Foundation Models and Prompts

Danna Gurari

University of Colorado Boulder Spring 2025



https://dannagurari.colorado.edu/course/neural-networks-and-deep-learning-spring-2025/

Review

- Last lecture:
 - Other data modalities
 - Internet-scale trained models
 - Scaling laws
 - Lab assignment 3: multimodal task
 - Programming tutorial
- Assignments (Canvas):
 - Lab assignment 3 due in 1 week
- Questions?

Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

Focus for 2010s: Closed-World Problems



Learning works well with "big" labeled datasets for the target task!

Kamath, Liu, and Whitaker. Deep Learning for NLP and Speech Recognition. 2019.

Open Problems: Beyond Closed-World Setting



https://arxiv.org/pdf/2210.09263.pdf

Open Problems: Beyond Closed-World Setting

Vocabulary / Language

Open vocabulary and Zero-shot:

generalize to task with no labeled training data for the target task (e.g., novel categories), where the former problem permits annotations with novel category (for a different task)

> Closed-world/Closed-set: assumes all target classes are available at training



running. while shirt



border collie. person. running, dog, while shirt standing/sitting

mask-wearing food flowers textures

Open world/In the wild for different tasks (e.g., detection): succeed for all categories, whether seen or not seen during training

person, dog Visual Content Training **Out-of-domain/Robustness Testing:** Domain same content observed differently

Open set classification/Out-of-distribution Detection:

predict whether a sample is drawn from the distribution observed at training time

https://arxiv.org/pdf/2210.09263.pdf

Beyond "Big" Data: Few-Shot Learning



Perrett et al. Use Your Head: Improving Long-Tail Video Recognition. CVPR 2023.

Beyond "Big" Data: Few-Shot Learning Intuition



Lake et al, 2013, 2015

Given one example per category, identify the category of the query

https://www.youtube.com/watch?v=9j4iH9TPTd8

Beyond "Big" Data: Zero-Shot Learning Intuition



What is this?

How many examples do you think you would need to see to recognize one of these?

Beyond "Big" Data: Zero-Shot Learning Intuition



Could see 0 examples if you knew the object fuses a person on top with a horse on the bottom

Beyond "Big" Data: Zero-Shot Learning Intuition



Could see 0 examples of a zebra if you knew it looks like a horse with black and white stripes

Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

New Paradigm:

Foundation Models Can Generalize Beyond Closed-World Settings With Limited Training Data

Definition of "Foundation Model"



Coined in 2021, it references the recent paradigm shift to develop a single model that can implicitly support many downstream tasks.

Bommasani et al. On the Opportunities and Risks of Foundation Models. arXiv 2021.

Training and Evaluating Foundation Models



Evaluate with modern benchmark datasets for many:

1. Different tasks (e.g., object recognition, scene classification)

2. Different distributions of the same task (e.g., ImageNet vs data from blind people)

Bommasani et al. On the Opportunities and Risks of Foundation Models. arXiv 2021.

A Transition from Specialist to General-Purpose Models



Li et al. Multimodal foundation models: From specialists to general-purpose assistants. 2024.

Emergent Abilities from Foundation Models



In 2020, we observed a model can be used *as is for many downstream tasks* with just *prompting*!

https://docs.graphcore.ai/projects/bert-training/en/latest/bert.html

Pioneering Methods: Historical Context



Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

Starting Point: Better Performance from Increasing Model and Dataset Sizes (Scaling Law)



https://medium.com/building-the-metaverse/the-metaverse-and-artificial-intelligence-ai-577343895411

(GPT-3; Neurips 2020)

Language Models are Few-Shot Learners

Tom B. Brown*		Benjamin Mann*		Ryder* Me	anie Subbiah*
Jared Kaplan [†]	Prafulla Dha	riwal Arvin	d Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Aga	arwal Ariel I	Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ran	nesh Danie	el M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher He	sse Ma	rk Chen I	Eric Sigler	Mateusz Litwin	Scott Gray
Benjamin Chess		Jack	Clark	Christopher Berner	
Sam McCan	dlish	Alec Radford	Ilya Su	ıtskever I	Dario Amodei

