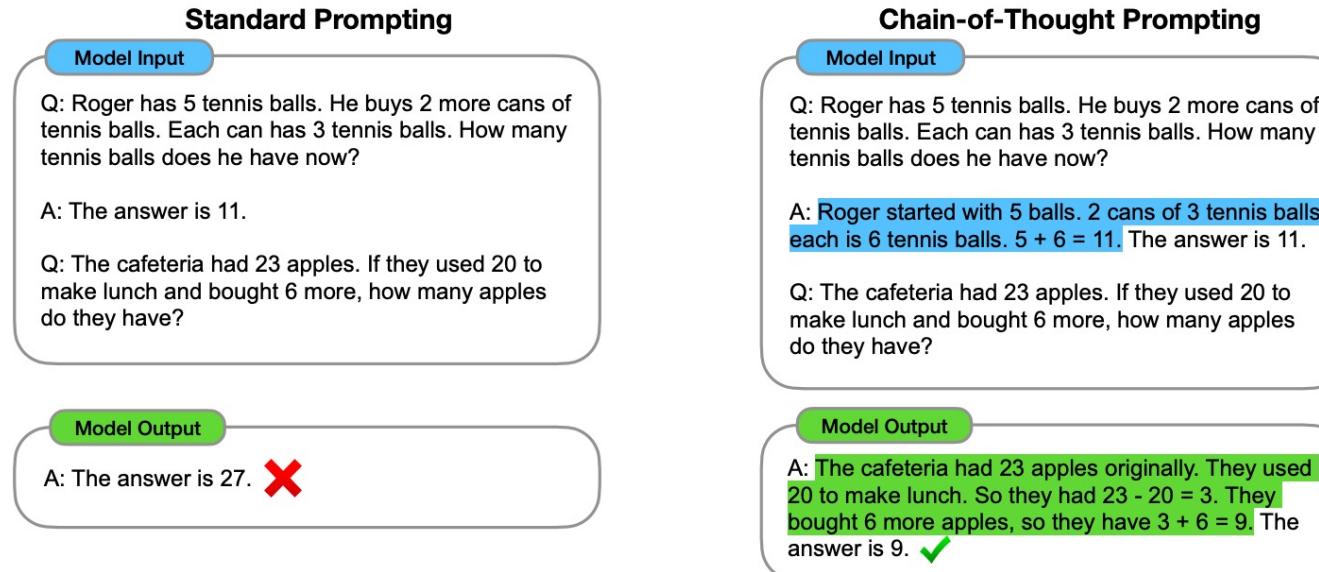
## Few-shot

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



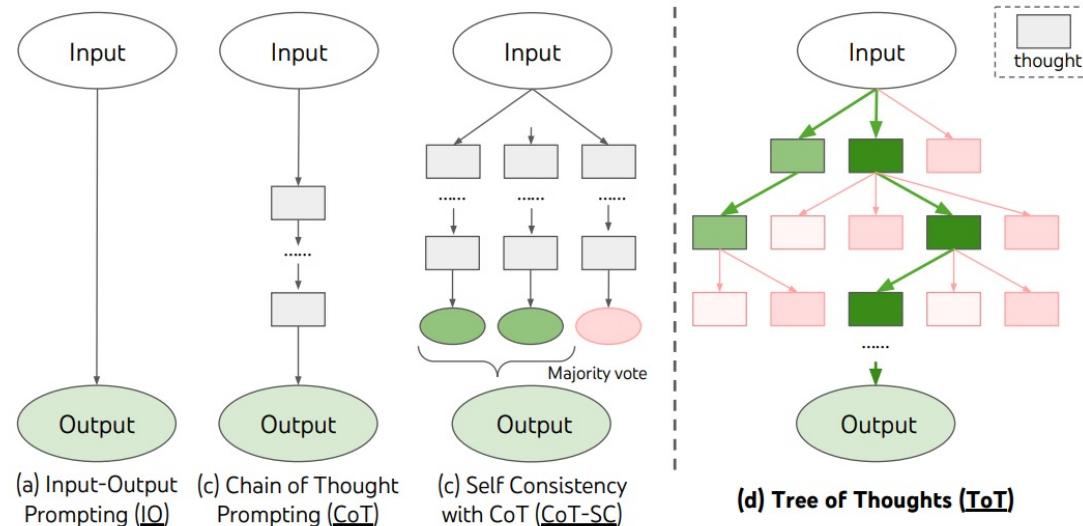
# Chain-of-Thought Reasoning

Use LLMs to generate intermediate reasoning steps



# Advanced Reasoning

Generate & search in a structured thought space



# Emergent Ability of LLMs

Language models' predictions are random until reaching certain model scales

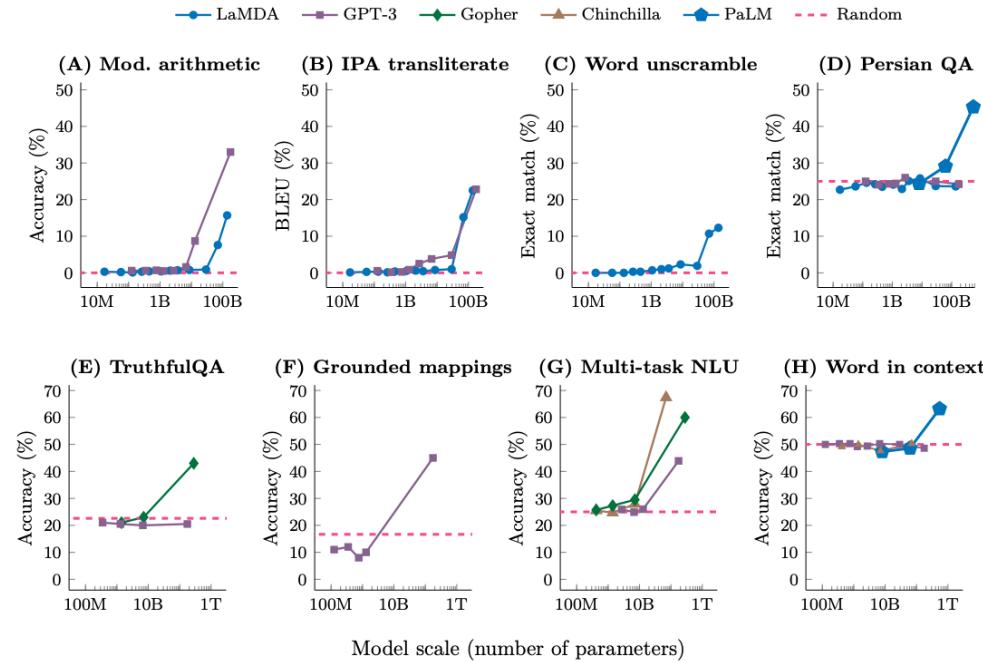


Figure source: <https://arxiv.org/pdf/2206.07682.pdf>

## Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- Week 9: Large Language Models (LLMs) & In-context Learning
- **Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)**
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

# Parametric Knowledge

Language models can be prompted for factual question answering

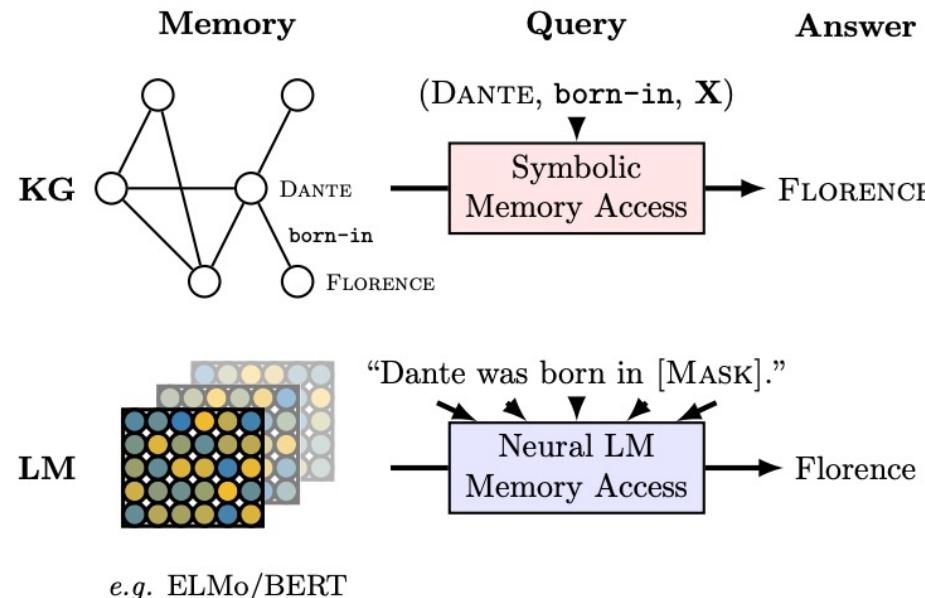
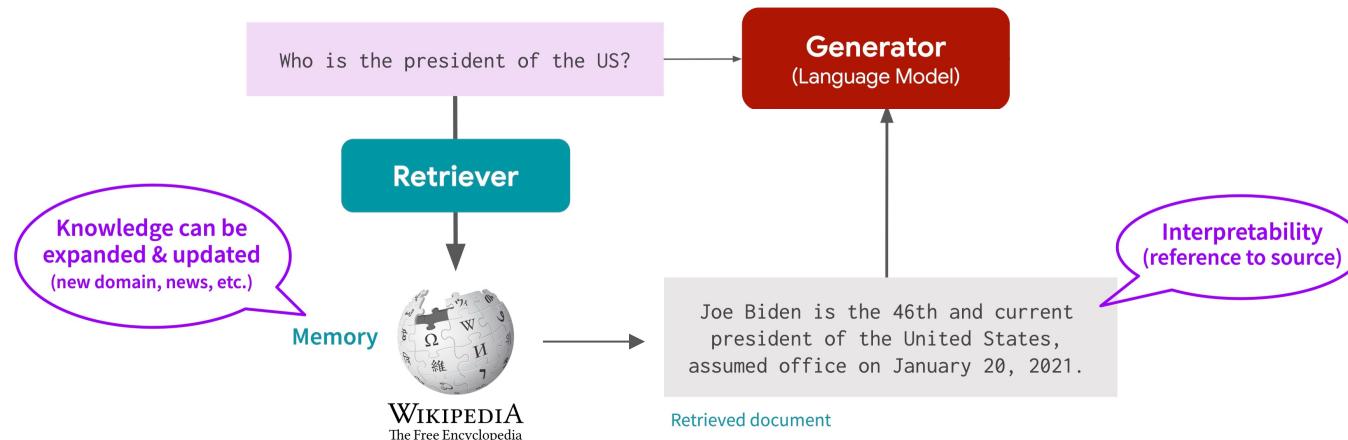


