



## State-of-the-art Seminar

# Understanding In-Context Learning

From the Architecture Perspective

Saqib Sarwar

August 16, 2025





# Outline

## Table of Contents

- In-Context Learning (ICL) - Motivation Example
- Pre-ICL Paradigm
- In-Context Learning
- Hypothesis 1: ICL as a Meta Optimizer
- Hypothesis 2: ICL as Bayesian Inference
- Hypothesis 3: ICL and Induction Heads
- ICL Across Architectures



# In-Context Learning

## Motivation Example

Answer in One Word.

"You have a right to perform your prescribed duties, but you are not entitled to the fruits of your actions." : Bhagavad Gita

"Love is patient, love is kind. It does not envy, it does not boast, it is not proud." : Corinthians

"The root of suffering is attachment." : Samyutta Nikaya

"And your Lord never forgets." : ?

Qur'an



# In-Context Learning

## Motivation Example

Answer in One Word.

"You have a right to perform your prescribed duties, but you are not entitled to the fruits of your actions." : Bhagavad Gita

"Love is patient, love is kind. It does not envy, it does not boast, it is not proud." : Corinthians

"The root of suffering is attachment." : Samyutta Nikaya

"And your Lord never forgets." : ?

Qur'an

"Love is patient, love is kind. It does not envy, it does not boast, it is not proud." : Christianity

"The root of suffering is attachment." : Buddhism

"And your Lord never forgets." : Islam

"You have a right to perform your prescribed duties, but you are not entitled to the fruits of your actions." : ?

Hinduism





# In-Context Learning

## Motivation Example

Answer in One Word.

"You have a right to perform your prescribed duties, but you are not entitled to the fruits of your actions." : Bhagavad Gita

"Love is patient, love is kind. It does not envy, it does not boast, it is not proud." : Corinthians

"The root of suffering is attachment." : Samyutta Nikaya

"And your Lord never forgets." : ?

Qur'an

"Love is patient, love is kind. It does not envy, it does not boast, it is not proud." : Christianity

"The root of suffering is attachment." : Buddhism

"And your Lord never forgets." : Islam

"You have a right to perform your prescribed duties, but you are not entitled to the fruits of your actions." : ?

Hinduism

"You have a right to perform your prescribed duties, but you are not entitled to the fruits of your actions." : Detachment

"And your Lord never forgets." : Omniscient

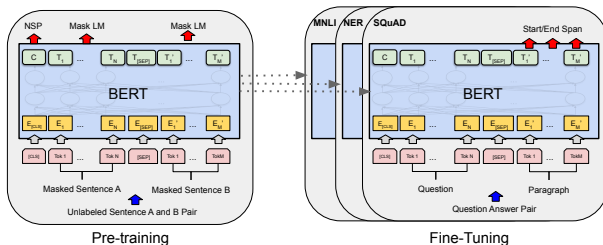
"Love is patient, love is kind. It does not envy, it does not boast, it is not proud." : Benevolent

"The root of suffering is attachment." : ?

Clinging



# The Pre-ICL Paradigm



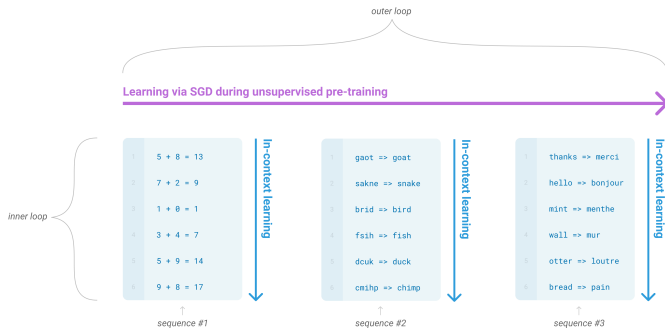
**Figure:** Pre-training and Supervised Fine-Tuning

1

<sup>1</sup>Devlin, J., et al. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186.



