

Updating Language Models

Joel Jang | MS Student @ KAIST | 02.11.2023

<https://joeljang.github.io/>

Table of Contents

Part 1 (~30 minutes)

- Towards Continual Knowledge Learning of Language Models [ICLR'22]
- TemporalWiki: A Lifelong Benchmark for Training and Evaluating Ever-Evolving Language Models [EMNLP'22]
- Knowledge Unlearning for Mitigating Privacy Risks in Language Models [*under review*]

Part 2 (~30 minutes)

- Exploring the Benefits of Training Expert Language Models over Instruction Tuning [*under review*]

Table of Contents

Part 1 (~30 minutes)

- Towards Continual Knowledge Learning of Language Models [ICLR'22]
- TemporalWiki: A Lifelong Benchmark for Training and Evaluating Ever-Evolving Language Models [EMNLP'22]
- Knowledge Unlearning for Mitigating Privacy Risks in Language Models [*under review*]

Part 2 (~30 minutes)

- Exploring the Benefits of Training Expert Language Models over Instruction Tuning [*under review*]

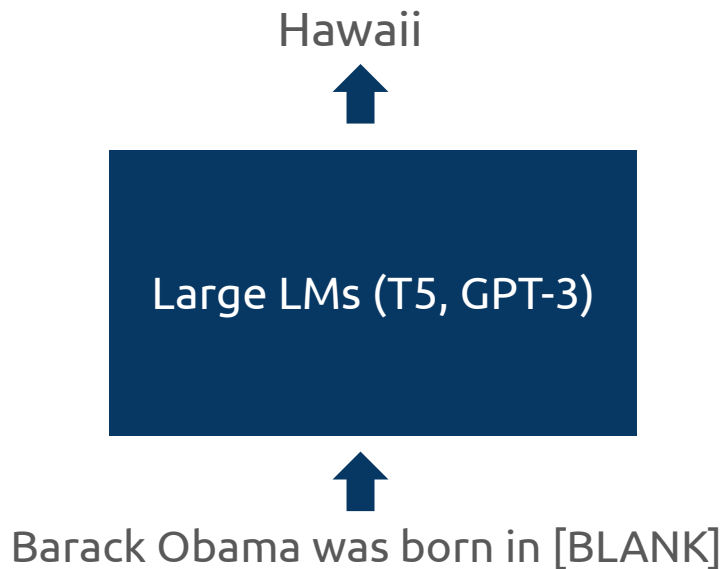
Towards Continual Knowledge Learning of Language Models [ICLR'22]

Joel Jang¹, Seonghyeon Ye¹, Sohee Yang¹, Joongbo Shin², Janghoon Han²,
Gyeunghun Kim², Stanley Choi², Minjoon Seo¹

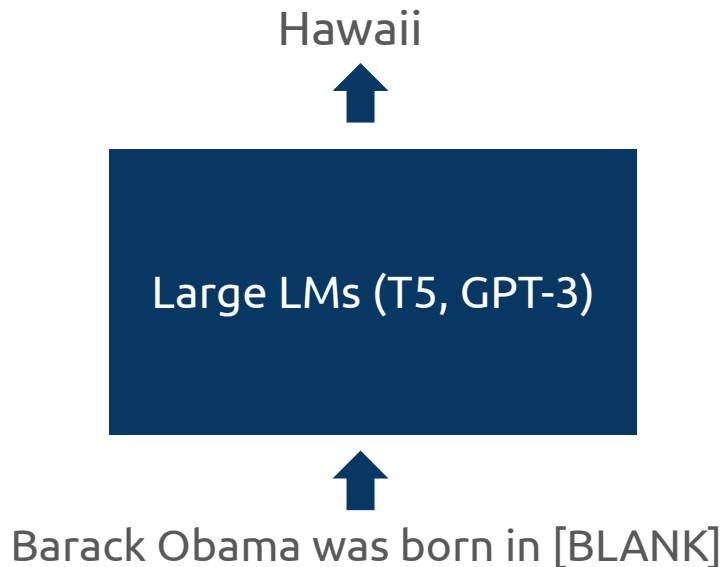
¹ **KAIST AI**
Graduate School of AI

²  **LG AI Research**

Motivation



Motivation



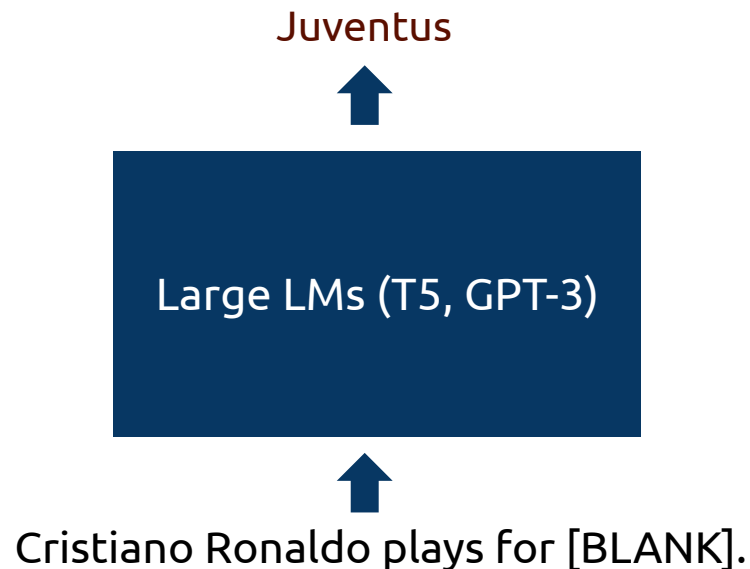
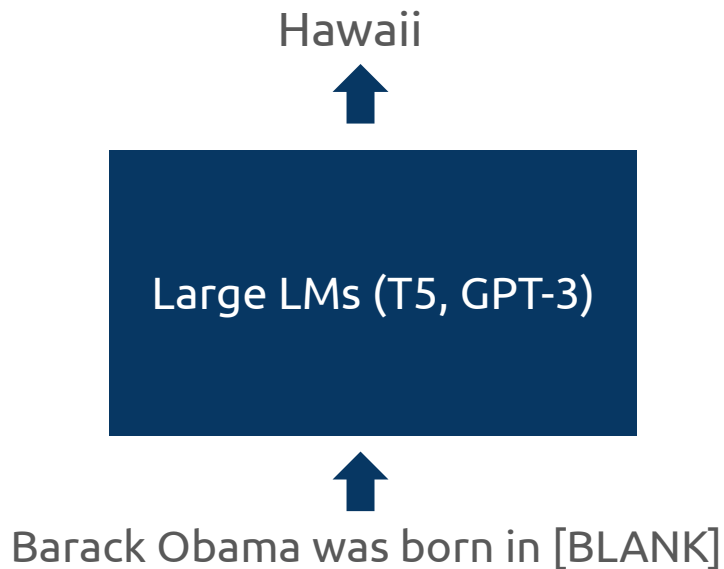
Open Domain Question Answering

