

Is In-Context Learning a Type of Error-Driven Learning Mechanism?



Yale Linguistics NAACL 2025

> Evidence from the Inverse Frequency Effects in Structural Priming

X ET VERI

In-context Learning (ICL) in LLMs

- *In-context* learning vs. *In-weights* learning
- ICL as a processing mechanism of LLMs;
- Is ICL functionally performing some error-driven learning?

In-context Learning (ICL) in LLMs

Structural Priming

- *In-context* learning vs. *In-weights* learning
- ICL as a processing mechanism of LLMs;
- Is ICL functionally performing some error-driven learning?

- Structural priming effect;
- The Inverse
 - Frequency Effect (IFE);
- Implicit learning, a type of error-driven learning, accounts for the IFE;;

In-context Learning (ICL) in LLMs

Structural Priming

- *In-context* learning vs. *In-weights* learning
- ICL as a processing mechanism of LLMs;
- Is ICL functionally performing some error-driven learning?

- Structural priming effect;
- The Inverse Frequency Effect (IFE);
- Implicit learning, a type of error-driven learning, accounts for the IFE;;

• Analogy between priming and ICL?

Current Study

- Do LLMs show the IFE?
- IFE as a diagnostic of error-driven learning?

In-context Learning (ICL) in LLMs

Structural Priming

Current Study

Discussion & Implications

- *In-context* learning vs. *In-weights* learning
- ICL as a processing mechanism of LLMs;
- Is ICL functionally performing some error-driven learning?

- Structural priming effect;
- The Inverse Frequency Effect (IFE);
- Implicit learning, a type of error-driven learning, accounts for the IFE;;

- Analogy between priming and ICL?
- Do LLMs show the IFE?
- IFE as a diagnostic of error-driven learning?

- Larger models show stronger IFE;
- There is an *implicit* error term involved in ICL;
- Humans and LLMs share a similar processing mechanism;

[Brown et al. 2020]

In-context Learning: having a demonstration (i.e. several <example, answer> pairs) of a (implicitly defined) task increases the model performance.

Language Models are Few-Shot Learners				
Tom B. Brown* Benjamin Mann* Nick Ryder* Melanie Subbiah*				
Jared Kaplan [†]	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher H	esse Mark Chen	Eric Sigler	Mateusz Litwin	Scott Gray
Benjamin Chess Jack Clark Christopher Berner		Berner		
Sam McCar	ndlish Alec R	adford Ilya Su	ıtskever D	ario Amodei
OpenAI				

[Brown et al. 2020]

In-context Learning: having a demonstration (i.e. several <example, answer> pairs) of a (implicitly defined) task increases the model performance.

Language Models are Few-Shot Learners							
Tom B. Brown* Benjamin Mann* Nick Ryder* Melanie Subbiah*							
Jared Kaplan [†]	Prafulla	Dhariwal	Arvind Neela	kantan	Pranav Shya	im G	irish Sastry
Amanda Askell	Sandhini	Agarwal	Ariel Herbert-	Voss	Gretchen Krue;	ger To	om Henighan
Rewon Child	Aditya	Ramesh	Daniel M. Zie	egler	Jeffrey Wu	Cleme	ens Winter
Christopher I	lesse	Mark Chen	Eric Sig	er	Mateusz Litwi	n S	cott Gray
Benjamin Chess Jack Clark Christopher B			ier Berne	r			
Sam McC	andlish	Alec Ra	adford	Ilya Sut	skever	Dario A	modei
OpenAI							

In-context Learning: having a demonstration (i.e. several <example, answer> pairs) of a (implicitly defined) task increases the model performance.

Language Models are Few-Shot Learners				
Tom B. Brown* Benjamin Mann* Nick Ryder* Melanie Subbiah*				
Jared Kaplan [†]	Prafulla Dhariwal	Arvind Neelakan	tan Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Vos	s Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher H	lesse Mark Che	n Eric Sigler	Mateusz Litwin	Scott Gray
Benjamin Chess Jack Clarl		Jack Clark	Christopher	Berner
Sam McCa	ndlish Alec	Radford Ily	a Sutskever 1	Dario Amodei
OpenAI				

- No gradient updates;
- Rapid: from a few examples;
- One-shot / Few-shot learning;

Finetuning



Finetuning



In-context Learning as implicitly performing gradient descent? — *in principle*, yes...

Finetuning



In-context Learning as implicitly performing gradient descent? — *in principle*, yes...

• ICL performs implicit Bayesian inference;

Finetuning



In-context Learning as implicitly performing gradient descent? — *in principle*, yes...

- ICL performs implicit Bayesian inference;
- ICL functionally performs gradient descent;

Finetuning



In-context Learning as implicitly performing gradient descent? — *in principle*, yes...

- ICL performs implicit Bayesian inference;
- ICL functionally performs gradient descent;
- ICL as a meta-optimization process equivalent to implicit fine-tuning;

Finetuning



In-context Learning as implicitly performing gradient descent? — *in principle*, yes...