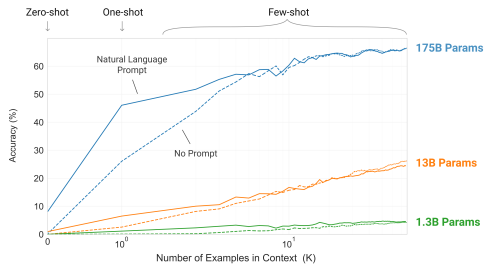
# In-Context Learning



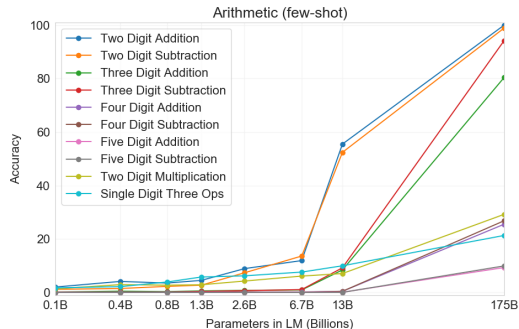
**Figure:** Language Model Meta Learning



# In-Context Learning: Scaling Effects



**Figure:** Word Scrambling and Manipulation Tasks



**Figure:** Arithmetic Tasks

3

<sup>3</sup>Brown, T., et al. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877–1901.

# Hypothesis 1: ICL as a Meta Optimizer

*Implicit Fine Tuning*





# Meta-ICL

	Meta-training	Inference
Task	$C$ meta-training tasks	An unseen <i>target</i> task
Data given	Training examples $\mathcal{T}_i = \{(x_j^i, y_j^i)\}_{j=1}^{N_i}, \forall i \in [1, C]$ ( $N_i \gg k$ )	Training examples $(x_1, y_1), \dots, (x_k, y_k)$ , Test input $x$
Objective	For each iteration, 1. Sample task $i \in [1, C]$ 2. Sample $k+1$ examples from $\mathcal{T}_i$ : $(x_1, y_1), \dots, (x_{k+1}, y_{k+1})$ 3. Maximize $P(y_{k+1} x_1, y_1, \dots, x_k, y_k, x_{k+1})$	$\operatorname{argmax}_{c \in C} P(c x_1, y_1, \dots, x_k, y_k, x)$

**Figure:** Meta-ICL Task

<sup>4</sup>Min, Sewon, et al. "MetalCL: Learning to Learn In Context." *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022, pp. 2791–2809.



## Meta-ICL

Meta-training		Inference
Task	$C$ meta-training tasks	An unseen <i>target</i> task
Data given	Training examples $\mathcal{T}_i = \{(x_j^i, y_j^i)\}_{j=1}^{N_i}, \forall i \in [1, C] \quad (N_i \gg k)$	Training examples $(x_1, y_1), \dots, (x_k, y_k)$ , Test input $x$
Objective	For each iteration, 1. Sample task $i \in [1, C]$ 2. Sample $k+1$ examples from $\mathcal{T}_i: (x_1, y_1), \dots, (x_{k+1}, y_{k+1})$ 3. Maximize $P(y_{k+1} x_1, y_1, \dots, x_k, y_k, x_{k+1})$	$\operatorname{argmax}_{c \in C} P(c x_1, y_1, \dots, x_k, y_k, x)$

**Figure: Meta-ICL Task**

Meta-train			Target	
Setting	# tasks	# examples	Setting	# tasks
HR	61	819,200	LR	26
Classification	43	384,022	Classification	20
Non-Classification	37	368,768		
QA	37	486,143	QA	22
Non-QA	33	521,342		
Non-NLI	55	463,579	NLI	8
Non-Paraphrase	59	496,106	Paraphrase	4

**Figure: Meta-ICL Experiments**

<sup>4</sup>Min, Sewon, et al. "MetalCL: Learning to Learn In Context." *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2022*, pp. 2791–2809.



# ICL as Meta Optimizer

- **Algorithmic Equivalence:** Transformers can simulate linear learners (GD, ridge, least-squares), transitioning to Bayesian estimators with depth/width. <sup>5</sup>

---

<sup>5</sup>Akyürek, Ege, et al. "What Learning Algorithm Is In-Context Learning? Investigations with Linear Models." *The Eleventh International Conference on Learning Representations*, 2023.