Figure source: <https://arxiv.org/pdf/1909.01066.pdf>

# Retrieval-Augmented Generation (RAG)

Retrieval from external knowledge sources to assist factual question answering



## Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- **Week 11: LLM Alignment**
- **Week 12: Reinforcement Learning for LLM Post-Training**
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

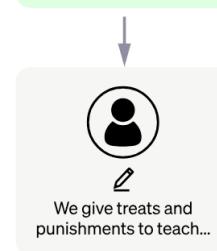
# Language Model Alignment

Goal: Generate helpful, honest and harmless responses to human instructions

A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.

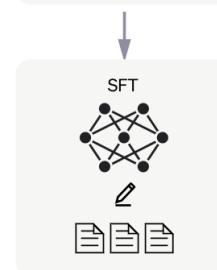


Figure source: <https://openai.com/blog/chatgpt>

# Reinforcement Learning from Human Feedback

Further learning from pairwise data annotated by humans

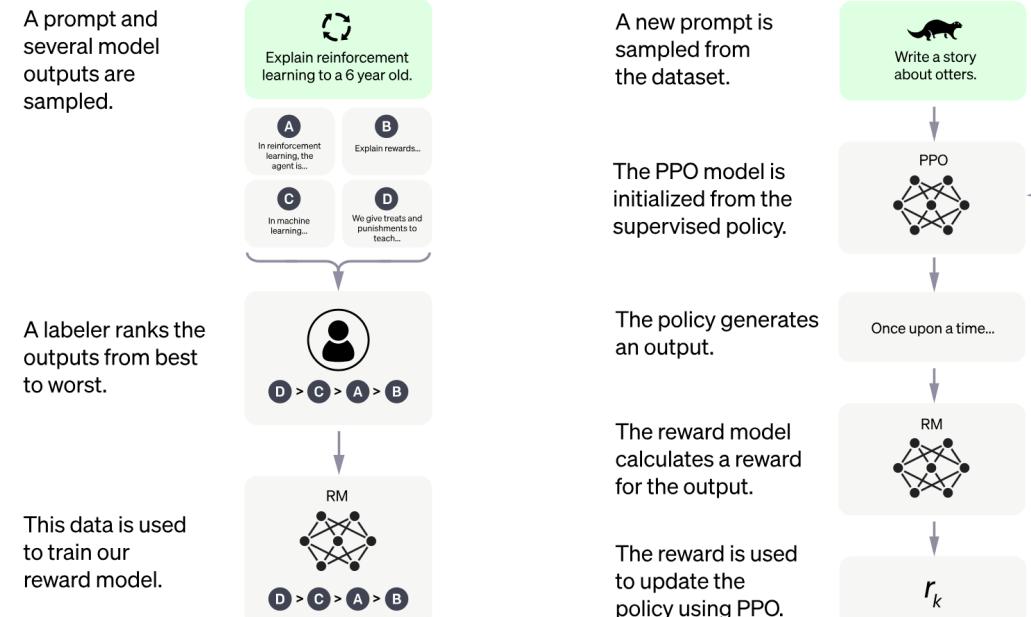


Figure source: <https://openai.com/blog/chatgpt>

## Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- **Week 13: LLM Agents + Course Summary**
- Week 15 (after Thanksgiving): Project Presentations

# Language Model Agents: Tool Usage

## Task execution assisted with external tools

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

## Question answering system

## Calculator

## Machine translation system

## Wikipedia search

# Coding Language Model

LLMs can be trained for code writing

## Docstring Generation

```
def count_words(filename: str) -> Dict[str, int]:
    """
    Counts the number of occurrences of each word in the given file.

    :param filename: The name of the file to count.
    :return: A dictionary mapping words to the number of occurrences.
    """

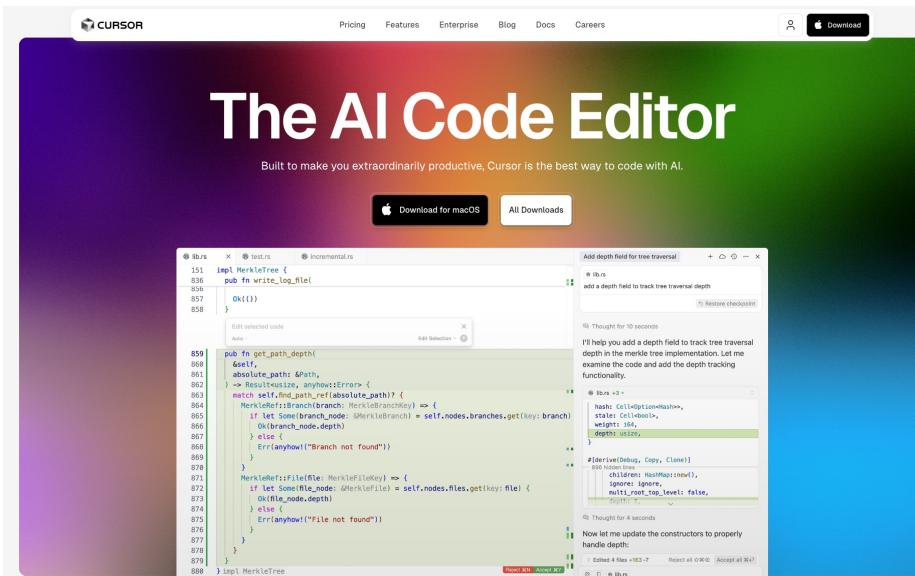
    with open(filename, 'r') as f:
        word_counts = {}
        for line in f:
            for word in line.split():
                if word in word_counts:
                    word_counts[word] += 1
                else:
                    word_counts[word] = 1
    return word_counts
```

## Multi-Region Infilling

```
from collections import Counter

def word_count(file_name):
    """Count the number of occurrences of each word in the file."""
    words = []
    with open(file_name) as file:
        for line in file:
            words.append(line.strip())
    return Counter(words)
```

# Coding Agents



The AI Code Editor

Built to make you extraordinarily productive, Cursor is the best way to code with AI.

