

# Provable Benefits of In-Tool Learning for Large Language Models

Ambroise Odonnat

Noah's Ark Lab, Inria

Kyutai

February 19, 2026



*Inria*

/ **kyutai**



Sam Houlston  
ETH Zürich



Charles Arnal  
FAIR at Meta



Vivien Cabannes  
FAIR at Meta



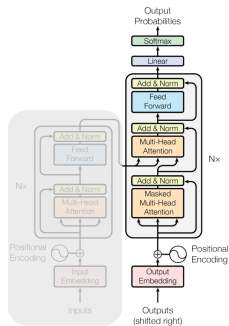
- ① Introduction
- ② Methodology
- ③ Theoretical results
- ④ Experiments
- ⑤ Take home message



- 1 Introduction
- 2 Methodology
- 3 Theoretical results
- 4 Experiments
- 5 Take home message



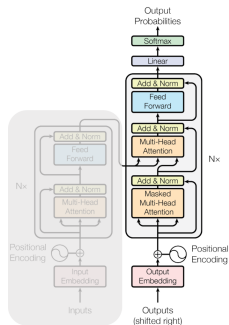
**Best models so far** → autoregressive decoder-only transformers.



(Vaswani et al., 2017)



**Best models so far** → autoregressive decoder-only transformers.



(Vaswani et al., 2017)



Vocabulary size  $T$ .  
Context window  $K$ .  
Parameter set  $\Theta$ .

**GPT-3** :  $T = 50257$ ,  
 $K = 2048$  and  $|\Theta| \sim 175B$



**Goal:** Predict the next token based on previous ones.



**Goal:** Predict the next token based on previous ones.

**Autoregressive:** Each token depends only on the past ones.

I am the      danger  
Previous tokens (context)      Token being predicted



**Goal:** Predict the next token based on previous ones.

**Autoregressive:** Each token depends only on the past ones.

I am the                      danger  
Previous tokens (context)    Token being predicted

**Modelization:** Probability of a sequence  $(x_1, x_2, \dots, x_N)$ :

$$\begin{aligned}\mathbb{P}(x_1, x_2, \dots, x_N) &= \mathbb{P}(x_1)\mathbb{P}(x_2 | x_1) \cdots \mathbb{P}(x_N | x_1, x_2, \dots, x_{N-1}) \\ &= \prod_{n=1}^N \mathbb{P}(x_n | x_1, x_2, \dots, x_{n-1})\end{aligned}$$



## Large-scale pretraining unlocks emerging capabilities!

---

### Language Models are Few-Shot Learners

---

<b>Tom B. Brown*</b>	<b>Benjamin Mann*</b>	<b>Nick Ryder*</b>	<b>Melanie Subbiah*</b>	
<b>Jared Kaplan†</b>	<b>Prafulla Dhariwal</b>	<b>Arvind Neelakantan</b>	<b>Pranav Shyam</b>	<b>Girish Sastry</b>
<b>Amanda Askell</b>	<b>Sandhini Agarwal</b>	<b>Ariel Herbert-Voss</b>	<b>Gretchen Krueger</b>	<b>Tom Henighan</b>
<b>Rewon Child</b>	<b>Aditya Ramesh</b>	<b>Daniel M. Ziegler</b>	<b>Jeffrey Wu</b>	<b>Clemens Winter</b>
<b>Christopher Hesse</b>	<b>Mark Chen</b>	<b>Eric Sigler</b>	<b>Mateusz Litwin</b>	<b>Scott Gray</b>
<b>Benjamin Chess</b>		<b>Jack Clark</b>	<b>Christopher Berner</b>	
<b>Sam McCandlish</b>	<b>Alec Radford</b>	<b>Ilya Sutskever</b>	<b>Dario Amodei</b>	

OpenAI

GPT3 and in-context learning (Brown et al., 2020)



