# In-Context-Learning For Few-Shot Adaptation of LLMs

Sunita Sarawagi IIT Bombay sunita@iitb.ac.in

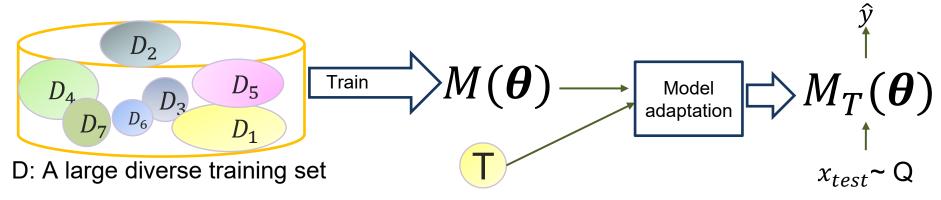
#### Few-shot Model Adaptation

Given a trained model  $M(\Theta)$  on a large diverse dataset D

User interested in a target task Q, gives a small labelled data

T: 
$$\{x_1 \ y_1 \ x_2 \ y_2 \ x_3 \ y_3 \ \dots \ x_k \ y_k\} \sim Q$$

Model adaptation = Find the best blend of  $M(\Theta)$  and T



D may contain a part  $D_i$  that is relevant to target task, but it is not explicitly identified.

#### Example: model adaptation

 $M(\boldsymbol{\theta})$ 

Target task

Conversation assistant



→ New users



Text-to-SQL system













→ Private pharma database



A large language model



Diverse topics and tasks

→ NER on physics documents



## The challenges of model adaptation

Accuracy: requires delicate balance.



- Robustness: Small T leads to high variance across T~Q
  - Cannot introduce additional instabilities, e.g. due to ordering
- Efficiency: Θ is large
  - Time of adaptation: Offline Vs Online (On-the-fly)
  - Test time: time for using after adaptation.

#### Model adaptation over the Ages

Thousands of paper under variants like domain adaptation, transfer-learning, few-shot learning...



Computer Science International Conference on Machine Learning

#### Deep Transfer Learning with Joint Adaptation Networks

Mingsheng Long Hanhua Zhu Jianmin Wang Michael I. Jordan

International Conference on Machine Learning · 21 May 2016

Chelsea Finn P. Abbeel S. Levine

#### Unsupervised Domain Adaptation by Backpropagation

Yaroslav Ganin V. Lempitsky Computer Science · <u>International Conference on Machine Learning</u>
26 September 2014

#### Model adaptation methods

- Fine-tuning and its variants
- Mixture of experts
- Task-vectors
- Matching-based methods

Each category explicitly planned for adaptation either when

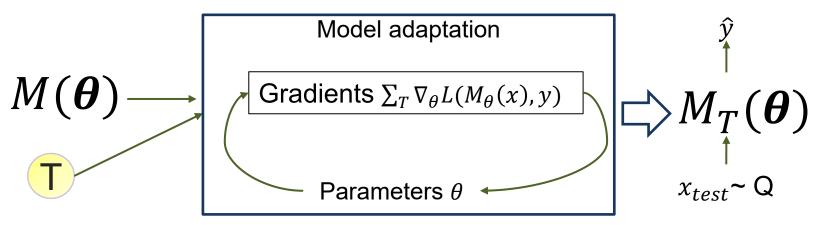
- Training  $M(\mathbf{\Theta})$ , or when
- Adapting with T

# Fine-tune parameters via gradient descent on T

Training of  $M(\theta)$ :