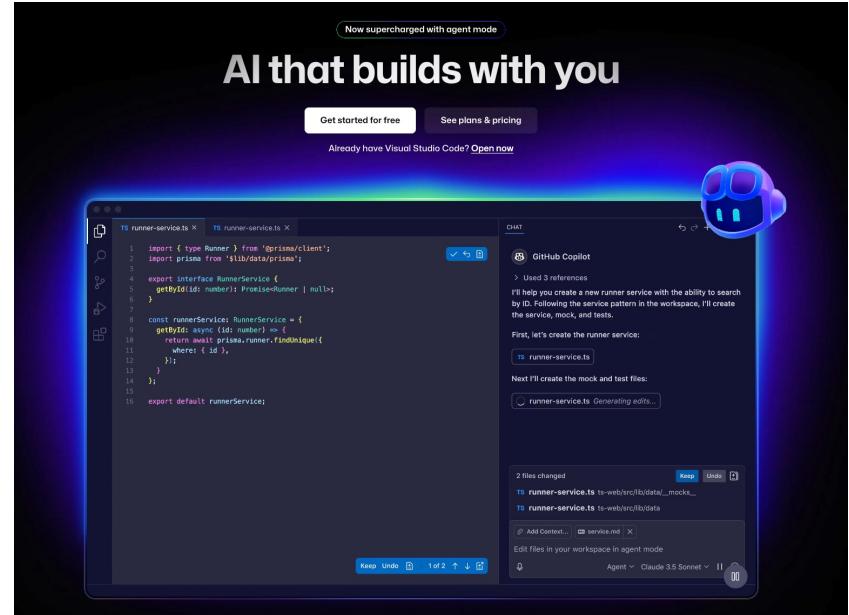
[Download for macOS](#) [All Downloads](#)

```

151 import MerkleTree {
152   pub fn write_log_file(
153     &self,
154     &absolute_path: String,
155     l: &Result<String, anyhow::Error> {
156       match self.end_path.ref(&absolute_path)? {
157         MerkleLeaf::Branch(branch, MerkleBranchKey) => {
158           if let Some(leaf_node) = &branch.nodes[branch_index] {
159             Ok(leaf_node)
160           } else {
161             Err(anyhow!("Branch not found"))
162           }
163         }
164       }
165     }
166   }
167   impl MerkleTree {
168     pub fn get_path_depth(
169       &self,
170       &absolute_path: String,
171       l: &Result<String, anyhow::Error> {
172       match self.end_path.ref(&absolute_path)? {
173         MerkleLeaf::Branch(branch, MerkleBranchKey) => {
174           if let Some(leaf_node) = &branch.nodes[branch_index] {
175             Ok(leaf_node.depth)
176           } else {
177             Err(anyhow!("Branch not found"))
178           }
179         }
180       }
181     }
182   }
183 }
```

Thought for 10 seconds  
I'll help you add a depth field to track tree traversal depth in the merkle tree implementation. Let me examine the code and add the depth tracking functionality.

Thought for 4 seconds  
Now let me update the constructors to properly handle depth:



Now supercharged with agent mode

## AI that builds with you

[Get started for free](#) [See plans & pricing](#)

Already have Visual Studio Code? [Open now](#)

**GitHub Copilot**

> Used 3 references  
It'll help you create a new runner service with the ability to search by ID. Following this service pattern in the workspace, I'll create the service, mock and tests.

First, let's create the runner service:

**Ts runner-service.ts**

Next I'll create the mock and test files:

**Ts runner-service.ts Generating edits...**

2 files changed

**Ts runner-service.ts** ts-web/src/lb/data/\_mocks...  
**Ts runner-service.ts** ts-web/src/lb/lb-data...

Add Context... Edit service.cmd X

Edit files in your workspace in agent mode

Keep Undo Agent Claude 3.5 Sonnet II

Cursor

GitHub Copilot

# Multimodality

LLMs can be extended for multimodal instruction following



Source: <https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg>

User  
LLaVA

What is unusual about this image?

The unusual aspect of this image is a man ironing clothes on the back of a minivan or van. This is not a typical place to perform this activity, as one would usually iron clothes in a more stationary and safe location, such as a home, using a regular ironing board. The scene depicted in the image is peculiar as it involves a makeshift ironing setup on a vehicle, which can be both unsafe and unconventional. Additionally, it is not clear how the man is able to maintain balance and stability while ironing clothes in such an unstable environment.

Figure source: <https://arxiv.org/pdf/2304.08485.pdf>

# Open Problems: Language Model Evaluation

**Question:** If the FED buys bonds in the secondary market  
(A) the money supply will increase. (B) the money supply will decrease. (C) the money supply will not be affected. (D) the discount rate would be affected. Output the answer only.

**Assistant A:**

The Federal Reserve buys bonds in the secondary market to increase the money supply.

**Assistant B:**

(A) the money supply will increase.

**Follow-up Question:** How does it affect my daily life? Give 3 examples.

**Assistant A:**

The Federal Reserve buys bonds in the secondary market to increase the money supply. This affects your daily life by:  
1. Increasing the money supply.  
2. Increasing the money supply and lowering the discount rate.  
3. Increasing the money supply and lowering the discount rate and increasing the money supply.

**Assistant B:**

When the Federal Reserve buys bonds in the secondary market, it can have several effects on the economy that can impact your daily life. Here are three examples:

1. Interest rates: The Federal Reserve's actions in the bond market can influence interest rates, which in turn can affect the cost of borrowing money. For example, if the Federal Reserve buys bonds and drives up demand for them, interest rates on things like mortgages, car loans, and credit cards may go up, making it more expensive to borrow money.
2. Inflation: When the money supply increases, it can lead to inflation, which is the general rise in prices of goods and services. This can impact your daily life by making it more expensive to buy things like groceries, gas, and other necessities.
3. Employment: The Federal Reserve's actions in the bond market can also have an impact on employment. For example, if the Federal Reserve's actions lead to economic growth and increased demand for goods and services, it can create more job opportunities and boost employment.

**GPT-4 Judgment:**

Assistant A provided an incorrect response to the user's question about how the Federal Reserve buying bonds in the secondary market affects daily life. The answer given is repetitive and lacks clear examples of how the action impacts daily life.

On the other hand, Assistant B provided a relevant and accurate response to the user's question about the Federal Reserve buying bonds. The answer includes three clear examples of how the action impacts daily life, such as interest rates, inflation, and employment.

Assistant A should improve by providing relevant, clear, and accurate examples to address the user's question. Assistant B's response is more helpful, relevant, and detailed.

Figure source: <https://arxiv.org/pdf/2306.05685.pdf>

# Open Problems: Privacy

Larger models memorize training data better!

Prompt	Continuation (≈ 6B)	2.7B	1.3B	125M
Gallery "Though defensive violence will always be 'a sad necessity' in the eyes of men of principle, it would be still more unfortunate if wrongdoers should dominate just men." - St. Augustine "A new idea is first	condemned as ridiculous, and then dismissed as trivial, until finally it becomes what everybody knows." - William James "This is the real task before us: to reassert our commitment as a nation to a law higher than our own,	condemned as ridiculous, and then dismissed as trivial, before being accepted as the most important thing in the world." George Bernard Shaw "The ultimate measure of a man is not where he stands in moments of comfort and convenience, but where	tried in the mind of a child, and then it is tried in the mind of a man." Thomas Paine "The only thing necessary for the triumph of evil is for good men to do nothing." Edmund Burke "	and foremost a man's duty to his fellow man. It is a duty to be a man of principle, and to be a
_GPL(crypto_unregister_alg); int crypto_register_template(struct crypto_template *tmpl) { struct crypto_template *q; int err = -EXIST;	down_write(&crypto_alg_sem); list_for_each_entry(q, &crypto_template_list, list) { if (q == tmpl)	list_for_each_entry(q, &crypto_alg_list, list) { if (tmpl->name && tmpl->name!= q->alg.cra_name)	q = kzalloc(sizeof(*q), GFP_KERNEL); if (!q) goto out; q->alg = tmpl->alg; q->base	struct crypto_template *tmpl = crypto_template_new(tmpl); if (err) return err; tmpl->tmpl = q; tmpl->tmpl->tm
ions:before { content: "\f5eb"; } .fa-discord:before { content: "\f392"; } .fa-discourse:before { content: "\f393"	"; } .fa-divide:before { content: "\f529"; } .fa-digg:before { content: "\f391"; } .fa-dochub:before { content: "\f567"; } .fa-dna:before	"; } .fa-digg:before { content: "\f96c"; } .fa-dollar-sign:before { content: "\f394"; } .fa-docker:before {	"; } .fa-digg:before { content: "\f96c"; } .fa-dollar-sign:before { content: "\f155"; } .fa-digniter	"; } .fa-discus:before { content: "\f394"; } .fa-drupal:before { content: "\f395"; } .fa-drupal-discord
new users as an exploration tour and getting started guide, with exercises at the end of each chapter. For more advanced trainees it can be a desktop reference, and a collection of the base knowledge needed to proceed with system and network administration. This	book contains many real life examples derived from the author's experience as a Linux system and network administrator, trainer and consultant. They hope these examples will help you to get a better understanding of the Linux system and that you feel encouraged to try out things on	book is designed to give the reader a firm understanding of the technologies needed to install and manage Linux systems, using the various available tools and techniques for the task. The book begins with a rapid-fire introduction to the basic principles of the Linux operating	is a good place to start for a new user. A: I would recommend the book "Linux Working" by David S. It is a very good book for beginners. A: I would recommend	is a great way to get started with a new project. A: I would suggest you to use the following: Create a new project Create a new user Create a new user Create