OpenAI

Architecture: Scaled-Up GPT-2 (aka GPT-3)

GPT-2 architecture with minor modifications; e.g.,

- More layers (96 vs 48) and so parameters (175B vs 1.5B)

Accepts up to 2,048 input tokens instead of 1,024 tokens.
Is the input size equivalent to approximately:
(a) 1/3 single-space page
(b) 1 single-space page
(c) 3 single-space pages

Is this model (a) discriminative or (b) generative?



Training Data: Scaled-Up from GPT-2

Dataset	Quantity (tokens)	Weight in training mix
Common Crawl (filtered)	410 billion	60%
WebText2	19 billion	22%
Books1	12 billion	8%
Books2	55 billion	8%
Wikipedia	3 billion	3%

(Noisy, lower quality Common Crawl data spanning 41 monthly crawls from 2016 to 2019, supplemented with 4 high-quality datasets)

Language composition: by word count, 93% English

How many GB for training? ~(a) 6, (b) 60, (c) 600, (d) 6,000 Could the entire dataset fit on your personal computer? - Not mine (it has ~250GB)

(web archives without HTML markup and non-text content https://commoncrawl.org/)



Brown et al. Language Models are Few-Shot Learners. Neurips 2020.

Training Protocol

- Same objective as GPT (predict next word)
- Higher quality datasets sampled more often (and so potential overfitting)
- Trained on V100 GPUs for 3.14E+23 flops (floating point operations)
- What do you think is the most accurate estimate for the training cost?
 (a) ~\$5,00
 - (b) ~\$50,000
 - (c) ~\$500,000
 - (d) ~\$5,000,000
- Required ~1.3 Gigawatt-hours of electricity (enough to power 121 homes in America for a year)

Two Key Emergent Abilities

- In-context learning: model completes document in format modeled by task-specific examples
- Chain-of-thought (CoT) prompting: model conveys its reasoning process when completing the task

Intuition: Can You Learn to Recognize Something With Zero/Few Examples?





https://www.youtube.co m/watch?v=9j4iH9TPTd8

So Can GPT-3 with Prompts: Instructions + Examples (Latter Called "In-Context Learning")

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



One-shot

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



Few-shot

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



(Typically uses maximum amount of examples that can fit in 2,048 sized context window: ~10 to 100)

Brown et al. Language Models are Few-Shot Learners. Neurips 2020.

Prompt Designed Per Dataset; e.g.,

$\texttt{Context} \ \rightarrow$	Q: What school did burne hogarth establish?
2	A:
Target Completion $ ightarrow$	School of Visual Arts

Figure G.35: Formatted dataset example for WebQA

Context \rightarrow Keinesfalls dürfen diese für den kommerziellen Gebrauch verwendet werden.

Target Completion \rightarrow In no case may they be used for commercial purposes.

Figure G.36: Formatted dataset example for $De \rightarrow En$. This is the format for one- and few-shot learning, for this and other langauge tasks, the format for zero-shot learning is "Q: What is the {language} translation of {sentence} A: {translation}."

Example Result: Fake News Generation

Title: United Methodists Agree to Historic Split Subtitle: Those who oppose gay marriage will form their own denomination Article: After two days of intense debate, the United Methodist Church has agreed to a historic split - one that is expected to end in the creation of a new denomination, one that will be "theologically and socially conservative," according to The Washington Post. The majority of delegates attending the church's annual General Conference in May voted to strengthen a ban on the ordination of LGBTQ clergy and to write new rules that will "discipline" clergy who officiate at same-sex weddings. But those who opposed these measures have a new plan: They say they will form a separate denomination by 2020, calling their church the Christian Methodist denomination.