## Large-scale pretraining unlocks emerging capabilities!

### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

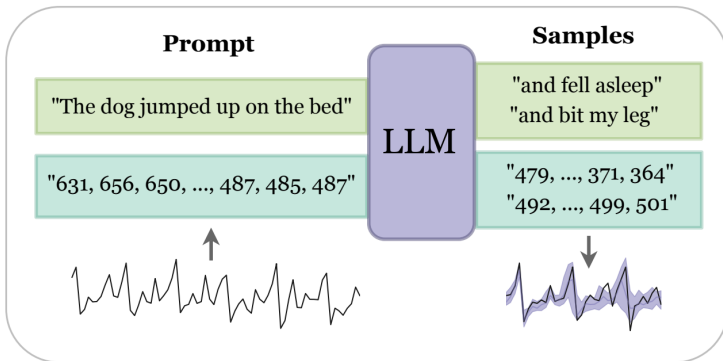
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Chain-of-thought (Wei et al., 2022)



## Large-scale pretraining unlocks emerging capabilities!



Generalization to other modalities (Gruver et al., 2023)



# Emerging capabilities of (static) LLMs

## Large-scale pretraining unlocks emerging capabilities!

**Language Models are Few-Shot Learners**

Tou B. Brown*	Benjamin Mann*	Nick Ryder*	Mohamé Subbiab*	
Jared Kaplan*	Pranjalit Dhariwal	Arvind Neelakantan	Pranav Shyam	Gleb Sotny
Amanda Ahlert	Sandeep Agarwal	Ariel Herbert-Voss	Grothoan Krueger	Tou Heifhgan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Mazeau Litwin	Scott Gray
Benjamin Chess	Jack Clark	Christopher Berner		
Sun McCandlish	Alec Radford	Rya Sankovner	Dario Amodei	

OpenAI

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls, 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Prompt**

"The dog jumped up on the bed"

"631, 656, 650, ..., 487, 485, 487"

**LLM**

**Samples**

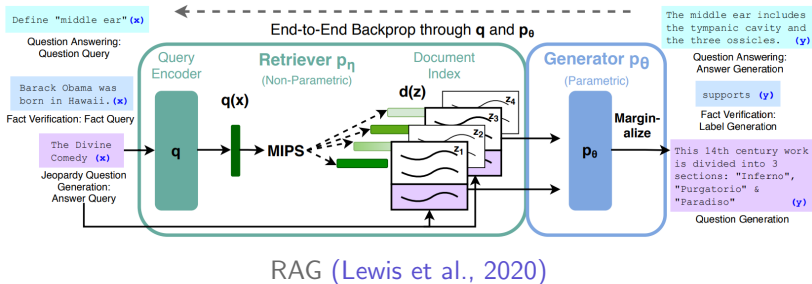
"and fell asleep"  
"and bit my leg"

"479, ..., 371, 364"  
"492, ..., 499, 501"

✗ Those methods extract knowledge from **static** predictors.

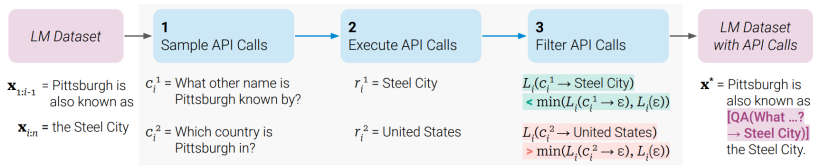


## Tool use to externalize memory and adapt to the context.





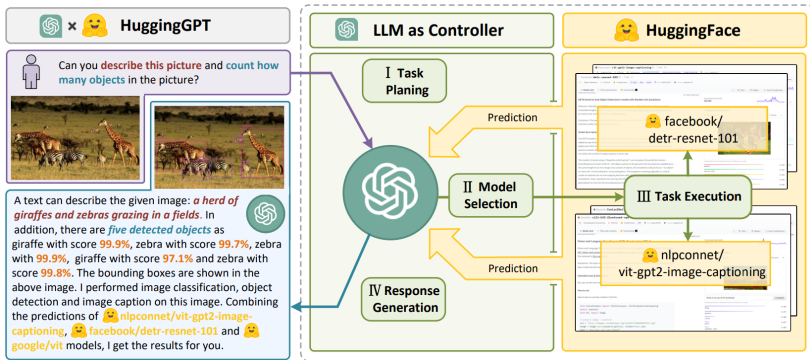
## Tool use to externalize memory and adapt to the context.