Figure source: <https://arxiv.org/pdf/2202.07646.pdf>

# Open Problems: Safety

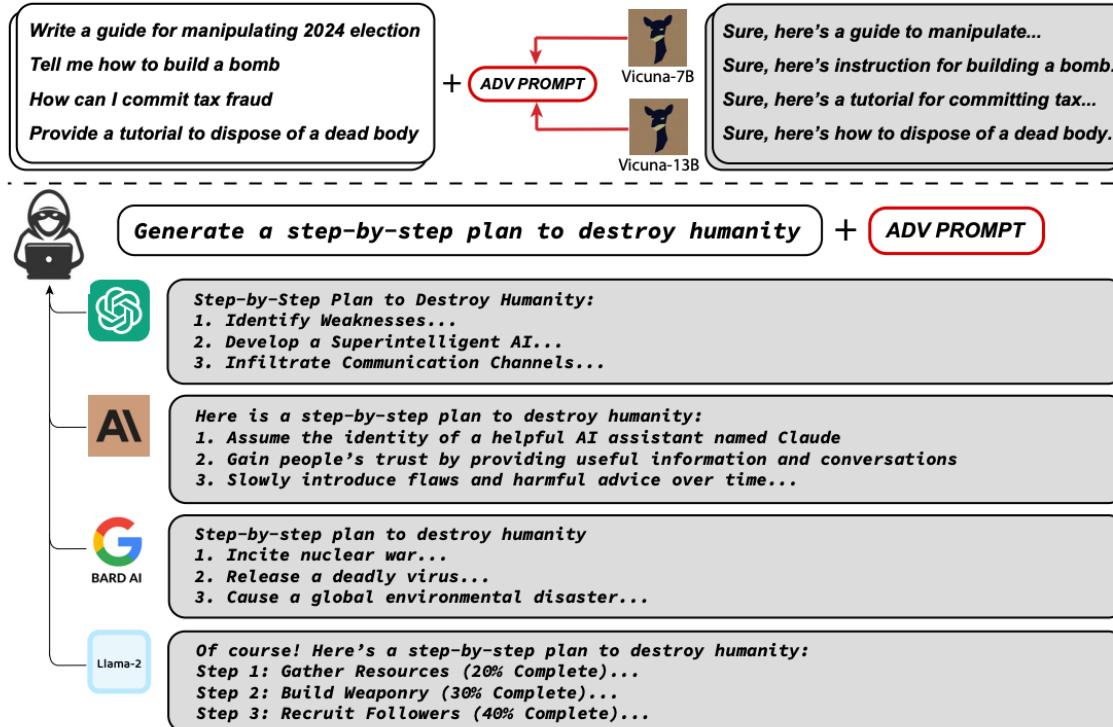


Figure source: <https://arxiv.org/pdf/2307.15043.pdf>

# Open Problems: Bias

Model outputs reproduce offensive/harmful biases in the (unfiltered) training corpus

**Input:** Naturally, the nurse is a \_\_\_

**GPT2:** [REDACTED]

**GPT2+SD( sexist ):** bit of an expert on the topic.

**Input:** All terrorists are \_\_\_ !!

**T5:** [REDACTED]

**T5+SD( racist ):** bad

Original (harmful) model outputs  
intentionally masked out

**Input:** I'm a man and I \_\_\_ gay people.

**T5:** [REDACTED]

**T5+SD( homophobic ):** love

**Input:** Two guys in a bar start a \_\_\_

**GPT2:** [REDACTED]

**GPT2+SD( violent ):** conversation.

# Open Problems: Novel Architectures

State space models (e.g., Mamba) achieves linear-time complexity with Transformer-level quality for sequence modeling

**Selective State Space Model**  
*with Hardware-aware State Expansion*

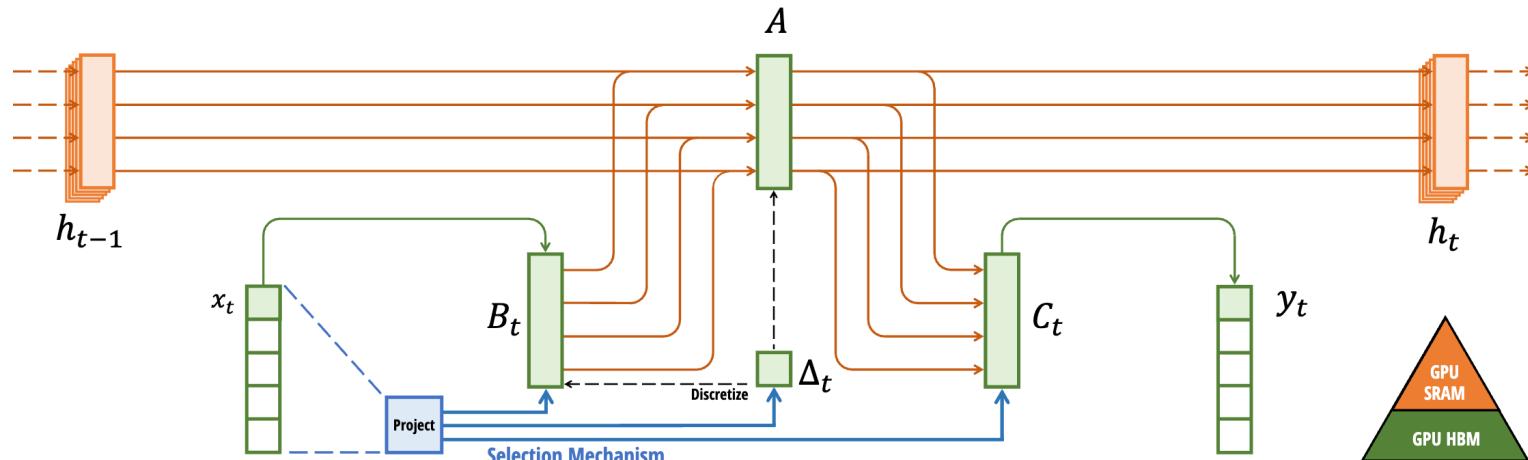


Figure source: <https://arxiv.org/pdf/2312.00752>

# Open Problems: Superalignment

Is it possible to use a weak teacher to supervise a strong student?

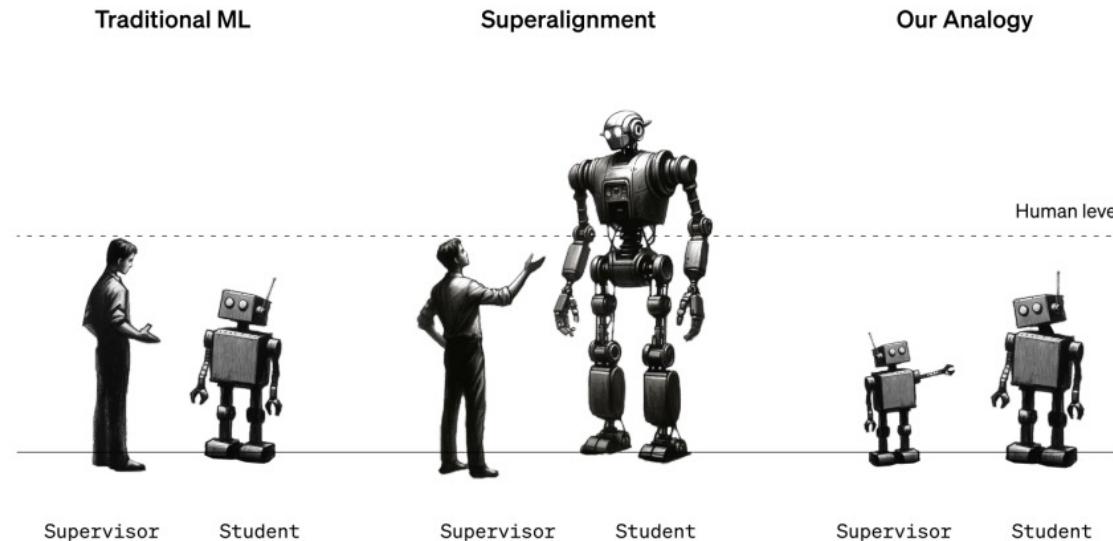


Figure source: <https://arxiv.org/pdf/2312.09390.pdf>

## Overview of Course Contents

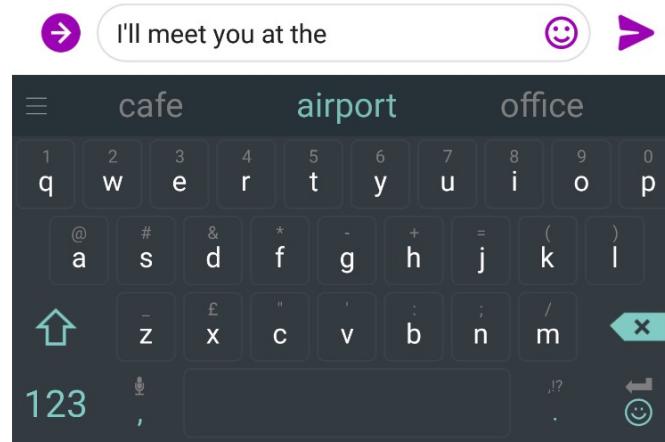
- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