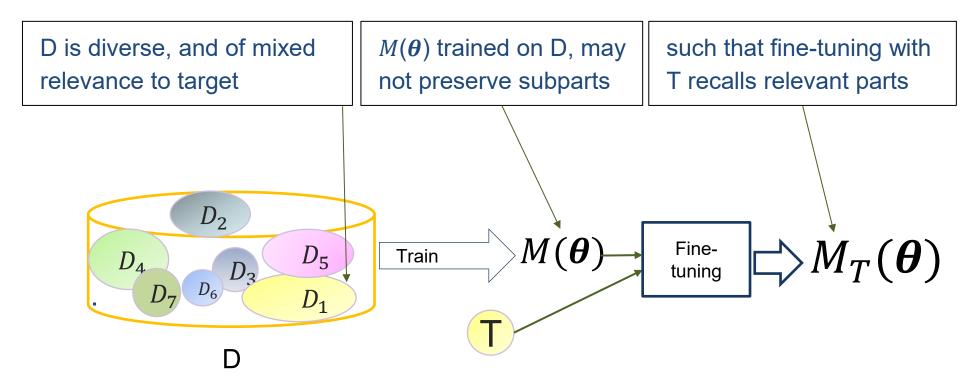
Normal loss on D

Adaptation with T: Update *θ* via gradient descent on T



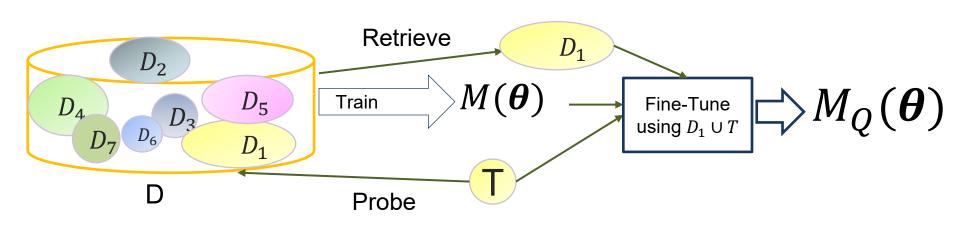
Explicit update of model parameters  $\Theta$  to minimize loss on target data T Generally provides good accuracy for modest sized T

# Does fine-tuning define the accuracy ceiling?

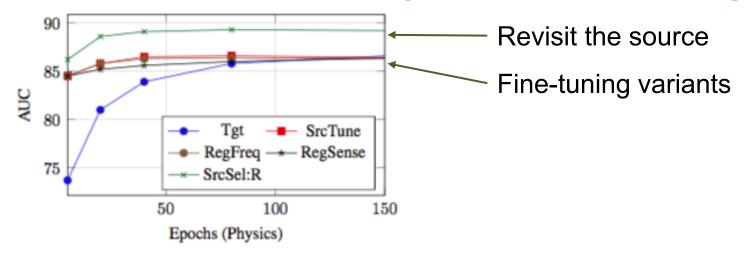


#### An accuracy ceiling: Revisit the source D

- Identify relevant  $D_1$  explicitly from D, Retrieval problem.
- Fine-tune with  $D_1 \cup T$ .



# Revisit the source for adapting word embeddings



Revisiting source significantly better than best fine-tuning.

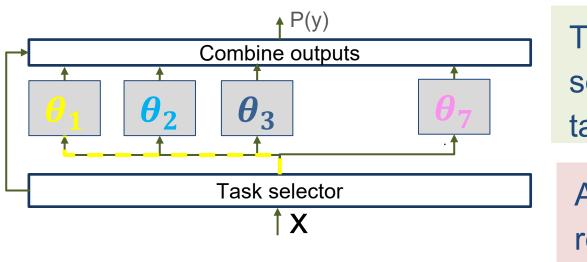
Not practical. Cannot revisit source

Useful intuition: important to maintain task identities for later specialization

Topic-Sensitive Attention on Generic Corpora Corrects Sense Bias in Pretrained Embeddings. Piratla, Sarawagi and Chakrabarti In *ACL*, 2019.

