

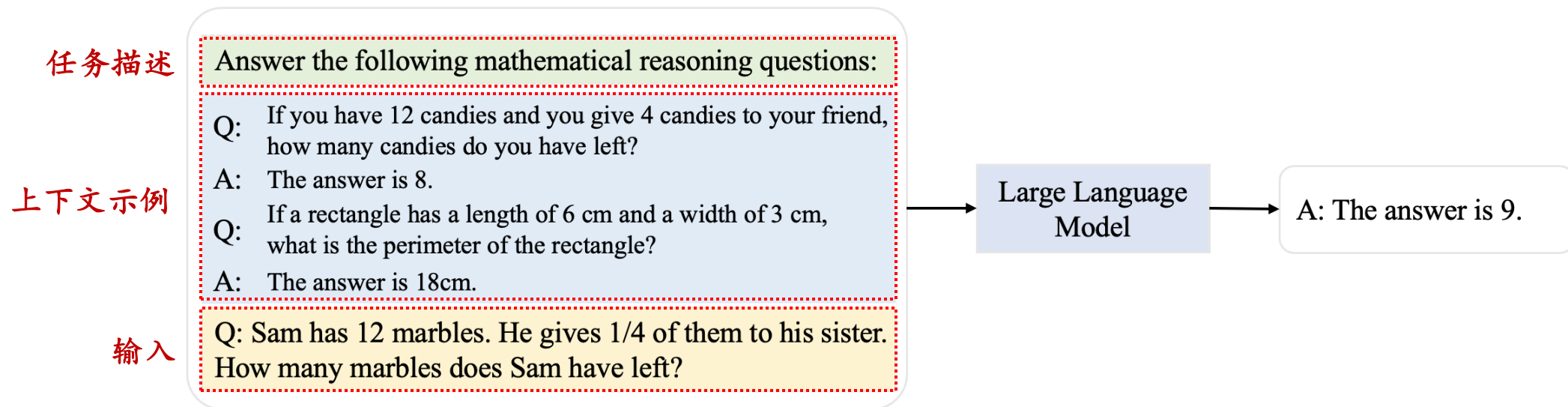
# 上下文学习

《大语言模型》编写团队：李军毅

## ➤ 定义

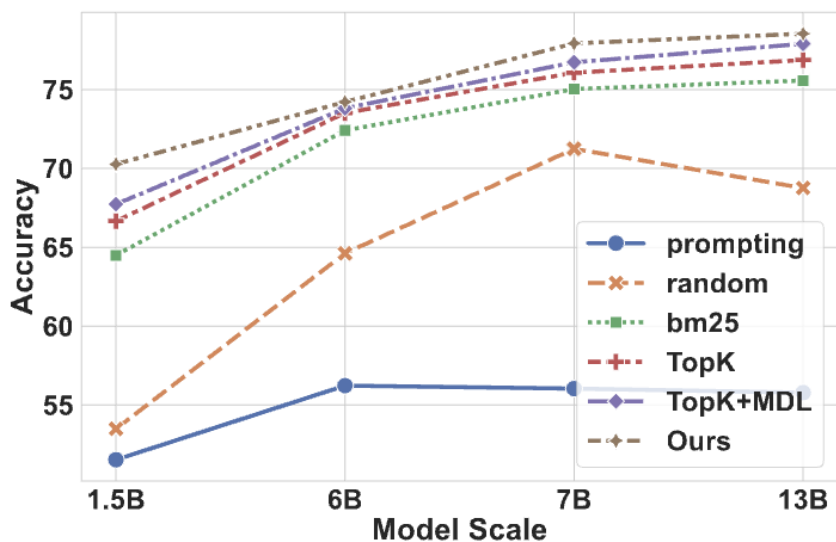
➤ 任务描述和 (或) 上下文示例所组成的自然语言文本作为提示

$$\underbrace{\text{LLM}}_{\text{大语言模型}} \left( \underbrace{I}_{\text{任务描述}}, \underbrace{f(x_1, y_1), \dots, f(x_k, y_k)}_{\text{示例}}, f(\underbrace{x_{k+1}}_{\text{输入}}, \underbrace{\quad}_{\text{答案}}) \right) \rightarrow \hat{y}_{k+1}.$$

