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

Zhenghao "Herbert" Zhou Yale University Hollymartin Lab @ University of Edinburgh December 4, 2024



Evidence from the Inverse Frequency Effects in Structural Priming













### 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);
- Two accounts: transient activation vs. implicit 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);
- Two accounts: transient activation vs. implicit learning;

• 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);
- Two accounts: transient activation vs. implicit learning

- 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* gradient component involved in ICL;
- Humans and LLMs share a similar processing mechanism;

Autoregressive Models

- GPT2 small: 117M
- GPT2 medium: 345M
- GPT2 large: 762M
- LLAMA2 7B: 7B
- LLAMA2 7B-chat: 7B
- LLAMA2 70B: 70B
- GPT3 davinci-002: 175B

*Predicting the next word!* 

Autoregressive Models

- GPT2 small: 117M
- GPT2 medium: 345M
- GPT2 large: 762M
- LLAMA2 7B: 7B
- LLAMA2 7B-chat: 7B
- LLAMA2 70B: 70B
- GPT3 davinci-002: 175B

#### output token *Predicting the next word!* Token probabilities (logits) Embeddings aardvark 0.19850038 Autoregressive Models 0.7089803 aarhus Decoder #12, Position #1 0.46333563 aaron Pick an output output vector token based on Х its probability GPT2 small: 117M (sample) The GPT2 medium: 345M -0.51006055 zvzzvva GPT2 large: 762M DECODER LLAMA2 7B: 7B ... LLAMA2 7B-chat: 7B DECODER LLAMA2 70B: 70B GPT3 davinci-002: 175B <S> 1 2 1024

https://jalammar.github.io/illustrated-gpt2/

#### output token *Predicting the next word!* Token probabilities (logits) Embeddings aardvark 0.19850038 **Autoregressive Models** 0.7089803 aarhus Decoder #12, Position #1 0.46333563 aaron Pick an output output vector token based on Х its probability (sample) GPT2 small: 117M The GPT2 medium: 345M -0.51006055 zvzzvva GPT2 large: 762M DECODER LLAMA2 7B: 7B ... LLAMA2 7B-chat: 7B DECODER LLAMA2 70B: 70B GPT3 davinci-002: 175B $\langle S \rangle$ 1 2 1024

**Parameters:** weights and biases; updated during training;  $\Rightarrow$  long-term memory~ish; **Activations:** temporarily variables generated and modified during generation (using LMs);

https://jalammar.github.io/illustrated-gpt2/

#### Scores *Predicting the next word!* (before softmax) Keys **Autoregressive Models** obey orders 0.11 0.00 0.81 0.79 robot must **Oueries** 0.50 0.30 robot must obey orders 0.19 0.48 robot must obey orders 0.53 0.98 0.95 0.14 must obey orders robot GPT2 small: 117M 0.86 0.38 0.90 0.81 robot must obev orders GPT2 medium: 345M GPT2 large: 762M LLAMA2 7B: 7B Scores Masked Scores LLAMA2 7B-chat: 7B (before softmax) (before softmax) LLAMA2 70B: 70B 0.11 0.00 0.81 0.79 0.11 -inf -inf -inf Apply Attention GPT3 davinci-002: 175B 0.19 0.50 0.30 0.48 Mask 0.19 0.50 -inf -inf 0.98 0.95 0.53 0.14 0.53 0.98 0.95 -inf 0.81 0.86 0.38 0.90 0.81 0.86 0.38 0.90

**Parameters:** weights and biases; updated during training;  $\Rightarrow$  long-term memory~ish; **Activations:** temporarily variables generated and modified during generation (using LMs);

https://jalammar.github.io/illustrated-gpt2/

Model

Model

Model



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



**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. Bro	wn* Benja	min Mann*	Nick R	yder* Me	lanie Subbiah*
Jared Kaplan <sup>†</sup>	Prafulla Dhariwa	Arvind Ne	elakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwa	Ariel Herbe	ert-Voss	Gretchen Krueger	Tom Henighan
<b>Rewon Child</b>	Aditya Ramesh	Daniel M.	Ziegler	Jeffrey Wu	Clemens Winter
Christopher He	esse Mark Cl	ien Eric S	Sigler	Mateusz Litwin	Scott Gray
Benja	min Chess	Jack Clar	·k	Christopher	Berner
Sam McCar	ndlish Ale	c Radford	Ilya Su	tskever I	Dario Amodei
OpenAI					

1 \* 1 = 23 \* 3 = 64 \* 2 = 67\*1=??