<sup>6</sup>Dai, Damai, et al. "Why Can GPT Learn In-Context? Language Models Implicitly Perform Gradient Descent as Meta-Optimizers." *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.

<sup>7</sup>Garg, Shivam, et al. "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes." *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*.





# ICL as Meta Optimizer

- **Algorithmic Equivalence:** Transformers can simulate linear learners (GD, ridge, least-squares), transitioning to Bayesian estimators with depth/width. <sup>5</sup>
- **Implicit Finetuning:** Transformer attention is dual to gradient descent. <sup>6</sup>  
⇒  $ICL \approx$  internal meta-gradient updates.

---

<sup>5</sup>Akyürek, Ege, et al. "What Learning Algorithm Is In-Context Learning? Investigations with Linear Models." *The Eleventh International Conference on Learning Representations*, 2023.

<sup>6</sup>Dai, Damai, et al. "Why Can GPT Learn In-Context? Language Models Implicitly Perform Gradient Descent as Meta-Optimizers." *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.

<sup>7</sup>Garg, Shivam, et al. "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes." *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*.



# ICL as Meta Optimizer

- **Algorithmic Equivalence:** Transformers can simulate linear learners (GD, ridge, least-squares), transitioning to Bayesian estimators with depth/width. <sup>5</sup>
- **Implicit Finetuning:** Transformer attention is dual to gradient descent. <sup>6</sup>  
⇒  $ICL \approx$  internal meta-gradient updates.
- **Function Class Generalization:** Transformers trained from scratch can in-context learn full function classes (linear, sparse linear, 2-layer MLPs, decision trees), even under distribution shift. <sup>7</sup>

---

<sup>5</sup>Akyürek, Ege, et al. "What Learning Algorithm Is In-Context Learning? Investigations with Linear Models." *The Eleventh International Conference on Learning Representations*, 2023.

<sup>6</sup>Dai, Damai, et al. "Why Can GPT Learn In-Context? Language Models Implicitly Perform Gradient Descent as Meta-Optimizers." *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.

<sup>7</sup>Garg, Shivam, et al. "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes." *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*.



# ICL as Meta Optimizer

- **Algorithmic Equivalence:** Transformers can simulate linear learners (GD, ridge, least-squares), transitioning to Bayesian estimators with depth/width. <sup>5</sup>
- **Implicit Finetuning:** Transformer attention is dual to gradient descent. <sup>6</sup>  
⇒  $ICL \approx$  internal meta-gradient updates.
- **Function Class Generalization:** Transformers trained from scratch can in-context learn full function classes (linear, sparse linear, 2-layer MLPs, decision trees), even under distribution shift. <sup>7</sup>

---

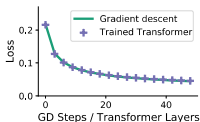
<sup>5</sup>Akyürek, Ege, et al. "What Learning Algorithm Is In-Context Learning? Investigations with Linear Models." *The Eleventh International Conference on Learning Representations*, 2023.

<sup>6</sup>Dai, Damai, et al. "Why Can GPT Learn In-Context? Language Models Implicitly Perform Gradient Descent as Meta-Optimizers." *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.

<sup>7</sup>Garg, Shivam, et al. "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes." *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*.



# ICL as Meta Optimizer



**Figure:** SGD-Transformer Equivalence

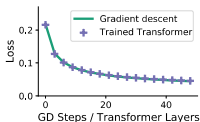
---

<sup>8</sup>von Oswald, Johannes, et al. "Transformers Learn In-Context by Gradient Descent." *Proceedings of the 40th International Conference on Machine Learning*, 2023.

<sup>9</sup>Wu, Shiguang, et al. "Why In-Context Learning Models are Good Few-Shot Learners?" *International Conference on Learning Representations*, 2025.