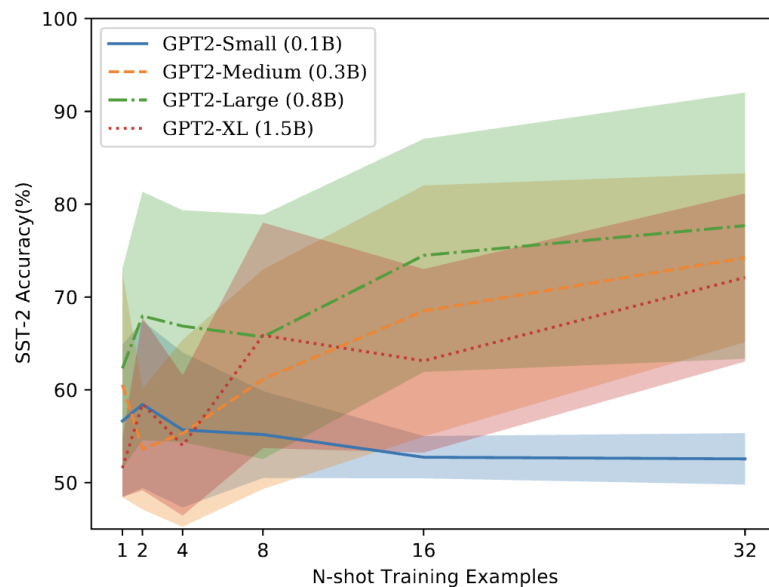


## ➤ 关键因素：上下文示例

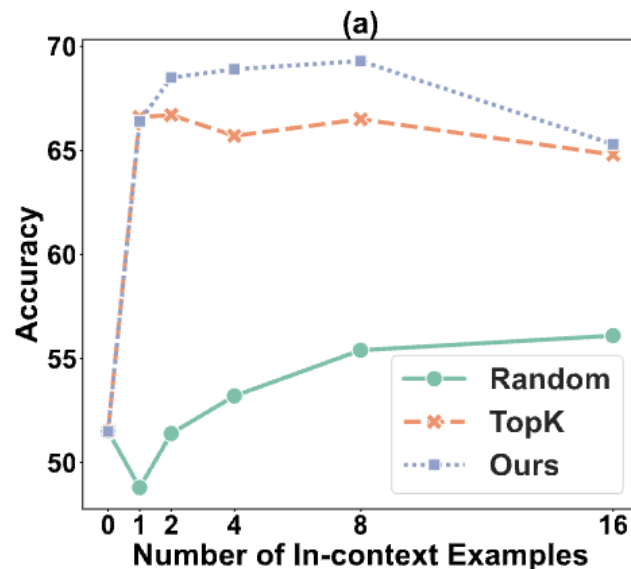
### ➤ 示例选择、示例顺序、示例数目



不同示例选择策略的性能差异



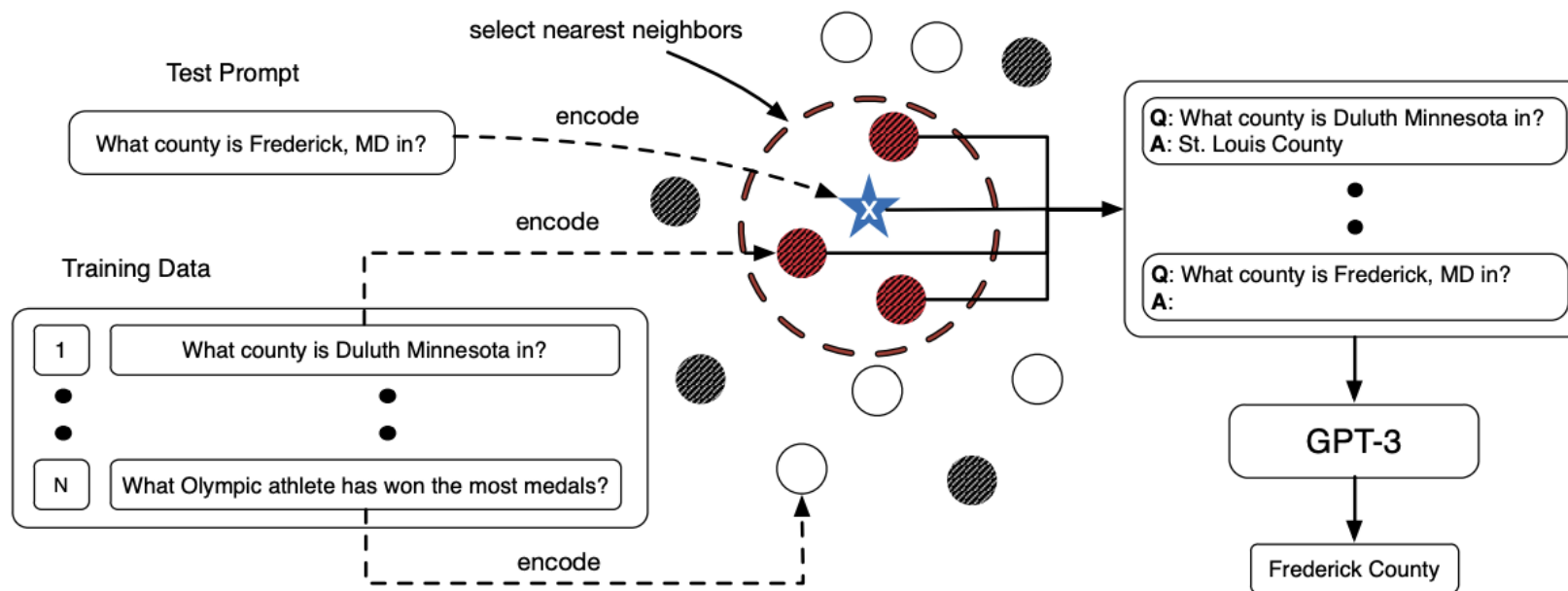
不同示例顺序导致性能方差大



不同示例数目的性能差异

## ➤ 基于相关度排序的方法

- 基于  $k$  近邻的相似度检索算法：根据候选样本与测试样本的嵌入语义相似度进行排序，并选择出最相关的  $k$  个示例



## ➤ 基于集合多样性的方法

➤ 针对特定任务选择出具有代表性的、信息覆盖性高的示例集合

➤ 如 MMR (Maximum Margin Relevance)、DPP (Determinantal Point Process) 算法

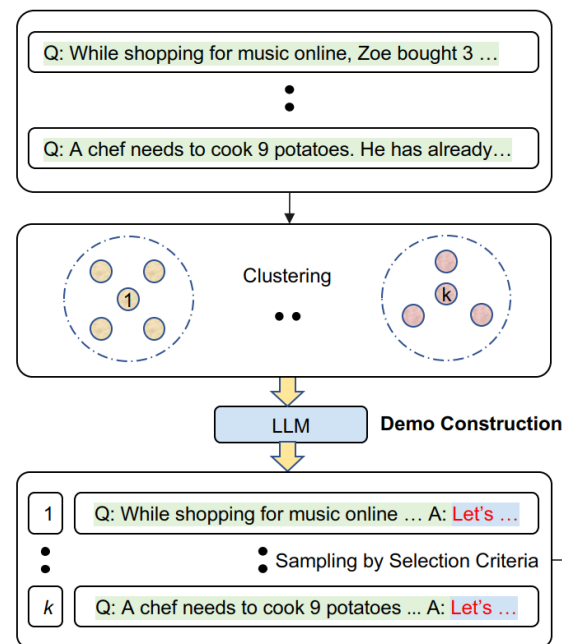
### Algorithm 1 MMR-Based Exemplar Selection

```
1: procedure MMRSELECT( $D, q, k, \mathcal{S}$ )  
   input: exemplar pool  $D = \{q_1 \dots q_n\}$ , test query  $q$ , num-  
   ber of shots  $m$  and similarity measurement  $\mathcal{S}$   
   output: selected exemplars  $T = \{q_1 \dots q_m\}$   
2:    $\mathbb{S} := [[\mathcal{S}(q_i, q_j)]]_{q_i, q_j \in D}$ ;  $\triangleright$  the pairwise similarity  
   between exemplars in  $D$   
3:    $\mathbb{Q} := [\mathcal{S}(q, q_i)]_{q_i \in D}$ ;  $\triangleright$  the similarity between query  
   and exemplars in  $T$   
4:    $T := \{\}$ ;  
5:   while  $|T| < k$  do  
6:      $\hat{q} := \text{Equation}(1)$ ;  $\triangleright$  get the next exemplar  
     based on Eq (1)  
7:      $T.\text{add}(\hat{q})$   
8:   return  $T$ ;
```

$$\arg \max_{q_j \in D/T} \lambda \mathcal{S}(q, q_j) - (1 - \lambda) \max_{q_i \in T} \mathcal{S}(q_j, q_i) \quad (1)$$