Fact Checking

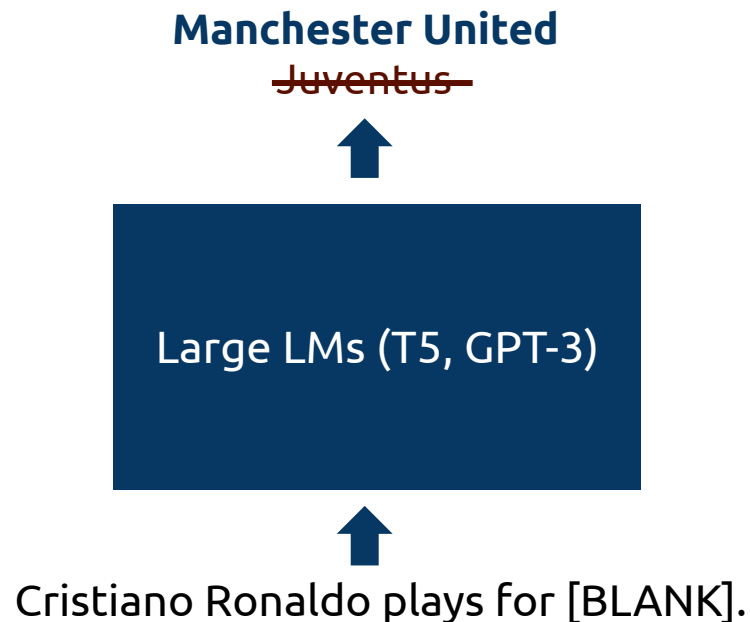
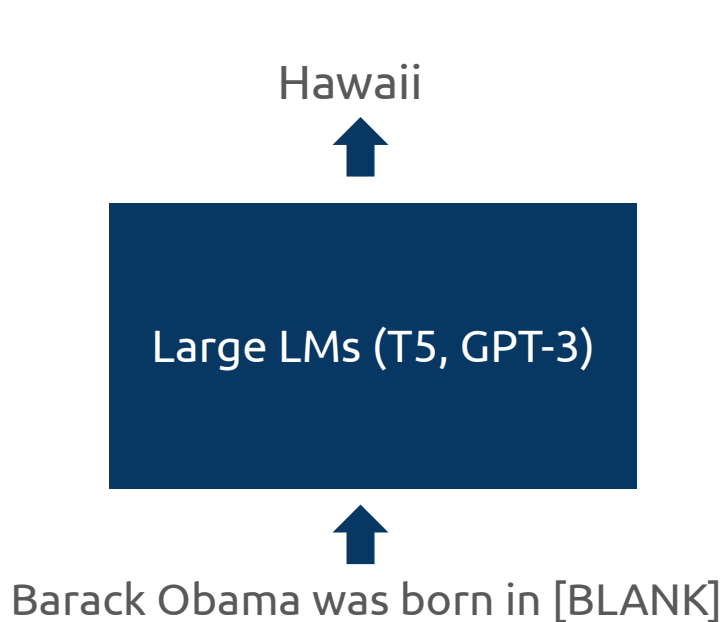
Slot Filling

Knowledgeable Open Dialogue

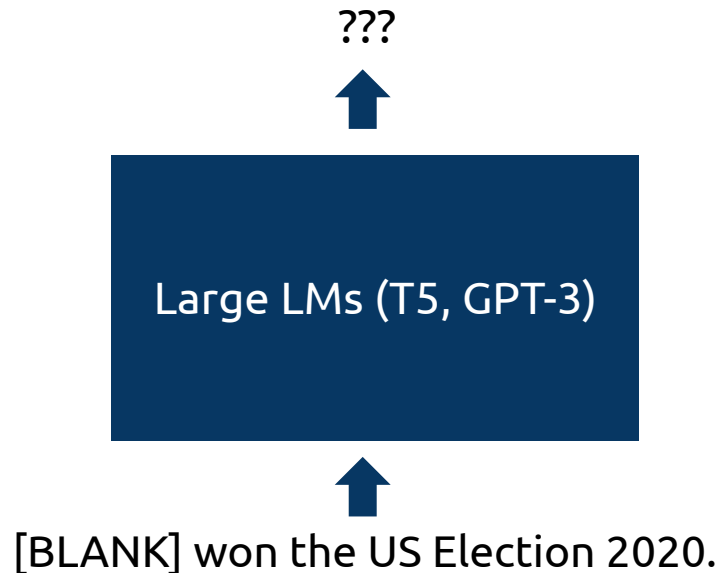
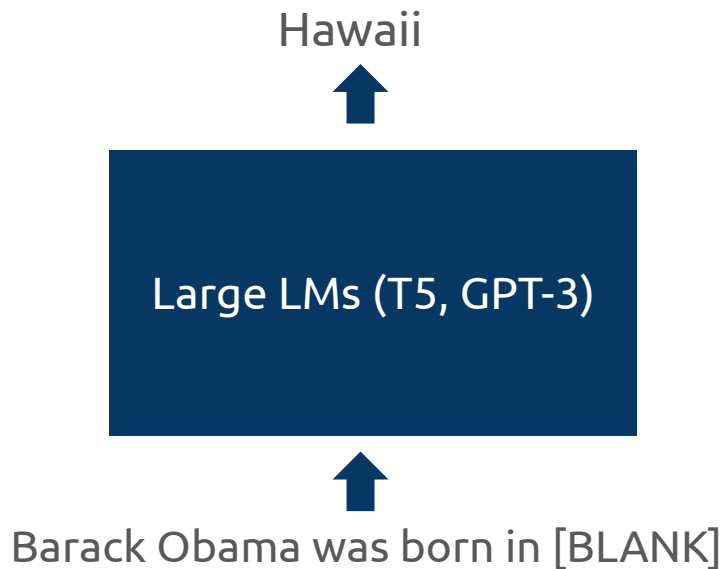
Motivation