# ICL as Meta Optimizer



**Figure:** SGD-Transformer Equivalence

**Meta-Learning View:** ICL acts as data-dependent meta-learning, distinct from gradient-/metric-/amortized meta-learners.  
⇒ Implicit algorithm is shaped by pretraining distribution.

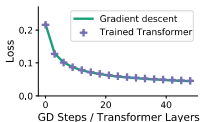
---

<sup>8</sup>von Oswald, Johannes, et al. "Transformers Learn In-Context by Gradient Descent." *Proceedings of the 40th International Conference on Machine Learning*, 2023.

<sup>9</sup>Wu, Shiguang, et al. "Why In-Context Learning Models are Good Few-Shot Learners?" *International Conference on Learning Representations*, 2025.

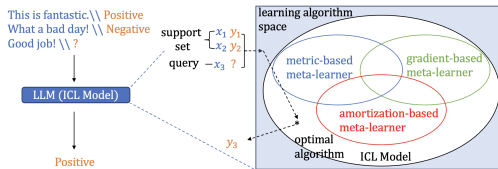


# ICL as Meta Optimizer



**Figure:** SGD-Transformer Equivalence

**Meta-Learning View:** ICL acts as data-dependent meta-learning, distinct from gradient-/metric-/amortized meta-learners.  
⇒ Implicit algorithm is shaped by pretraining distribution.



**Figure:** ICL as a Meta-Optimizer

8 9

<sup>8</sup>von Oswald, Johannes, et al. "Transformers Learn In-Context by Gradient Descent." *Proceedings of the 40th International Conference on Machine Learning*, 2023.

<sup>9</sup>Wu, Shiguang, et al. "Why In-Context Learning Models are Good Few-Shot Learners?" *International Conference on Learning Representations*, 2025.



## ICL as SGD?

- *Lingfeng et al* questioned the  $ICL \approx SGD$  approach.

---

<sup>10</sup>Shen, Lingfeng. "Do Pretrained Transformers Learn In-Context by Gradient Descent?" *Proceedings of the 41st International Conference on Machine Learning*, 2024.



## ICL as SGD?

- *Lingfeng et al* questioned the  $ICL \approx SGD$  approach.
- They questioned both the **Emergent equivalence** and the **Constructive Equivalence**

---

<sup>10</sup>Shen, Lingfeng. "Do Pretrained Transformers Learn In-Context by Gradient Descent?" *Proceedings of the 41st International Conference on Machine Learning*, 2024.





## ICL as SGD?

- *Lingfeng et al* questioned the  $ICL \approx SGD$  approach.
- They questioned both the **Emergent equivalence** and the **Constructive Equivalence**
- 

$$\underbrace{M_{\Theta_0}(\sigma_A \circ x_t) - M_{\Theta_0}(\sigma_B \circ x_t)}_{\text{The order sensitivity of ICL}} = \underbrace{M_{\Theta_{\sigma_A}}(x_t) - M_{\Theta_{\sigma_B}}(x_t)}_{\text{The order sensitivity of algorithm A}}$$

---

<sup>10</sup>Shen, Lingfeng. "Do Pretrained Transformers Learn In-Context by Gradient Descent?" *Proceedings of the 41st International Conference on Machine Learning*, 2024.



## ICL as SGD?

- *Lingfeng et al* questioned the  $ICL \approx SGD$  approach.
- They questioned both the **Emergent equivalence** and the **Constructive Equivalence**
- 

$$\underbrace{M_{\Theta_0}(\sigma_A \circ x_t) - M_{\Theta_0}(\sigma_B \circ x_t)}_{\text{The order sensitivity of ICL}} = \underbrace{M_{\Theta_{\sigma_A}}(x_t) - M_{\Theta_{\sigma_B}}(x_t)}_{\text{The order sensitivity of algorithm A}}$$

- Across all metrics(Accuracy, Cosine similarity among others), ICL and GD show inconsistent behavior.

---

<sup>10</sup>Shen, Lingfeng. "Do Pretrained Transformers Learn In-Context by Gradient Descent?" *Proceedings of the 41st International Conference on Machine Learning*, 2024.



