# 大模型技术基础

《大语言模型》编写团队：赵鑫

# 大语言模型

➤ 定义：通常是指具有超大规模参数的预训练语言模型

➤ 架构：主要为 Transformer解码器架构

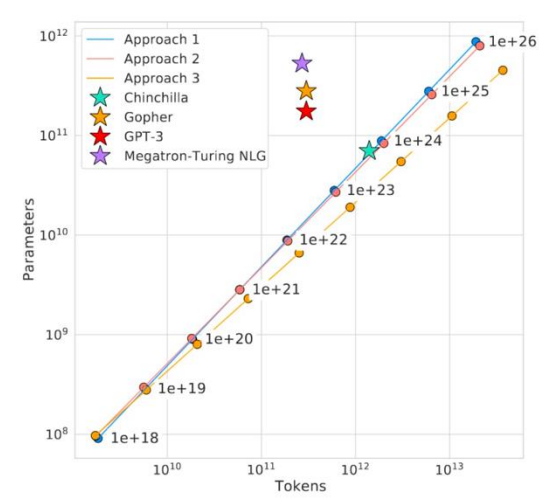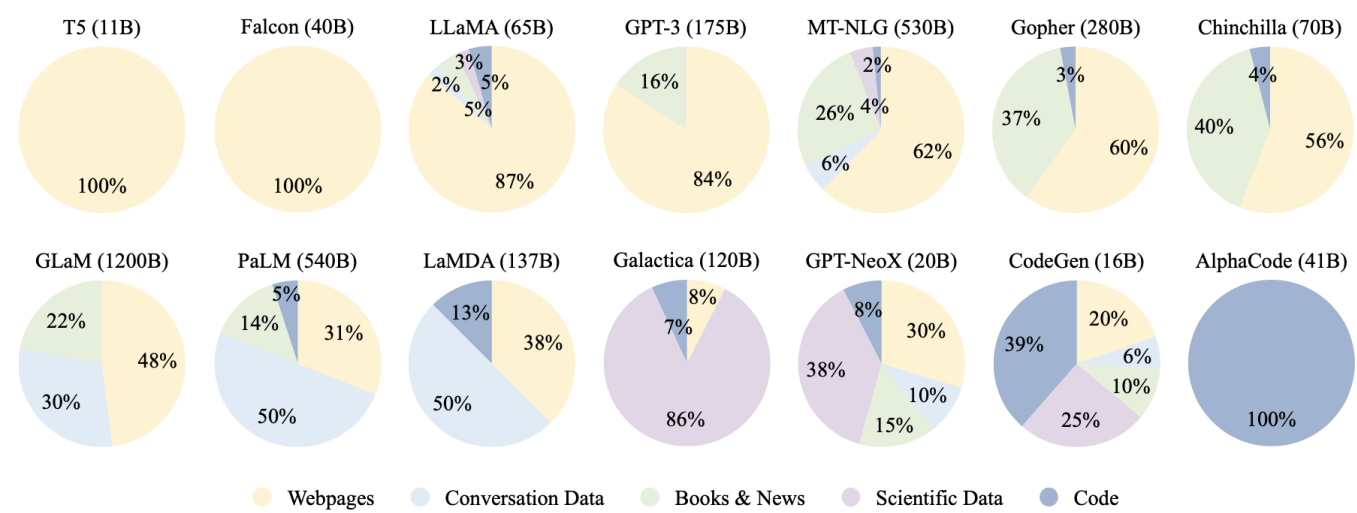➤ 训练：预训练（base model）、后训练（instruct model）



大语言模型，2024

《大语言模型》教材课件

# 大语言模型

➤ 定义：通常是指具有超大规模参数的预训练语言模型

➤ 架构：主要为 Transformer解码器架构

➤ 训练：预训练（base model）、后训练（instruct model）

| 对比方面 | 预训练 (Pre-training) | 后训练 (Post-training) |
|---|---|---|
| 核心目标 | 建立模型基础能力 | 将基座模型适配到具体应用场景 |
| 数据资源 | 数万亿词元的自然语言文本 | 数十万、数百万到数千万指令数据 |
| 所需算力 | 耗费百卡、千卡甚至万卡算力数月时间 | 耗费数十卡、数百卡数天到数十天时间 |
| 使用方式 | 通常为few-shot提示 | 可以直接进行zero-shot使用 |