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

Translate English to French:	$\longleftarrow$ task description
sea otter => loutre de mer	$\longleftarrow$ example
cheese =>	← prompt
	Translate English to French: sea otter => loutre de mer cheese =>

#### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



### Model

#### One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

1	Translate English to French:	← task description
	sea otter => loutre de mer	←— example
	cheese =>	←— prompt



#### **In-context** Learning:

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

#### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



#### One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

1	Translate English to French:	← task description
	sea otter => loutre de mer	←— example
	cheese =>	← prompt



#### **In-context** Learning:

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

#### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.





#### Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



#### One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

1	Translate English to French:	$\longleftarrow$ task description
	sea otter => loutre de mer	←— example
	cheese =>	$\longleftarrow$ prompt

#### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.





#### In-context Learning:

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

### In-*weights* Learning (Fine-tuning):

- Gradient-based;
- Slow: need many examples;
- Standard supervised learning;



The model is trained via repeated gradient updates using a large corpus of example tasks.

Model



### Open Question: ICL ≈(functional) Gradient Descent?

#### Finetuning



### Open Question: ICL ≈(functional) 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;

E.g. [Xie et al. 2022, von Oswald et al. 2022, Dai et al. 2023]

### Open Question: ICL ≈(functional) 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;

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

E.g. [Xie et al. 2022, von Oswald et al. 2022, Dai et al. 2023]

### Interim Summary 1

### 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?

Phenomenon

**Linguistic Adaptation:** the linguistic knowledge representations that are used for language processing **change** in response to language input.

E.g. [Bock 1986, Chang 2012]

Phenomenon

**Linguistic Adaptation:** the linguistic knowledge representations that are used for language processing **change** in response to language input.



E.g. [Bock 1986, Chang 2012]

Phenomenon

**Linguistic Adaptation:** the linguistic knowledge representations that are used for language processing **change** in response to language input.

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

Phenomenon

**Linguistic Adaptation:** the linguistic knowledge representations that are used for language processing **change** in response to language input.

**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.

E.g. [Bock 1986, Chang 2012]

Phenomenon

**Linguistic Adaptation:** the linguistic knowledge representations that are used for language processing **change** in response to language input.

**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 (IFE)

Phenomenon

**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 (IFE)

Phenomenon

**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
Phenomenon

**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]

Phenomenon

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



Phenomenon

**Prime in DO Structure** 

**Target in PD Structure** 

**Unprimed log probability = -50** 

A doctor bought a chief a plate. <u>The secretary drew the card for the band.</u> A doctor designed a chief a plate. <u>The secretary drew the card for the band.</u>

 $T_{\rm PD}|P_{\rm DO}|$ 

Phenomenon

Prime in DO Structure

**Target in PD Structure** 

**Unprimed log probability = -50** 

A doctor bought a chief a plate. <u>The secretary drew the card for the band.</u> A doctor designed a chief a plate. <u>The secretary drew the card for the band.</u>

 $T_{\rm PD}|P_{\rm DO}|$ 

#### Prime Verb

Bring Buy Find Draw Design

Phenomenon

**Prime in DO Structure** 

Target in PD Structure

**Unprimed log probability = -50** 

A doctor bought a chief a plate. <u>The secretary drew the card for the band.</u> A doctor designed a chief a plate. <u>The secretary drew the card for the band.</u>

 $T_{\rm PD}|P_{\rm DO}|$ 

<u>Prime Verb</u>	<u>Verb PD Bias</u>	Primed log probability (priming magnitude)
Bring	0.23	<b>-50.5</b> (0.5)
<b>Buy</b>	0.27	<b>-51.1</b> (1.1)
<b>Find</b>	0.41	-51.8 (1.8)
Draw	0.52	<b>-52.2</b> (2.2)
<b>Design</b>	0.77	<b>-52.9</b> (2.9)

Phenomenon

**Prime in DO Structure** 

\_\_\_\_\_

**Target in PD Structure** 

**Unprimed log probability = -50** 

A doctor bought a chief a plate. <u>The secretary drew the card for the band.</u> A doctor designed a chief a plate. <u>The secretary drew the card for the band.</u>

 $T_{
m PD}|P_{
m DO}|$ 



#### Transient **Activation**





[Pickering & Branigan 1998]

#### Transient **Activation**





[Pickering & Branigan 1998]

#### Transient Activation









#### Theory 2



#### Theory 2



#### Theory 2



## Interim Summary 2

error-driven

learning?