与问题相关但与  
集合示例相似度低

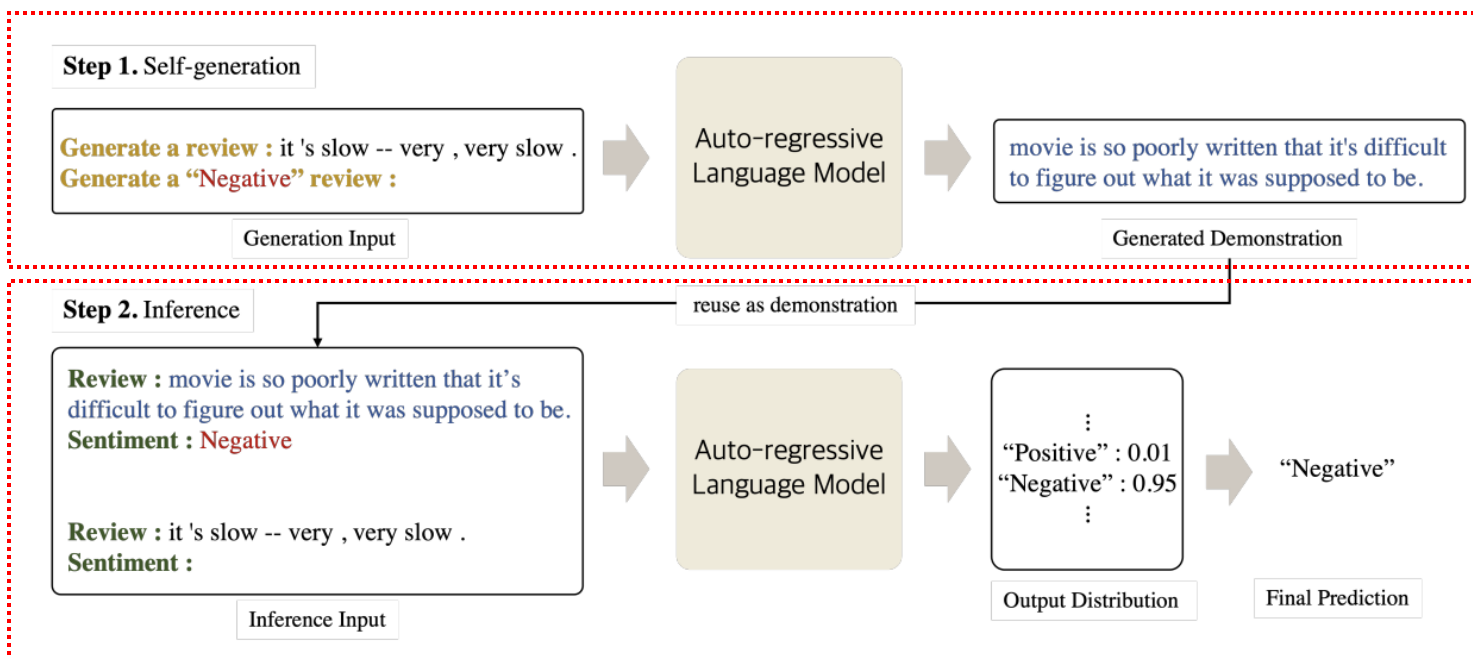
将样例进行聚类  
选取多样化的样例



# 示例选择

## ➤ 基于大语言模型的方法

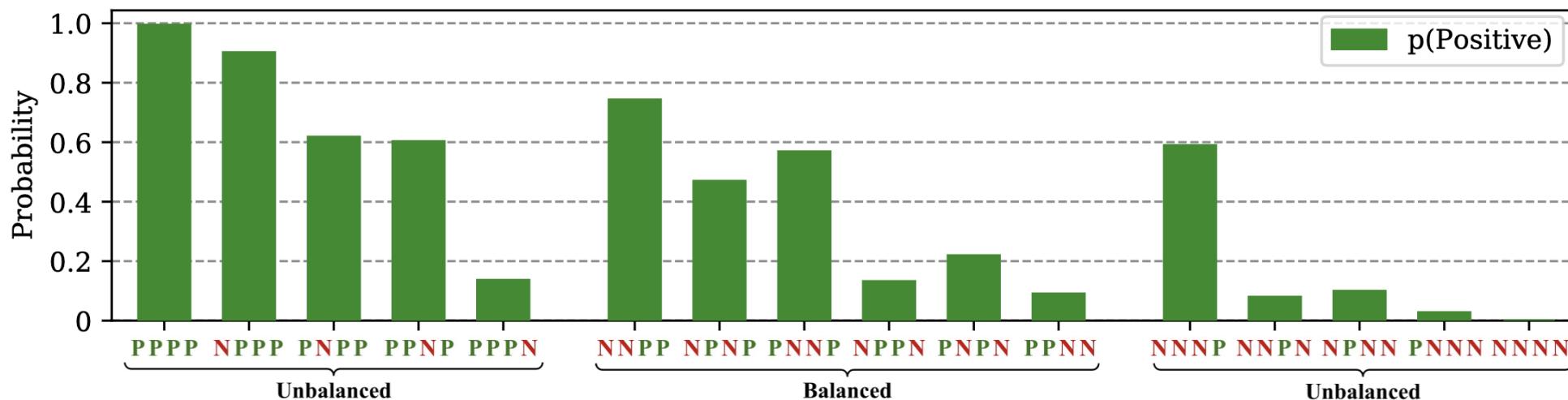
- 根据加入示例前后大模型性能的增益来评估示例选择的有效性
- 直接利用大模型生成与测试样本相关的上下文示例



生成与输入评论相关且为  
“消极”的评论作为示例

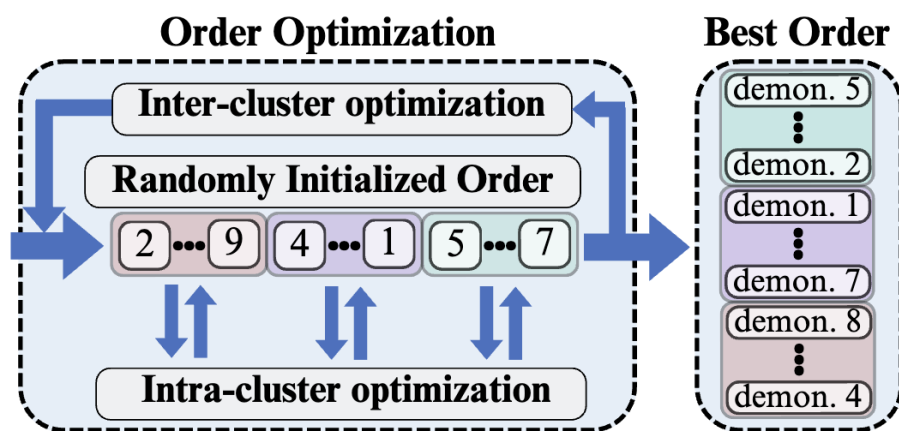
基于大模型生成的评论示  
例预测输入评论的情感

- 基于相关度的示例顺序
  - 大模型在预测时更倾向于提示末端的示例 (更靠近测试输入)