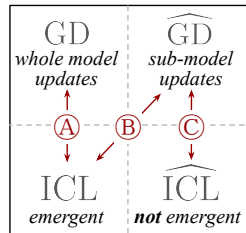
## ICL as SGD?

- *Lingfeng et al* questioned the  $ICL \approx SGD$  approach.
- They questioned both the **Emergent equivalence** and the **Constructive Equivalence**
- 

$$\underbrace{M_{\Theta_0}(\sigma_A \circ x_t) - M_{\Theta_0}(\sigma_B \circ x_t)}_{\text{The order sensitivity of ICL}} = \underbrace{M_{\Theta_{\sigma_A}}(x_t) - M_{\Theta_{\sigma_B}}(x_t)}_{\text{The order sensitivity of algorithm A}}$$

- Across all metrics (Accuracy, Cosine similarity among others), ICL and GD show inconsistent behavior.

10



**Figure:** ICL GD Equivalence

<sup>10</sup>Shen, Lingfeng. "Do Pretrained Transformers Learn In-Context by Gradient Descent?" *Proceedings of the 41st International Conference on Machine Learning*, 2024.

## Hypothesis 2: ICL as Bayesian Inference





## ICL and Bayesian Inference

$$p(\text{output} \mid \text{prompt}) = \int_{\text{concept}} p(\text{output} \mid \text{concept}, \text{prompt}) p(\text{concept} \mid \text{prompt}) d(\text{concept})$$

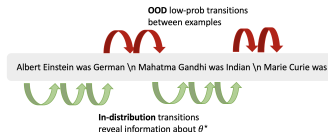
---

<sup>11</sup>Ahuja, Kabir, et al. "In-Context Learning through the Bayesian Prism." *The Twelfth International Conference on Learning Representations*, 2024.



# ICL and Bayesian Inference

$$p(\text{output} \mid \text{prompt}) = \int_{\text{concept}} p(\text{output} \mid \text{concept}, \text{prompt}) p(\text{concept} \mid \text{prompt}) d(\text{concept})$$



**Figure:** Signal and the OOD <sup>a</sup>

---

<sup>a</sup>Xie, Sang Michael, et al. "An Explanation of In-context Learning as Implicit Bayesian Inference." *The Tenth International Conference on Learning Representations, ICLR 2022*, 2022.

---

<sup>11</sup>Ahuja, Kabir, et al. "In-Context Learning through the Bayesian Prism." *The Twelfth International Conference on Learning Representations*, 2024.



# ICL and Bayesian Inference

$$p(\text{output} \mid \text{prompt}) = \int_{\text{concept}} p(\text{output} \mid \text{concept}, \text{prompt}) p(\text{concept} \mid \text{prompt}) d(\text{concept})$$

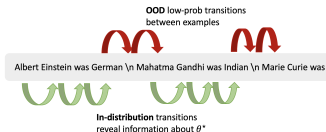
## Hierarchical Meta ICL Setup

$$c \sim \pi = (\pi_1, \dots, \pi_M),$$

$$f \sim p(f \mid c),$$

$$(x_i, y_i) : x_i \sim p(x), y_i = f(x_i), i = 1, \dots, N,$$

$x_{N+1}$ : query input,  $y_{N+1}$ : to predict.



**Figure:** Signal and the OOD <sup>a</sup>

## Bayes optimal inference:

$$\Pr(c \mid \{(x_i, y_i)\}_{i=1}^N) \propto \pi_c \prod_{i=1}^N \Pr(y_i \mid x_i, c),$$

$$\Pr(y_{N+1} \mid x_{N+1}, \text{context}) = \sum_{c=1}^M \Pr(c \mid \text{context}) \mathbb{E}_{f \sim p(\cdot \mid c)} [\Pr(y_{N+1} \mid x_{N+1}, f)].$$

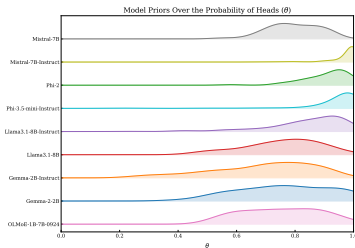
<sup>a</sup>Xie, Sang Michael, et al. "An Explanation of In-context Learning as Implicit Bayesian Inference." *The Tenth International Conference on Learning Representations, ICLR 2022, 2022*.