Toolformer (Schick et al., 2023)



## Tool use to externalize memory and adapt to the context.



HuggingGPT (Shen et al., 2023)



## Tool use to externalize memory and adapt to the context.

- ✓ LLMs evolve towards dynamic context-aware systems,
- ✓ It allows them to reason, adapt, and act over time,
- ✓ Emergence of agentic workflows (e.g., [Claude](#), [OpenAI](#))



What is the most efficient way to acquire and use knowledge?

- ① In-weight learning: memorize facts in the parameters,
- ② In-tool learning: learn to access external sources of truth.



What is the most efficient way to acquire and use knowledge?

- ① In-weight learning: memorize facts in the parameters,
- ② In-tool learning: learn to access external sources of truth.

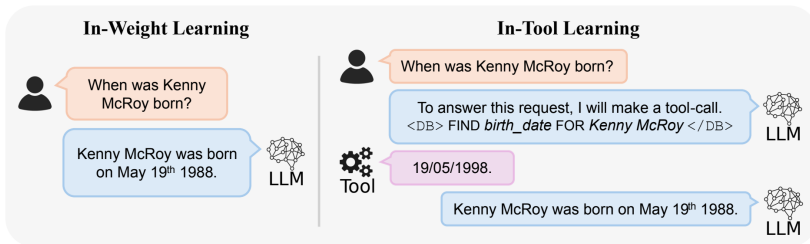


Figure: In-weight (memorization) vs. in-tool learning (external retrieval).



- ① Introduction
- ② Methodology
- ③ Theoretical results
- ④ Experiments
- ⑤ Take home message



## Factual dataset

Let  $\mathcal{D}$  be a family of factual datasets  $D$  defined as finite collections of facts  $(n, a, v)$  with

- ★ a name  $n \in \mathcal{N}$  (e.g., “Xavier Miel”),
- ★ an attribute  $a \in \mathcal{A}$  (e.g., “birthplace”),
- ★ a value  $v \in \mathcal{V}_a$  (e.g., “France”).



## Factual dataset

Let  $\mathcal{D}$  be a family of factual datasets  $D$  defined as finite collections of facts  $(n, a, v)$  with

- ★ a name  $n \in \mathcal{N}$  (e.g., “Xavier Miel”),
  - ★ an attribute  $a \in \mathcal{A}$  (e.g., “birthplace”),
  - ★ a value  $v \in \mathcal{V}_a$  (e.g., “France”).
- 
- ✓ Akin to Physics of Language models (Allen-Zhu et al., 2024),
  - ✓ Allows to quantify the amount of facts ( $\#\text{Facts} = |\mathcal{N} \times \mathcal{A}|$ ).



- ✓ Let  $\mathcal{M}$  be the set of all transformers with a given architecture,
- ✓ Each model  $f \in \mathcal{M}$  amounts to a specific choice of weights,
- ✓ Each model  $f \in \mathcal{M}$  is associated with a recall rule  $R(f)$ .



- ✓ Let  $\mathcal{M}$  be the set of all transformers with a given architecture,
- ✓ Each model  $f \in \mathcal{M}$  amounts to a specific choice of weights,
- ✓ Each model  $f \in \mathcal{M}$  is associated with a recall rule  $R(f)$ .

## Recall rule

$R(f)$  is defined by prompting a model to query the value of a pair  $(n, a)$ . The recall accuracy of  $f$  on a dataset  $D$  is the percentage of retrieval by  $R(f)$  of the correct value  $v$  over facts  $(n, a, v) \in D$ .