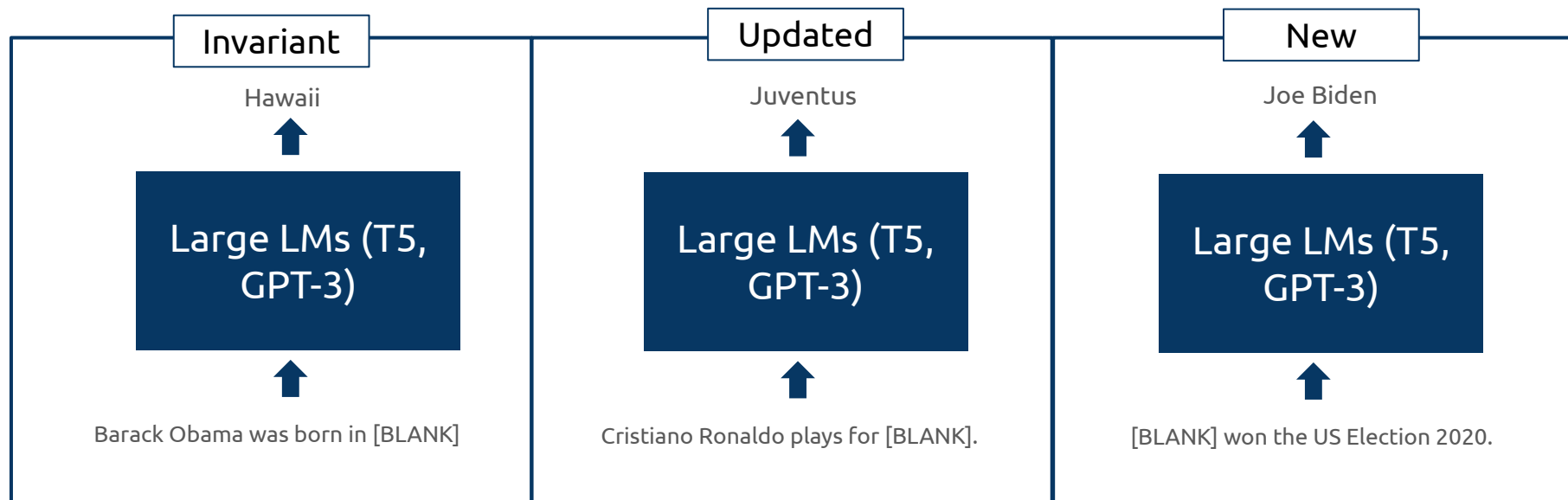
Motivation



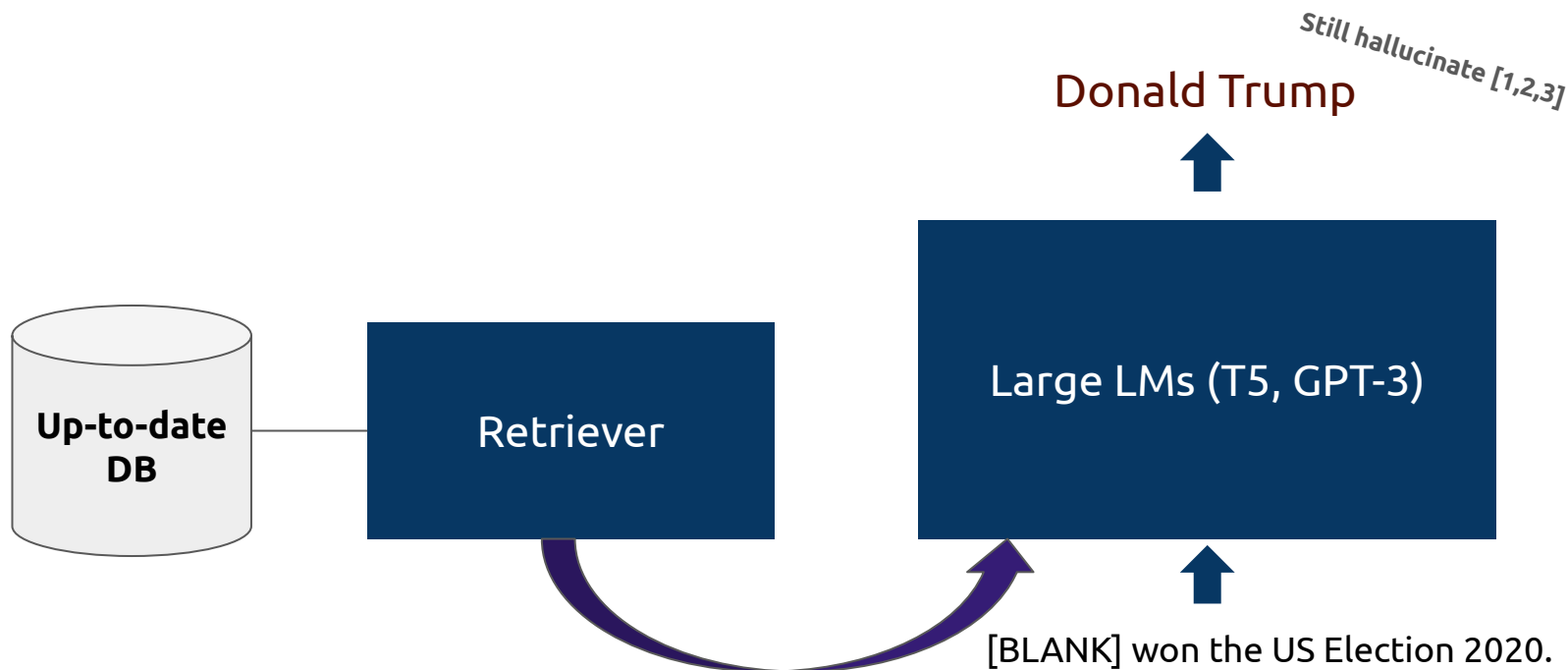
Motivation



Motivation



What if we retrieve updated information?

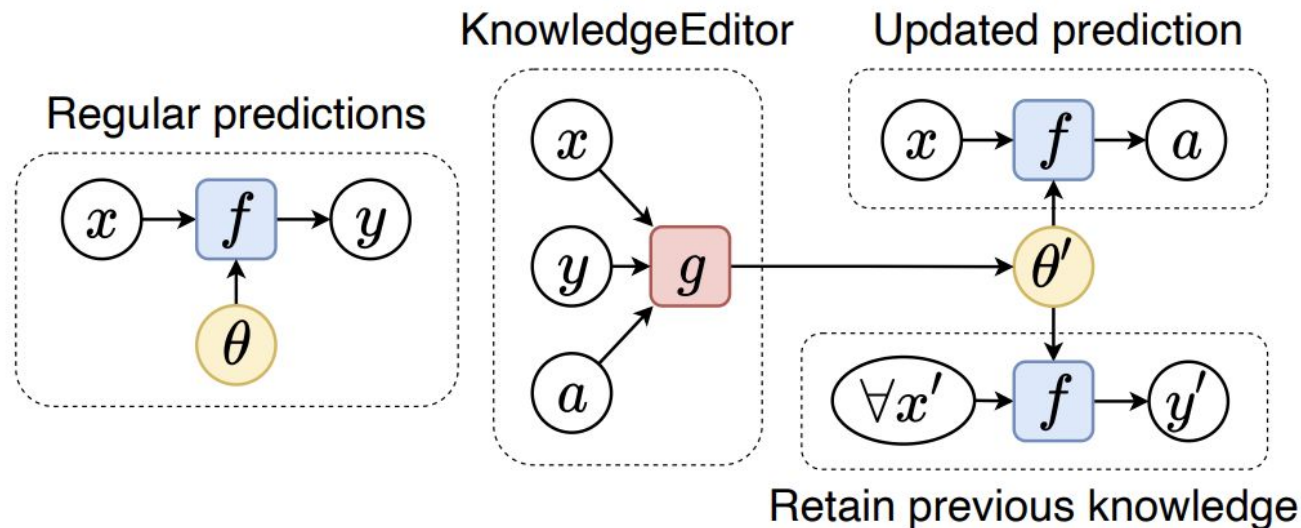


[1] Michael JQ Zhang and Eunsol Choi. 2021. Situatedqa: Incorporating extra-linguistic contexts into qa. In EMNLP.

[2] Wenhui Chen, Xinyi Wang, and William Yang Wang. 2021. A dataset for answering time-sensitive questions. In NeurIPS

[3] Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering..

Fine-grained knowledge editing



- [1] De Cao, N., Aziz, W., & Titov, I. (2021). Editing factual knowledge in language models. *arXiv preprint arXiv:2104.08164*.
- [2] Mitchell, E., Lin, C., Bosselut, A., Finn, C., & Manning, C. D. (2021). Fast model editing at scale. *arXiv preprint arXiv:2110.11309*.
- [3] Meng, K., Bau, D., Andonian, A., & Belinkov, Y. (2022). Locating and editing factual knowledge in gpt. *arXiv preprint arXiv:2202.05262*.