- ICL performs implicit Bayesian inference;
- ICL functionally performs gradient descent;
- ICL as a meta-optimization process equivalent to implicit fine-tuning;

Current Case Study: Is there an error-based learning process in the forward pass? testing with off-the-shelf LLMs and natural language!

Structural Priming: speakers tend to reuse the syntactic structures they have recently encountered during production or comprehension.

Structural Priming: speakers tend to reuse the syntactic structures they have recently encountered during production or comprehension.

Our focus: Double Object (DO) vs. Prepositional Dative (PD) for ditransitive predicates.

Structural Priming: speakers tend to reuse the syntactic structures they have recently encountered during production or comprehension.

Our focus: Double Object (DO) vs. Prepositional Dative (PD) for ditransitive predicates.

• DO: Alice sent Bob a letter.

E.g. [Bock 1986, Chang 2012]

Structural Priming: speakers tend to reuse the syntactic structures they have recently encountered during production or comprehension.

Our focus: Double Object (DO) vs. Prepositional Dative (PD) for ditransitive predicates.

- DO: Alice sent Bob a letter.
- PD: Alice sent a letter to Bob.

E.g. [Bock 1986, Chang 2012]

Inverse Frequency Effect: the less preferred (lower frequency) syntactic structure causes a stronger priming effect than the more preferred (higher frequency) structural alternative.

Inverse Frequency Effect: the less preferred (lower frequency) syntactic structure causes a stronger priming effect than the more preferred (higher frequency) structural alternative.



Verb Bias: buy is biased towards DO design towards PD

Inverse Frequency Effect: the less preferred (lower frequency) syntactic structure causes a stronger priming effect than the more preferred (higher frequency) structural alternative.



E.g. [Jaeger & Snider 2007]

Inverse Frequency Effect: the less preferred (lower frequency) syntactic structure causes a stronger priming effect than the more preferred (higher frequency) structural alternative.



ICL as Structural Priming?

- Instead of viewing "A terrible movie. → negative" as an input-output pair, we can view it as a single "sentence" with a particular *structure*:
- Structure = movie review + arrow + sentiment label
 A terrible movie. → *negative*

ICL as Structural Priming?

- Instead of viewing "A terrible movie. → negative" as an input-output pair, we can view it as a single "sentence" with a particular *structure*:
- Structure = movie review + arrow + sentiment label *A terrible movie.* \rightarrow *negative*
- Framed this way, in-context learning is structural priming!

ICL as Structural Priming?

- Instead of viewing "A terrible movie. → negative" as an input-output pair, we can view it as a single "sentence" with a particular *structure*:
- Structure = movie review + arrow + sentiment label *A terrible movie.* \rightarrow *negative*
- Framed this way, in-context learning is structural priming!

 \Rightarrow The IFE as a diagnostic of the error-driven learning mechanism in ICL!











Assumption from Priming Theories: only some error-driven learning mechanism could lead to the IFE.

- Fine-tuning Mode: (with weight update) IFE 🔽
- Concatenation Mode: (no weight update) IFE ?



Corpus

- 22 ditransitive verbs;
- 50 target sentences per verb;
- For each target sentence, pair it with a prime sentence with each prime verb;

22 x 50 (target sentences) x 21 (prime sentences) = 23100 <prime, target> pairs

Dataset adapted from Sinclair et al. 2022

Corpus

- 22 ditransitive verbs;
- 50 target sentences per verb;
- For each target sentence, pair it with a prime sentence with each prime verb;

22 x 50 (target sentences) x 21 (prime sentences) = 23100 <prime, target> pairs

Each <prime, target> pair \Rightarrow 4 structural combinations \Rightarrow **92400 trials.**

DO prime + DO target	DO prime + PD target
PD prime + DO target	PD prime + PD target

Corpus

- 22 ditransitive verbs;
- 50 target sentences per verb;
- For each target sentence, pair it with a prime sentence with each prime verb;

22 x 50 (target sentences) x 21 (prime sentences) = 23100 <prime, target> pairs

Each <prime, target> pair \Rightarrow 4 structural combinations \Rightarrow **92400 trials.**

DO prime + DO target	DO prime + PD target
PD prime + DO target	PD prime + PD target

Prime: A professor promised a student a letter. **Target**: The secretary drew the card for the boss.

Dataset adapted from Sinclair et al. 2022

Quantifying Verb Biases and the IFE

Verb Bias of verb V on structure X:

$$bias(V, ext{PD}) = rac{1}{|\mathcal{S}_V|} \sum_{t_{ ext{PD}} \in \mathcal{S}_V} rac{\mathcal{P}(t_{ ext{PD}})}{\mathcal{P}(t_{ ext{PD}}) + \mathcal{P}(t_{ ext{DO}})}$$

Quantifying Verb Biases and the IFE

Verb Bias of verb V on structure X:



Quantifying Verb Biases and the IFE

Verb Bias of verb V on structure X:

$$bias(V, ext{PD}) = rac{1}{|\mathcal{S}_V|} \sum_{t_{ ext{PD}} \in \mathcal{S}_V} rac{\mathcal{P}(t_{ ext{PD}})}{\mathcal{P}(t_{ ext{PD}}) + \mathcal{P}(t_{ ext{DO}})}$$