In-context Learning (ICL) in LLMs		Structural Priming	Current Study	Discussion & Implications
<ul> <li><i>In-context</i> learning vs. <i>In-weights</i> learning</li> <li>ICL as a processing mechanism of LLMs;</li> <li>Is ICL functionally</li> </ul>	•	Structural priming effect; The Inverse Frequency Effect (IFE); Two accounts:		

vs. implicit learning;

# ICL as 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!

**Data Distributional Properties Drive Emergent In-Context Learning in Transformers** Stephanie C.Y. Chan **Adam Santoro** Andrew K. Lampinen Jane X. Wang DeepMind DeepMind DeepMind DeepMind Aaditya K. Singh Pierre H. Richemond James L. McClelland University College London DeepMind DeepMind, Stanford University Felix Hill DeepMind

**Burstiness as a distributional property!** 

Data Di	stributiona	al Prop	erties Dr	ive
Emergent In-	Context Le	earning	g in Tran	sformers
Stephanie C.Y. Chan	Adam Santoro	Andrew	<b>K. Lampinen</b>	Jane X. Wang
DeepMind	DeepMind	De	eepMind	DeepMind
Aaditya K. Singh	Pierre H. Ric	<b>:hemond</b>	James L.	McClelland
University College Londor	DeepM	ind	DeepMind, St	anford University
	<b>Felix</b> Deep	Hill Mind		

**Burstiness as a distributional property!** 

**Compositional Structures!** 

#### **Rethinking the Role of Demonstrations:** What Makes In-Context Learning Work?

Sewon Min <sup>1,2</sup>	Xinxi Lyu <sup>1</sup>	Ari Holtzma	$\mathbf{an}^1$	Mikel Artetxe <sup>2</sup>
Mike Lewis <sup>2</sup>	Hannaneh	Hajishirzi <sup>1,3</sup>	Luke	Zettlemoyer <sup>1,2</sup>

Data Dis	stributional	Prop	erties Dr	ive
Emergent In-(	Context Lea	rning	g in Trans	sformers
Stephanie C.Y. Chan DeepMind	Adam Santoro	Andrew Dec	<b>K. Lampinen</b> epMind	Jane X. Wang DeepMind
Aaditya K. Singh	Pierre H. Riche	mond	James L.	McClelland
University College London	DeepMind		DeepMind, St	anford University
	Felix Hi DeepMi	<b>ill</b> nd		

**Burstiness as a distributional property!** 

**Compositional Structures!** 

Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?

Sewon Min<sup>1,2</sup> Xinxi Lyu<sup>1</sup> Ari Holtzman<sup>1</sup> Mikel Artetxe<sup>2</sup> Mike Lewis<sup>2</sup> Hannaneh Hajishirzi<sup>1,3</sup> Luke Zettlemoyer<sup>1,2</sup>

**Parallel Structures in Pre-training Data Yield In-Context Learning** 

**Parallelism!** 

Yanda Chen<sup>1</sup> Chen Zhao<sup>2,3</sup> Zhou Yu<sup>1</sup> Kathleen McKeown<sup>1</sup> He He<sup>2</sup>

Context Lear	ning in Tran	sformers
Adam Santoro An	ndrew K. Lampinen	Jane X. Wang
DeepMind	DeepMind	DeepMind
Pierre H. Richem	ond James I	McClelland
DeepMind	DeepMind, S	Stanford University
	Adam Santoro Ar DeepMind Pierre H. Richem DeepMind	Adam Santoro       Andrew K. Lampinen         DeepMind       DeepMind         Pierre H. Richemond       James I         DeepMind       DeepMind, S

**Burstiness as a distributional property!** 

**Compositional Structures!** 

Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?