#### A practical approximation: mixture of experts

Model: mixture of experts with a task selector



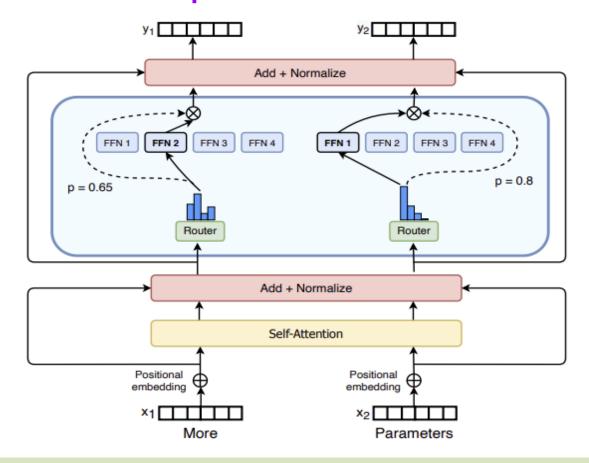
Training of  $M(\theta)$ : Task selector supervised by tasks in D

Adaptation with T: recognize the most relevant expert

Separate parameters for each task may be too much.

Multi-Source Domain Adaptation with Mixture of Experts. EMNLP 2018. Guo, Shah, Barzilay

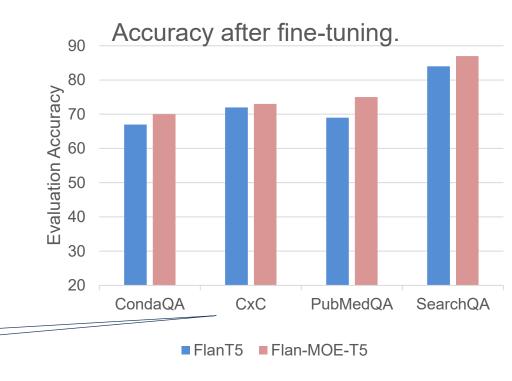
#### Modern Mixture of Experts: Switch Transformers



# Greater accuracy gains on fine-tuning MoEs

Base model: Flan T5

MoE variant: Flan-MoE-T5

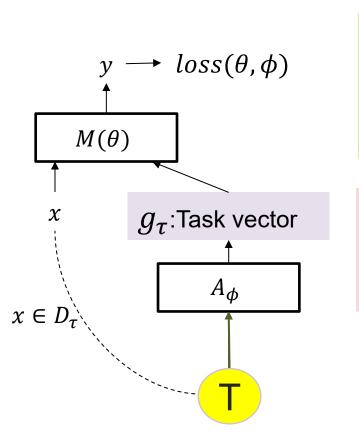


Various NLP tasks

Mixture-of-Experts Meets Instruction Tuning: A Winning Combination for Large Language Models. ICLR 2024. Shen, Hou, Zhou, Du, Longpre, Wei, Chung, Zoph, Fedus, Chen, Tu Vu, Yuexin Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vincent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell, Denny Zhou

#### Adaptation with Task vectors

Another way to maintain task identities



#### Training:

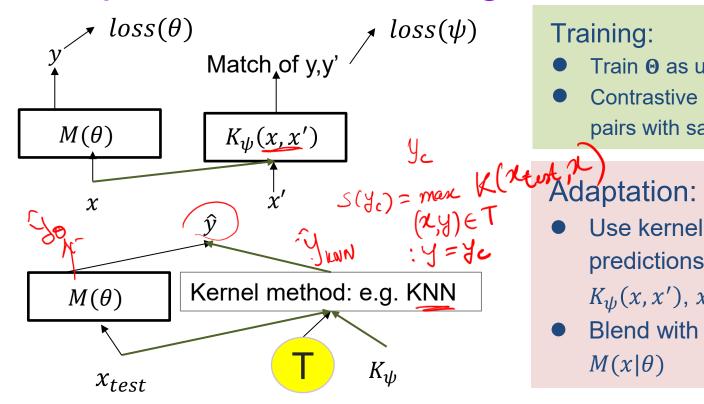
- Train  $\Theta$ ,  $\phi$  jointly
- Requires task boundaries

Adaptation:

Use T to extract  $g_t = A_{\phi}(T)$ 

 $\hat{y} = M(x_{test}, g_t)$ 

# Adaptation with Matching Networks



#### **Training:**

- Train O as usual
- Contrastive  $K_{tt}(x, x')$  so example pairs with same labels are close

- Use kernel method to get predictions from T using  $K_{\psi}(x, x'), x' \in T.$
- Blend with predictions from  $M(x|\theta)$
- Labeled Memory Networks for Online Model Adaptation. In AAAI, 2018. Shiv Shankar and Sunita Sarawagi
- Speech-enriched Memory for Inference-time Adaptation of ASR Models to Word Dictionaries Mittal, Sarawagi, Jyothi. EMNLP 2023.
- Conditional Tree Matching for Inference-Time Adaptation of Tree Prediction Models. Varma, Awasthi, Sarawagi ICML 2023

# An example matching-based adaptation for Text2SQL

#### Text to SQL

#### Schema

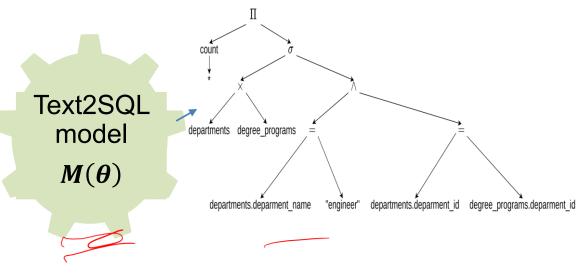
Department		
<u>ld</u>	Name	Description
14	Informatik	
11	Software system	
21	Engineer	

Degree programs		
<u>ld</u>	Name	Dept id
1	PhD	14
2	Masters	21
3	Bachelors	21

#### Question

How many degrees does the engineering department offer?

# Relational Algebra Tree representation of SQL



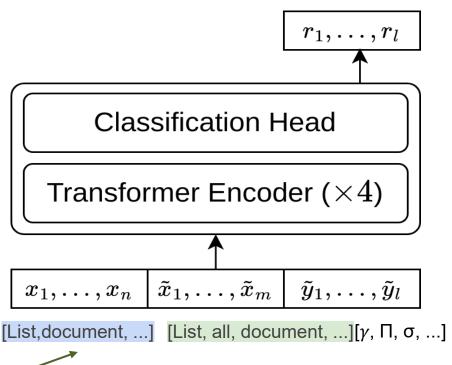
X

# Designing matching kernel: K((x,y),(x',y'))Two parts since y is structured.

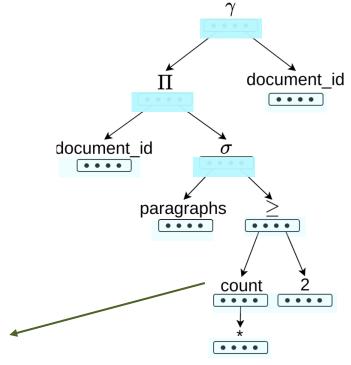
- Node-level reasoning of the relevance of specific subtrees of labelled examples to current question
  - Fine-grained sharing.
- Similarity based on alignment of respective trees -> structurally informed

#### Fine-grained relevance scoring

Relevance score each subtree



Wrong nodes marked irrelevant



User Question  $\boldsymbol{\mathcal{X}}$ 

Labeled Question  $oldsymbol{x}'$ , Tree  $oldsymbol{y}'$  from T

#### Relevance weighted candidate tree $y_r$ labeled tree $y'_{x}$ document id document id document id document id paragraphs paragraphs . . . . count count count . . . .