  - 根据示例与测试样本的相似度排序，相似度越高则越靠近测试样本

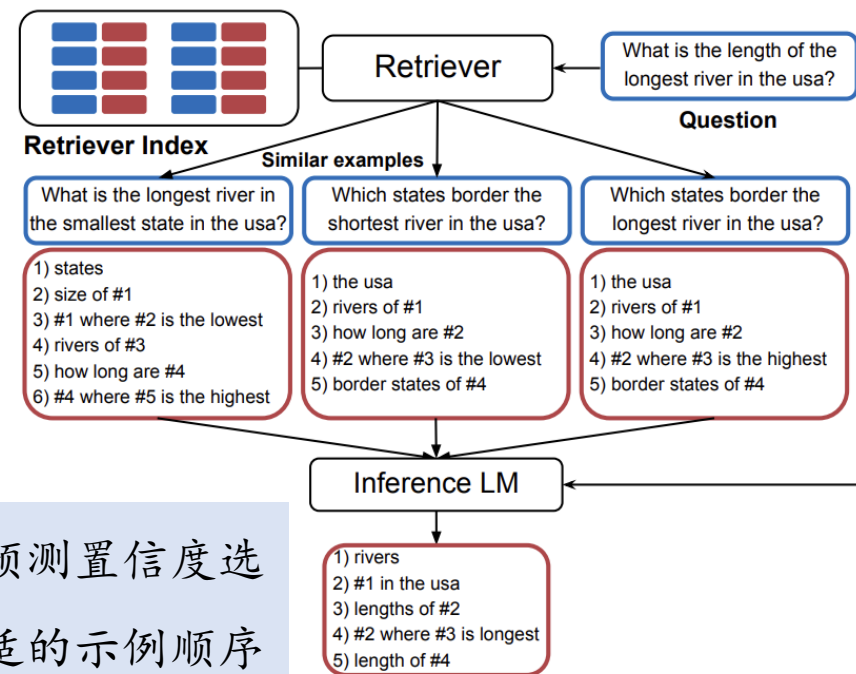


测试样本的预测结果与末端示例的标签高度正相关

- 基于任务性能的示例顺序
  - 选择少量数据测试大模型基于不同示例顺序的任务性能
  - 大模型预测置信度也可作为评估示例顺序的指标



根据任务性能优化示例顺序，  
选择最佳的示例顺序



根据预测置信度选  
择合适的示例顺序



# 上下文学习的工作机制

## ➤ 大模型执行上下文学习的两种范式

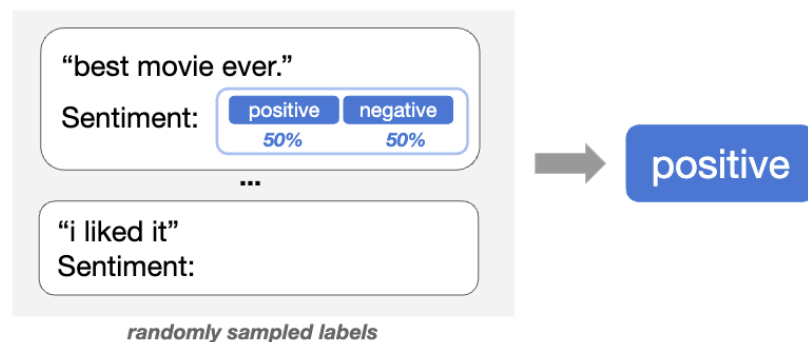
### ➤ 任务识别 (Task Recognition)