Sewon Min<sup>1,2</sup> Xinxi Lyu<sup>1</sup> Ari Holtzman<sup>1</sup> Mikel Artetxe<sup>2</sup> Mike Lewis<sup>2</sup> Hannaneh Hajishirzi<sup>1,3</sup> Luke Zettlemoyer<sup>1,2</sup>

**Parallel Structures in Pre-training Data Yield In-Context Learning** 

**Parallelism!** 

Yanda Chen<sup>1</sup> Chen Zhao<sup>2,3</sup> Zhou Yu<sup>1</sup> Kathleen McKeown<sup>1</sup> He He<sup>2</sup>

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

#### Current Approach: Overview



### Current Approach: Overview



**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;

# Corpus

#### Simulations

- 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

# 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

Simulations

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

# $t_{ m PD}|p_{ m PD},t_{ m PD}|p_{ m DO},t_{ m DO}|p_{ m PD},t_{ m DO}|p_{ m DO}$

**Prime**: A doctor brought a plate to a chief. **Target**: The secretary drew the card for the band.

# 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

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

Simulations

66



Simulations

67





















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

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

#### Results

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



#### Results

**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 and run the model — without weight update;



**Concatenation Mode:** concatenating the prime and target sentences and run the model — without weight update;





**Concatenation Mode:** concatenating the prime and target sentences and run the model — **without weight update**;



*Larger* models show more significant inverse frequency effects!



**Concatenation Mode:** concatenating the prime and target sentences and run the model — without weight update;





**Concatenation Mode:** concatenating the prime and target sentences and run the model — without weight update;



*Larger* models show more significant inverse frequency effects!

## Interim Summary 3

In-context Learning (ICL) in LLMs		Structural Priming		Current Study	Discussion & Implications
<ul> <li><i>In-context</i> learning vs. <i>In-weights</i> learning</li> <li>ICL as a processing</li> </ul>	•	Structural priming effect; The Inverse Frequency Effect	•	Analogy between priming and ICL; Do LLMs show the IFE? – yes, with	

- mechanism of LLMs;
- Is ICL functionally performing some error-driven learning?
- (IFE);
- Two accounts: transient activation vs. implicit learning;
- various degrees!
- IFE as a diagnostic of error-driven learning!

## Reasoning

[Assumption] Psycholinguistic Theories: only some kinds of gradient-based, error-driven learning mechanism could give rise to the IFE.

#### [Connecting Priming to ICL]

Having a prime sentence as the prompt in the context window conditions the probability distribution over the target sentence via ICL.

[Previous Works on Priming in LLMs] LLMs do show human-like behaviors for standard structural priming. **[Core Hypothesis]** If ICL involves a *gradient* component when processing the context, then LLMs with strong ICL capabilities should show the IFE.

[Previous Works on LLMs' sizes and ICL] The ICL capability scales with LLMs' size. [Additional Hypothesis] Because larger models have stronger ICL capabilities, they should display a more significant IFE.

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, supporting the hypothesis that [LLMs is functionally performing gradient descent] \*\*\*

# Thanks for Listening! Q&A Session 🙋 🧕

#### 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;
- ICL is functionally equivalent to Gradient Descent?

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

- 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* gradient component involved in ICL;
- Humans and LLMs share a similar processing mechanism!

## Open Discussion: activation vs. learning?

#### **Enriching Transient Activation to simulate implicit learning?**



**Dynamic Field Theory** 



**Gradient Symbolic Computation** 

#### On-going Work 1: Task / Function Vectors?



				Co	ntext =	÷					
0	0.33	1		1.1	0.71	0.11	-0.47	0.28	-0.13		
Ч	0.093	0.62		0.99	1	0.048	1	0.76	0.65	- 2	
2	0.12	0.93		0.27	0.32	0.0097	0.79	0.57	0.98		
m ·	0.16	0.69		0.47	0.52	1.4	1.2		0.93		
4	0.096	0.69		0.27	0.42	0.64	0.81	0.75	0.66	- 1	
<u>د</u>	0.19	0.81		0.19	0.32	0.022	0.49	1.2	0.17		
, o .	0.21	1		0.25	0.25	-0.6	0.42	1.2	-0.1		
- 7	0.17		2.4	1.1	0.32	-0.25	1.1	1.2	-0.25	- C	)
8	0.19	1.4	2.1	0.95	0.41	-0.28	0.83	1.4	-1.3		
о. С	0.12	1.9	0.98	1.1	0.24	0.016	1.2	1.3	-1.8		
10	0.12		-0.12	1.3	0.047	0.19	1.2	0.89	-2.7		-1
Ħ	0.031	1.8	-0.76	1.1	0.0027	0.15	0.8	0.28	-2.6		
12	0.033		-1	1.1	-0.14	-0.0041	0.72	0.45			
13	-0.061	1.1	-1.4	1	-0.57	-0.029	0.52	-0.22	-2.4		-2
14	-0.042	0.29		0.92	-1.1	0.066	0.27	-0.28	-2.5		~
15	0.19	1.8	-0.81	1	-0.081	-0.2	0.52				
	- A -	princess -	brought -	u D	- chicken -	ţ,	a d	mother -			
			loker	DOSITIO	n where	EV IS INIE	ected				