Alignment based relevance-weighted K((x, y), (x', y'))

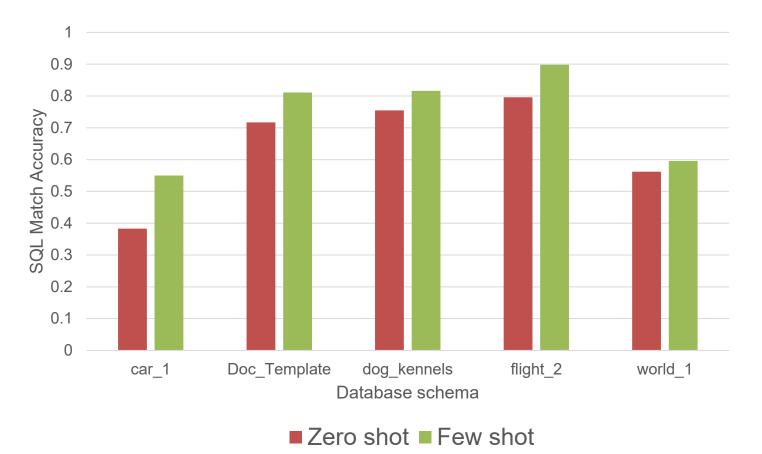
**X** Contextualized

 $\boldsymbol{x}'$  Contextualized

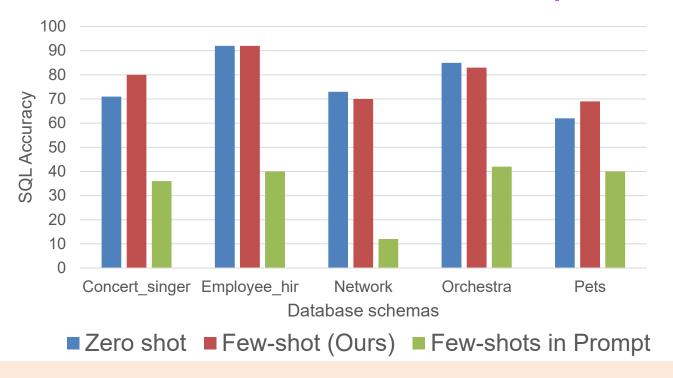
Conditional Tree Matching for Inference-Time Adaptation of Tree Prediction Models. Harshit Varma, Abhijeet Awasthi and Sunita Sarawagi In *ICML* 2023.

. . . .

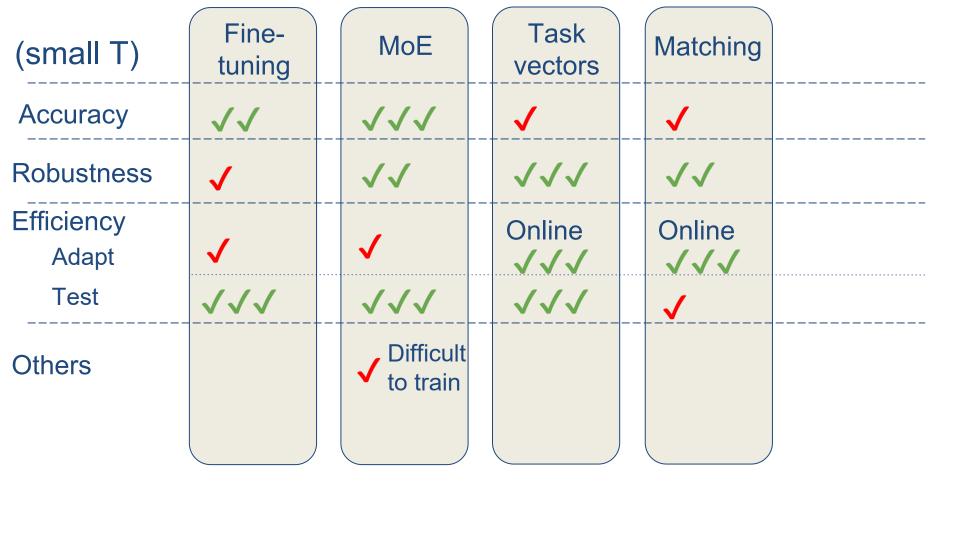
#### Accuracy after adaptation with few examples



#### Greater robustness to irrelevant examples



Reasoning about fine-grained relevance of few-shots is important for robustness



#### Model adaptation methods

Fine-tuning and its variants

Mixture of experts and task-vectors

Matching / memory augmented methods.