- 使用预训练阶段习得的知识解决任务
- 例如，预测评论的情感

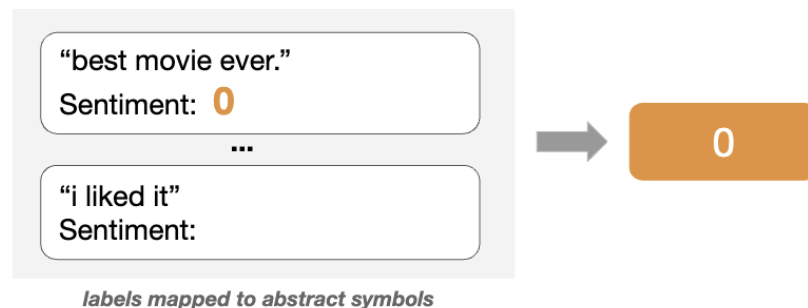
### ➤ 任务学习 (Task Learning)

- 从示例中学习如何解决新任务
- 例如，预测评论的情感并转换为 0/1

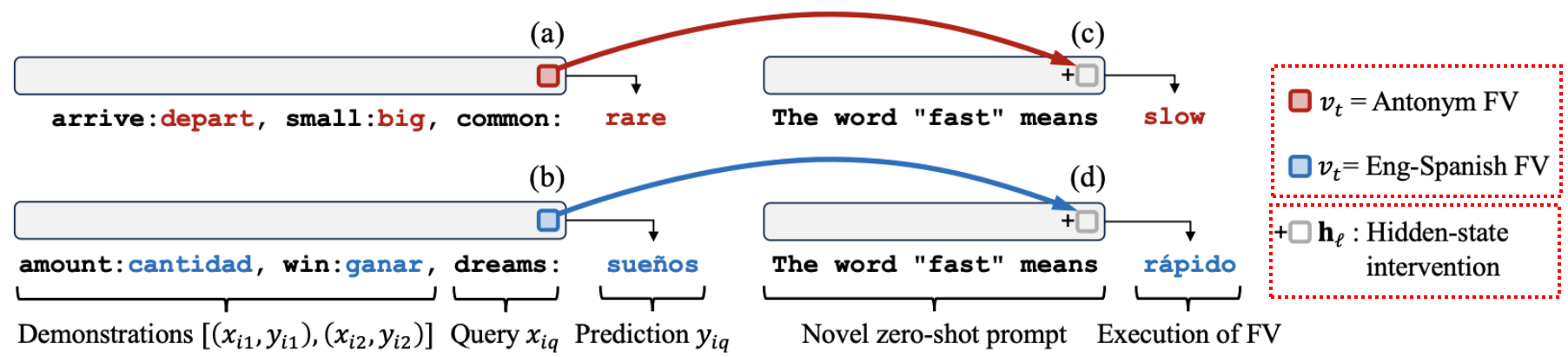
**Random**  
Task Recognition



**Abstract**  
Task Learning



- 任务识别
  - 模型在预训练阶段编码表征任务信息的隐向量 (学习的先验知识), 可用于在推理阶段解决任务
  - 以预测反义词和西班牙语为例



模拟预训练阶段  
学习的隐向量  
加入扰动向量  
稳健性良好

# 上下文学习的工作机制

## ➤ 任务学习

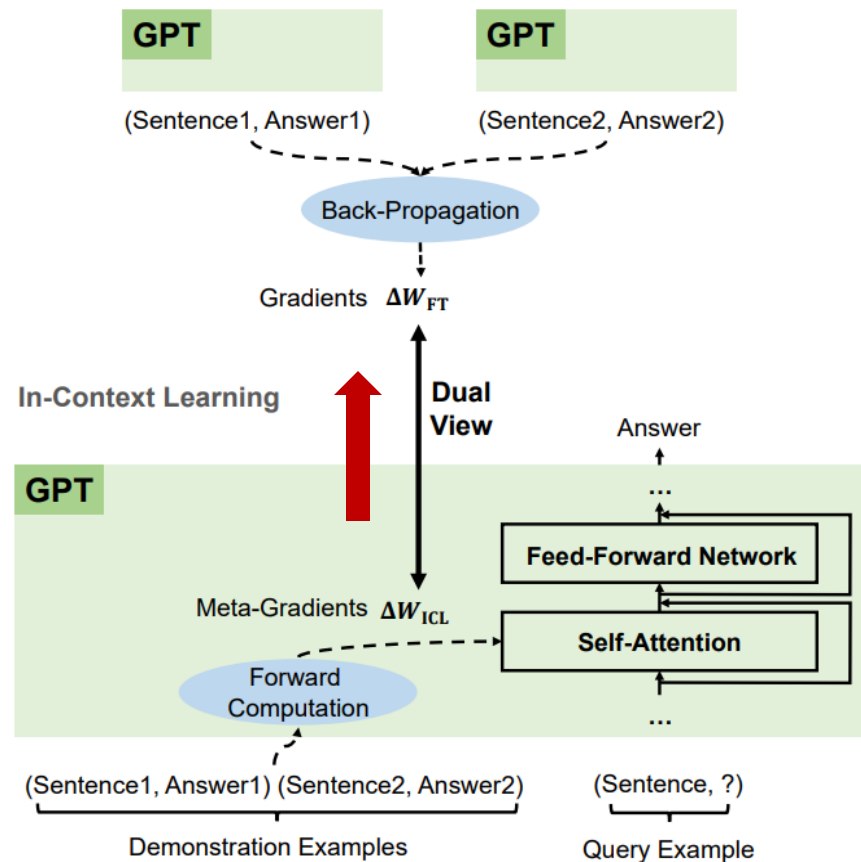
### ➤ 上下文学习是微调的对偶形式

➤ 注意力机制 → 梯度下降

➤ LinearAttn ( $V, K, q$ )