IFE: the priming effect for verb V in DO form on PD targets

$$PrimeBias(ext{PD}| ext{DO},V) = rac{1}{|T_{ ext{PD}}|\cdot|P_{ ext{DO}}^V|}\sum_{t_{ ext{PD}}\in T_{ ext{PD}}}\sum_{p_{ ext{DO}}^V\in P_{ ext{DO}}^V}rac{\mathcal{P}(t_{ ext{PD}}|p_{ ext{DO}}^V)}{\mathcal{P}(t_{ ext{DO}}|p_{ ext{DO}}^V)+\mathcal{P}(t_{ ext{PD}}|p_{ ext{DO}}^V)}$$

 $\bar{\mathcal{P}}(T_{\mathrm{PD}}|P_{\mathrm{PD}}^V)$ **Increasing PD Biases**













- **IFE:** double negative slopes;
- **Standard Priming:** PD-PD has higher intercept than DO-PD;

Fine-tuning Mode: fine-tuning the model with the prime sentence and use the updated model to run the target sentence — with weight update;

Fine-tuning Mode: fine-tuning the model with the prime sentence and use the updated model to run the target sentence — with weight update;



Fine-tuning Mode: fine-tuning the model with the prime sentence and use the updated model to run the target sentence — with weight update;



Even *GPT2-small* shows significant inverse frequency effects!

Concatenation Mode: concatenating the prime and target sentences as an ICL sequence and run the model — without weight update;

Concatenation Mode: concatenating the prime and target sentences as an ICL sequence and run the model — without weight update;



Concatenation Mode: concatenating the prime and target sentences as an ICL sequence and run the model — without weight update;



Concatenation Mode: concatenating the prime and target sentences as an ICL sequence and run the model — without weight update;



55

We used the IFE as a diagnostic on the error-driven nature of ICL as a processing mechanism of LLMs.

We used the IFE as a diagnostic on the error-driven nature of ICL as a processing mechanism of LLMs.

• Generalizing beyond standard notion of ICL, connecting priming with prompting

We used the IFE as a diagnostic on the error-driven nature of ICL as a processing mechanism of LLMs.

- Generalizing beyond standard notion of ICL, connecting priming with prompting
- Larger LLMs show more significant IFE 🤒 🧐

We used the IFE as a diagnostic on the error–driven nature of ICL as a processing mechanism of LLMs.

- Generalizing beyond standard notion of ICL, connecting priming with prompting
- Larger LLMs show more significant IFE 🤒 🧐
- At least in the case of priming, error-driven learning happens in ICL!***

Thanks for Listening!





In-context Learning (ICL) in LLMs

Structural Priming

Current Study

Discussion & Implications

- In-context learning vs. In-weights learning
- ICL as a processing mechanism of LLMs;
- Is ICL functionally performing some error-driven learning?

- Structural priming effect;
- The Inverse Frequency Effect (IFE);
- Two accounts: transient activation vs. implicit learning;

- Structural priming effect;
- The Inverse Frequency Effect (IFE);
- Implicit learning, a type of error-driven learning, accounts for the IFE;;

- Larger models show stronger IFE;
- There is an *implicit* gradient component involved in ICL;
- Humans and LLMs share a similar processing mechanism!

Selected Reference

- Sarah Bernolet and Robert J. Hartsuiker. 2010. Does verb bias modulate syntactic priming? *Cognition*, 114(3):455–461.
- J. Kathryn Bock. 1986. Syntactic persistence in language production. *Cognitive Psychology*, 18(3):355–387.
- Kathryn Bock and Carol A Miller. 1991. Broken agreement. *Cognitive Psychology*, 23(1):45–93.
- Laurel Brehm, Pyeong Whan Cho, Paul Smolensky, and Matthew A. Goldrick. 2022. PIPS: A Parallel Planning Model of Sentence Production. *Cognitive Science*, 46(2):e13079.
- Pyeong Whan Cho, Matthew Goldrick, Richard L. Lewis, and Paul Smolensky. 2018. Dynamic encoding of structural uncertainty in gradient symbols. In *Proceedings of the 8th Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2018)*, pages 19–28, Salt Lake City, Utah. Association for Computational Linguistics.
- Pyeong Whan Cho, Matthew Goldrick, and Paul Smolensky. 2020. Parallel parsing in a Gradient Symbolic Computation parser.
- John Hale and Paul Smolensky. 2006. Harmonic gram- mars and harmonic parsers for formal languages. *Smolensky and Legendre*, pages 393–416.
- T. Florian Jaeger and Neal Snider. 2007. Implicit Learning and Syntactic Persistence: Surprisal and Cumulativity. *University of Rochester Working Papers in the Language Sciences*, 3:26–44.
- Martin J. Pickering and Holly P. Branigan. 1998. The representation of verbs: Evidence from syntactic priming in language production. *Journal of Memory and Language*, 39(4):633–651.
- Eunkyung Yi, Jean-Pierre Koenig, and Douglas Roland. 2019. Semantic similarity to high-frequency verbs affects syntactic frame selection. *Cognitive Linguistics*, 30(3):601–628.