#### On-going Work 1: Task / Function Vectors cont.



(a) Input: "Italy, Russia, China, Japan, France"

FV	Task	<b>Expected Output</b>
V <sub>AC</sub>	First-Copy	Italy
V <sub>AD</sub>	First-Capital	Rome
$v_{BC}$	Last-Copy	France
$v_{BD}^{*}$	Last-Capital	Paris



## On-going Work 2: IFE in other ICL "Tasks"?

- If the IFE is a diagnostic for the error-driven nature of adaptation, then we could apply it to any other ICL tasks;
  - 1. Other structural priming instances:
    - a. Active-passive;
    - b. Complementizer (*that*) priming;
  - 2. Non-linguistic tasks where we can sort demonstration difficulty w.r.t baseline performance;
    - a. Two-digit {addition, multiplication}: 10 x 10 is easier than 31 x 67;
    - b. Country-capital mapping: more commonly known <country, capital> pairs leads to lower surprisal than less known pairs;

#### More Intricate Pattern: Pronoun vs. NoPronoun



**Pronoun Corpus:** replacing all indirect object with a pronoun;

#### Pronoun vs. NoPronoun: results

		Intercept		Slope		95% CI of Slopes		
Model	Pronoun Obj?	PD-PD	DO-PD	PD-PD	DO-PD	PD-PD	DO-PD	
GPT2-small	True	0.370	0.278	0.011	-0.007	(-0.049, 0.072)	(-0.058, 0.044)	
GPT2-small	False	0.746	0.653	0.014	0.006	(-0.042, 0.07)	(-0.06, 0.072)	
GPT2-medium	True	0.351	0.256	-0.013	-0.026	(-0.078, 0.053)	(-0.073, 0.022)	
GPT2-medium	False	0.748	0.590	-0.023	-0.035	(-0.079, 0.032)	(-0.125, 0.054)	
GPT2-large	True	0.330	0.241	0.011	-0.037+	(-0.043, 0.065)	(-0.09, 0.015)	
GPT2-large	False	0.698	0.487	-0.003	-0.02	(-0.062, 0.055)	(-0.098, 0.058)	
Llama-7b	True	0.392	0.229	-0.02	-0.086****	(-0.067, 0.026)	(-0.126, -0.045)	
Llama-7b	False	0.807	0.627	-0.026	-0.111+	(-0.102, 0.049)	(-0.279, 0.057)	
Llama-7b-chat	True	0.413	0.263	-0.012	-0.095****	(-0.067, 0.043)	(-0.146, -0.044)	
Llama-7b-chat	False	0.788	0.605	-0.013	-0.102	(-0.115, 0.089)	(-0.289, 0.085)	
Llama-13b	True	0.434	0.256	-0.059**	-0.099****	(-0.113, -0.005)	(-0.134, -0.063)	
Llama-13b	False	0.859	0.685	-0.066+	-0.177*	(-0.163, 0.031)	(-0.388, 0.033)	
GPT3-davinci-002	True	0.403	0.223	-0.078****	-0.078****	(-0.121, -0.035)	(-0.114, -0.043)	
GPT3-davinci-002	False	0.851	0.632	-0.064+	-0.145*	(-0.153, 0.025)	(-0.301, 0.012)	

**Pattern:** more significant IFE results in the WithPronoun condition than the NoPronoun condition;

**Interpretation**: ICL capability correlates with data scale – more WithPronoun sentences in the training data!

#### Future Direction: ERP as evidence of learning signal?

# The N400 ERP component reflects an error-based implicit learning signal during language comprehension



Predictive coding: A theory under which the functional role of neural systems centers on prediction and learning from prediction error.