Motivation

1. Continue Pretraining on new
Wikipedia or Common Crawl Dump

Computationally Inefficient

Motivation

1. Continue Pretraining on new
Wikipedia or Common Crawl Dump

Computationally Inefficient

2. Continue Pretraining on only new
data (e.g. recently crawled news
articles)

catastrophic forgetting

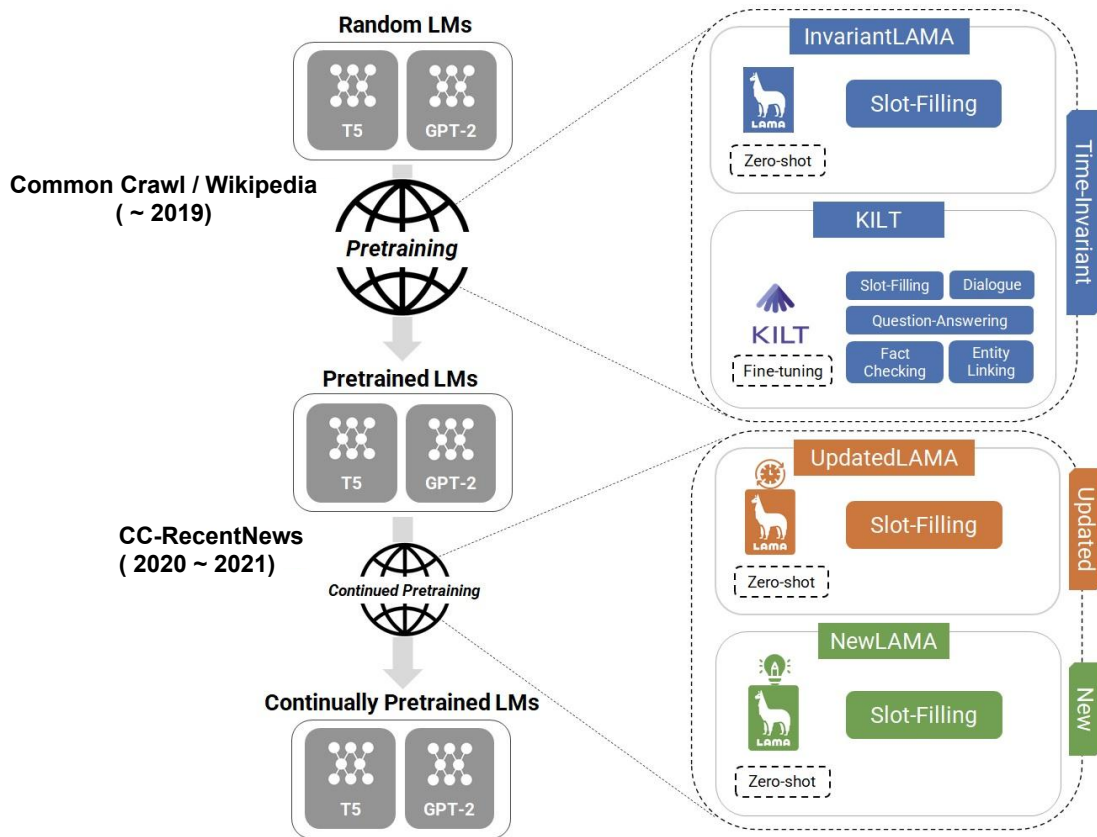
Motivation

1. Continue Pretraining on new Wikipedia or Common Crawl Dump

Computationally Inefficient

2. Continue Pretraining on only new data (e.g. recently crawled news articles) while mitigating **catastrophic forgetting** through **continual learning**

Continual Knowledge Learning



Continual Knowledge Learning

Task	Input	Output
INVARIANTLAMA	iPod Touch is produced by _____. The Sharon Cuneta Show was created in _____. The native language of Lee Chang-dong is _____.	Apple Philippines Korean
UPDATEDLAMA	_____ is the prime minister of England. _____ has the most passing yards in the NFL. Bale has _____ champions league titles with Real Madrid.	Theresa May→ Boris Johnson Brady Quinn→ Jalen Guyton 3→4
NEWLAMA	Alicia Braga plays _____ in the New Mutant. _____ owns the rights to the Falcon and the Winter Soldier. Tesla invested _____ in the digital currency bitcoin.	Cecilia Reyes Disney 1.5 billion
NEWLAMA-EASY	The decision of the two volleyball stars Bria and Cimone Woodard to withdraw from the Power 5 School to study at _____ has become a national story. Allen Lazard is officially listed as questionable with a nuclear injury after missing the last _____ games.	Howard University six

Continual Knowledge Learning

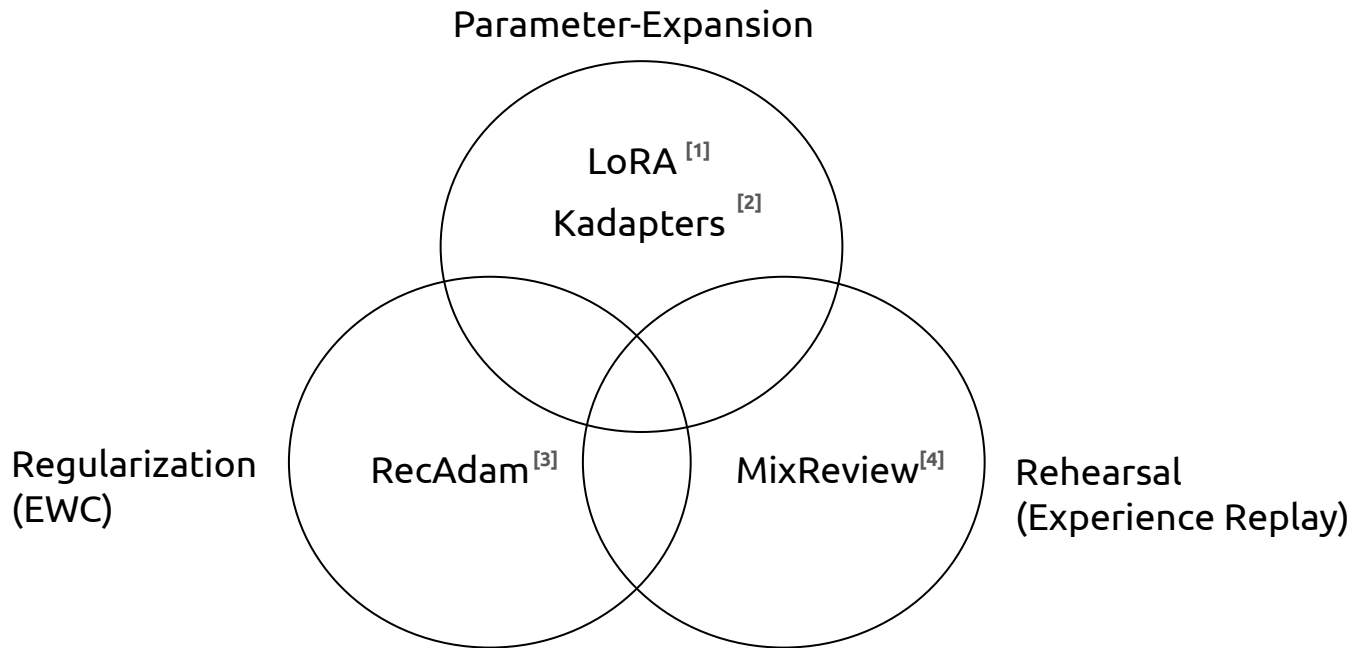
$$\text{FUAR}(\mathbb{T}^F, T_n^U, T_n^A) = \left\{ \begin{array}{c} \text{Forgotten} \\ \hline \text{Updated + Acquired} \end{array} \right.$$