$$\begin{aligned}\mathcal{F}(\mathbf{x}) &= (W_0 + \Delta W) \mathbf{x} \\ &= W_0 \mathbf{x} + \Delta W \mathbf{x} \\ &= W_0 \mathbf{x} + \sum_i (\mathbf{e}_i \otimes \mathbf{x}_i'^T) \mathbf{x} \\ &= W_0 \mathbf{x} + \sum_i \mathbf{e}_i (\mathbf{x}_i'^T \mathbf{x}) \\ &= W_0 \mathbf{x} + \text{LinearAttn}(E, X', \mathbf{x})\end{aligned}$$

Finetuning



# 上下文学习的工作机制

## ➤ 任务学习

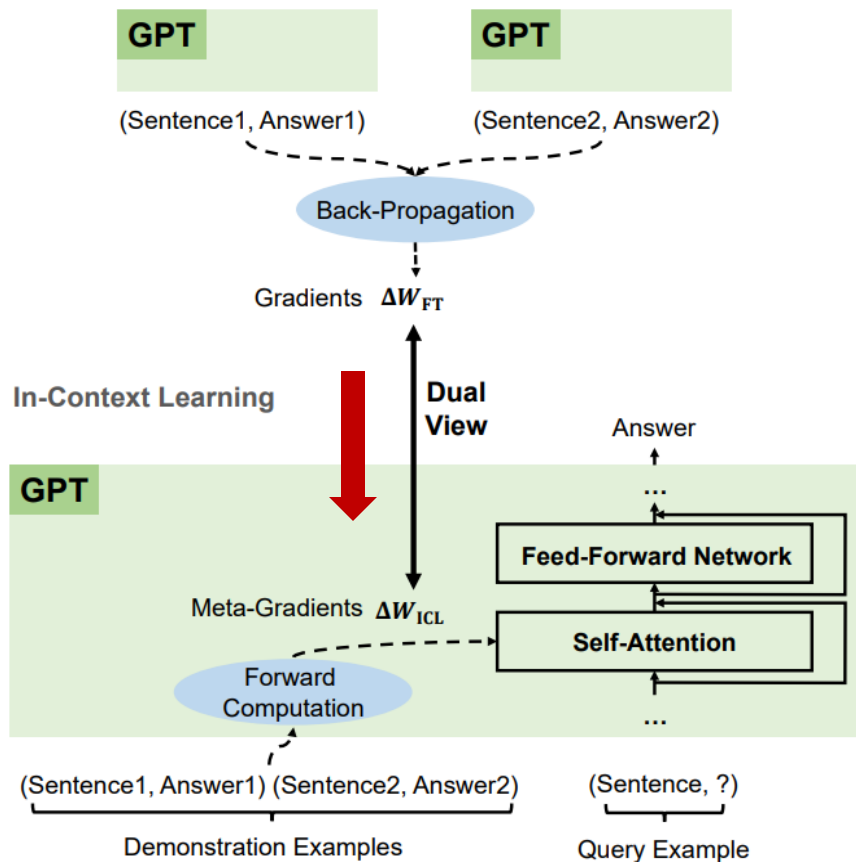
### ➤ 上下文学习是微调的对偶形式

#### ➤ 梯度下降 → 注意力机制

$$\begin{aligned}\mathcal{F}_{\text{ICL}}(\mathbf{q}) &= \text{Attn}(V, K, \mathbf{q}) \\ &\approx W_V[X'; X] (W_K[X'; X])^T \mathbf{q} \\ &= \boxed{W_V X (W_K X)^T \mathbf{q}} + W_V X' (W_K X')^T \mathbf{q} \\ &= \tilde{\mathcal{F}}_{\text{ICL}}(\mathbf{q}). \quad \text{W}_{\text{ZSL}}: \text{no demonstration}\end{aligned}$$

$$\begin{aligned}\tilde{\mathcal{F}}_{\text{ICL}}(\mathbf{q}) &= W_{\text{ZSL}} \mathbf{q} + W_V X' (W_K X')^T \mathbf{q} \\ &= W_{\text{ZSL}} \mathbf{q} + \text{LinearAttn}(W_V X', W_K X', \mathbf{q}) \\ &= W_{\text{ZSL}} \mathbf{q} + \sum_i W_V \mathbf{x}'_i \left( (W_K \mathbf{x}'_i)^T \mathbf{q} \right) \\ &= W_{\text{ZSL}} \mathbf{q} + \sum_i \left( W_V \mathbf{x}'_i \otimes (W_K \mathbf{x}'_i)^T \right) \mathbf{q} \\ &= W_{\text{ZSL}} \mathbf{q} + \Delta W_{\text{ICL}} \mathbf{q} \\ &= (W_{\text{ZSL}} + \Delta W_{\text{ICL}}) \mathbf{q}.\end{aligned}$$