# Future Direction: ERP as evidence of learning signal, cont.

SPECIAL ISSUE: Cognitive Computational Neuroscience of Language

#### Strong Prediction: Language Model Surprisal Explains Multiple N400 Effects

James A. Michaelov<sup>1</sup>, Megan D. Bardolph<sup>1</sup>, Cyma K. Van Petten<sup>2</sup>, Benjamin K. Bergen<sup>1</sup>, and Seana Coulson<sup>1</sup>

<sup>1</sup>Department of Cognitive Science, University of California, San Diego, La Jolla, CA, USA <sup>2</sup>Department of Psychology, Binghamton University, State University of New York, Binghamton, NY, USA

Keywords: distributional semantics, ERPs, N400, neural language models, predictive coding



#### Neural language model gradients predict event-related brain potentials

Stefan L. Frank Centre for Language Studies Radboud University Nijmegen, the Netherlands stefan.frank@ru.nl

## **Additional Slides**



**Existing Line of research:** treating <u>LMs as psycholinguistic subjects</u> and probing whether they have acquired the linguistic knowledge that humans used for sentence processing.



# **Existing Line of research:** treating <u>LMs as psycholinguistic subjects</u> and probing whether they have acquired the linguistic knowledge that humans used for sentence processing.

The Emergence of Number and Syntax Units in LSTM Language Models

Yair Lakretz Cognitive Neuroimaging Unit NeuroSpin center 91191, Gif-sur-Yvette, France yair.lakretz@gmail.com

**Theo Desbordes** Facebook AI Research Paris, France tdesbordes@fb.com

Stanislas Dehaene Cognitive Neuroimaging Unit NeuroSpin center 91191, Gif-sur-Yvette, France stanislas.dehaene@gmail.com German Kruszewski Facebook AI Research Paris, France germank@gmail.com

Dieuwke Hupkes ILLC, University of Amsterdam Amsterdam, Netherlands d.hupkes@uva.nl

> Marco Baroni Facebook AI Research Paris, France mbaroni@fb.com



#### What do RNN Language Models Learn about Filler–Gap Dependencies?

Ethan Wilcox<sup>1</sup>, Roger Levy<sup>2</sup>, Takashi Morita<sup>3,4</sup>, and Richard Futrell<sup>5</sup>

<sup>1</sup>Department of Linguistics, Harvard University, wilcoxeg@g.harvard.edu
<sup>2</sup>Department of Brain and Cognitive Sciences, MIT, rplevy@mit.edu
<sup>3</sup>Primate Research Institute, Kyoto University, tmorita@alum.mit.edu
<sup>4</sup>Department of Linguistics and Philosophy, MIT
<sup>5</sup>Department of Language Science, UC Irvine, rfutrell@uci.edu

Model

**Existing Line of research:** treating <u>LMs as psycholinguistic subjects</u> and probing whether they have acquired the linguistic knowledge that humans used for sentence processing.

#### When a sentence does not introduce a discourse entity, Transformer-based models still sometimes refer to it

Sebastian Schuster Center for Data Science Department of Linguistics New York University schuster@nyu.edu

#### Tal Linzen

Center for Data Science Department of Linguistics New York University linzen@nyu.edu

Model



#### Neural Language Models as Psycholinguistic Subjects: Representations of Syntactic State

Richard Futrell<sup>1</sup>, Ethan Wilcox<sup>2</sup>, Takashi Morita<sup>3,4</sup>, Peng Qian<sup>5</sup>, Miguel Ballesteros<sup>6</sup>, and Roger Levy<sup>5</sup>

<sup>1</sup>Department of Language Science, UC Irvine, rfutrell@uci.edu
<sup>2</sup>Department of Linguistics, Harvard University, wilcoxeg@g.harvard.edu
<sup>3</sup>Primate Research Institute, Kyoto University, tmorita@alum.mit.edu
<sup>4</sup>Department of Linguistics and Philosophy, MIT
<sup>5</sup>Department of Brain and Cognitive Sciences, MIT, {pqian, rplevy}@mit.edu
<sup>6</sup>IBM Research, MIT-IBM Watson AI Lab, miguel.ballesteros@ibm.com