11

<sup>11</sup>Ahuja, Kabir, et al. "In-Context Learning through the Bayesian Prism." *The Twelfth International Conference on Learning Representations, 2024*.



# ICL and Bayesian Inference



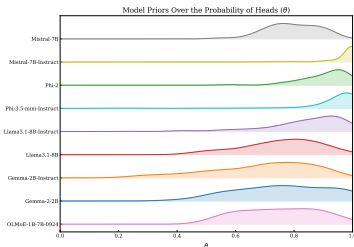
**Figure:** Coin Prior

<sup>12</sup>Gupta, Ritwik, et al. "Enough Coin Flips Can Make LLMs Act Bayesian." *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2025.

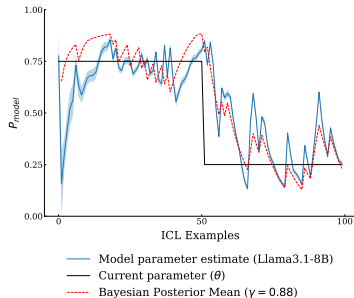




# ICL and Bayesian Inference



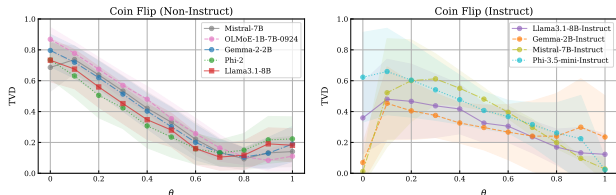
**Figure:** Coin Prior



**Figure:** Posterior



# ICL and Bayesian Inference



**Figure:** Biased Coin Instruct vs Non-Instruct

1. LLMs have biased *priors*.
2. Initial predictions diverge from ground truth due to these.
3. Explicit biasing (using prompts) improves only Instruct LLMs.
4. ICL helps remove the bias, similar to *Bayesian Updates*.

13

<sup>13</sup>Gupta, Ritwik, et al. "Enough Coin Flips Can Make LLMs Act Bayesian." *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2025.

## Hypothesis 3: ICL and Induction Heads





# Induction Heads

•

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.

---

<sup>14</sup>Olsson, Catherine, et al. "In-context Learning and Induction Heads." *CoRR*, *abs/2209.11895*, 2022.



# Induction Heads

- 

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.

- They **emerge** during training alongside ICL ability.

---

<sup>14</sup>Olsson, Catherine, et al. "In-context Learning and Induction Heads." *CoRR*, *abs/2209.11895*, 2022.



# Induction Heads

- 

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.

- They **emerge** during training alongside ICL ability.
- Core mechanism behind ICL.

---

<sup>14</sup>Olsson, Catherine, et al. "In-context Learning and Induction Heads." *CoRR*, *abs/2209.11895*, 2022.



# Induction Heads

- 

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.

- They **emerge** during training alongside ICL ability.
- Core mechanism behind ICL.
- Ablation shows they are **causal** for ICL in small models.

---

<sup>14</sup>Olsson, Catherine, et al. "In-context Learning and Induction Heads." *CoRR*, *abs/2209.11895*, 2022.



# Induction Heads

•

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.

- They **emerge** during training alongside ICL ability.
- Core mechanism behind ICL.
- Ablation shows they are **causal** for ICL in small models.
- Found across various model sizes and tasks.

---

<sup>14</sup>Olsson, Catherine, et al. "In-context Learning and Induction Heads." *CoRR*, *abs/2209.11895*, 2022.





# Induction Heads

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.



**Figure:** Induction Heads

- They **emerge** during training alongside ICL ability.
- Core mechanism behind ICL.
- Ablation shows they are **causal** for ICL in small models.
- Found across various model sizes and tasks.