*此部分算力估计为一个大致估计，需要根据模型大小、数据数量、训练框架等多方面因素确定

《大语言模型》教材课件

➢ 大语言模型预训练（Pre-training）

　　➢ 使用与下游任务无关的大规模数据进行模型参数的初始训练

　　　　➢ 基于Transformer解码器架构，进行下一个词预测

　　　　➢ 数据数量、数据质量都非常关键

➢ **大语言模型后训练（Post-Training）**

  ➢ **指令微调（Instruction Tuning）**
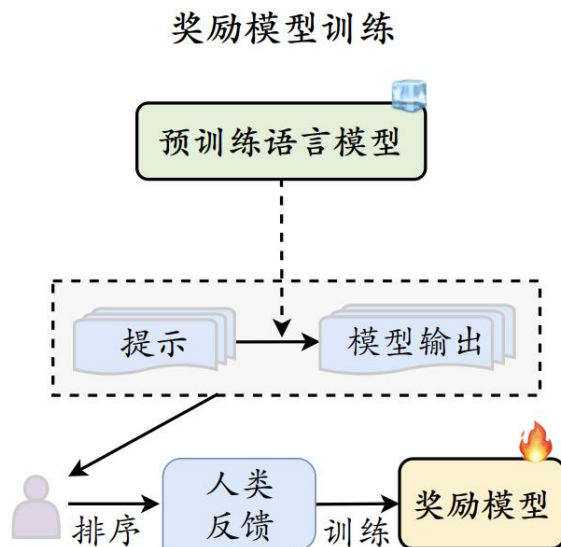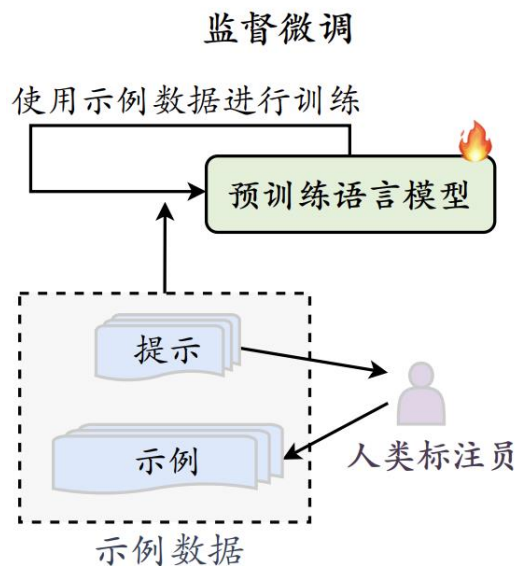
   ➢ 使用输入与输出配对的指令数据对于模型进行微调

   ➢ 提升模型通过问答形式进行任务求解的能力

# 大语言模型构建概览

➢ 大语言模型后训练（Post-Training）

➢ 人类对齐（Human Alignment）

➢ 将大语言模型与人类的期望、需求以及价值观对齐

➢ 基于人类反馈的强化学习对齐方法（RLHF）

# 大语言模型构建概览

➢ 大模型的研发已经成为一项系统工程



GPT-4 Technical Report

《大语言模型》教材课件

# 扩展定律

➢ 什么是扩展定律

　➢ 通过扩展参数规模、数据规模和计算算力，大语言模型的能力会出现显著提升

　➢ 扩展定律在本次大模型浪潮中起到了重要作用



Test Loss vs Compute (PF-days, non-embedding): $L = (C_{min}/2.3 \cdot 10^8)^{-0.050}$

Test Loss vs Dataset Size (tokens): $L = (D/5.4 \cdot 10^{13})^{-0.095}$

Test Loss vs Parameters (non-embedding): $L = (N/8.8 \cdot 10^{13})^{-0.076}$

Scaling Laws for Neural Language Models, Arxiv, 2020

# 扩展定律

➢ KM扩展定律

  ➢ OpenAI 团队建立了神经语言模型性能与参数规模（$N$）、数据规模（$D$）和计算算力（$C$）之间的幂律关系

$$L(N) = \left(\frac{N_c}{N}\right)^{\alpha_N}, \quad \alpha_N \sim 0.076, N_c \sim 8.8 \times 10^{13}$$

$$L(D) = \left(\frac{D_c}{D}\right)^{\alpha_D}, \quad \alpha_D \sim 0.095, D_c \sim 5.4 \times 10^{13}$$

$$L(C) = \left(\frac{C_c}{C}\right)^{\alpha_C}, \quad \alpha_C \sim 0.050, C_c \sim 3.1 \times 10^{8}$$