Model

## Probing Classifiers

Pro	noun_la	arge_bin	_label_t	arget_T	hreeLay	er	s . 1.0	Pr	onoun_	large_b	in_label_	target_	OneLaye	er	- 1.0
layer0 -	0.50	0.50	1.00	1.00	1.00		- 1.0	layer0 -	0.50	0.50	1.00	1.00	1.00		- 1.0
layer3 -	0.50	0.50	1.00	1.00	1.00			layer3 -	0.50	0.50	1.00	1.00	1.00		
layer6 -	0.50	0.50	1.00	1.00	1.00		- 0.9	layer6 -	0.50	0.50	1.00	1.00	1.00		- 0.9
layer9 -	0.50	0.50	1.00	1.00	1.00			layer9 -	0.50	0.50	1.00	1.00	1.00		
layer12 -	0.50	0.50	1.00	1.00	1.00			layer12 -	0.50	0.50	1.00	1.00	1.00		
layer15 -	0.50	0.50	1.00	1.00	1.00		- 0.8	layer15 -	0.50	0.50	1.00	1.00	1.00		- 0.8
layer18 -	0.50	0.50	1.00	1.00	1.00			layer18 -	0.50	0.50	1.00	1.00	1.00		
layer21 -	0.50	0.50	1.00	1.00	1.00		- 0.7	layer21 -	0.50	0.50	1.00	1.00	1.00		- 0.7
layer24 -	0.50	0.50	1.00	1.00	1.00			layer24 -	0.50	0.50	1.00	1.00	1.00		
layer27 -	0.50	0.50	1.00	1.00	1.00			layer27 -	0.50	0.50	1.00	1.00	1.00		
layer30 -	0.50	0.50	1.00	1.00	1.00		- 0.6	layer30 -	0.50	0.50	1.00	1.00	1.00		- 0.6
layer33 -	0.50	0.50	1.00	1.00	1.00			layer33 -	0.50	0.50	1.00	1.00	1.00		
final -	0.50	0.50	1.00	1.00	1.00			final -	0.50	0.50	1.00	1.00	1.00		
	Period	v	DP1_mid	DP1_end	DP1_after		- 0.5	-	Period	v	DP1_mid	DP1_end	DP1_after		- 0.5

	layer30 -	1.50	1.50	0.50	0.50	0.47			1
	layer27 -	1.50	1.50	0.52	0.49	0.49			1
	layer24 -	1.50	1.50	0.50	0.49	0.46	-	0.8	1
	layer21 -	1.50	1.50	0.58	0.48	0.43			1
	layer18 -	1.50	1.50	0.57	0.47	0.45	-	1.0	1
	layer15 -	1.50	1.50	0.53	0.49	0.51			1
).5	layer12 -	1.50	1.50	0.55	0.48	0.49			1
	layer9 -	1.50	1.50	0.56	0.51	0.54	-	1.2	
	layer6 -	1.50	1.50	0.62	0.53	0.54			
).6	layer3 -	1.50	1.50	0.59	0.58	0.71	-	1.4	
	layer0 -	1.50	1.50	0.70	0.75	0.70			
	Pro	noun lai	ae cor	it label	large T	hreeLav	rers		
).7									

		Pr	onoun_	large_co	ont_labe	l_large_	OneLaye	er	
		layer0 -	1.50	1.50	0.68	0.84	0.59		
-	1.4	layer3 -	1.50	1.50	0.55	0.56	0.50		- 1.4
		layer6 -	1.50	1.50	0.54	0.57	0.47		
-	1.2	layer9 -	1.50	1.50	0.55	0.48	0.49		- 1.2
		layer12 -	1.50	1.50	0.50	0.45	0.48		
		layer15 -	1.50	1.50	0.52	0.45	0.47		
1	1.0	layer18 -	1.50	1.50	0.46	0.45	0.44		- 1.0
		layer21 -	1.50	1.50	0.50	0.45	0.46		
-	0.8	layer24 -	1.50	1.50	0.49	0.44	0.43		- 0.8
		layer27 -	1.50	1.50	0.48	0.43	0.48		
		layer30 -	1.50	1.50	0.48	0.46	0.44		- 0.4
	0.0	layer33 -	1.50	1.50	0.51	0.44	0.46		0.0
		final -	1.50	1.50	0.49	0.44	0.44		

Period V DP1\_mid DP1\_end DP1\_after

#### Verb Bias: With Pronoun vs. No Pronoun

1.0 PD Biases for Each Prime Verb 0.8 0.6 0.4 0.2 0.0 St State State to the state of the state CISTON DESIGN **Prime Verbs** 

No Pronoun

With Pronoun



#### GPT2-small doesn't show IFE in TA mode



# Explanations?





**Implicit Learning Mode:** GPT2-small does have the capability to capture the IFE!

Presence vs. absence of Error-driven weight update mechanism! There is some mechanism that functionally performs gradient update in LLMs!

**Transient Activation Mode:** GPT2-small doesn't show the IFE, while larger models show stronger IFEs.



# Thanks for Listening! Q&A Session 🙋



#### 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;

- 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* 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.