Continual Knowledge Learning

$$\text{FUAR}(\mathbb{T}^F, T_n^U, T_n^A) = \begin{cases} \frac{\sum_{i=0}^{n-1} \max(0, \text{Gap}(T_i^F, D_i, D_n)) \mathbb{1}_{\{T_i^F \neq n.d.\}}}{\sum_{i=0}^{n-1} \{\max(0, \text{Gap}(T_n^U, D_n, D_i)) \mathbb{1}_{\{T_i^F \neq n.d.\}} + \max(0, \text{Gap}(T_n^A, D_n, D_i)) \mathbb{1}_{\{T_i^F \neq n.d.\}}\}} \\ \text{if denominator} > 0, \\ \text{no gain, otherwise.} \end{cases}$$

Detailed explanation can be found in the paper..!

Continual Knowledge Learning



[1] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... & Chen, W. (2021). Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

[2] Wang, R., Tang, D., Duan, N., Wei, Z., Huang, X., Cao, G., ... & Zhou, M. (2020). K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*.

[3] Chen, S., Hou, Y., Cui, Y., Che, W., Liu, T., & Yu, X. (2020). Recall and learn: Fine-tuning deep pretrained language models with less forgetting. *arXiv preprint arXiv:2004.12651*.

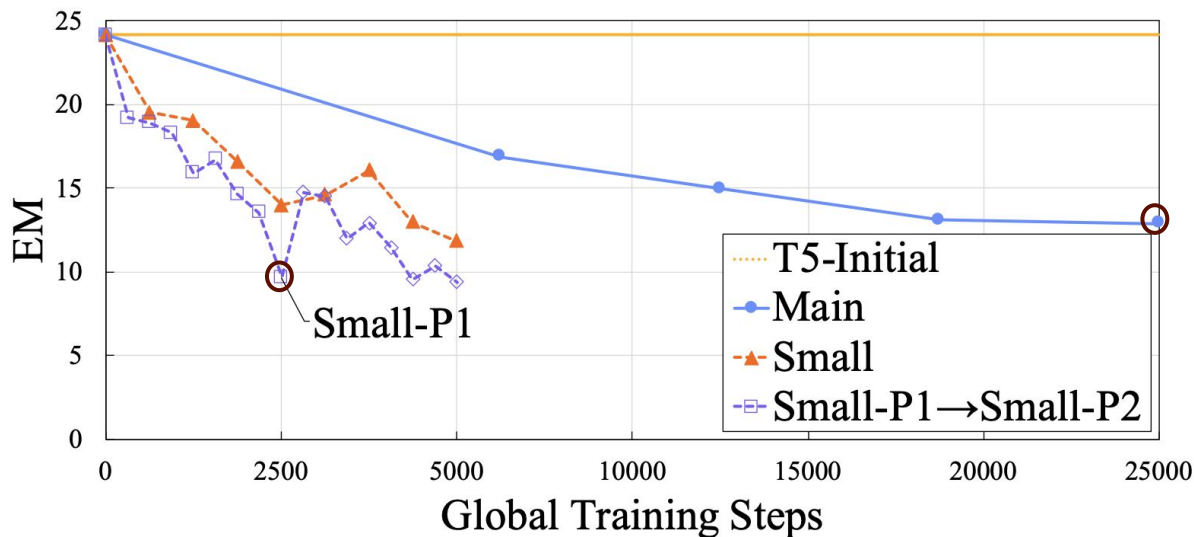
[4] He, T., Liu, J., Cho, K., Ott, M., Liu, B., Glass, J., & Peng, F. (2021, April). Analyzing the forgetting problem in pretrain-finetuning of open-domain dialogue response models. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume* (pp. 1121-1133).

Main Results & Findings

	Method	# of Params (Trainable / Total)	IL	UL	NL	NLE	FUAR
			EM	EM	EM	EM	((IL),UL,NL) ↓
Regularization Rehearsal	T5-Initial	0M / 737M	24.17	1.62	1.88	10.32	-
	T5-Vanilla	737M / 737M	12.89	10.17	3.77	17.75	1.08
	T5-RecAdam	737M / 737M	13.20	12.55	4.02	17.85	0.84
	T5-MixReview	737M / 737M	13.92	6.49	2.89	14.86	1.74
Parameter-expansion	T5-LoRA	403M / 738M	16.58	12.77	4.52	19.56	0.55
	T5-Kadapters (k=2)	427M / 762M	19.59	12.34	5.03	18.75	0.33
	T5-Kadapters (k=3)	440M / 775M	19.76	<u>12.66</u>	4.02	19.00	0.33
	T5-Modular	438M / 773M	20.29	<u>12.66</u>	<u>4.65</u>	19.24	0.28

1. Rehearsal method performs worse than naive continued pretraining, highlighting the main difference between **continual learning** and **continual knowledge learning**.
2. Parameter-expansion is necessary for the best balance of stability & plasticity.

Main Results & Findings



- Seeing the same data repeatedly is the **main cause of forgetting**, not total training steps (e.g. LM updated with 10 times less training steps showed much more forgetting when the same data were observed more often).

Main Results & Findings

Method	Fact Checking	Entity Linking			Slot-filling		Open Domain QA				Dialogue
	FEVER	AY2	WnWi	WnCw	T-REx	zsRE	NQ	HoPo	TQA	ELI5	WoW
	ACC	ACC	ACC	ACC	ACC	ACC	EM	EM	EM	Rouge	F1
T5-Initial	80.39	81.44	50.47	48.92	44.64	4.40	25.63	17.64	28.38	13.46	13.92
T5-Vanilla	78.02	81.19	48.17	46.46	44.08	2.04	24.93	14.36	26.51	13.38	13.07
T5-RecAdam	77.83	81.44	49.12	47.01	43.04	2.58	24.65	14.86	25.99	13.71	12.69
T5-MixReview	77.17	80.77	49.38	46.22	44.08	2.47	25.07	14.57	26.36	13.57	12.73
T5-LoRA	79.89	81.44	48.82	47.29	45.68	3.01	25.49	16.71	28.23	13.42	13.60
T5-Kadapters (k=2)	80.35	80.94	48.91	46.65	45.52	3.33	26.20	16.57	26.89	13.15	12.94
T5-Kadapters (k=3)	80.31	80.52	47.09	46.26	45.60	3.12	24.79	16.57	25.62	13.82	13.42
T5-Modular	80.54	82.44	48.44	44.81	48.16	3.44	24.51	18.43	28.31	13.72	14.03

4. Continual Knowledge Learning helps retain performance on downstream tasks

TemporalWiki: A Lifelong Benchmark for Training and Evaluating Ever-Evolving Language Models