- ✓ Let  $\mathcal{M}$  be the set of all transformers with a given architecture,
- ✓ Each model  $f \in \mathcal{M}$  amounts to a specific choice of weights,
- ✓ Each model  $f \in \mathcal{M}$  is associated with a recall rule  $R(f)$ .

## Learnability

We say that the model class  $\mathcal{M}$  *can solve* the recall task on the family of datasets  $\mathcal{D}$  if, for each dataset  $D \in \mathcal{D}$ , there exists a model  $f \in \mathcal{M}$  that achieves a perfect recall accuracy on  $D$ .



**In-weight: Train the model to generate the answer directly.**



**In-weight: Train the model to generate the answer directly.**

- ✓ Templates  $\phi_i$  for queries,  $\psi_i$  for answers,
- ✓ Query  $Q = \phi_1(a) \circ \phi_2(n) \circ \phi_3(a)$ ,
- ✓ Answer  $A = \psi_1(n) \circ \psi_2(a) \circ \psi_3(v)$ .



**In-weight: Train the model to generate the answer directly.**

- ✓ Templates  $\phi_i$  for queries,  $\psi_i$  for answers,
- ✓ Query  $Q = \phi_1(a) \circ \phi_2(n) \circ \phi_3(a)$ ,
- ✓ Answer  $A = \psi_1(n) \circ \psi_2(a) \circ \psi_3(v)$ .

$$Q = \underbrace{\text{Where was}}_{\phi_1(a)} \underbrace{\text{Xavier Miel}}_{\phi_2(n)} \underbrace{\text{born?}}_{\phi_3(a)}$$
$$A = \underbrace{\text{Xavier Miel}}_{\psi_1(n)} \underbrace{\text{was born in}}_{\psi_2(a)} \underbrace{\text{France}}_{\psi_3(v)}$$



**In-tool: Train the model to use a tool to retrieve the answer.**



## In-tool: Train the model to use a tool to retrieve the answer.

- ✓ Templates  $\phi_i$  for queries,  $\psi_i$  for answers,  $\chi_i$  **for tool calls**,
- ✓ Query  $Q = \phi_1(a) \circ \phi_2(n) \circ \phi_3(a)$ ,
- ✓ **Tool**  $T = \chi_1(a) \circ \chi_2(n)$ ,
- ✓ Answer  $A = \psi_1(n) \circ \psi_2(a) \circ \psi_3(v)$ .



**In-tool: Train the model to use a tool to retrieve the answer.**

- ✓ Templates  $\phi_i$  for queries,  $\psi_i$  for answers,  $\chi_i$  **for tool calls**,
- ✓ Query  $Q = \phi_1(a) \circ \phi_2(n) \circ \phi_3(a)$ ,
- ✓ **Tool**  $T = \chi_1(a) \circ \chi_2(n)$ ,
- ✓ Answer  $A = \psi_1(n) \circ \psi_2(a) \circ \psi_3(v)$ .

$T = \underbrace{\text{To answer this request, I will make a tool-call. <DB> FIND birthplace}}_{\chi_1(a)}$   
 $\underbrace{\text{FOR Xavier Miel </DB>}}_{\chi_2(n)}$

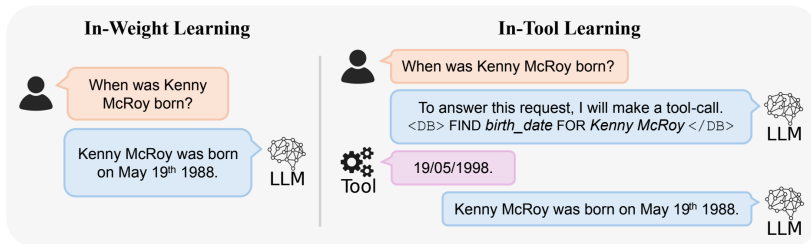


Figure: In-weight (memorization) vs. in-tool learning (external retrieval).