The Post notes that the denomination, which claims 12.5 million members, was in the early 20th century the "largest Protestant denomination in the U.S.," but that it has been shrinking in recent decades. The new split will be the second in the church's history. The first occurred in 1968, when roughly 10 percent of the denomination left to form the Evangelical United Brethren Church. The Post notes that the proposed split "comes at a critical time for the church, which has been losing members for years," which has been "pushed toward the brink of a schism over the role of LGBTQ people in the church." Gay marriage is not the only issue that has divided the church. In 2016, the denomination was split over ordination of transgender clergy, with the North Pacific regional conference voting to ban them from serving as clergy, and the South Pacific regional conference voting to allow them.

Results: Smooth Scaling with Model Capacity

What is the trend for zero-shot performance? (recall zero-shot example)





Results: Smooth Scaling with Model Capacity



Results: Smooth Scaling with Model Capacity



Summary of GPT-3 Findings

Tested on 10s of NLP datasets, showing strong performance overall and occasionally state-of-the-art performance (outperforming fine-tuned methods)!

Brown et al. Language Models are Few-Shot Learners. Neurips 2020.

Two Key Emergent Abilities

- In-context learning: model completes document in format modeled by task-specific examples
- Chain-of-thought (CoT) prompting: model conveys its reasoning process when completing the task

Intuition: How Much Reasoning Do You Do?

- What is 2 + 2?
 - 4
- What is 24 x 14?
 - My reasoning process:
 - 24 x 10 = 240
 - 24 x (14-10) = 24 x 4 = 96
 - 240+96 = 336 (final answer!)
 - Did you follow the same or a different reasoning path to generate the answer?
- We can also encourage a model to reason/think before responding
(Pioneering paper; Neurips 2022)

Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

Jason Wei Xuezhi Wang Dale Schuurmans Maarten Bosma

Brian Ichter Fei Xia Ed H. Chi Quoc V. Le Denny Zhou

Google Research, Brain Team {jasonwei,dennyzhou}@google.com

CoT Prompting

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output A: The answer is 27.

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9.

Guides model to show intermediate reasoning steps

Wei et al. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. Neurips 2022.

CoT Prompting: Why Does it Work?

СоТ	Originally, Leah had 32 chocolates and her sister had 42. So in total they had $32 + 42 =$ 74. After eating 35, they had 74 - $35 = 39$ pieces left in total. The answer is 39.	Julie is reading a 120-page book. Yesterday, she read 12 pages and today, she read 24 pages. So she read a total of $12 + 24 = 36$ pages. Now she has $120 - 36 = 84$ pages left. Since she wants to read half of the remaining pages, she should read $84 / 2 = 42$ pages. The answer is 42.	Performance can still improve with invalid	
Invalid Reasoning	Originally, Leah had 32 chocolates and her sister had 42. So her sister had 42 - $32 = 10$ chocolates more than Leah has. After eating 35, since $10 + 35 =$ 45, they had 45 - 6 = 39 pieces	Yesterday, Julie read 12 pages. Today, she read $12 * 2 = 24$ pages. So she read a total of 12 + 24 = 36 pages. Now she needs to read 120 - 36 = 84 more pages. She wants to read half of the remaining pages tomorrow, so she needs to read $84 / 2 = 42$ pages tomorrow.	examples, suggesting benefit of examples is revealing target format	

The answer is 42. \checkmark

left in total. The answer is 39.

Wang et al. Towards Understanding Chain-of-Thought Prompting: An Empirical Study of What Matters. arXiv 2022.

CoT Prompting Enhancement: Self-Consistency



Wang et al. Self-consistency Improves Chain of Thought Reasoning in Language Models. ICLR 2023.

CoT Prompting Enhancement: Self-Consistency

Introduces diversity by randomly sampling, at each time step, the next token from top *K* most probable tokens (from softmax layer)



Wang et al. Self-consistency Improves Chain of Thought Reasoning in Language Models. ICLR 2023.

CoT Variant

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis					
balls. Each can has 3 tennis balls. How many tennis balls does					
he have now?					