Each category explicitly planned for adaptation either when training  $M(\boldsymbol{\Theta})$  or adapting to T

## **Emergent Model Adaptation in LLMs**

 Solution to model adaptation just emerges in the form of in-context learning!

 Models trained with next-token prediction loss on huge natural document collections can be adapted just with few labelled examples in-context..

# Adaptation via In-Context Learning in LLMs

Predicted label  $\hat{y} \sim M_T(x_{test}|\theta)$ 

**Training: Default** 

Adaptation: Just a forward pass.

## Examples

LLM: 
$$M(\theta)$$

2!3 = 6, 3!3 = 9, 5!4 = 20, 6!7 =

63

LLM:  $M(\theta)$ 

2!3 = 6, 3!3 = 9, 5!4 = 20, 9!7 =

## Why does ICL work?

 Hypothesis 1: Transformers implement gradient descent algorithm over IC examples

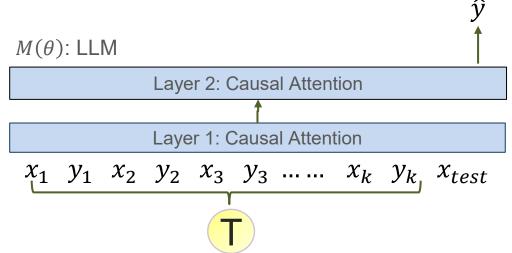
 Hypothesis 2: IC examples recognize tasks in pre-training e.g. via task vectors

 Hypothesis 3: Self-attention implements matching-based adaptation via induction heads We observe a parallel between these hypothesis and age-old methods of model adaptation.

# Hypothesis 1: Transformers implement gradient descent over label loss during ICL

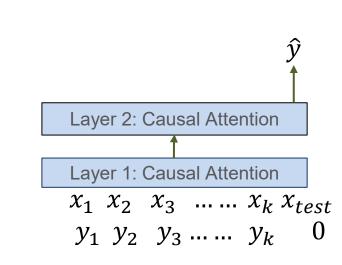
But how can GD be implemented in one forward pass over T?

- No obvious loss on IC examples, leave alone gradients on loss?
- $\circ$  No  $\theta$  parameters are updated!



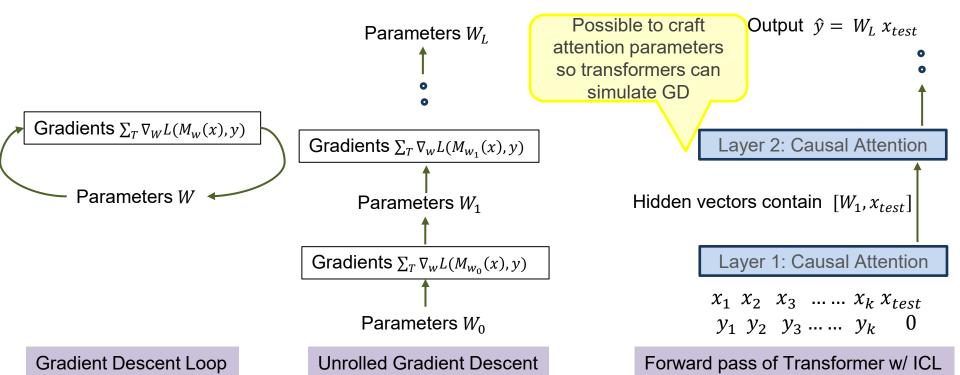
#### Assume

- Input stacked as [xi,yi]
- Task is linear regression
  - $\circ$   $x_i$  is a real vector
  - $\circ y_i = w_{\tau} x_i$
- Loss is square loss
  - $\circ \min_{w} \sum_{(x_i, y_i) \in T} (y_i wx_i)^2$
- Transformer attention is linear: no softmax.