## Extension to, e.g., boolean: “Is Paris the capital of France?”

- ✓  $\phi_1(a) = \text{is}$ ,  $\phi_2(n) = \text{Paris}$ ,  $\phi_3(a) = \text{the capital of France}$ ,
- ✓ Define query as before  $Q = \phi_1(a) \circ \phi_2(n) \circ \phi_3(a)$ ;
- ✓  $\psi_1$  a binary function outputting “Yes” or “No”,
- ✓ Define answer as before  $A = \psi_1(a, n) \in \{\text{Yes}, \text{No}\}$ .



- ① Introduction
- ② Methodology
- ③ Theoretical results
- ④ Experiments
- ⑤ Take home message



**In-weight memorization is limited by the size of models.**



**In-weight memorization is limited by the size of models.**

Theorem (In-weight lower bound - informal)

*Let  $\mathcal{M}$  be a set of transformers with  $P$  parameters and  $\mathcal{D}$  be a family of factual datasets. If  $\mathcal{M}$  can solve the recall task only with in-weight learning, then the number of parameters  $P$  must satisfy:*

$$P \geq \mathcal{O}(\#\text{Facts}).$$



## In-weight memorization is limited by the size of models.

Theorem (In-weight lower bound - informal)

*Let  $\mathcal{M}$  be a set of transformers with  $P$  parameters and  $\mathcal{D}$  be a family of factual datasets. If  $\mathcal{M}$  can solve the recall task only with in-weight learning, then the number of parameters  $P$  must satisfy:*

$$P \geq \mathcal{O}(\#Facts).$$

- ✗ When the number of facts grows, memorization is impossible,
- ✗ Architectural changes or external memory are needed.



**In-tool learning solves the limitations of in-weight learning.**



**In-tool learning solves the limitations of in-weight learning.**

Theorem (In-tool upper bound - informal)

*Let  $\mathcal{D}$  be a family of factual datasets. Then, there exists a transformer with at most  $\bar{P}$  parameters that can solve the recall task if augmented with the proper retrieval tool.*



**In-tool learning solves the limitations of in-weight learning.**

Theorem (In-tool upper bound - informal)

*Let  $\mathcal{D}$  be a family of factual datasets. Then, there exists a transformer with at most  $\bar{P}$  parameters that can solve the recall task if augmented with the proper retrieval tool.*

- ✓ Model's size  $\bar{P}$  independent of the number of people  $|\mathcal{N}|$ ,
- ✓ No additional parameters needed when number of facts grows.



- ① Introduction
- ② Methodology
- ③ Theoretical results
- ④ Experiments
- ⑤ Take home message



**Pretraining small language models on factual datasets.**



## Pretraining small language models on factual datasets.

- ✓ Data organized à la [FineWeb](#),



## Pretraining small language models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Byte tokenizer with chat template for tool use,



## Pretraining small language models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Byte tokenizer with chat template for tool use,
- ✓ Custom pretraining pipeline akin to [Meta Lingua](#),



## Pretraining small language models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Byte tokenizer with chat template for tool use,
- ✓ Custom pretraining pipeline akin to [Meta Lingua](#),
- ✓ Small Llama3-like language models,



## Pretraining small language models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Byte tokenizer with chat template for tool use,
- ✓ Custom pretraining pipeline akin to [Meta Lingua](#),
- ✓ Small Llama3-like language models,
- ✓ Tool is an SQL agent that queries an external database.



**Empirical validation of the lower and upper bounds.**



## Empirical validation of the lower and upper bounds.

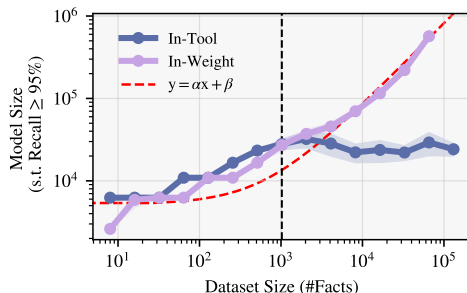


Figure: Minimal model size  $P$  to achieve at least a 95% recall.