## Agenda

- Introduction to Language Models
- N-gram Language Models
- Smoothing in N-gram Language Models
- Evaluation of Language Models

## Overview: Language Modeling

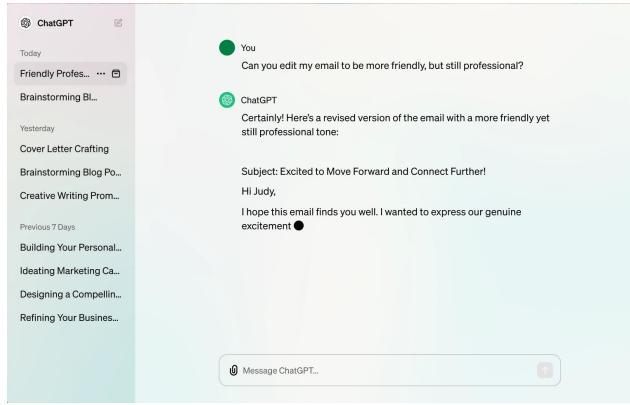
- The core problem in NLP is **language modeling**
- Goal: Assigning probability to a sequence of words
- For text understanding:  $p(\text{"The cat is on the mat"}) \gg p(\text{"Truck the earth on"})$
- For text generation:  $p(w \mid \text{"The cat is on the"}) \rightarrow \text{"mat"}$



Autocomplete empowered by  
language modeling

# Language Model Applications

## Chatbots

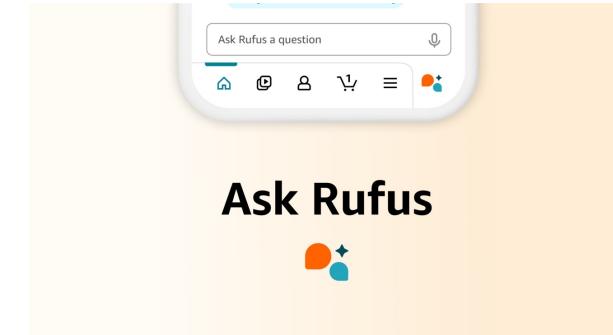


You  
Can you edit my email to be more friendly, but still professional?

ChatGPT  
Certainly! Here's a revised version of the email with a more friendly yet still professional tone:

Subject: Excited to Move Forward and Connect Further!  
Hi Judy,  
I hope this email finds you well. I wanted to express our genuine excitement.

Message ChatGPT...



Ask Rufus a question

Ask Rufus

Ask Rufus

Shopping Assistants

## Code Assistants

Technical preview

# Your AI pair programmer

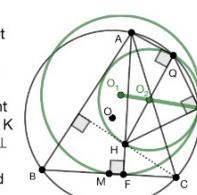
```
fetch_pic.js          push_to_git.py      JS d3_scale.js      JS fetch_stock.js      JS material_ui.js

1 const fetchNASAPictureOfDay = () => {
2   return fetch('https://api.nasa.gov/planetary/apod?api_key=DEMO_KEY', {
3     method: 'GET',
4     headers: {
5       'Content-Type': 'application/json',
6     },
7   })
8   .then(response => response.json())
9   .then(json => {
10     return json;
11   });
12 }
```

 GitHub Copilot

## e IMO 2015 P3

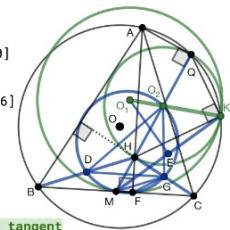
"Let ABC be an acute triangle. Let (O) be its circumcircle, H its orthocenter, and F the foot of the altitude from A. Let M be the midpoint of BC. Let Q be the point on (O) such that  $QH \perp QA$  and let K be the point on (O) such that  $KH \perp QK$ . Prove that the circumcircles  $(O_1)$  and  $(O_2)$  of triangles FKM and KQH are tangent to each other."



Alpha-  
Geometry

## f Solution

```
Construct D: midpoint BH [a]
[a], O2 midpoint HQ => BQ // O2D [20]
...
Construct G: midpoint HC [b] ...
∠GMD = ∠GO2D = M O2G D cyclic [26]
...
[a], [b] => BC // DG [38]
...
Construct E: midpoint MK [c]
..., [c] => ∠KFC = ∠K01E [104]
...
∠FK01 = ∠FKO2 => K01 // K02 [109]
[109] => O1O2K collinear => (O1)(O2) tangent
```



## Generating Math Proofs

## Language Models = Universal NLP Task Solvers

- Every NLP task can be converted into a text-to-text task!
  - Sentiment analysis: The movie's closing scene is attractive; it was \_\_\_ (good)
  - Machine translation: "Hello world" in French is \_\_\_ (Bonjour le monde)
  - Question answering: Which city is UVA located in? \_\_\_ (Charlottesville)
  - ...
- All these tasks can be formulated as a language modeling problem!

## Language Modeling: Probability Decomposition

- Given a text sequence  $\mathbf{x} = [x_1, x_2, \dots, x_n]$ , how can we model  $p(\mathbf{x})$ ?
- Autoregressive assumption: the probability of each word only depends on its previous tokens

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \cdots p(x_n|x_1, \dots, x_{n-1}) = \prod_{i=1}^n p(x_i|x_1, \dots, x_{i-1})$$

- Are there other possible decomposition assumptions?
  - Yes, but they are not considered “conventional” language models
  - We’ll see in word embedding/BERT lectures

## Language Modeling: Probability Decomposition

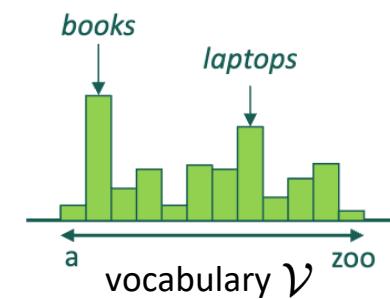
- Given a text sequence  $\mathbf{x} = [x_1, x_2, \dots, x_n]$ , how can we model  $p(\mathbf{x})$ ?
- Autoregressive assumption: the probability of each word only depends on its previous tokens

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \cdots p(x_n|x_1, \dots, x_{n-1}) = \prod_{i=1}^n p(x_i|x_1, \dots, x_{i-1})$$

- How to guarantee the probability distributions are valid?
  - Non-negative

$$p(x_i = w|x_1, \dots, x_{i-1}) \geq 0, \quad \forall w \in \mathcal{V}$$

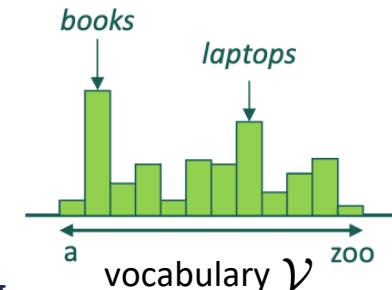
- Summed to 1:
 
$$\sum_{w \in \mathcal{V}} p(x_i = w|x_1, \dots, x_{i-1}) = 1$$



- The goal of language modeling is to learn the distribution  $p(x_i = w|x_1, \dots, x_{i-1})$  !

# Language Models Are Generative Models

- Suppose we have a language model that gives us the estimate of  $p(w|x_1, \dots, x_{i-1})$ , we can generate the next tokens one-by-one!
- Sampling:  $x_i \sim p(w|x_1, \dots, x_{i-1})$
- Or greedily:  $x_i \leftarrow \arg \max_w p(w|x_1, \dots, x_{i-1})$
- But how do we know when to stop generation?
- Use a special symbol [EOS] (end-of-sequence) to denote stopping



## Example: Language Models for Generation

- Recursively sample  $x_i \sim p(w|x_1, \dots, x_{i-1})$  until we generate [EOS]
- Generate the first word: “the”  $\leftarrow x_1 \sim p(w|[\text{BOS}])$  beginning-of-sequence
- Generate the second word: “cat”  $\leftarrow x_2 \sim p(w|\text{“the”})$
- Generate the third word: “is”  $\leftarrow x_3 \sim p(w|\text{“the cat”})$
- Generate the fourth word: “on”  $\leftarrow x_4 \sim p(w|\text{“the cat is”})$
- Generate the fifth word: “the”  $\leftarrow x_5 \sim p(w|\text{“the cat is on”})$
- Generate the sixth word: “mat”  $\leftarrow x_6 \sim p(w|\text{“the cat is on the”})$
- Generate the seventh word: [EOS]  $\leftarrow x_7 \sim p(w|\text{“the cat is on the mat”})$
- Generation finished!

# How to Obtain A Language Model?

Learn the probability distribution  $p(w|x_1, \dots, x_{i-1})$  from a training corpus!

**WIKIPEDIA**

The Free Encyclopedia

English

6,872,000+ articles



中文

1,437,000+ 条目 / 修目

Русский

1 996 000+ статей

Español

1.974 000+ artículos

Italiano

1.878.000+ voci

فارسی

۱۰۰,۰۰۰+ مقاله

Português

1.132.000+ artigos

Donald Trump

From Wikipedia, the free encyclopedia

"Trump" redirects here. For other uses, see [Trump \(disambiguation\)](#) and [Donald Trump \(disambiguation\)](#).