What can transformers learn in-context? a case study of simple function classes. NeurIPS 2022. Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant.

#### Each layer implements a gradient step during forward pass of ICL



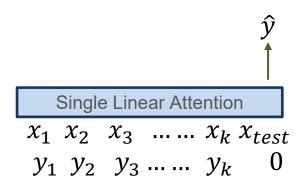
- Transformers learn in-context by gradient descent. Johannes Von Oswald et al. ICML 2023
- What learning algorithm is in-context learning? investigations with linear models. Ekin Akyurek et al. ICLR 2023

# But can we learn these parameters that implement gradient descent using standard next-token loss?

#### Proved for special cases

#### Special case I: One-layer linear Transformer

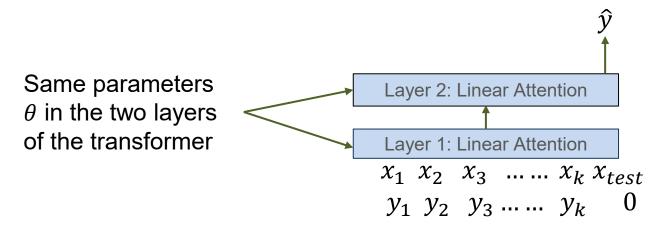
 Optima of training loss leads to transformer parameters that implement single pre-conditioned gradient descent step.



- Trained transformers learn linear models in-context. JMLR 2024. Ruigi Zhang, Spencer Frei, Peter L. Bartlett
- Transformers learn to implement preconditioned gradient descent for in-context learning. NeurIPS 2023.
   Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, Suvrit Sra

# Proved for special cases

Special case II: Linear Multi-layer looped Transformer

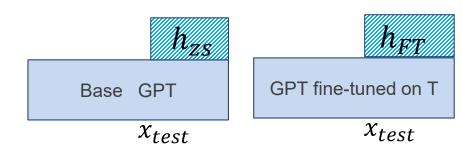


But these results are for numerical regression tasks, and that too on special transformers trained for this task. Do not generalize to text and NLP

Can looped transformers learn to implement multi-step gradient descent for in-context learning? ICML 2024. Khashayar Gatmiry, Nikunj Saunshi, Sashank J. Reddi, Stefanie Jegelka, and Sanjiv Kumar

#### What about text and pre-trained LLMs?

Dai et al (2023) support GD hypothesis for ICL for text and pre-trained LLMs



$$h_{FT} - h_{ZS} \sim h_{ICL} - h_{ZS}$$
 Base GPT with ICL 
$$x_1 \quad y_1 \quad x_2 \quad y_2 \quad x_3 \quad y_3 \quad \dots \quad x_k \quad y_k \quad x_{test}$$

Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. ACL 2023. Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei.

# ICL performs Gradient descent hypothesis is refuted..

Even an untrained transformer has high similarity between corresponding vectors while showing no capability for ICL

In-context learning and gradient descent revisited. NAACL 2024. G Deutch, N Magar, T Natan, G Dar.

Others: Li et al ACL 2024, Shen et al ICML 2024

# Hypothesis 2:

ICL identifies task from the pre-training set

#### ICL does task selection

LLMs provide the correct answer by recognizing this as a multiplication problem. They do not really learn to do multiplication from these examples.

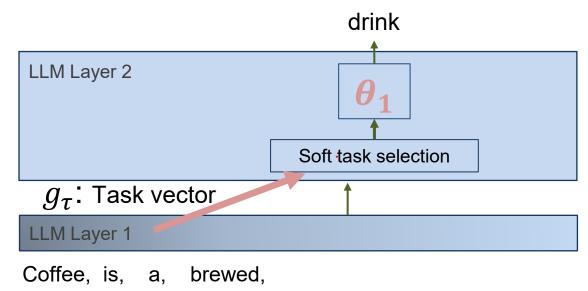
# How does task selection capability emerge in LLMs after pre-training with next token loss?