Loss vs Model and Dataset Size

Params: 708M, 302M, 85M, 3M, 25M, 393.2K
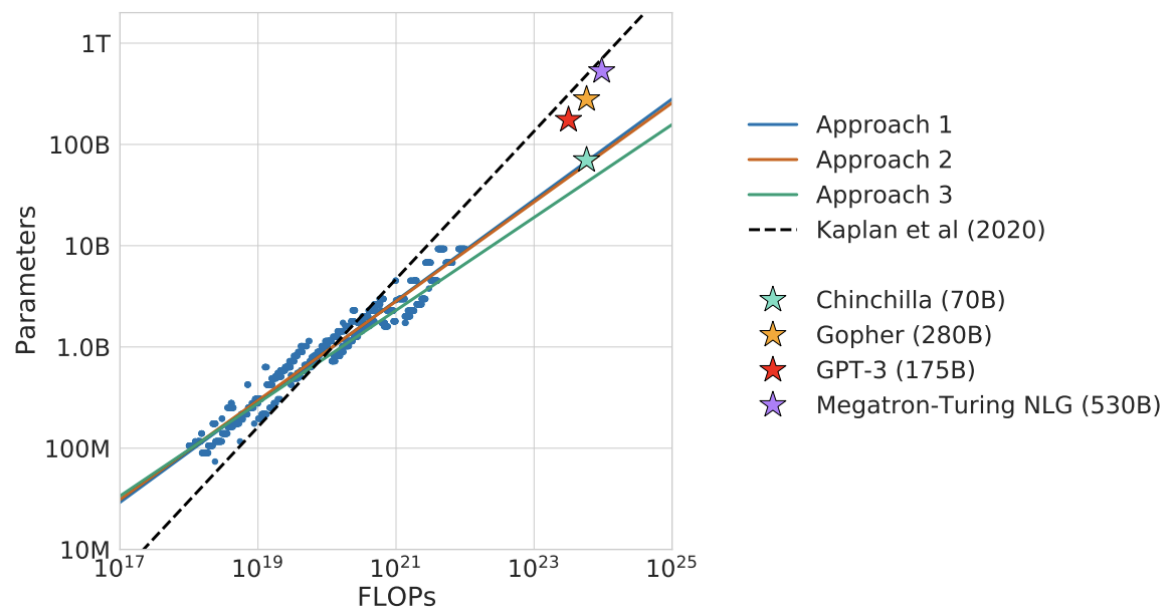
# 扩展定律

➤ Chinchilla扩展定律

  ➤ DeepMind 团队于 2022 年提出了另一种形式的扩展定律,旨在指导大语言模型充分利用给定的算力资源优化训练

$$L(N, D) = E + \frac{A}{N^\alpha} + \frac{B}{D^\beta},$$

$$N_{\text{opt}}(C) = G\left(\frac{C}{6}\right)^a, \quad D_{\text{opt}}(C) = G^{-1}\left(\frac{C}{6}\right)^b,$$

> ## 深入讨论

> > ## 模型的语言建模损失可以进行下述分解

$$L(x) = \underbrace{L_\infty}_{\text{不可约损失}} + \underbrace{\left(\frac{x_0}{x}\right)^{\alpha_x}}_{\text{可约损失}}$$

可约损失：真实分布和模型分布之间KL散度，可通过优化减少

不可约损失：真实数据分布的熵，无法通过优化减少

> > ## 扩展定律可能存在边际效益递减

> > > 随着模型参数、数据数量的扩展，模型性能增益将逐渐减小

> > > 目前开放数据已经接近枯竭，难以支持扩展定律的持续推进

# 扩展定律

➢ 深入讨论

➢ 可预测的扩展（Predictable Scaling）

➢ 使用小模型性能去预估大模型的性能，或帮助超参数选择

➢ 训练过程中使用模型早期性能来预估后续性能



(a) Batch size scaling curve

(b) Learning rate scaling curve

# 涌现能力

➢ 什么是涌现能力

  ➢ 原始论文定义："在小型模型中不存在、但在大模型中出现的能力"