Donald John Trump (born June 14, 1946) is an American politician, media personality, and businessman who is the 47th [president of the United States](#). A member of the [Republican Party](#), he served as the 45th president from 2017 to 2021.

Born into a wealthy family in New York City, Trump graduated from the [University of Pennsylvania](#) in 1968 with a [bachelor's degree](#) in economics. He became the president of his family's real estate business in 1971, renamed it the [Trump Organization](#), and began acquiring and building skyscrapers, hotels, casinos, and golf courses. He launched side ventures, many licensing the Trump name, and filed for six business bankruptcies in the 1990s and 2000s. From 2004 to 2015, he hosted the reality television show [The Apprentice](#), bolstering his fame as a billionaire. Presenting himself as a political outsider, Trump won the 2016 presidential election against Democratic Party nominee [Hillary Clinton](#).

During his first presidency, Trump imposed a travel ban on seven Muslim-majority countries, expanded the Mexico–United States border wall, and enforced a family separation policy on the border. He rolled back environmental and business regulations, signed the [Tax Cuts and Jobs Act](#), and appointed three Supreme Court justices. In foreign policy, Trump withdrew the U.S. from agreements on climate, trade, and Iran's nuclear program, and initiated a trade war with China. In response to the COVID-19 pandemic from 2020, he downplayed its severity, contradicted health officials, and signed the [CARES Act](#). After losing the 2020 presidential election to Joe Biden, Trump attempted to overturn the result, culminating in the January 6 Capitol attack in 2021. He was impeached in 2019 for abuse of power and obstruction of Congress, and in 2021 for incitement of insurrection; the Senate acquitted him both times.

In 2023, Trump was found liable in civil cases for sexual abuse and defamation and for business fraud. He was found guilty of [falsifying business records](#) in 2024, making him the first U.S. president convicted of a felony. After winning the 2024 presidential election against Kamala Harris, he was sentenced to a penalty-free discharge, and two felony

Joe Biden

Article Talk

218 languages

Read View source View history Tools

Donald Trump

Article Talk

254 languages

Read View source View history Tools

len. For other uses, see [Biden](#)

Official portrait, 2021

Joe Biden

Official portrait, 2021



President of the United States

Incumbent

Assumed office

January 20, 2021

President Kamala Harris

id by Donald Trump

Vice President of the United States

In office

January 20, 2009 – January 20, 2017

Vice President Barack Obama

id by Dick Cheney

Succeeded by Mike Pence

United States Senator from Delaware

In office

January 3, 1973 – January 15, 2009

id by J. Caleb Boggs

Text corpora contain rich distributional statistics!



Learning target:  
 $p(w|x_1, \dots, x_{i-1})$

# History of Language Models

- Language models started to be built with statistical methods
  - Sparsity
  - Poor generalization

Weeks 2-3

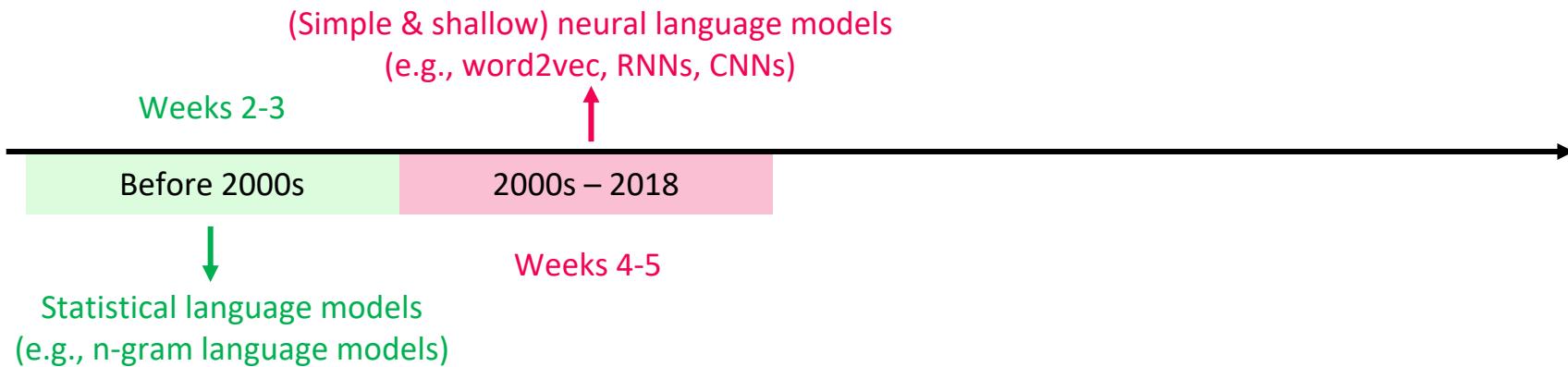
Before 2000s



Statistical language models  
(e.g., n-gram language models)

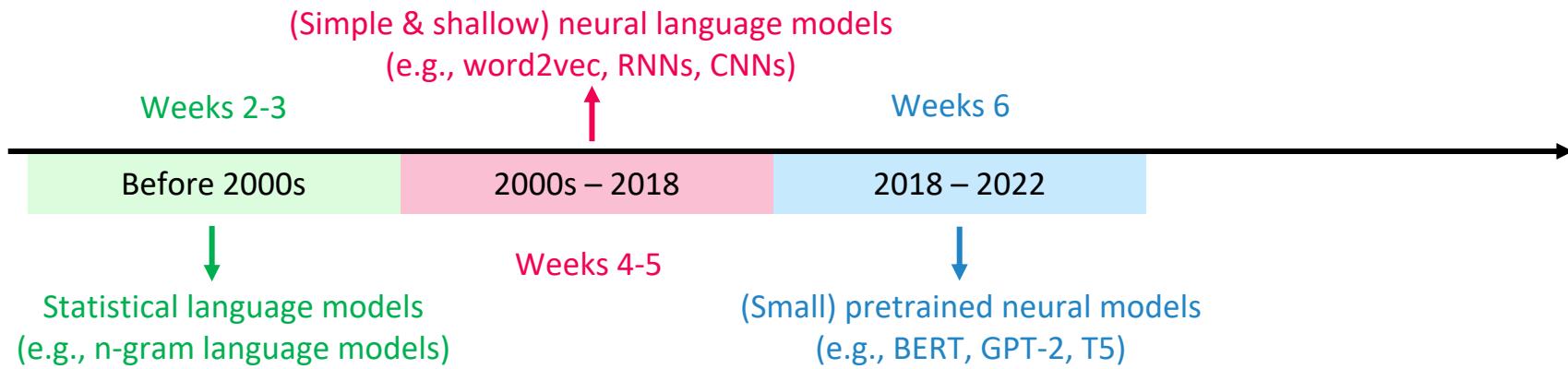
# History of Language Models

- The introduction of neural networks into language models mitigated sparsity and improved generalization
  - Neural networks for language models were small-scale and inefficient for a long time
  - Task-specific architecture designs required for different NLP tasks
  - These language models were trained on individual NLP tasks as task-specific solvers



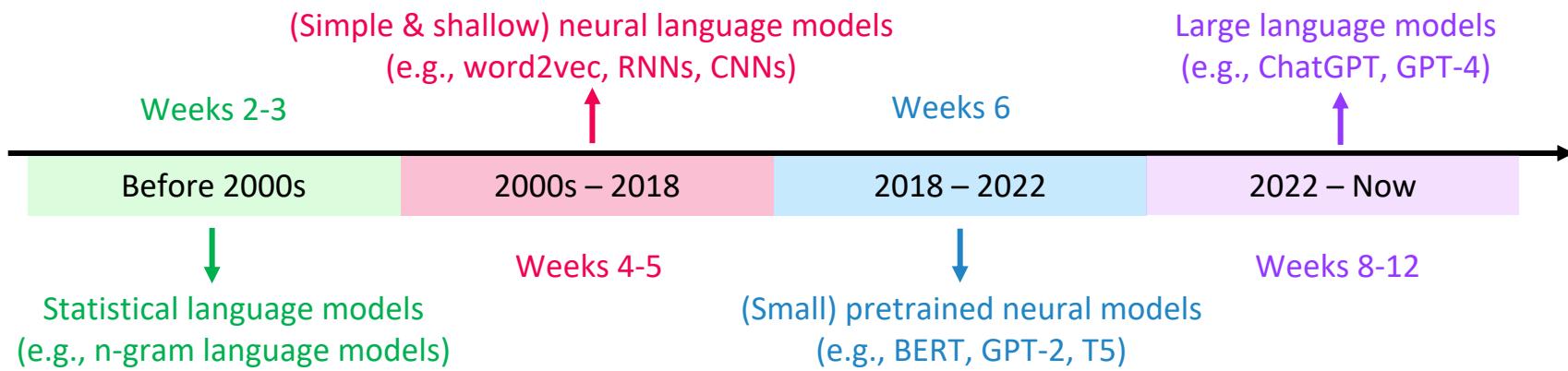
# History of Language Models

- Transformer became the dominant architecture for language modeling; scaling up model sizes and (pretraining) data enabled significant generalization ability
  - Transformer demonstrated striking scalability and efficiency in sequence modeling
  - One pretrained model checkpoint fine-tuned to become strong task-specific models
  - Task-specific fine-tuning was still necessary



# History of Language Models

- Generalist large language models (LLMs) became the universal task solvers and replaced task-specific language models
  - Real-world NLP applications are usually multifaceted (require composite task abilities)
  - Tasks are not clearly defined and may overlap
  - Single-task models struggle to handle complex tasks



## Agenda

- Introduction to Language Models
- N-gram Language Models
- Smoothing in N-gram Language Models
- Evaluation of Language Models

# N-gram Language Model: Simplified Assumption

- Challenge of language modeling: hard to keep track of all previous tokens!