## Empirical validation of the lower and upper bounds.

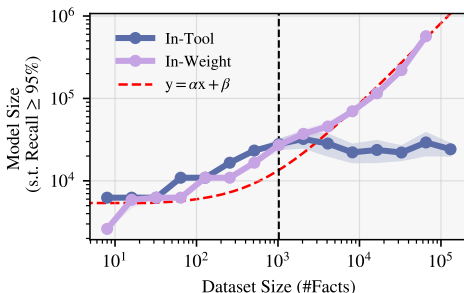


Figure: Minimal model size  $P$  to achieve at least a 95% recall.

- ✗ In-weight scales linearly with the number of facts,
- ✓ In-tool remains constant after a critical dataset size.



**In-tool models transition from memorization to rule learning.**



In-tool models transition from memorization to rule learning.

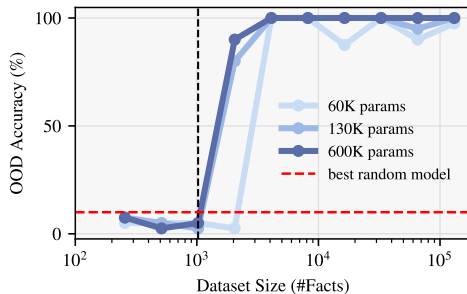


Figure: In-tool recall accuracy on out-of-distribution databases.



In-tool models transition from memorization to rule learning.

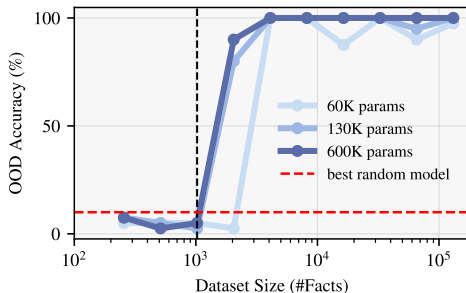


Figure: In-tool recall accuracy on out-of-distribution databases.

- ✗ Models start to memorize facts, similar to in-weight learning,
- ✓ After the transition, models learn to formulate in-tool queries.

# What counts as a fact?



**In-weight memorization is easier with interdependent facts.**



In-weight memorization is easier with interdependent facts.

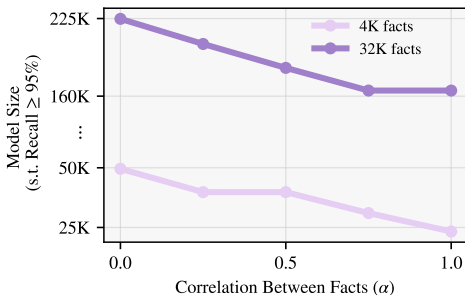


Figure: Minimal model size  $P$  to achieve at least a 95% recall.



In-weight memorization is easier with interdependent facts.

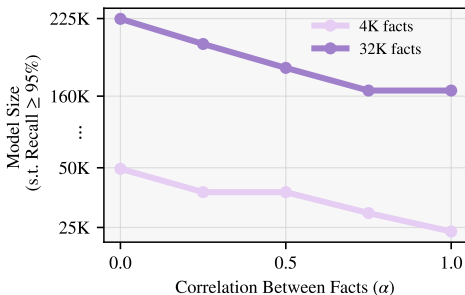


Figure: Minimal model size  $P$  to achieve at least a 95% recall.

- ✓ Correlation breaks independence between triplets  $(n, a, v)$ ,
- ✓ Reduced number of “effective” facts, easier to store.