A: The answer is 11.

A:

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

(Output) The answer is 8. X

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there? A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4.

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: Let's think step by step.

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.

Models perform better even when just asked to reason slowly

(Examples from GPT-

3) Kojima et al. Large Language Models are Zero-Shot Reasoners. Neurips 2022.

These Findings Helped Inspire a New Era of Foundation Models and Prompts

- In-context learning: model completes document in format modeled by task-specific examples
- Chain-of-thought (CoT) prompting: model conveys its reasoning process when completing the task



November 2022: The Release of ChatGPT Changed the World's Expectations from Al...

Bing Introducing ChatGPT

🕼 OpenAl

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

API ~

Research ~

ChatGPT ~

Safetv

Try ChatGPT 7 Read about ChatGPT Plus

https://openai.com/blog/chatgpt



And helped NVIDIA, OpenAI's GPU provider, get much richer ⁽²⁾ (market value of \$4.3 TRILLION in Jan 2025)

(Excellent Summary of Prompts; ACM Surveys 2023)

Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing

Pengfei Liu Carnegie Mellon University pliu3@cs.cmu.edu

Zhengbao Jiang Carnegie Mellon University zhengbaj@cs.cmu.edu Weizhe Yuan Carnegie Mellon University weizhey@cs.cmu.edu Jinlan Fu National University of Singapore jinlanjonna@gmail.com

Hiroaki Hayashi Carnegie Mellon University hiroakih@cs.cmu.edu Graham Neubig Carnegie Mellon University gneubig@cs.cmu.edu

(Another Excellent Read; arXiv 2023)

A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT

Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C. Schmidt Department of Computer Science Vanderbilt University, Tennessee

e.g., "The Persona" Pattern

- "From now on, act as a security reviewer. Pay close attention to the security details of any code that we look at. Provide outputs that a security reviewer would regarding the code."
- "You are going to pretend to be a Linux terminal for a computer that has been compromised by an attacker. When
 I type in a command, you are going to output the corresponding text that the Linux terminal would produce."

To Help With Prompt Engineering, Many Companies Now Sell Prompts; e.g.,

PromptBase Search Prompts			Q	Marketplace Ger	nerate Hire I	ogin Sell	
(2) AI Models	ChatGPT Midjourney	v Marketing 🗸	Shop Other	✓ Free Prom	pt Free Gen	erator 🗸	login→ 🎦
Mid	Se PROMPT Se	arch Prompts, @authors or #	tags	Q		Marketplace L	Login Register
Sta							
Pro <mark>Prompt M</mark>	Unleash the power of Artific	cial Intelligence					
save or Business	Prompt A	ttack				ATT F	
Find Super Prompt,	Your #1 F	Prompt			Rojon		
Prompt. +10k Prom	Marketpla	ace			MIDJOURNEY		© дрт
Sell a Prompt	PromptAttack is a marketplace w	here you can purchase and s	ell high-quality prompts th	at Jagged Cu @midrun	t Out Punk Posters \$2.99	Research Paper S @laxman1986	Summarizer \$2.99
	generate optimal stunning results	while also reducing your API	expenses.				Studio Quality Product Fruit @midrun \$2.99

Prompt Engineering Work Also Helped Inspire Rise of "Reasoning" Models; e.g.,

- Sep 2024: OpenAl's 01, world's first reasoning model
- Nov 2024: Alibaba's QWQ, or Questions with Qwen (model open-source; extends Meta's Llama)
- Dec 2024: Google's Gemini Flash Thinking and OpenAl's 03
- Jan 2025: DeepSeek's R1 (model open-source)

To Help With Leveraging These Huge Models, Companies Now Sell Access to Them; e.g.,



Programming Tutorials

- Everley's programming tutorial for Hugging Face Inference Playground
- Ch 6 of this book:



Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

Motivating Observation

- Foundation models achieve better performance for NLP tasks when provided "in-context" examples.
 - i.e., [Task description, Examples, Prompt]
 - e.g., "Translate English to Spanish. Computer -> Computadora. Vision ->
- Idea: Use in-context few-shot learning for image-based prompts

Novel Idea: Image Inpainting



Designed to adapt to any "image-to-image translation" task by using the model as is (e.g., no fine-tuning required)

Bar et al. Visual Prompting via Image Inpainting. Neurips 2022.