  ➢ 模型扩展到一定规模时，特定任务性能突然出现显著跃升趋势，远超随机水平

# 涌现能力

➤ 涌现能力可能部分归因于评测设置

  ➤ 本课程定义其为"**代表性能力**"，并不区分是否在小模型中存在



Are Emergent Abilities of Large Language Models a Mirage?, NIPS 2023

# 涌现能力

➢ 代表性能力

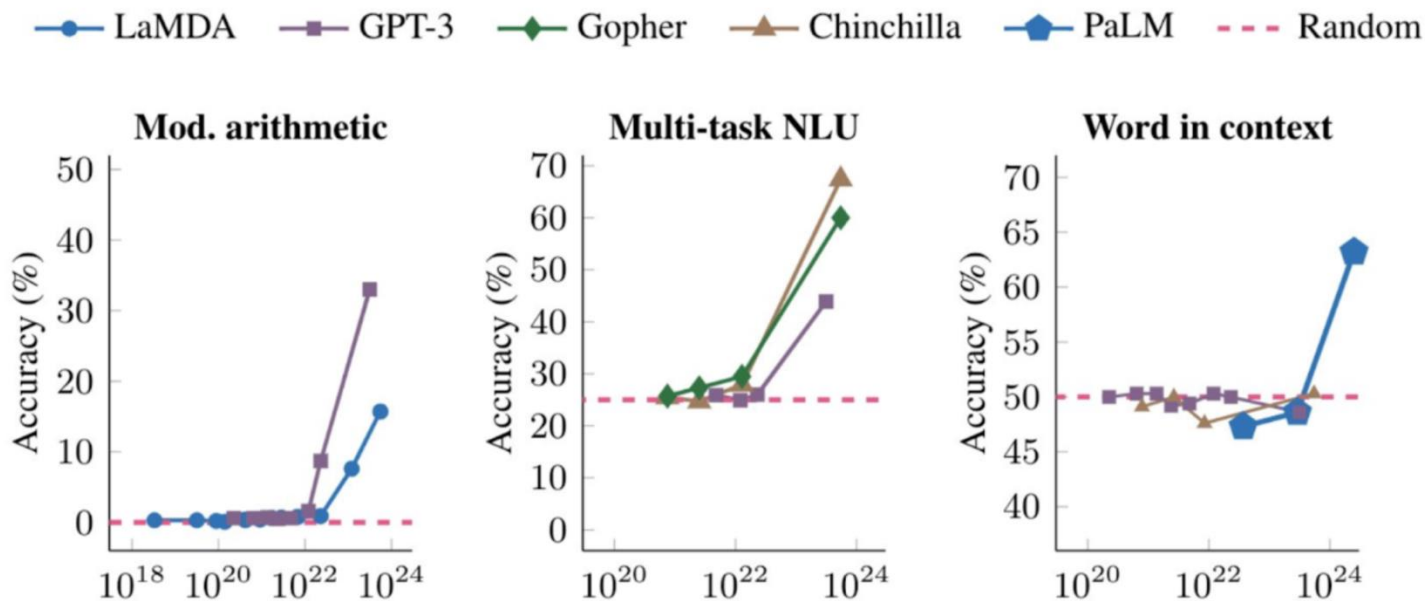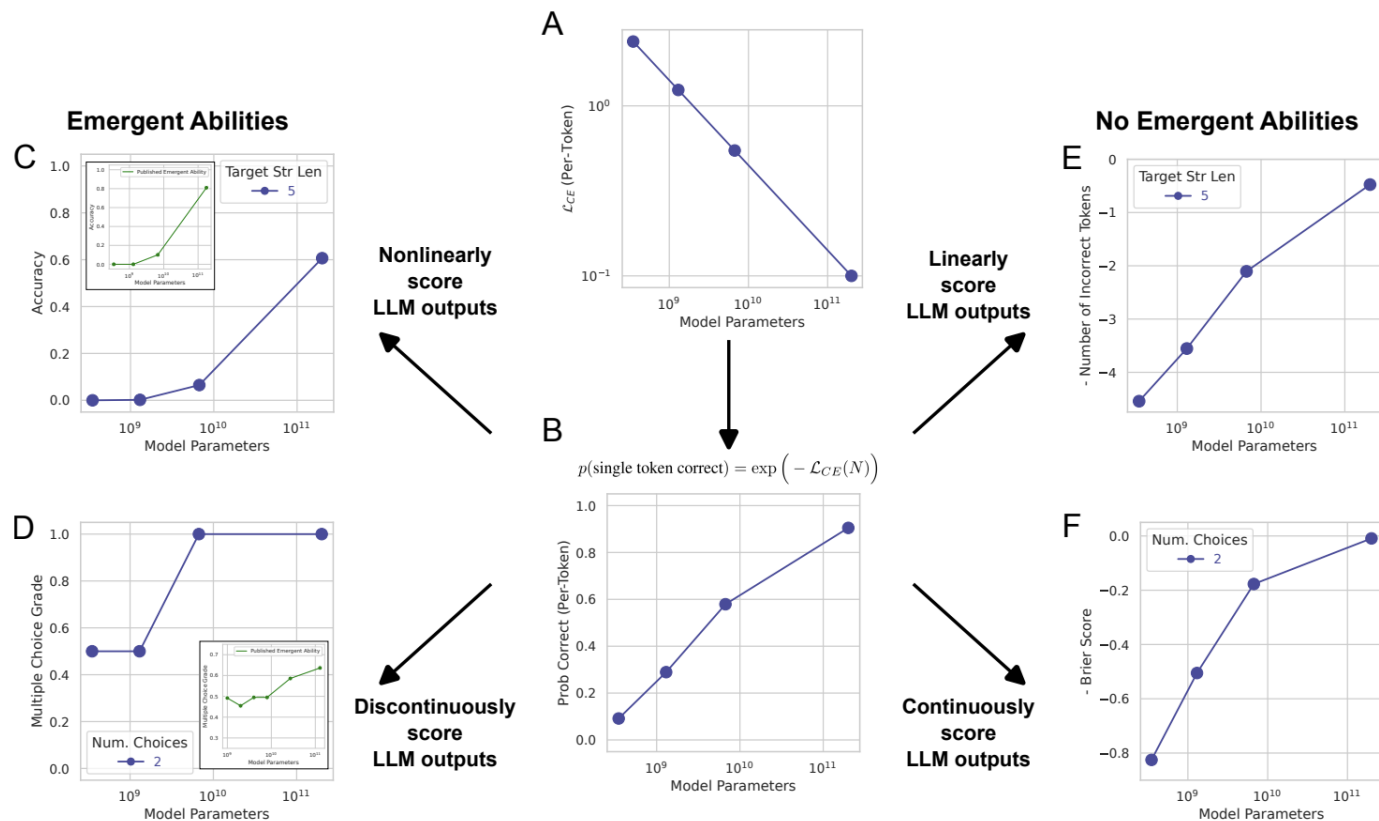  ➢ 指令遵循（Instruction Following）

