# Question Answering & Factual Knowledge

**Yu Meng**

University of Virginia

yumeng5@virginia.edu

Oct 25, 2024

# 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 and Neural Language Models

- Week 6-7: Language Modeling with Transformers (Pretraining + Fine-tuning)

- Week 8: Large Language Models (LLMs) & In-context Learning

- Week 9-10: Reasoning, Knowledge, and Retrieval-Augmented Generation (RAG)

- Week 11: LLM Alignment

- Week 12: Language Agents

- Week 13: Recap + Future of NLP

- Week 15 (after Thanksgiving): Project Presentations

# (Recap) Reasoning: Overview

- **Reasoning** (rough definition): perform deductive, inductive, commonsense, or logical reasoning via generating or analyzing text with language models

- Deductive reasoning: draw specific conclusions from general principles or premises
  - E.g.: "All humans are mortal" + "Socrates is a human" => "Socrates is mortal"

- Inductive reasoning: make generalizations based on specific observations
  - E.g.: "The sun has risen in the east every day" => "The sun will rise in the east tomorrow"

- Commonsense reasoning: rely on world knowledge or commonsense understanding to make predictions or answer questions
  - E.g.: "If I drop a ball, what will happen?" => "It will fall"

- Mathematical/logical reasoning: follow specific rules or procedures to arrive at a correct answer
  - E.g.: "If 3 apples cost $6, how much do 5 apples cost?" => "$10"

# (Recap) Standard vs. CoT Prompting

CoT prompting is especially effective for large models

# (Recap) CoT Can Be Triggered Zero-shot

Just add "Let's think step by step" at the beginning of the answer

### (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.

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 answer is 8. ✗

### (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. ✓

### (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 ✗

### (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.* ✓

Paper: https://arxiv.org/pdf/2205.11916

# (Recap) CoT Demo

CoT prompt: How many 'r' letters are there in the following word: strawberry? Let's think step by step

No-CoT prompt: How many 'r' letters are there in the following word: strawberry? Answer without reasoning steps

llama-3.1-70b-instruct

🔍 Expand to see the descriptions of 75 models ◀

💬 Scroll down and start chatting

> How many 'r' letters are there in the following word: strawberry? Answer without reasoning steps

There are 2 'r' letters in the word "strawberry". **Wrong result**

llama-3.1-70b-instruct

🔍 Expand to see the descriptions of 75 models ◀

💬 Scroll down and start chatting

> How many 'r' letters are there in the following word: strawberry? Let's think step by step

**Step 1: Write down the word**
The word is: strawberry

**Step 2: Break down the word into individual letters**
Here are the individual letters: s-t-r-a-w-b-e-r-r-y

**Step 3: Count the number of 'r' letters**
Let's count the 'r' letters:

1. r (first 'r')

2. r (second 'r')

3. r (third 'r')

There are 3 'r' letters in the word "strawberry". **Correct result**

Figure source: https://lmarena.ai/?model=llama-3.1-70b-instruct

# (Recap) Self-consistency CoT

Intuition: if multiple different ways of thinking lead to the same answer, one has greater confidence that the final answer is correct

Paper: https://arxiv.org/pdf/2203.11171

Join at
**slido.com**
**#1886 206**

# Agenda

- Question Answering (QA)

- Parametric Knowledge in LMs

- Hallucination

- Non-parametric Knowledge

# Introduction to Question Answering

- **Question Answering (QA):** build systems that can automatically answer questions posed by humans in natural language

- Categorization by application domain: closed-domain vs. open-domain QA

- **Closed-domain** QA: answer questions within a specific domain
    - Example: medical, legal, technical fields
    - Models are trained on specialized knowledge to be highly accurate within their domain

- **Open-domain** QA: answer questions from any domain
    - Typically rely on vast (external) knowledge sources like the web or large text corpora
    - Most LLM applications consider open-domain QA settings

# Introduction to Question Answering

- **Question Answering (QA):** build systems that can automatically answer questions posed by humans in natural language

- Categorization by modeling approach: extractive vs. abstractive QA

- **Extractive** QA: output a span of text extracted directly from a given context
  - A natural language understanding task (reading comprehension)
  - Example: context: "The human brain contains approximately 86 billion neurons" Q: "How many neurons are in the human brain?" A: "86 billion"
  - Can be done with encoder-only LMs (e.g., BERT)