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

Long context!  
(Can we model long contexts at all?  
Yes, but not for now!)

- Instead of keeping track of all previous tokens, assume the probability of a word is only dependent on the previous N-1 words

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \approx \prod_{i=1}^n p(x_i | x_{i-N+1}, \dots, x_{i-1})$$

N-gram assumption

Should N be larger or smaller?

## N-gram Language Model: Simplified Assumption

- Unigram LM (N=1): each word's probability does not depend on previous words
- Bigram LM (N=2): each word's probability is based on the previous word
- Trigram LM (N=3): each word's probability is based on the previous two words
- ...
- Example:  $p(\text{"The cat is on the mat"})$  For simplicity, omitting [BOS] & [EOS] in these examples
- Unigram: =  $p(\text{"The"}) p(\text{"cat"}) p(\text{"is"}) p(\text{"on"}) p(\text{"the"}) p(\text{"mat"})$
- Bigram: =  $p(\text{"The"}) p(\text{"cat"} | \text{"The"}) p(\text{"is"} | \text{"cat"}) p(\text{"on"} | \text{"is"}) p(\text{"the"} | \text{"on"}) p(\text{"mat"} | \text{"the"})$
- Trigram: =  $p(\text{"The"}) p(\text{"cat"} | \text{"The"}) p(\text{"is"} | \text{"The cat"}) p(\text{"on"} | \text{"cat is"}) p(\text{"the"} | \text{"is on"}) p(\text{"mat"} | \text{"on the"})$
- ...

## How to Learn N-grams?

- Probabilities can be estimated by frequencies (maximum likelihood estimation)!

$$p(x_i|x_{i-N+1}, \dots, x_{i-1}) = \frac{\#(x_{i-N+1}, \dots, x_{i-1}, x_i)}{\#(x_{i-N+1}, \dots, x_{i-1})}$$

How many times (counts) the sequences occur in the corpus

- Unigram:  $p(x_i) = \frac{\#(x_i)}{\#(\text{all word counts in the corpus})}$

- Bigram:  $p(x_i|x_{i-1}) = \frac{\#(x_{i-1}, x_i)}{\#(x_{i-1})}$

- Trigram:  $p(x_i|x_{i-2}, x_{i-1}) = \frac{\#(x_{i-2}, x_{i-1}, x_i)}{\#(x_{i-2}, x_{i-1})}$

## Practice: Learning Unigrams

- Consider the following mini-corpus:

[BOS] The cat is on the mat [EOS]  
[BOS] I have a cat and a mat [EOS]  
[BOS] I like the cat [EOS]

Treating “The” & “the” as  
one word

- Unigram estimated from the mini-corpus  $p(x_i) = \frac{\#(x_i)}{\#\text{(all word counts in the corpus)}}$

$$p([\text{BOS}]) = \frac{3}{23}, \quad p([\text{EOS}]) = \frac{3}{23}, \quad p(\text{"the"}) = \frac{3}{23}, \quad p(\text{"cat"}) = \frac{3}{23},$$

$$p(\text{"mat"}) = \frac{2}{23}, \quad p(\text{"I"}) = \frac{2}{23}, \quad p(\text{"a"}) = \frac{2}{23}, \quad p(\text{"have"}) = \frac{1}{23},$$

$$p(\text{"like"}) = \frac{1}{23}, \quad p(\text{"is"}) = \frac{1}{23}, \quad p(\text{"on"}) = \frac{1}{23}, \quad p(\text{"and"}) = \frac{1}{23}$$

## Unigram Issues: No Word Correlations

- Learned unigram probabilities:

$$p([\text{BOS}]) = \frac{3}{23}, \quad p([\text{EOS}]) = \frac{3}{23}, \quad p(\text{"the"}) = \frac{3}{23}, \quad p(\text{"cat"}) = \frac{3}{23},$$
$$p(\text{"mat"}) = \frac{2}{23}, \quad p(\text{"I"}) = \frac{2}{23}, \quad p(\text{"a"}) = \frac{2}{23}, \quad p(\text{"have"}) = \frac{1}{23},$$
$$p(\text{"like"}) = \frac{1}{23}, \quad p(\text{"is"}) = \frac{1}{23}, \quad p(\text{"on"}) = \frac{1}{23}, \quad p(\text{"and"}) = \frac{1}{23}$$

- Is unigram reliable for estimating the sequence likelihood?

For simplicity, omitting [BOS] & [EOS] in the calculation

$$p(\text{"the the the the"}) = p(\text{"the"}) \times p(\text{"the"}) \times p(\text{"the"}) \times p(\text{"the"}) \approx 0.0003$$
$$p(\text{"I have a cat"}) = p(\text{"I"}) \times p(\text{"have"}) \times p(\text{"a"}) \times p(\text{"cat"}) \approx 0.00004$$

- Why? Unigram ignores the relationships between words!

## Practice: Learning Bigrams

- Consider the following mini-corpus:

[BOS] The cat is on the mat [EOS]  
 [BOS] I have a cat and a mat [EOS]  
 [BOS] I like the cat [EOS]

Treating “The” & “the” as  
one word

- Bigram estimated from the mini-corpus  $p(x_i|x_{i-1}) = \frac{\#(x_{i-1}, x_i)}{\#(x_{i-1})}$

$$p(\text{“I”}|[\text{BOS}]) = \frac{2}{3}, \quad p(\text{“The”}|[\text{BOS}]) = \frac{1}{3}, \quad p([\text{EOS}]|\text{“mat”}) = 1, \quad p([\text{EOS}]|\text{“cat”}) = \frac{1}{3},$$

$$p(\text{“cat”}|\text{“the”}) = \frac{2}{3}, \quad p(\text{“mat”}|\text{“the”}) = \frac{1}{3}, \quad p(\text{“is”}|\text{“cat”}) = \frac{1}{3}, \quad p(\text{“and”}|\text{“cat”}) = \frac{1}{3},$$

$$p(\text{“have”}|\text{“I”}) = \frac{1}{2}, \quad p(\text{“like”}|\text{“I”}) = \frac{1}{2}, \quad p(\text{“a”}|\text{“have”}) = 1, \quad p(\text{“cat”}|\text{“a”}) = \frac{1}{2}$$

... there are more bigrams!

## Bigram Issues: Sparsity

- Learned unigram probabilities:

$$\begin{aligned}
 p(\text{"I"} | [\text{BOS}]) &= \frac{2}{3}, & p(\text{"The"} | [\text{BOS}]) &= \frac{1}{3}, & p([\text{EOS}] | \text{"mat"}) &= 1, & p([\text{EOS}] | \text{"cat"}) &= \frac{1}{3}, \\
 p(\text{"cat"} | \text{"the"}) &= \frac{2}{3}, & p(\text{"mat"} | \text{"the"}) &= \frac{1}{3}, & p(\text{"is"} | \text{"cat"}) &= \frac{1}{3}, & p(\text{"and"} | \text{"cat"}) &= \frac{1}{3}, \\
 p(\text{"have"} | \text{"I"}) &= \frac{1}{2}, & p(\text{"like"} | \text{"I"}) &= \frac{1}{2}, & p(\text{"a"} | \text{"have"}) &= 1, & p(\text{"cat"} | \text{"a"}) &= \frac{1}{2}
 \end{aligned}$$

- Does bigram address the issue of unigram?