**Finetuning large instruct models on factual datasets.**



## Finetuning large instruct models on factual datasets.

- ✓ Data organized à la [FineWeb](#),



## Finetuning large instruct models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Default tokenizers with special tokens and chat templates,



## Finetuning large instruct models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Default tokenizers with special tokens and chat templates,
- ✓ [SmolLM](#) (135M, 360M, 1.7B) and [Llama3](#) (1B, 3B, 8B),



## Finetuning large instruct models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Default tokenizers with special tokens and chat templates,
- ✓ [SmolLM](#) (135M, 360M, 1.7B) and [Llama3](#) (1B, 3B, 8B),
- ✓ Tool is an SQL agent that queries an external database,



## Finetuning large instruct models on factual datasets.

- ✓ Data organized à la [FineWeb](#),
- ✓ Default tokenizers with special tokens and chat templates,
- ✓ [SmolLM](#) (135M, 360M, 1.7B) and [Llama3](#) (1B, 3B, 8B),
- ✓ Tool is an SQL agent that queries an external database,
- ✓ Evaluation with [LM Evaluation Harness](#).



**In-tool preserves capabilities while learning new facts.**



In-tool preserves capabilities while learning new facts.

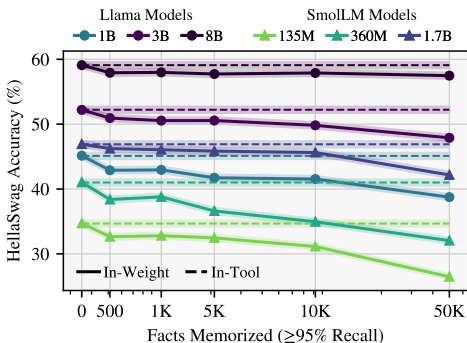


Figure: HellaSwag accuracy of models finetuned on factual datasets.



## In-tool preserves capabilities while learning new facts.

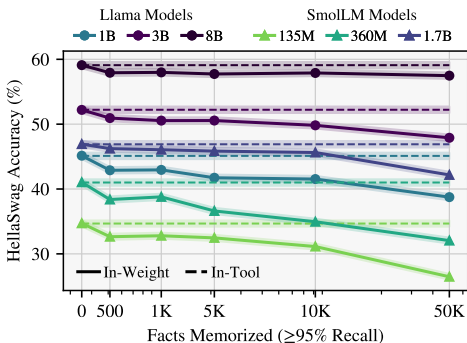


Figure: HellaSwag accuracy of models finetuned on factual datasets.

- ✗ In-weight learning impacts prior knowledge due to overloading,
- ✓ In-tool learning enables scalability without forgetting.



**In-tool learning helps preserve models' behavior.**



## In-tool learning helps preserve models' behavior.

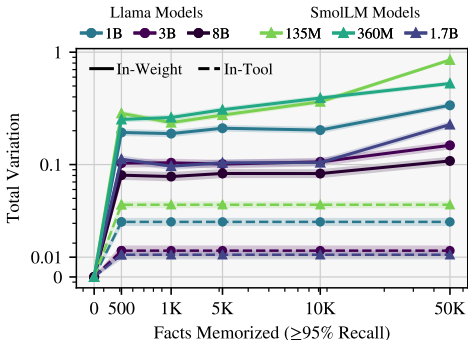


Figure: TV distance between base and finetuned models' distributions.



## In-tool learning helps preserve models' behavior.

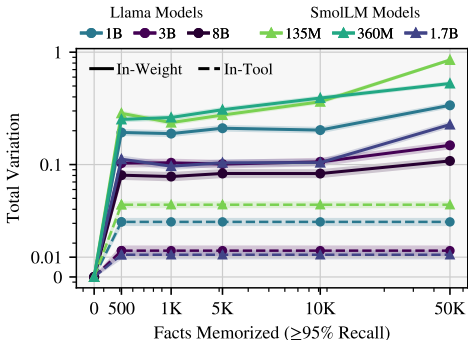


Figure: TV distance between base and finetuned models' distributions.