- **Abstractive** QA: synthesize the answer in its own words (rephrasing/summarizing)
  - Example: context: "Albert Einstein published his theory of special relativity which introduced the famous equation $E=mc^2$, which relates energy (E) to mass (m) and the speed of light (c)" Q: "What did Einstein contribute to physics?" A: "Einstein made significant contributions to the theory of special relativity which established the relationship between energy and mass"
  - Need to use a generative LM (e.g., GPT)

# Introduction to Question Answering

- **Question Answering (QA):** build systems that can automatically answer questions posed by humans in natural language

- Categorization by access to external source: closed-book vs. open-book QA

- **Closed-book** QA: answer questions without access to any external information
  - Accuracy depends heavily on how well the training data covered the relevant information
  - Similar to a human answering a question from memory without looking anything up

- **Open-book** QA: can access external knowledge source to answer the questions
  - Typically using retrieval from reliable external sources that contain
  - Similar to a human answering a question by looking it up in a book or online resource

# Agenda

- Question Answering (QA)

- Parametric Knowledge in LMs

- Hallucination

- Non-parametric Knowledge

# Prompting LMs: Parametric Knowledge

- LMs have learned from a lot of facts in their pretraining data

- LMs can be directly prompted to generate answers to factoid questions (Closed-book QA setting)

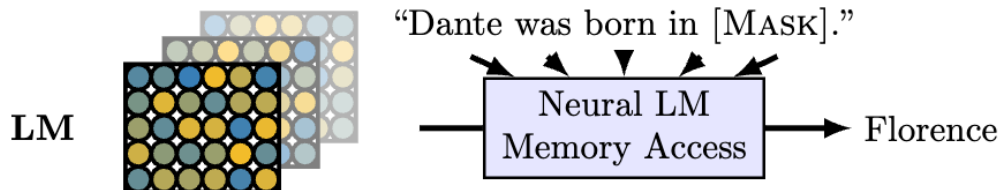- Example:

$P(w$ Q: Who wrote the book ''The Origin of Species"?  A: $)$  prompt

- Since prompting LLMs only relies on the information stored within the parameters of the model itself, this kind of knowledge is called **parametric knowledge**

14/35

# Language Model as Knowledge Bases

- **Acquisition**: LM's knowledge is derived from the vast amount of pretraining data

- **Access**: information is accessed through natural language prompts

- **Update/maintenance**: re-training/fine-tuning the model with new data

- **Pros**:
  - Handle a wide range of natural language queries with contextual understanding
  - Generalize to unseen queries not seen during training

- **Cons**:
  - May produce incorrect/outdated information
  - Lack interpretability/transparency



Paper: https://arxiv.org/pdf/1909.01066

# (Real) Knowledge Bases

- **Acquisition**: manually constructed by human annotators

- **Access**: information is accessed through queries in specific formats

- **Update/maintenance**: adding/modifying/deleting entries (incrementally) by humans

- **Pros**:
  - Precise & verifiable

- **Cons**:
  - Not able to handle natural language
  - Require massive human efforts to construct & maintain



Knowledge Graph
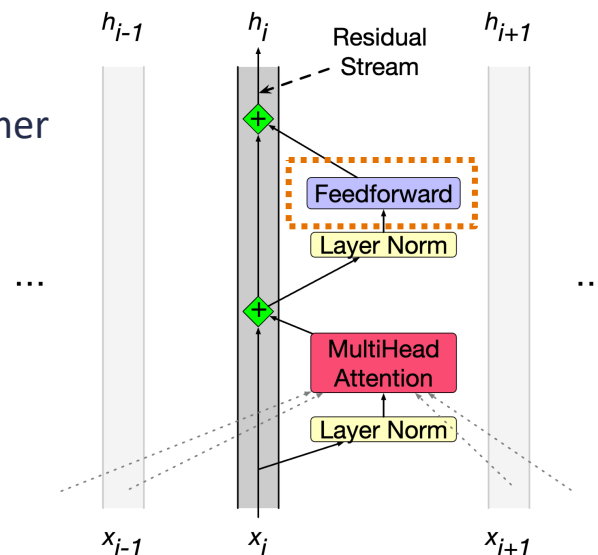
Paper: https://arxiv.org/pdf/1909.01066