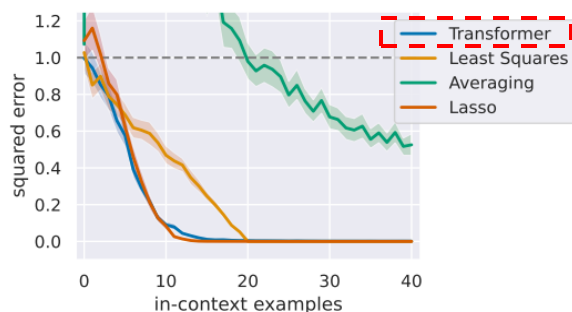
#### Finetuning



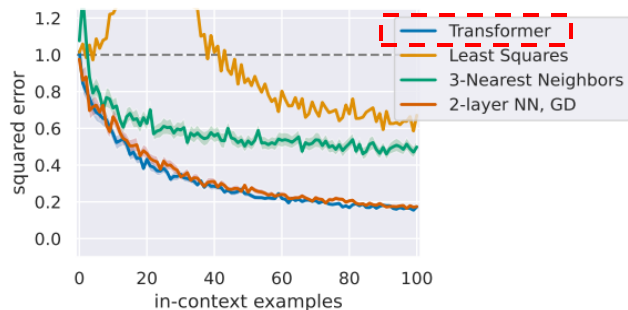
# 上下文学习的工作机制

## ➤ 任务学习

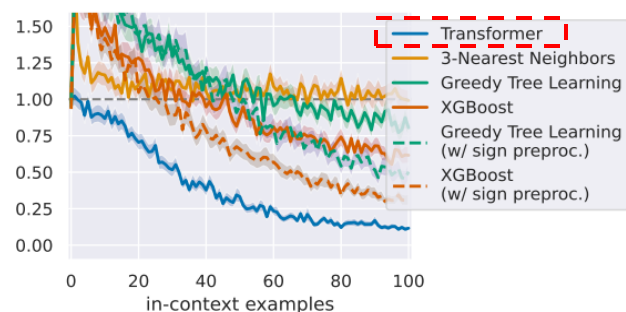
### ➤ 可以通过上下文学习学会复杂函数的输入输出映射关系



(a) Sparse linear functions



(c) 2-layer NN



(b) Decision trees



(d) 2-layer NN, eval on linear functions

- 在**预训练阶段**，大模型通过其内部参数能够学习隐式的函数关系
- 在**上下文学习阶段**，借助给定示例，大模型能够通过这个隐式模型模拟出诸如决策树等复杂学习算法

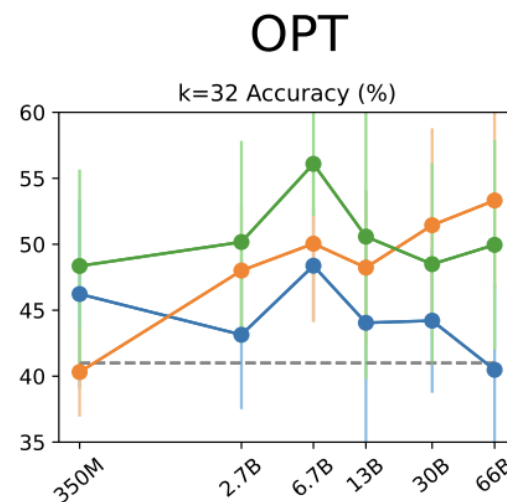
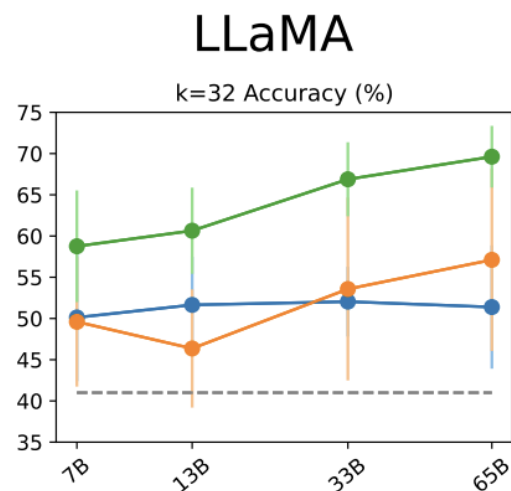
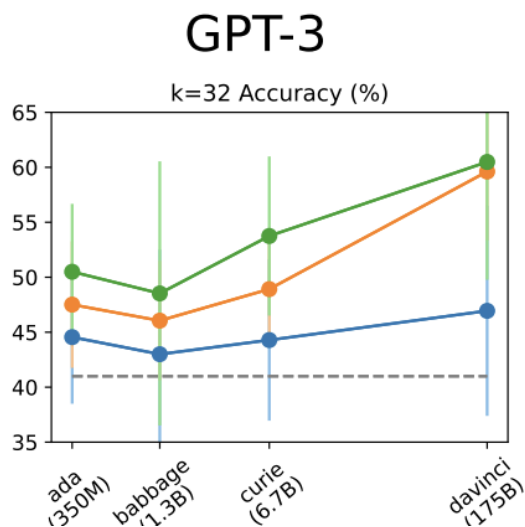
# 上下文学习的工作机制