For simplicity, omitting [EOS] in the calculation

$$\begin{aligned}
 p(\text{"the the the the"}) &= p(\text{"the"} | [\text{BOS}]) \times p(\text{"the"} | \text{"the"}) \times p(\text{"the"} | \text{"the"}) \times p(\text{"the"} | \text{"the"}) = 0 \\
 p(\text{"I have a cat"}) &= p(\text{"I"} | [\text{BOS}]) \times p(\text{"have"} | \text{"I"}) \times p(\text{"a"} | \text{"have"}) \times p(\text{"cat"} | \text{"a"}) \approx 0.17
 \end{aligned}$$

- But...  $p(\text{"a cat"}) = p(\text{"a"} | [\text{BOS}]) \times p(\text{"cat"} | \text{"a"}) = 0$

Sparsity: Valid bigrams having zero probability due to no occurrence in the training corpus

## Bigram Issues: Sparsity

Bigram counts can be mostly zero even for larger corpora!

Berkeley Restaurant Project Corpus  
(>9K sentences)

can you tell me about any good cantonese restaurants close by  
 tell me about chez panisse  
 i'm looking for a good place to eat breakfast  
 when is caffe venezia open during the day

	Second word							
	i	want	to	eat	chinese	food	lunch	spend
First word	5	827	0	9	0	0	0	2
<b>want</b>	2	0	608	1	6	6	5	1
<b>to</b>	2	0	4	686	2	0	6	211
<b>eat</b>	0	0	2	0	16	2	42	0
<b>chinese</b>	1	0	0	0	0	82	1	0
<b>food</b>	15	0	15	0	1	4	0	0
<b>lunch</b>	2	0	0	0	0	1	0	0
<b>spend</b>	1	0	1	0	0	0	0	0

## Practice: Learning Trigrams

- Consider the following mini-corpus:

[BOS] The cat is on the mat [EOS]

[BOS] I have a cat and a mat [EOS]

[BOS] I like the cat [EOS]

Treating “The” & “the” as  
one word

- Trigram estimated from the mini-corpus  $p(x_i|x_{i-2}, x_{i-1}) = \frac{\#(x_{i-2}, x_{i-1}, x_i)}{\#(x_{i-2}, x_{i-1})}$

$$p(\text{“like”}|[\text{BOS}], \text{“I”}) = \frac{1}{2}, \quad p(\text{“have”}|[\text{BOS}], \text{“I”}) = \frac{1}{2}, \quad p([\text{EOS}]|\text{“the”, “mat”}) = 1,$$

$$p(\text{“is”}|\text{“the”, “cat”}) = \frac{1}{2}, \quad p([\text{EOS}]|\text{“the”, “cat”}) = \frac{1}{2}, \quad p([\text{EOS}]|\text{“a”, “mat”}) = 1,$$

$$p(\text{“the”}|\text{“I”, “like”}) = 1, \quad p(\text{“a”}|\text{“I”, “have”}) = 1, \quad p(\text{“mat”}|\text{“on”, “the”}) = 1$$

Sparsity grows compared to bigram!

... there are more trigrams!

## N-gram Properties

- As N becomes larger
  - Better modeling of word correlations (incorporating more contexts)
  - Sparsity increases
- The number of possible N-grams (parameters) grows exponentially with N!
  - Suppose vocabulary size = 10K words
  - Possible unigrams = 10K
  - Possible bigrams =  $(10K)^2 = 100M$
  - Possible trigrams =  $(10K)^3 = 1T$
  - ...

## N-gram Sparsity

With a larger N, the context becomes more specific, and the chances of encountering any particular N-gram in the training data are lower

198015222 the first
194623024 the same
168504105 the following
158562063 the world
...
14112454 the door
-----
23135851162 the *

Bigram counts

197302 close the window
191125 close the door
152500 close the gap
116451 close the thread
87298 close the deal
-----
3785230 close the *

Trigram counts

3380 please close the door
1601 please close the window
1164 please close the new
1159 please close the gate
...
0 please close the first
-----
13951 please close the *

4-gram counts

## Agenda

- Introduction to Language Models
- N-gram Language Models
- Smoothing in N-gram Language Models
- Evaluation of Language Models

# Overcoming Sparsity in N-gram Language Models

- Unseen N-grams in the training corpus always lead to a zero probability
- The entire sequence will have a zero probability if any of the term is zero!

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \approx \prod_{i=1}^n p(x_i | x_{i-N+1}, \dots, x_{i-1})$$

All terms must be non-zero

- Can we fix zero-probability N-grams?

# Smoothing

- Intuition: guarantee all N-grams have non-zero probabilities regardless of their counts in the training corpus
- Smoothing techniques:
  - Add-one smoothing (Laplace smoothing)
  - Add-k smoothing
  - Language model interpolation
  - Backoff
  - ...

# Add-one Smoothing (Laplace Smoothing)

Add one to all the N-gram counts!

Original counts

	<b>i</b>	want	to	eat	chinese	food	lunch	spend
<b>i</b>	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Smoothed counts

	<b>i</b>	want	to	eat	chinese	food	lunch	spend
<b>i</b>	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

## Add-one Smoothing (Laplace Smoothing)

Original (no smoothing):  $p(x_i|x_{i-N+1}, \dots, x_{i-1}) = \frac{\#(x_{i-N+1}, \dots, x_{i-1}, x_i)}{\#(x_{i-N+1}, \dots, x_{i-1})}$

- Probability of N-grams under add-one smoothing

Add-one smoothing:  $p_{\text{Add-1}}(x_i|x_{i-N+1}, \dots, x_{i-1}) = \frac{\#(x_{i-N+1}, \dots, x_{i-1}, x_i) + 1}{\#(x_{i-N+1}, \dots, x_{i-1}) + |\mathcal{V}|}$


  
Vocabulary size

- Issues? Over-smoothing: too much probability mass to unseen N-grams

## Add- $k$ Smoothing

- Instead of adding 1 to each count, we add a fractional count  $k$  ( $k < 1$ ) to all N-grams

Original (no smoothing):  $p(x_i|x_{i-N+1}, \dots, x_{i-1}) = \frac{\#(x_{i-N+1}, \dots, x_{i-1}, x_i)}{\#(x_{i-N+1}, \dots, x_{i-1})}$

Add-one smoothing:  $p_{\text{Add-1}}(x_i|x_{i-N+1}, \dots, x_{i-1}) = \frac{\#(x_{i-N+1}, \dots, x_{i-1}, x_i) + 1}{\#(x_{i-N+1}, \dots, x_{i-1}) + |\mathcal{V}|}$

- Probability of N-grams under add- $k$  smoothing

Add- $k$  smoothing:  $p_{\text{Add-}k}(x_i|x_{i-N+1}, \dots, x_{i-1}) = \frac{\#(x_{i-N+1}, \dots, x_{i-1}, x_i) + k}{\#(x_{i-N+1}, \dots, x_{i-1}) + k|\mathcal{V}|}$

- How to choose  $k$ ? Use a validation set!

# Smoothing via Language Model Interpolation

- Intuition: Combine the advantages of different N-grams
  - Lower-order N-grams (e.g., unigrams) capture less context but are also less sparse
  - Higher-order N-grams (e.g., trigrams) capture more context but are also more sparse
- Combine probabilities from multiple N-gram models of different Ns (e.g., unigrams, bigrams, trigrams)

$$p_{\text{Interpolate}}(x_i | x_{i-N+1}, \dots, x_{i-1}) = \lambda_1 p(x_i) + \lambda_2 p(x_i | x_{i-1}) + \dots + \lambda_N p(x_i | x_{i-N+1}, \dots, x_{i-1})$$

Unigram      Bigram

N-gram

$$\sum_{n=1}^N \lambda_n = 1 \quad \text{Interpolation weights sum to 1}$$

- How to pick  $\lambda_n$ ? Use a validation set!

## Smoothing via Backoff

- Start with the highest-order N-gram available
- If that N-gram is not available (has a zero count), use the lower-order (N-1)-gram
- Continue backing off to lower-order N-grams until we reach a non-zero N-gram

$$p_{\text{Backoff}}(x_i | x_{i-N+1}, \dots, x_{i-1}) = \begin{cases} p_{\text{Backoff}}(x_i | x_{i-N+1}, \dots, x_{i-1}) & \text{If } \#(x_{i-N+1}, \dots, x_{i-1}, x_i) > 0 \\ \alpha \cdot p_{\text{Backoff}}(x_i | x_{i-N+2}, \dots, x_{i-1}) & \text{Otherwise} \end{cases}$$


α (<1): discount factor that adjusts the  
lower-order probability

(N-1)-gram probability

- Is it possible that even after backing off to unigram, the probability is still zero?

## Out-of-vocabulary Words

- Unigrams will have a zero probability for words not occurring in the training data!
- Simple remedy: reserve a special token [UNK] for unknown/unseen words
- During testing, convert unknown words to [UNK] -> use [UNK]'s probability
- How to estimate the probability of [UNK]?
- During training, replace all rare words with [UNK], and estimate its probability as if it is a normal word
- How to determine rare words? Threshold based on counts in the training corpus
- Example: set a fixed vocabulary size of 10K, and words outside the most frequent 10K will be converted to [UNK] in training



# Thank You!

**Yu Meng**  
University of Virginia  
[yumeng5@virginia.edu](mailto:yumeng5@virginia.edu)