# Feedforward Parameters in Transformer

- FFN in Transformer is a 2-layer network (one hidden layer, two weight matrices)

$$\mathrm{FFN}(\boldsymbol{x}_i) = \mathrm{ReLU}(\boldsymbol{x}_i \boldsymbol{W}_1) \boldsymbol{W}_2$$

- FFN constitutes ~2/3 of the total parameters of Transformer



Figure source: https://web.stanford.edu/~jurafsky/slp3/9.pdf

# Feedforward Parameters Are Neural Memories

Viewing FFN as key-value memories

$$\mathrm{FFN}(\boldsymbol{x}_i) = \mathrm{ReLU}(\boldsymbol{x}_i \boldsymbol{W}_1) \boldsymbol{W}_2$$

$$\boldsymbol{x}_i \in \mathbb{R}^{d_1}$$

$$\mathrm{FFN}(\boldsymbol{x}_i) = \mathrm{ReLU}(\boldsymbol{x}_i \boldsymbol{K}) \boldsymbol{V}$$

Rectified linear unit
(ReLU)

$$\boldsymbol{K} \in \mathbb{R}^{d_1 \times d_2} \qquad \boldsymbol{V} \in \mathbb{R}^{d_2 \times d_1}$$

key vectors (column vectors in $\boldsymbol{K}$) act as
**pattern detectors** over the input sequence

value vectors (row vectors in $\boldsymbol{V}$) represent
**distributions over the output vocabulary**

$$\mathrm{FFN}(\boldsymbol{x}_i) = \sum_{j=1}^{d_2} \mathrm{ReLU}(\boldsymbol{x}_i \cdot \boldsymbol{k}_j) \boldsymbol{v}_j$$

weights of value vectors

Paper: https://arxiv.org/pdf/2012.14913

# **Memory Keys Correspond to Input Patterns**

Each individual key vector corresponds to a specific pattern over the input prefix

| Key | Pattern | Example trigger prefixes |
|---|---|---|
| $\mathbf{k}^1_{449}$ | Ends with *"substitutes"* (shallow) | *At the meeting, Elton said that "for artistic reasons there could be no substitutes* <br> *In German service, they were used as substitutes* <br> *Two weeks later, he came off the substitutes* |
| $\mathbf{k}^6_{2546}$ | Military, ends with *"base"/"bases"* (shallow + semantic) | *On 1 April the SRSG authorised the SADF to leave their bases* <br> *Aircraft from all four carriers attacked the Australian base* <br> *Bombers flying missions to Rabaul and other Japanese bases* |
| $\mathbf{k}^{10}_{2997}$ | a *"part of"* relation (semantic) | *In June 2012 she was named as one of the team that competed* <br> *He was also a part of the Indian delegation* <br> *Toy Story is also among the top ten in the BFI list of the 50 films you should* |
| $\mathbf{k}^{13}_{2989}$ | Ends with a time range (semantic) | *Worldwide, most tornadoes occur in the late afternoon, between 3 pm and 7* <br> *Weekend tolls are in effect from 7:00 pm Friday until* <br> *The building is open to the public seven days a week, from 11:00 am to* |
| $\mathbf{k}^{16}_{1935}$ | TV shows (semantic) | *Time shifting viewing added 57 percent to the episode's* <br> *The first season set that the episode was included in was as part of the* <br> *From the original NBC daytime version , archived* |