## ➤ 任务识别和任务学习

- 大语言模型能够同时展现出任务识别和任务学习两种能力
- 这两种能力的强弱与模型的规模紧密相关

规模较小的模型就能  
具备任务识别能力

较大规模的模型具备  
更强的任务学习能力



● 任务识别    ● 任务学习

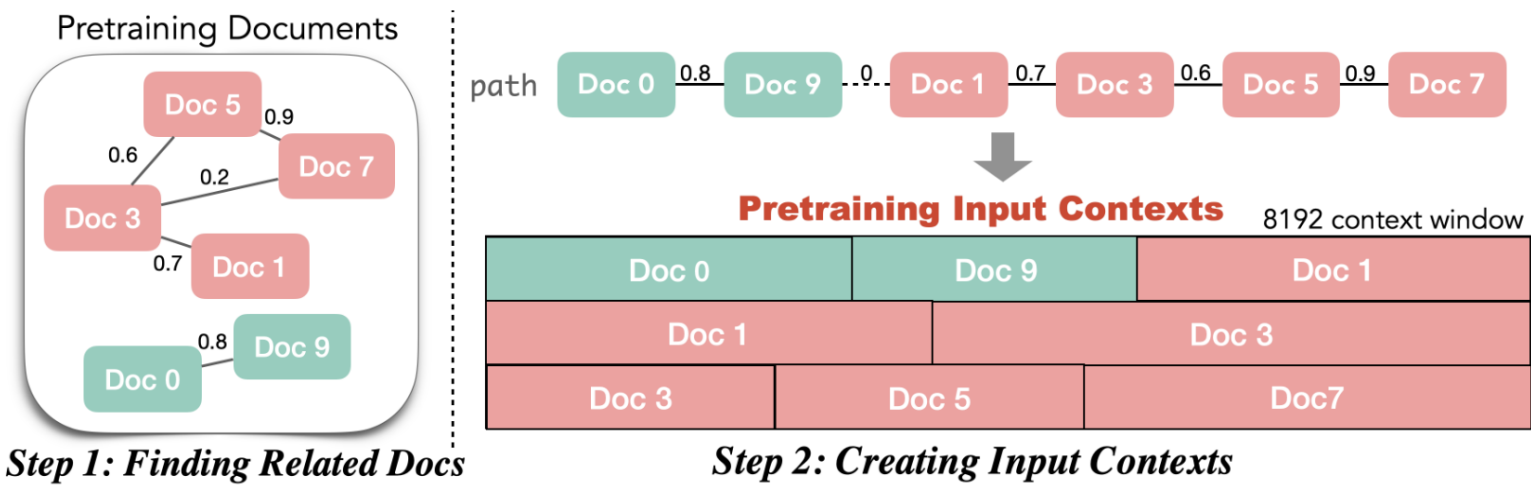
- 训练任务
  - 通过设计专门的训练任务进行继续预训练或微调也能获得上下文学习能力
  - MetaICL: 元训练任务训练模型根据具体任务示例和输入预测相应的输出

Method	HR→LR	Class →Class	non-Class →Class	QA →QA	non-QA →QA	non-NLI →NLI	non-Para →Para
All target tasks							
0-shot	34.8	34.2	34.2	40.2	40.2	25.5	34.2
PMI 0-shot	35.1	33.8	33.8	40.2	40.2	27.9	39.2
Channel 0-shot	36.5	37.3	37.3	38.7	38.7	33.9	39.5
In-context	38.2/35.3	37.4/33.9	37.4/33.9	40.1/38.7	40.1/38.7	34.0/28.3	33.7/33.1
PMI In-context	39.2/33.7	38.8/30.0	38.8/30.0	40.3/38.8	40.3/38.8	33.0/28.0	38.6/33.4
Channel In-context	43.1/38.5	46.3/40.3	46.3/40.3	40.8/38.1	40.8/38.1	39.9/34.8	45.4/40.9
Multi-task 0-shot	35.6	37.3	36.8	45.7	36.0	40.7	30.6
Channel Multi-task 0-shot	38.8	40.9	42.2	42.1	36.4	36.8	35.1
MetaICL	43.3/41.7	43.4/39.9	38.1/31.8	<b>46.0</b> /44.8	38.5/36.8	49.0/44.8	33.1/33.1
Channel MetaICL	<b>49.1</b> /46.8	<b>50.7</b> /48.0	<b>50.6</b> /48.1	44.9/43.5	<b>41.9</b> /40.5	<b>54.6</b> /51.9	<b>52.2</b> /50.3
Fine-tune	46.4/40.0	50.7/44.0	50.7/44.0	41.8/39.1	41.8/39.1	44.3/32.8	54.7/48.9
Fine-tune w/ meta-train	52.0/47.9	53.5/48.5	51.2/44.9	46.7/44.5	41.8/39.5	57.0/44.6	53.7/46.9

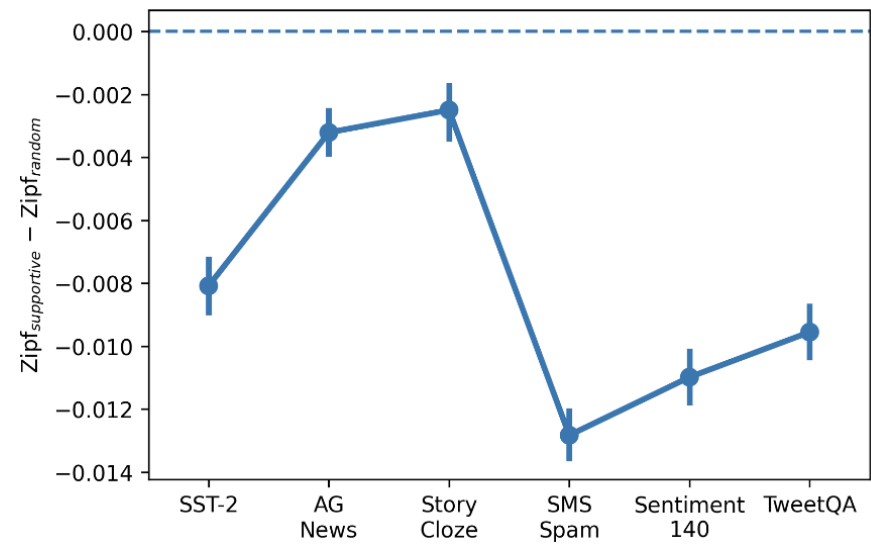
# 预训练对上下文学习的影响



- 预训练数据
- 预训练数据的多样性和长程依赖关系，以及高密度低频长尾词汇



根据训练数据相关性对文档排序依次训练



Zipfian系数显著低于随机训练数据





谢谢