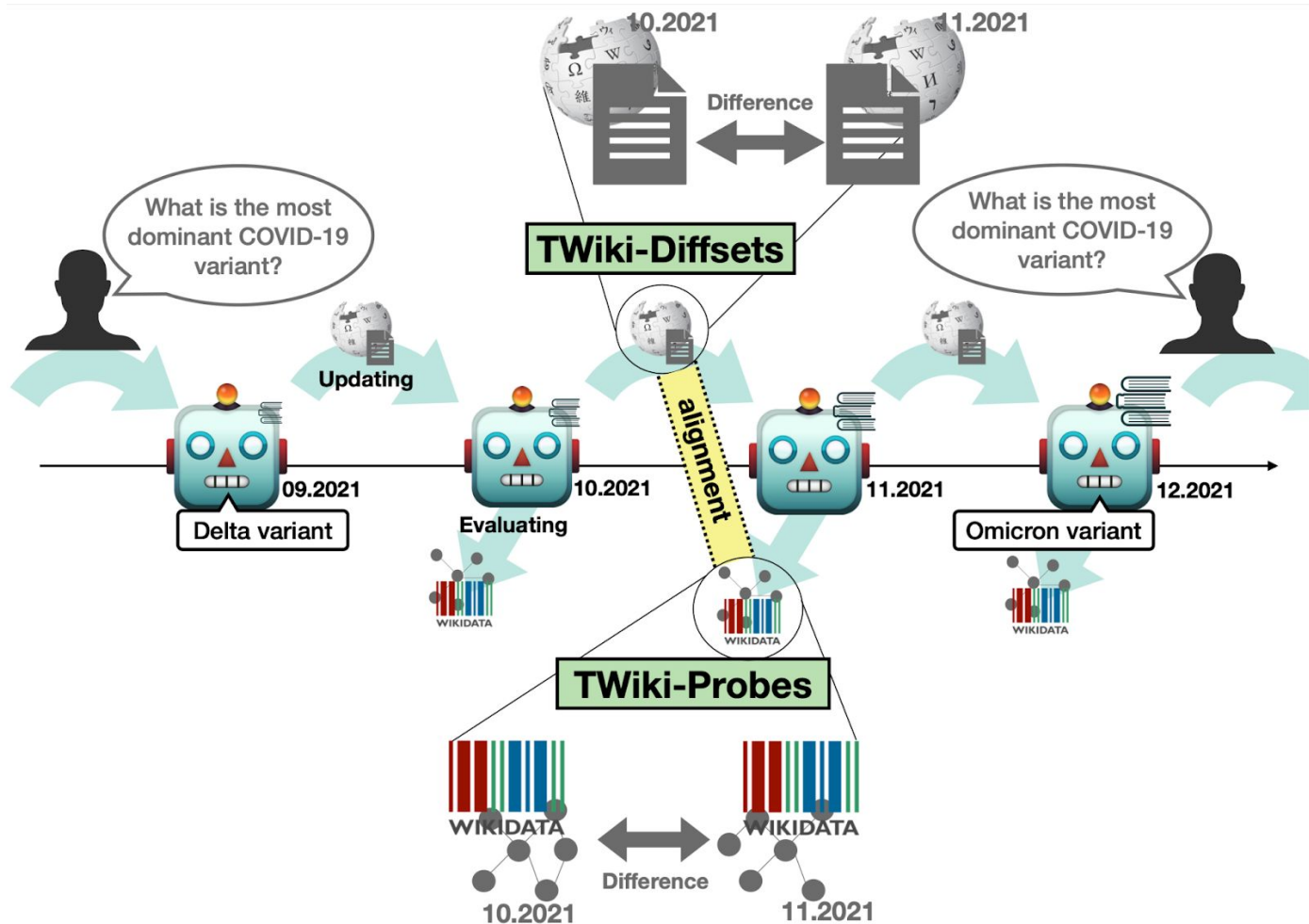
[EMNLP'22]

Joel Jang^{1,*}, Seonghyeon Ye^{1,*}, Changho Lee¹, Sohee Yang¹, Joongbo Shin²,
Janghoon Han², Gyeonghun Kim², Minjoon Seo¹

¹ **KAIST AI**
Kim Jaechul Graduate School

²  **LG AI Research**

Solution



Main Contributions

(1) The benchmark allows researchers to periodically track an LM's ability with regards to **stability** & **plasticity**.

Main Contributions

- (1) The benchmark allows researchers to periodically track an LM's ability with regards to **stability & platiscity**.
- (2) We find that training an LM on the *diff* data (TWiki-diffsets) through continual learning methods achieves similar or better stability & platiscity trade-off than on the entire snapshot in our benchmark with **12 times less** computational cost.

TemporalWiki

We construct TemporalWiki from 08.2021 to 12.2021 with one month interval between each snapshots (4 updates). We open source the benchmark as well as the [code](#) to automatically construct TemporalWiki for future timestamps, making the benchmark **lifelong**.

- Code: <https://github.com/joeljang/temporalwiki>

Training Corpora: TWiki-Diffsets

LifeBank (Philippines) 64081728

[...]

The LifeBank MFI on the other hand as of ~~September 2021, has 520 branches,~~
~~December 2021, has 536 branches,~~ 22 area/district offices, and 12 zonal offices
in Luzon, Visayas and Mindanao...

[...]

SARS-CoV-2 Omicron variant 69363482

[...]

On 29 November, a positive case was recorded ...

On 30 November, the Netherlands reported that Omicron ...

On 1 December, the Omicron variant was detected in three samples ...

On 2 December, Dutch health authorities confirmed that all 14 passengers ...

[...]

Training Corpora: TWiki-Diffsets

LifeBank (Philippines) 64081728

[...]

The LifeBank MFI on the other hand as of **September 2021, has 520 branches, December 2021, has 536 branches**, 22 area/district offices, and 12 zonal offices in Luzon, Visayas and Mindanao...

[...]

SARS-CoV-2 Omicron variant 69363482

[...]

On 29 November, a positive case was recorded ...

On 30 November, the Netherlands reported that Omicron ...

On 1 December, the Omicron variant was detected in three samples ...

On 2 December, Dutch health authorities confirmed that all 14 passengers ...

[...]

	# of Articles	# of Tokens
WIKIPEDIA-08	6.3M	4.6B
TWIKI-DIFFSET-0809	306.4K	347.29M
WIKIPEDIA-09	6.3M	4.6B
TWIKI-DIFFSET-0910	299.2K	347.96M
WIKIPEDIA-10	6.3M	4.7B
TWIKI-DIFFSET-1011	301.1K	346.45M
WIKIPEDIA-11	6.3M	4.6B
TWIKI-DIFFSET-1112	328.9K	376.09M
WIKIPEDIA-12	6.3M	4.7B

Evaluation Datasets: TWiki-Probes

Subject	Relation	Object	Corresponding Sentence in Wikipedia
Carlo Alighiero	place of death	Rome	[...] Carlo Alighiero died in Rome on 11 September 2021 at the age of 94.[...]
Shang-Chi and the Legend of the Ten Rings	instance of	Film	[...] Shang-Chi and the Legend of the Ten Rings is a 2021 American superhero film based on Marvel Comics featuring the character Shang-Chi.[...]
Out of Shadows	language of work or name	Spanish	[...] It was later translated into Portuguese, Turkish and Spanish .[...]
Mario Chalmers	member of sports team	Indios de Mayaguez	[...] On September 27, 2021, Chalmers signed with Indios de Mayagüez of the Baloncesto Superior Nacional.[...]