Paper: https://arxiv.org/pdf/2012.14913
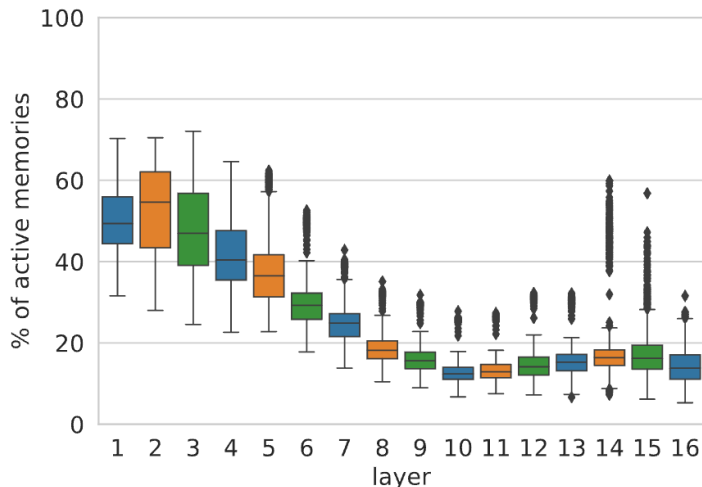
# Memory Values Correspond to Output Tokens

Each value vector (roughly) matches a predicted token distribution

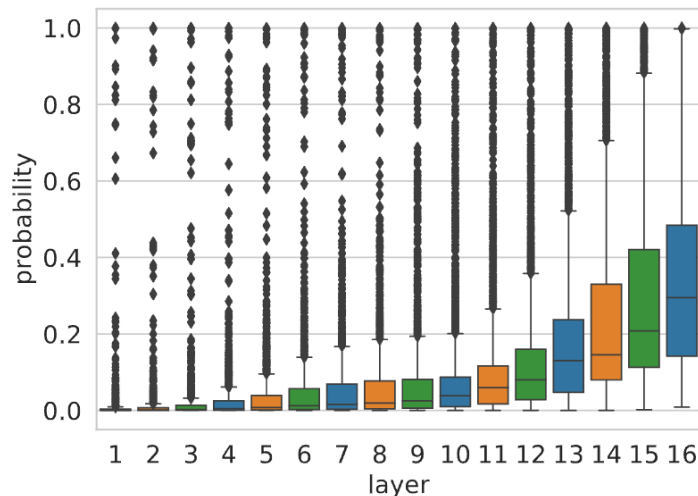| Value | Prediction | Trigger example |
|---|---|---|
| $\mathbf{v}_{222}^{15}$ | each | But when bees and wasps resemble *each* |
| $\mathbf{v}_{752}^{16}$ | played | Her first role was in Vijay Lalwani's psychological thriller Karthik Calling Karthik, where Padukone was cast as the supportive girlfriend of a depressed man (*played* |
| $\mathbf{v}_{2601}^{13}$ | extratropical | Most of the winter precipitation is the result of synoptic scale, low pressure weather systems (large scale storms such as *extratropical* |
| $\mathbf{v}_{881}^{15}$ | part | Comet served only briefly with the fleet, owing in large *part* |
| $\mathbf{v}_{2070}^{16}$ | line | Sailing from Lorient in October 1805 with one ship of the *line* |
| $\mathbf{v}_{3186}^{12}$ | jail | On May 11, 2011, four days after scoring 6 touchdowns for the Slaughter, Grady was sentenced to twenty days in *jail* |

Paper: https://arxiv.org/pdf/2012.14913

# Memory Aggregation

- "Active" memories (memory vectors with non-zero coefficients) are typically sparse
- The residual connection sequentially refines token prediction from layer to layer



Fraction of active memory
units across layers

Output token probability is
gradually refined across layers

Paper: https://arxiv.org/pdf/2012.14913

# Further Reading on LLM Parametric Knowledge

- How Much Knowledge Can You Pack Into the Parameters of a Language Model? [Roberts et al., 2020]

- Extracting Training Data from Large Language Models [Carlini et al., 2021]

- Locating and Editing Factual Associations in GPT [Meng et al., 2022]

# Agenda

- Question Answering (QA)

- Parametric Knowledge in LMs

- Hallucination

- Non-parametric Knowledge

# Hallucination

- **Hallucination**: LM generates information that is factually incorrect, misleading, or fabricated, even though it may sound plausible or convincing

- Why does hallucination happen?
  - Limited knowledge: LLMs are trained on finite datasets, which don't have access to all possible information; when asked about topics outside their training data, they may generate plausible-sounding but incorrect responses
  - Overgeneralization: LLMs may apply patterns they've learned from one context to another where they don't apply, leading to incorrect conclusions
  - Lack of common sense: While LLMs can process and generate human-like text, they often lack the ability to apply commonsense reasoning to their outputs
  - …

# Hallucination Examples

- (**Limited knowledge**) Q: "What were the main features of the iPhone 15 Pro Max?"
  LLM (trained before 2023): "The iPhone 15 Pro Max features a revolutionary holographic display, quantum computing chip, and telepathic user interface."