    ➢ 大语言模型能够按照自然语言指令来执行对应的任务

> 代表性能力

> 上下文学习（In-context Learning）

> 在提示中为语言模型提供自然语言指令和任务示例，无需显式梯度更新就能为测试样本生成预期输出



Answer the following mathematical reasoning questions:

Q: If you have 12 candies and you give 4 candies to your friend, how many candies do you have left?
A: The answer is 8.
Q: If a rectangle has a length of 6 cm and a width of 3 cm, what is the perimeter of the rectangle?
A: The answer is 18cm.

$N$x

Q: Sam has 12 marbles. He gives 1/4 of them to his sister. How many marbles does Sam have left?



Language Models are Few-shot Learners, NIPS 2020

# 涌现能力

➢ 代表性能力

➢ 逐步推理（Step-by-step Reasoning）

➢ 在提示中引入任务相关的中间推理步骤来加强复杂任务的求解，从而获得更可靠的答案

**上下文学习**

Answer the following mathematical reasoning questions:

Nx

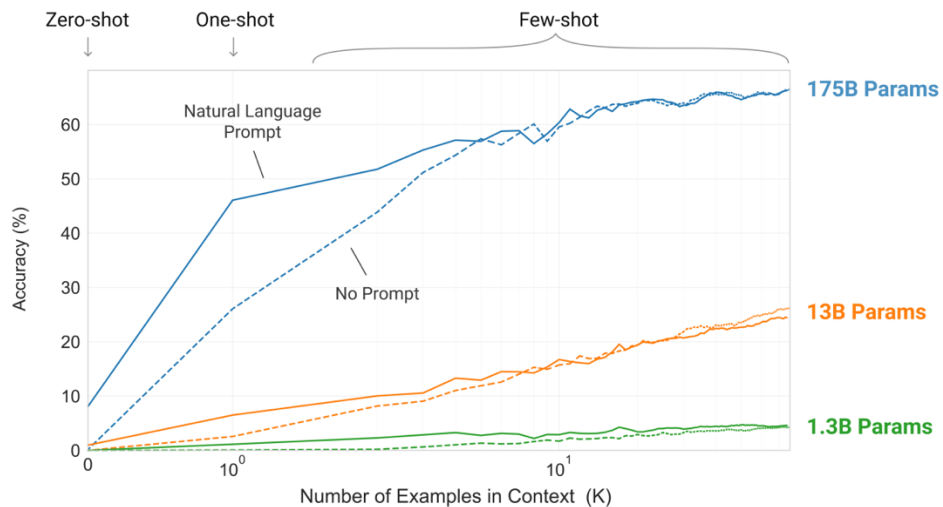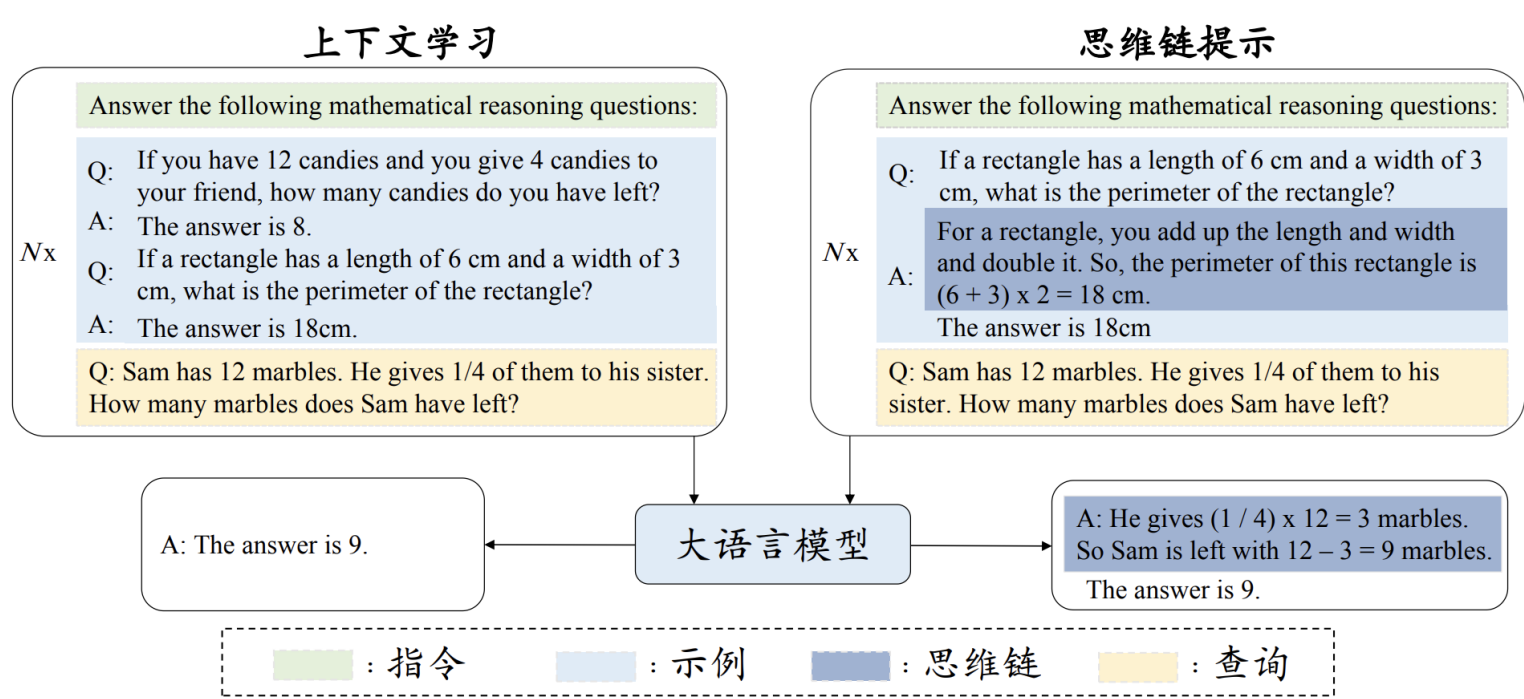Q: If you have 12 candies and you give 4 candies to your friend, how many candies do you have left?
A: The answer is 8.
Q: If a rectangle has a length of 6 cm and a width of 3 cm, what is the perimeter of the rectangle?
A: The answer is 18cm.