Approach: Prompting for Image Inpainting



Bar et al. Visual Prompting via Image Inpainting. Neurips 2022.

Approach

Extended in 2023 to standard vision benchmarks:



Training: Masked Image Modeling



Uses self-supervised learning such that the model predict values in masked out patches

Uses standard vision benchmarks for each evaluated task

Experimental Results

(Used in prompt the best performing example-per pair per task from all examples in the training dataset)

Model achieves state-of-the-art performance on depth estimation for NYUv2 dataset and outperforms other generalist models on several more tasks.

Qualitative Results: In-Domain Results



Qualitative Results: In-Domain Results



Qualitative Results: In-Domain Results



Qualitative Results: Open-Vocabulary Results (i.e., Categories Not Seen at Training)

Shows in-context examples, prompts, and predictions for keypoint detection, object segmentation, and instance segmentation



Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

(Pioneering LMM; Neurips 2022)

Flamingo: a Visual Language Model for Few-Shot Learning

Jean-Baptiste Alayra	c ^{*,‡} Jeff Donahu	e* Pauline Luc*	Antoine Miech*	
Iain Barr [†] Yana	Hasson [†] Karel Ler	nc [†] Arthur Mensch	[†] Katie Millican [†]	
Malcolm Reynolds †	Roman Ring † Eliza	a Rutherford [†] Serka	n Cabi Tengda Han	
Zhitao Gong	Sina Samangooei	Marianne Monteiro	Jacob Menick	
Sebastian Borgeaud	Andrew Brock	Aida Nematzadeh	Sahand Sharifzadeh	
Mikolaj Binkowski	Ricardo Barreira	Oriol Vinyals	Andrew Zisserman	

Karen Simonyan*,‡

* Equal contributions, ordered alphabetically, [†] Equal contributions, ordered alphabetically, [‡] Equal senior contributions

DeepMind



Figure 1: Development timeline of Multimodal models grouped in four proposed architecture types.

Wadekar et al. The Evolution of Multimodal Model Architectures. arXiv 2024.

Key Challenge: Interleaved Image-Text Input



Achieved Strong Zero-Shot Performance with Multimodal Prompts for Many Tasks; e.g.,



Achieved Strong Zero-Shot Performance with Multimodal Prompts for Many Tasks; e.g.,



Results: Larger Models Perform Better



What are trends for different model sizes and numbers of "shots"?

Largest model achieved state-of-the-art on 6 of 16 tasks, outperforming fine-tuned models

LMMs Extending Beyond Language + Vision

- e.g., PandaGPT: text, image, video, audio, and sensors for depth (3D), thermal (infrared radiation), or inertial measurement units (IMU) of motion and position; connects objects with sound, 3D shape, temperature, and motion
- e.g., SpeechGPT: language and speech
- e.g., NExT-GPT: text, images, videos, and audio

Class Discussion: What are Risks of Using Foundation Models (Today, Generative AI)?

• Model biases/limitations can affect all downstream models

Scaling Laws Still In Play: Training Costs (and Emergent Abilities) Are Soaring



"Existential" threat for those without "deep pockets"; how to make a contribution?

https://epoch.ai/blog/how-much-does-it-cost-to-train-frontier-ai-models
Next Two Lectures: How Can You Contribute to Neural Networks and Deep Learning in this "Scaling Law" Era?



Today's Topics

- Motivation
- Foundation Models
- Prompting "Large Language Models" for NLP
- Prompting "Large Vision Models" for Computer Vision
- Prompting "Large Multimodal Models"