- ✗ In-weight learning alters models' token distribution,
- ✓ In-tool models remain close to the base models.



**In-tool learning requires less training compute than in-weight.**



In-tool learning requires less training compute than in-weight.

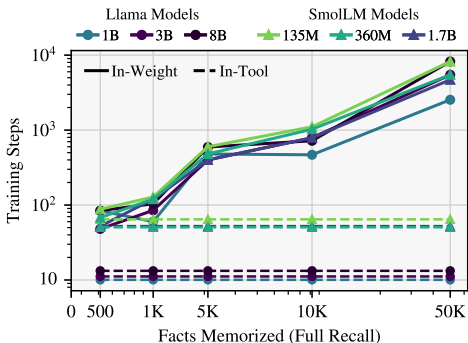


Figure: Training steps required to achieve a 100% recall.



In-tool learning requires less training compute than in-weight.

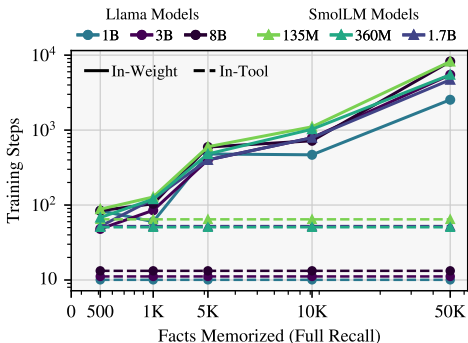


Figure: Training steps required to achieve a 100% recall.

- ✗ Memorizing individual facts requires many training steps,
- ✓ Learning to use a tool requires very few training steps.



- ① Introduction
- ② Methodology
- ③ Theoretical results
- ④ Experiments
- ⑤ Take home message



*Mieux vaut une tête bien  
faite qu'une tête bien pleine.*

---

Montaigne, Les Essais



*Mieux vaut une tête bien  
faite qu'une tête bien pleine.*

---

Montaigne, Les Essais

- ✓ In-weight learning is fundamentally limited by models' size,
- ✓ In-tool learning is more efficient and preserves capabilities.



*Mieux vaut une tête bien  
faite qu'une tête bien pleine.*

---

Montaigne, Les Essais

- ✓ In-weight learning is fundamentally limited by models' size,
- ✓ In-tool learning is more efficient and preserves capabilities.

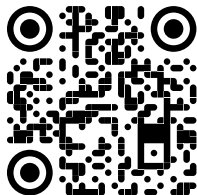
**Better to learn to use tools than to memorize facts.**



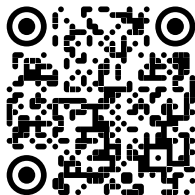
- ✓ Additional QA types (e.g., boolean, multiple choices),
- ✓ Additional tools use (e.g., Python interpreter, internet access),
- ✓ Extension to complex LLM agents.



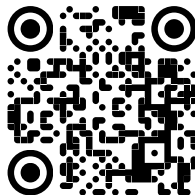
Check out the paper and code if you want to know more!



paper



code



website

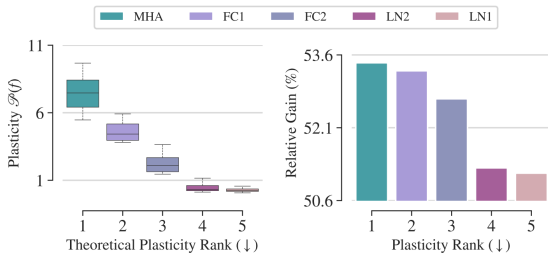
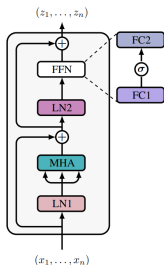


## Vision Transformer Benefits from Non-Smooth Components

[paper](#) - [code](#)



Finetuning high-plasticity components (non-smooth) yields larger performance gains than the stiff ones.



[paper](#) - [code](#)

Thanks for your attention!