Q: Sam has 12 marbles. He gives 1/4 of them to his sister. How many marbles does Sam have left?

**思维链提示**

Answer the following mathematical reasoning questions:

Nx

Q: If a rectangle has a length of 6 cm and a width of 3 cm, what is the perimeter of the rectangle?
A: For a rectangle, you add up the length and width and double it. So, the perimeter of this rectangle is (6 + 3) x 2 = 18 cm.
The answer is 18cm

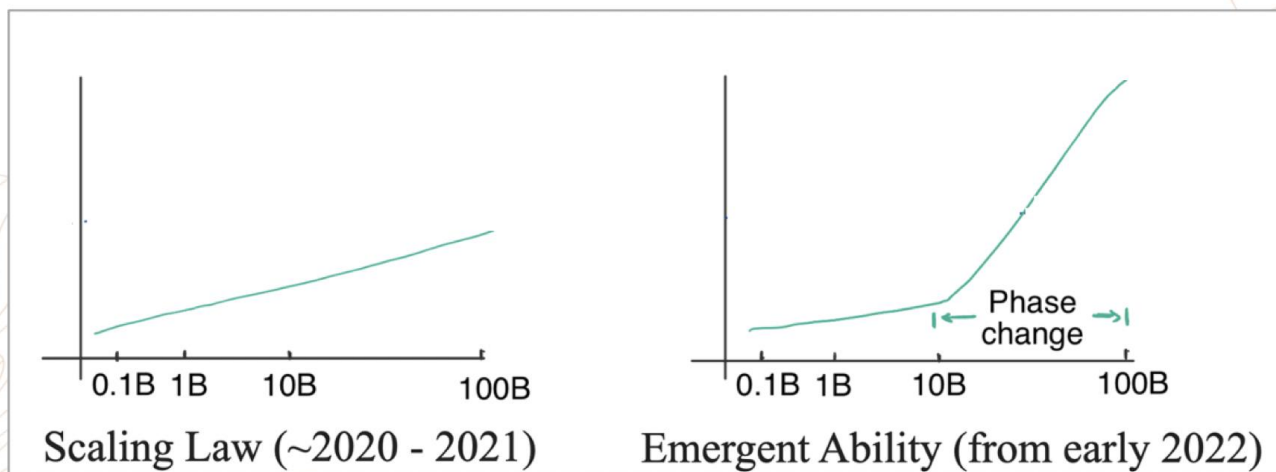Q: Sam has 12 marbles. He gives 1/4 of them to his sister. How many marbles does Sam have left?

A: The answer is 9.

**大语言模型**

A: He gives (1 / 4) x 12 = 3 marbles.
So Sam is left with 12 − 3 = 9 marbles.
The answer is 9.

☐ :指令   ☐ :示例   ☐ :思维链   ☐ :查询

# 涌现能力

➤ 涌现能力与扩展定律的关系

　　➤ 涌现能力和扩展定律是两种描述规模效应的度量方法

Scaling law describes a *predictable* increase pattern, but with *diminishing return*



Scaling Law (~2020 - 2021)

Emergent Ability (from early 2022)

Phase change

Model scaling is the key to the emergence of strong abilities

可以理解为是一种较为平滑的多任务损失平均（LM loss）

Emergent abilities transcend the scaling law, making the increase *unpredictable but profitable*

非平滑的、某种特定能力或任务的性能跃升 (Task loss)

《大语言模型》教材课件

# 总结

➢ 大模型核心技术

  ➢ 规模扩展：扩展定律奠定了早期大模型的技术路线，产生了巨大的性能提升

  ➢ 数据工程：数据数量、数据质量以及配制方法极其关键

  ➢ 高效预训练：需要建立可预测、可扩展的大规模训练架构

  ➢ 能力激发：预训练后可以通过微调、对齐、提示工程等技术进行能力激活

  ➢ 人类对齐：需要设计对齐技术减少模型使用风险，并进一步提升模型性能

  ➢ 工具使用：使用外部工具加强模型的弱点，拓展其能力范围

谢谢