<sup>14</sup>Olsson, Catherine, et al. "In-context Learning and Induction Heads." *CoRR*, abs/2209.11895, 2022.



# Induction Heads

$$[A^*][B^*] \dots [A] \rightarrow [B]$$

where  $A^* \approx A$  and  $B^* \approx B$  are similar in some space.

- They **emerge** during training alongside ICL ability.
- Core mechanism behind ICL.
- Ablation shows they are **causal** for ICL in small models.
- Found across various model sizes and tasks.



Figure: Induction Heads

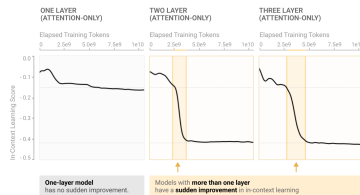
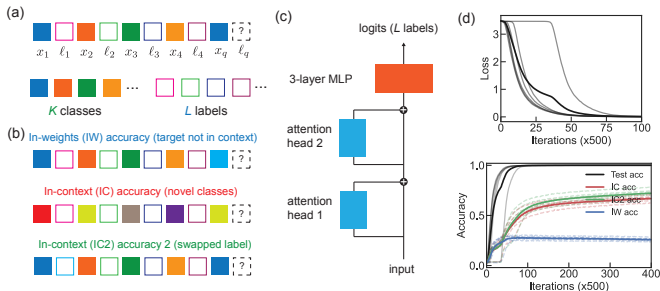


Figure: Abrupt Loss Transition



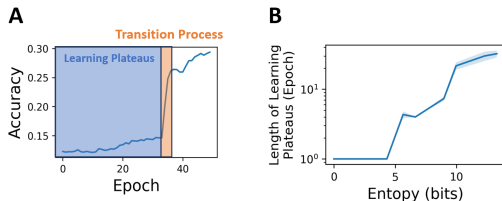
# Learning Plateau's and Abrupt Switching of ICL



**Figure:** Interpreting In-Context Classification Task



# Learning Plateau's and Abrupt Switching of ICL



**Figure:** Plateau in ICL

⇒ **Burstiness, Large Vocabulary Size, Skewed Classes** and **High Diversity** within Class promote ICL. <sup>a</sup>

⇒ Decompose Representation from parameters (W) and parameters + context (C).

⇒ Transferring Embeddings and Initial layers eliminates plateaus. <sup>b</sup>

<sup>a</sup>Reddy, Gautam. "The Mechanistic Basis of Data Dependence and Abrupt Learning in an In-Context Classification Task." *International Conference on Learning Representations*, 2024.

<sup>b</sup>Fu, Jingwen, et al. "Breaking through the Learning Plateaus of In-context Learning in Transformer." *Proceedings of the 41st International Conference on Machine Learning*, 2024.



## ICL Across Architectures

- *Lee et al* conduct an empirical study comparing ICL performance across diverse model architectures such as **CNNs, RNNs, Transformers** and **SSMs**.

---

<sup>16</sup>Lee, Ivan, et al. "Is Attention Required for ICL? Exploring the Relationship Between Model Architecture and In-Context Learning Ability." *The Twelfth International Conference on Learning Representations*, 2024.

<sup>17</sup>Park, Jongho, et al. "Can Mamba Learn How to Learn? A Comparative Study on In-Context Learning Tasks." *arXiv preprint arXiv:2402.04248*, 2024.



## ICL Across Architectures

- *Lee et al* conduct an empirical study comparing ICL performance across diverse model architectures such as **CNNs, RNNs, Transformers** and **SSMs**.
- All considered architectures achieve ICL, and some attention alternatives not only match but even surpass transformers.

---




<sup>16</sup>Lee, Ivan, et al. "Is Attention Required for ICL? Exploring the Relationship Between Model Architecture and In-Context Learning Ability." *The Twelfth International Conference on Learning Representations*, 2024.

<sup>17</sup>Park, Jongho, et al. "Can Mamba Learn How to Learn? A Comparative Study on In-Context Learning Tasks." *arXiv preprint arXiv:2402.04248*, 2024.