#### ICL does Task Selection

Pretraining documents: mix of related topics



Prefix of sentence recognizes topic as latent vectors  $g_{\tau} \rightarrow$  soft select task  $\rightarrow$  recall words from related documents seen earlier.



- An explanation of in-context learning as implicit Bayesian inference. ICLR 2022. Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma
- The learnability of in-context learning. NeurIPS 2023. Noam Wies, Yoav Levine, and Amnon Shashua

#### ICL goes beyond task selection

#### Large LLMs do learn novel input-output mappings

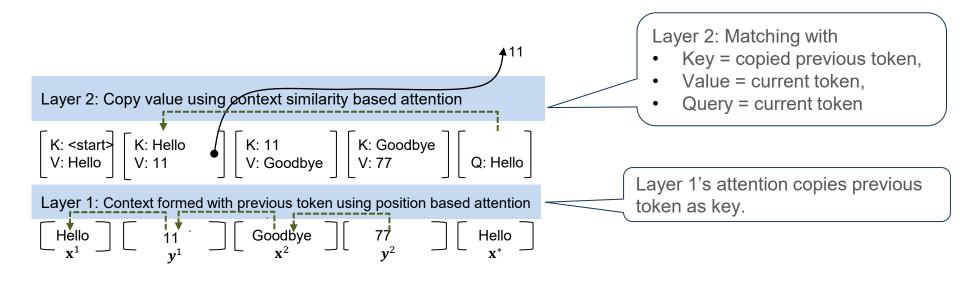
- Replacing positive/negative with foo/bar or flipping labels works for sentiment classification
- Unseen key-value associative mappings are learned
- Totally synthetic formal Markovian languages are learned

- What do language models learn in context? the structured task hypothesis. ACL 2024. Jiaoda Li, Yifan Hou, Mrinmaya Sachan, and Ryan Cotterell.
- Why larger language models do in-context learning differently? ICLR 2024. Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang

# Hypothesis 3: Self-attention implements matching-based adaptation

#### ICL= Matching with Induction circuits

Induction circuits created over two layers



In-context learning and induction heads. Olsson et al 2022

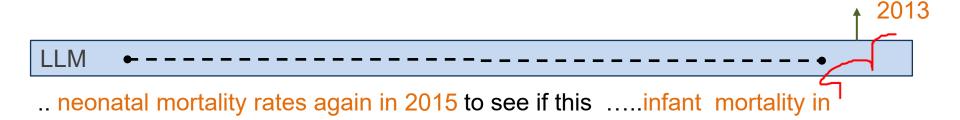
The mechanistic basis of data dependence and abrupt learning in an in-context classification task. ICLR 2024. Gautam Reddy

But why do transformer parameters orient to create in-context induction heads during next-token loss pre-training?

## Pre-training data contains repeated phrases

Corpus contains co-occurring phrases in similar templates.

care. We will **estimate** infant and neonatal mortality rates again in **2015** to see if this trend continues and, if so, to assess how it can **be** reversed. Infant mortality in **2013** was

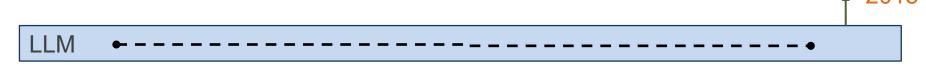


Parallel structures in pre-training data yield in-context learning. Chen et al. In ACL, 2024. Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, He He

#### Pre-training data contains repeated phrases

Destroy repetitions from context during training causes ICL to drop by 50% as against 2% random drops

care. We will **estimate** infant and neonatal mortality rates again in 2015 to see if this trend continues and, if so, to assess how it can be reversed. Infant mortality in 2013 was



.. strawberries ran December crazy to see if this .....infant mortality in

Random tokens

Parallel structures in pre-training data yield in-context learning. Chen et al. In ACL, 2024. Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, He He