- (**Overgeneralization**) Q: "How do you form the past tense in Japanese?"
  LLM: "In Japanese, you typically add '-ed' to the end of verbs to form the past tense, just like in English." (incorrect)

- (**Lack of common sense**) Q: "How many tennis balls can fit in a typical smartphone?"
  LLM: "Approximately 15-20 tennis balls can fit in a typical smartphone, depending on the model and screen size."

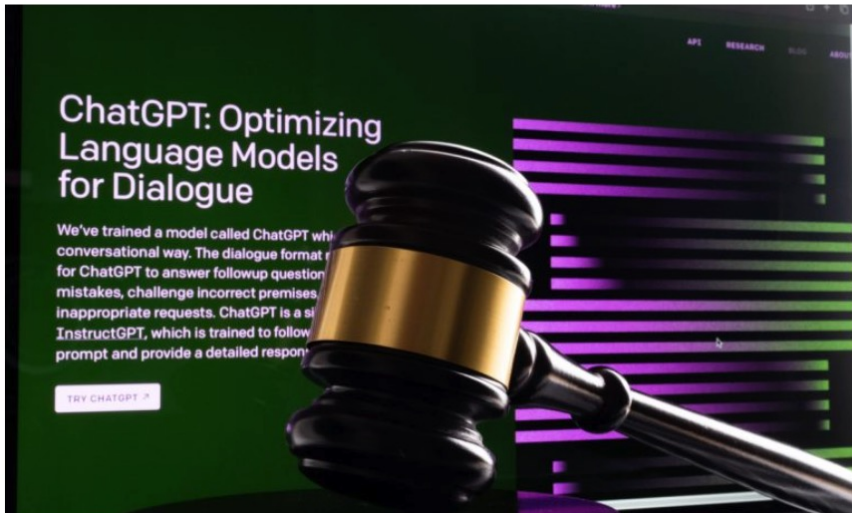what's your knowledge cutoff date?

My knowledge cutoff date is October 2023. This means I don't have information on events or developments that have occurred after that time. How can I assist you with your question?

Knowledge cutoff date of GPT4o

# Concerns About Hallucination

Still a concerning issue in modern LLMs!

## Attorneys Face Sanctions After Citing Case Law 'Hallucinated' by ChatGPT

BY **PYMNTS** | MAY 30, 2023

**ChatGPT: Optimizing Language Models for Dialogue**

We've trained a model called ChatGPT whi... conversational way. The dialogue format ... for ChatGPT to answer followup question... mistakes, challenge incorrect premises,... inappropriate requests. ChatGPT is a si... InstructGPT, which is trained to follow ... prompt and provide a detailed respon...

TRY CHATGPT ↗

Figure source: https://www.pymnts.com/artificial-intelligence-2/2023/attorneys-face-sanctions-after-citing-information-hallucinated-by-chatgpt/

# Further Reading on Hallucination

- LLM Lies: Hallucinations are not Bugs, but Features as Adversarial Examples [Yao et al., 2023]

- Towards Mitigating Hallucination in Large Language Models via Self-Reflection [Ji et al., 2023]

- Hallucination is Inevitable: An Innate Limitation of Large Language Models [Xu et al., 2024]

# Agenda

- Question Answering (QA)