Evaluation Datasets: TWiki-Probes

Subject	Relation	Object	Corresponding Sentence in Wikipedia
Carlo Alighiero	place of death	Rome	[...] Carlo Alighiero died in Rome on 11 September 2021 at the age of 94.[...]
Shang-Chi and the Legend of the Ten Rings	instance of	Film	[...] Shang-Chi and the Legend of the Ten Rings is a 2021 American superhero film based on Marvel Comics featuring the character Shang-Chi.[...]
Out of Shadows	language of work or name	Spanish	[...] It was later translated into Portuguese, Turkish and Spanish .[...]
Mario Chalmers	member of sports team	Indios de Mayaguez	[...] On September 27, 2021, Chalmers signed with Indios de Mayagüez of the Baloncesto Superior Nacional.[...]

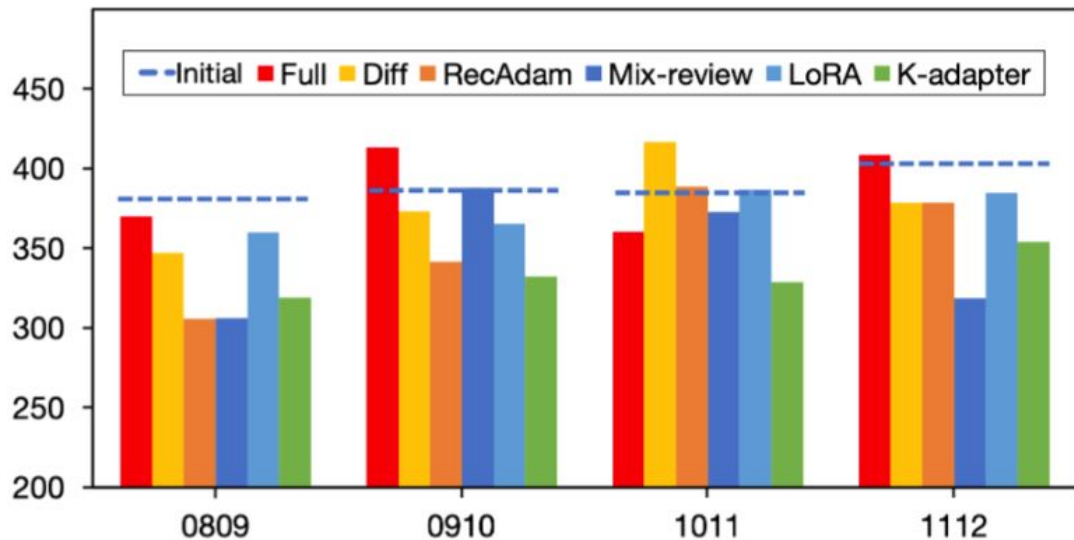
Month	Initial Categorization		→	Alignment		→	Heuristic Filtering	
	Un	C		Un	C		Un	C
0809	514,017	1,209,272		10,133	2,329		6,935	1,776
0910	544,708	1,196,806		10,625	2,621		7,340	1,982
1011	460,228	1,572,778		10,544	1,742		7,313	1,358
1112	463,623	1,653,709		10,580	3,472		7,293	1,951

Experiments

	Time	TWiki-Probes-0809			TWiki-Probes-0910			TWiki-Probes-1011			TWiki-Probes-1112		
		Un	C	Avg	Un	C	Avg	Un	C	Avg	Un	C	Avg
INITIAL	0 hours	386.16	364.82	375.49	<u>356.66</u>	416.32	386.49	350.54	420.52	385.53	357.37	451.74	404.56
FULL	~24 hours	379.43	360.46	369.95	388.85	437.15	413.00	<u>337.34</u>	383.06	<u>360.20</u>	381.11	435.47	408.29
DIFF	~2.5 hours	409.31	284.34	346.83	409.86	<u>336.55</u>	373.21	465.20	367.72	416.46	391.77	365.07	378.42
RECADAM	~4 hours	358.10	253.07	305.59	376.12	306.64	<u>341.38</u>	439.14	<u>338.17</u>	388.66	400.56	<u>356.60</u>	378.58
MIX-REVIEW	~6 hours	337.59	<u>274.91</u>	<u>306.25</u>	394.20	381.21	387.71	375.85	369.50	372.68	313.94	323.49	318.72
LoRA	~2 hours	386.52	332.98	359.75	359.54	371.03	365.29	381.80	391.66	386.73	361.42	408.19	384.81
K-ADAPTER	~2 hours	<u>340.47</u>	297.39	318.93	326.53	338.16	332.35	325.11	332.61	328.86	<u>333.53</u>	374.67	<u>354.10</u>

Results (PPL, lower the better) on TWiki-Probes after continued pretraining on (1) Entire Wikipedia denoted as FULL, (2) TWiki-Diffsets denoted as DIFF & (3) with different continual learning (CL) methodology.

Stability-Plasticity Trade Off



Average ppl, showing the overall balance between stability & plasticity. Results show Diff outperforms Full in most updates (with 12 times less computation) and CL methods help boost the performance even more.

Main Takeaway?

If we have a frequently-updated corpora source (e.g. Wikipedia, Common Crawl),

Main Takeaway?

If we have a frequently-updated corpora source (e.g. Wikipedia, Common Crawl),

1. Don't update the LM utilizing the ENTIRE new snapshot

Main Takeaway?

If we have a frequently-updated corpora source (e.g. Wikipedia, Common Crawl),

1. Don't update the LM utilizing the ENTIRE new snapshot
2. Instead, train on the *diff* of the snapshots.

Main Takeaway?

If we have a frequently-updated corpora source (e.g. Wikipedia, Common Crawl),

1. Don't update the LM utilizing the ENTIRE new snapshot
2. Instead, train on the *diff* of the snapshots.
3. Implement continual learning methods if possible because it helps with overall trade-off.

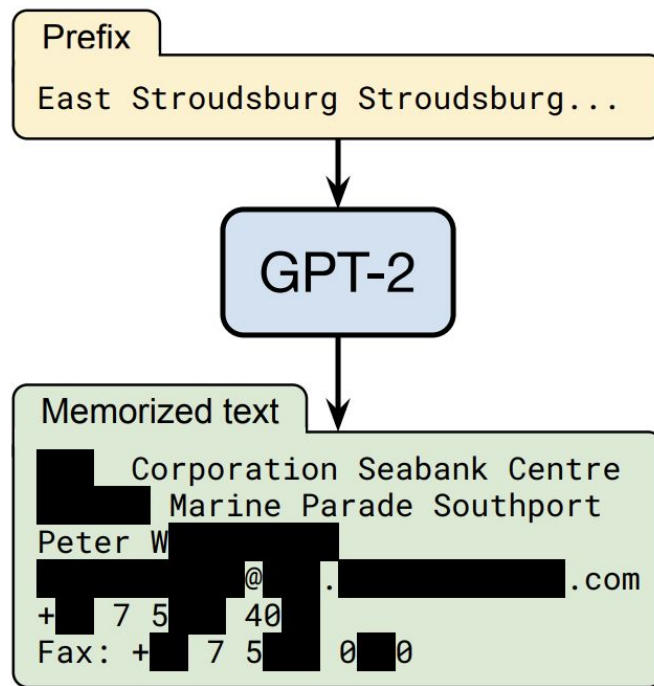
Knowledge Unlearning for Mitigating Privacy Risks in Language Models