## ICL Across Architectures

- *Lee et al* conduct an empirical study comparing ICL performance across diverse model architectures such as **CNNs, RNNs, Transformers** and **SSTMs**.
- All considered architectures achieve ICL, and some attention alternatives not only match but even surpass transformers.

Task	Prompt	Target
Associative Recall	a, 1, b, 3, c, 2, b	3
Linear Regression	$x_1, y_1, x_2, y_2, x_3, y_3, x_4$	$y_4$ $\exists \mathbf{w}$ such that $\forall i, y_i = \mathbf{x}_i \cdot \mathbf{w}$
Multiclass Classification	$x_1, b, x_2, a, x_3, a, x_4$	b $x_1, x_4 \sim \mathcal{N}(y_b, I_d)$ $x_2, x_3 \sim \mathcal{N}(y_a, I_d)$
Image Classification		4 bursty training prompt
		2 non-bursty training prompt
		0 evaluation prompt
Language Modeling	Colorless green ideas sleep	furiously

16 17

<sup>16</sup>Lee, Ivan, et al. "Is Attention Required for ICL? Exploring the Relationship Between Model Architecture and In-Context Learning Ability." *The Twelfth International Conference on Learning Representations*, 2024.

<sup>17</sup>Park, Jongho, et al. "Can Mamba Learn How to Learn? A Comparative Study on In-Context Learning Tasks." *arXiv preprint arXiv:2402.04248*, 2024.



# ICL Across Architectures

- *Tong et al* test the ICL capabilities in **MLPs** and **MLP-Mixer** architectures.

---

<sup>18</sup>Tong, William L., and Cengiz Pehlevan. "MLPs Learn In-Context on Regression and Classification Tasks." *The Thirteenth International Conference on Learning Representations*, 2025.





# ICL Across Architectures

- *Tong et al* test the ICL capabilities in **MLPs** and **MLP-Mixer** architectures.
- Experimentally, they observe that given sufficient compute, MLP based architectures perform comparably with the transformers.

---

<sup>18</sup>Tong, William L., and Cengiz Pehlevan. "MLPs Learn In-Context on Regression and Classification Tasks." *The Thirteenth International Conference on Learning Representations*, 2025.



## ICL Across Architectures

- *Tong et al* test the ICL capabilities in **MLPs** and **MLP-Mixer** architectures.
- Experimentally, they observe that given sufficient compute, MLP based architectures perform comparably with the transformers.
- Also, for tasks which align with the architecture bias, MLP based architectures even outperform the transformers.

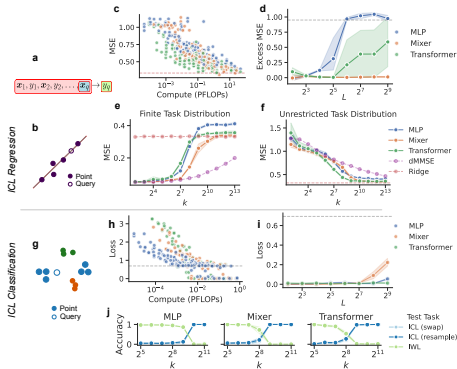
---

<sup>18</sup>Tong, William L., and Cengiz Pehlevan. "MLPs Learn In-Context on Regression and Classification Tasks." *The Thirteenth International Conference on Learning Representations*, 2025.



# ICL Across Architectures

- *Tong et al* test the ICL capabilities in **MLPs** and **MLP-Mixer** architectures.
- Experimentally, they observe that given sufficient compute, MLP based architectures perform comparably with the transformers.
- Also, for tasks which align with the architecture bias, MLP based architectures even outperform the transformers.



18

<sup>18</sup>Tong, William L., and Cengiz Pehlevan. "MLPs Learn In-Context on Regression and Classification Tasks." *The Thirteenth International Conference on Learning Representations*, 2025.

# Thank you!

This presentation is typeset using the  $\LaTeX$  beamer package with the  
UIC Presentation Template:

<https://www.overleaf.com/latex/templates/uic-presentation-template/dgjbtyvtgqcg>