- Parametric Knowledge in LMs

- Hallucination
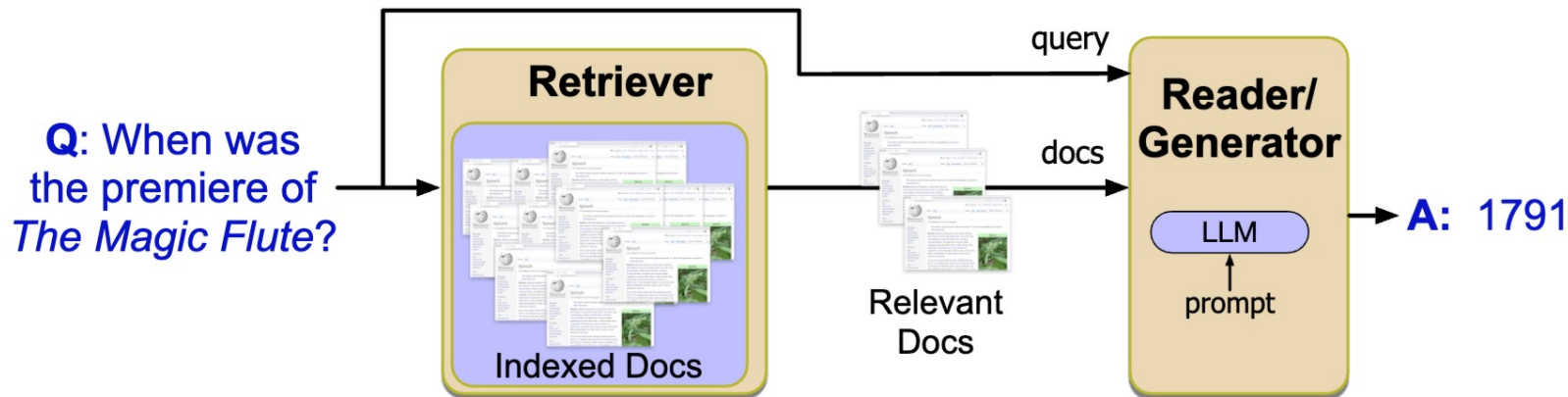
- **Non-parametric Knowledge**

# Non-parametric Knowledge

- **Non-parametric knowledge**: (external) information not stored in the model's parameters but can be accessed or retrieved when needed

- Examples:
    - External knowledge bases/graphs
    - Pretraining corpora
    - User-provided documents/passages

- Non-parametric knowledge is typically used to **augment** parametric knowledge (typically via **retrieval**) for more accurate factoid question answering

- Benefits of **non-parametric knowledge**
    - Incorporate more information without increasing model size
    - Easier updates and modifications to the knowledge base
    - Improve model interpretability

# Overview: Retrieval-Augmented Generation

- Use a retriever to obtain relevant documents to the query from an external text collection

- Use LLMs to generate answers given the documents and a prompt



Figure source: https://web.stanford.edu/~jurafsky/slp3/14.pdf

# Sparse vs. Dense Retrieval

- **Sparse** retrieval: based on traditional information retrieval (IR) techniques where the representations of documents and queries are sparse (most vector values are zero)
  - Example: TF-IDF
  - Pros: simple and interpretable
  - Cons: lack semantic understanding

- **Dense** retrieval: encode documents and queries into dense vectors (embeddings) using deep neural networks
  - Example: BERT-based encoding methods
  - Pros: semantic & contextualized understanding
  - Cons: computationally more expensive and less interpretable

# Summary: Scaling Up LLMs & Emergent Abilities

- Scaling up LLMs involves three factors
  - Add more parameters (scaling up model size)
  - Add more data (scaling up dataset size)
  - Train for more iterations (scaling up compute)

- The language modeling loss exhibits a power-law relationship with each of the three factors

- Given a specific compute, an optimal model size can be derived

- Emergent ability: a model has random performance until a certain scale, after which performance increases notably

# Summary: Reasoning with LLMs

- Reasoning: perform deductive, inductive, commonsense, or logical reasoning by generating text with language models

- Chain-of-thought (CoT): the model breaks down complex problems into a step-by-step reasoning process Instead of directly answering the question

- CoT can be either zero-shot or few-shot

- Self-consistency CoT involves generating multiple reasoning paths and deriving the final answer with majority voting

# Summary: Question Answering & Knowledge

- Question answering (QA) can be categorized in different ways:
  - Closed-domain vs. open-domain QA
  - Extractive vs. abstractive QA
  - Closed-book vs. open-book QA

- LMs can be used as knowledge bases via prompting

- Feedforward network parameters can store factual knowledge (parametric knowledge)

- LMs can hallucinate, which can be mitigated by retrieving non-parametric knowledge

# Thank You!

**Yu Meng**
University of Virginia
yumeng5@virginia.edu