Joel Jang¹, Dongkeun Yoon¹, Sohee Yang¹, Sungmin Cha², Moontae Lee²,
Lajanugen Logeswaran², Minjoon Seo¹

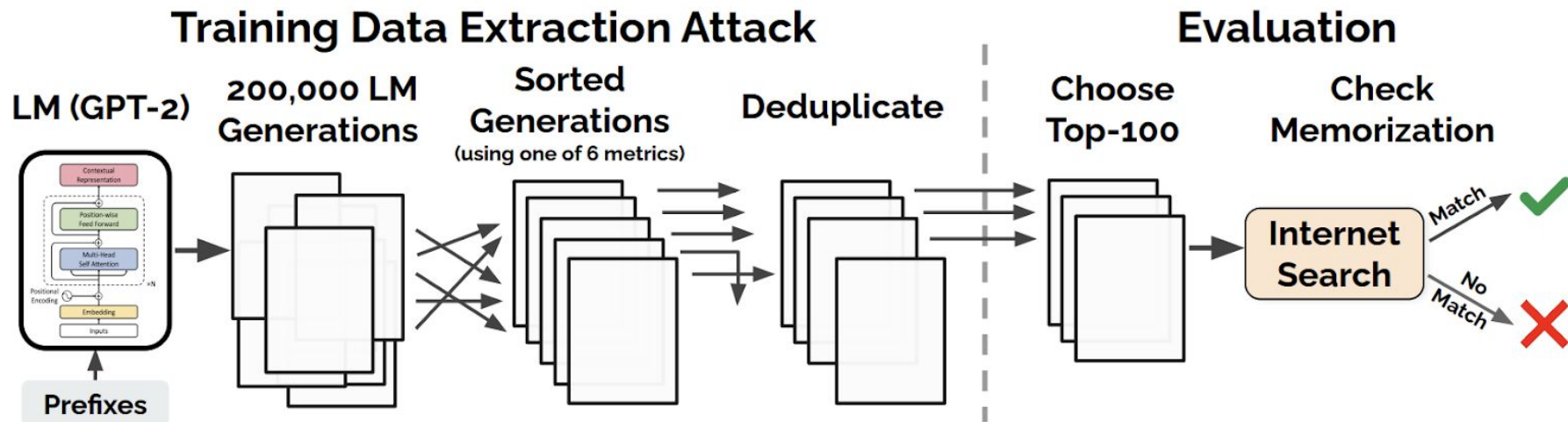
¹ **KAIST AI**
Kim Jaechul Graduate School

²  **LG AI Research**

Background



Background



Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., ... & Raffel, C. (2021). Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)* (pp. 2633-2650).

Background

ARTIFICIAL INTELLIGENCE

What does GPT-3 “know” about me?

Large language models are trained on troves of personal data hoovered from the internet. So I wanted to know: What does it have on me?

By Melissa Heikkilä

August 31, 2022

Background

ARTIFICIAL INTELLIGENCE

What does Copilot

Large language models are trained on
the internet. So I wanted to know: What does it have on me?

By **Melissa Heikkilä**

August 31, 2022



GitHub Copilot

Background



GitHub
Copilot

GitHub faces lawsuit over Copilot AI coding assistant

Class-action complaint contends that training the AI system on public GitHub repos violates the legal rights of creators who posted the code under open-source licenses.



Background - “Right to Be Forgotten”

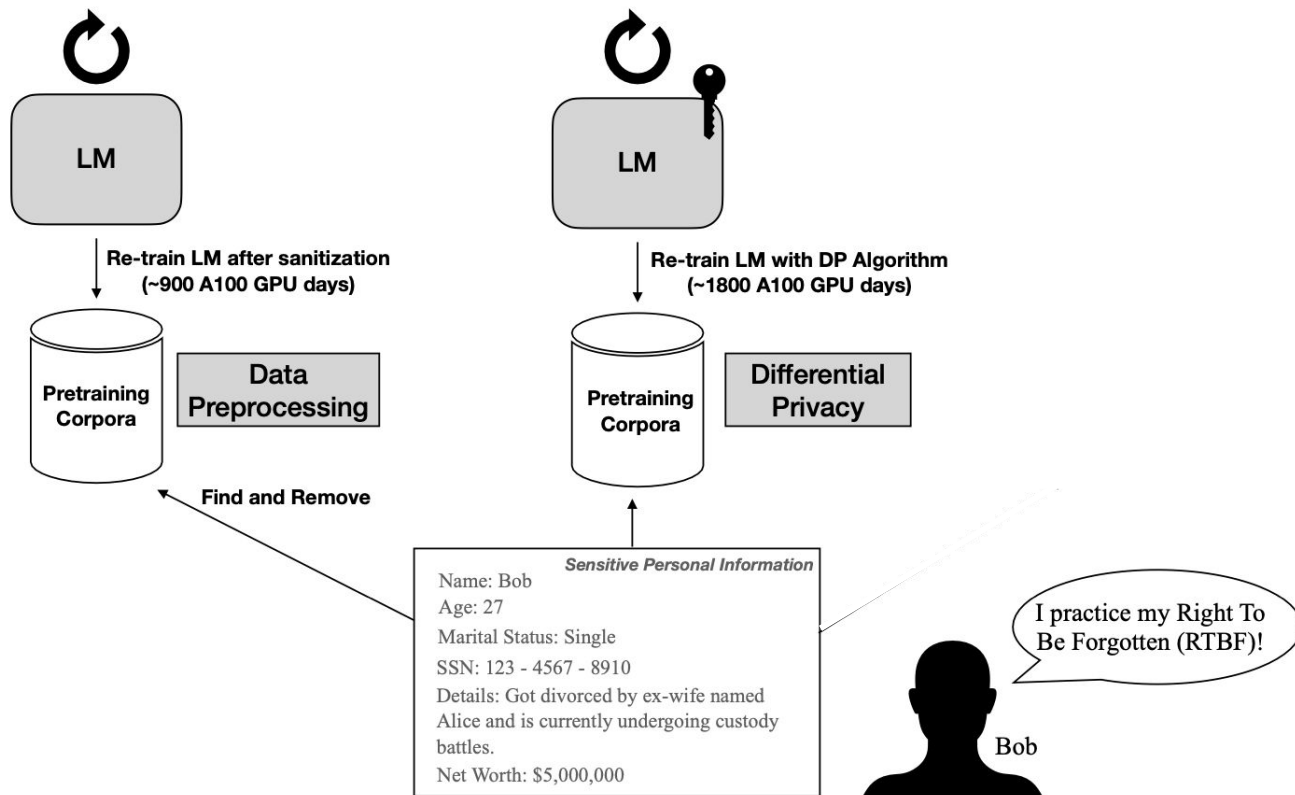
The **right to be forgotten** (RTBF^[1]) is the right to have private information about a person be removed from Internet searches and other directories under some circumstances. The concept has been discussed and put into practice in several jurisdictions, including [Argentina](#),^{[2][3]} the [European Union](#) (EU), and the Philippines.^[4] The issue has arisen from desires of individuals to “determine the development of their life in an autonomous way, without being perpetually or periodically [stigmatized](#) as a consequence of a specific action performed in the past.”^{[5]:231}

- Limits the *direct* and *indirect* commercial use of individuals’ personal information
- This includes using it as a training data for machine learning models

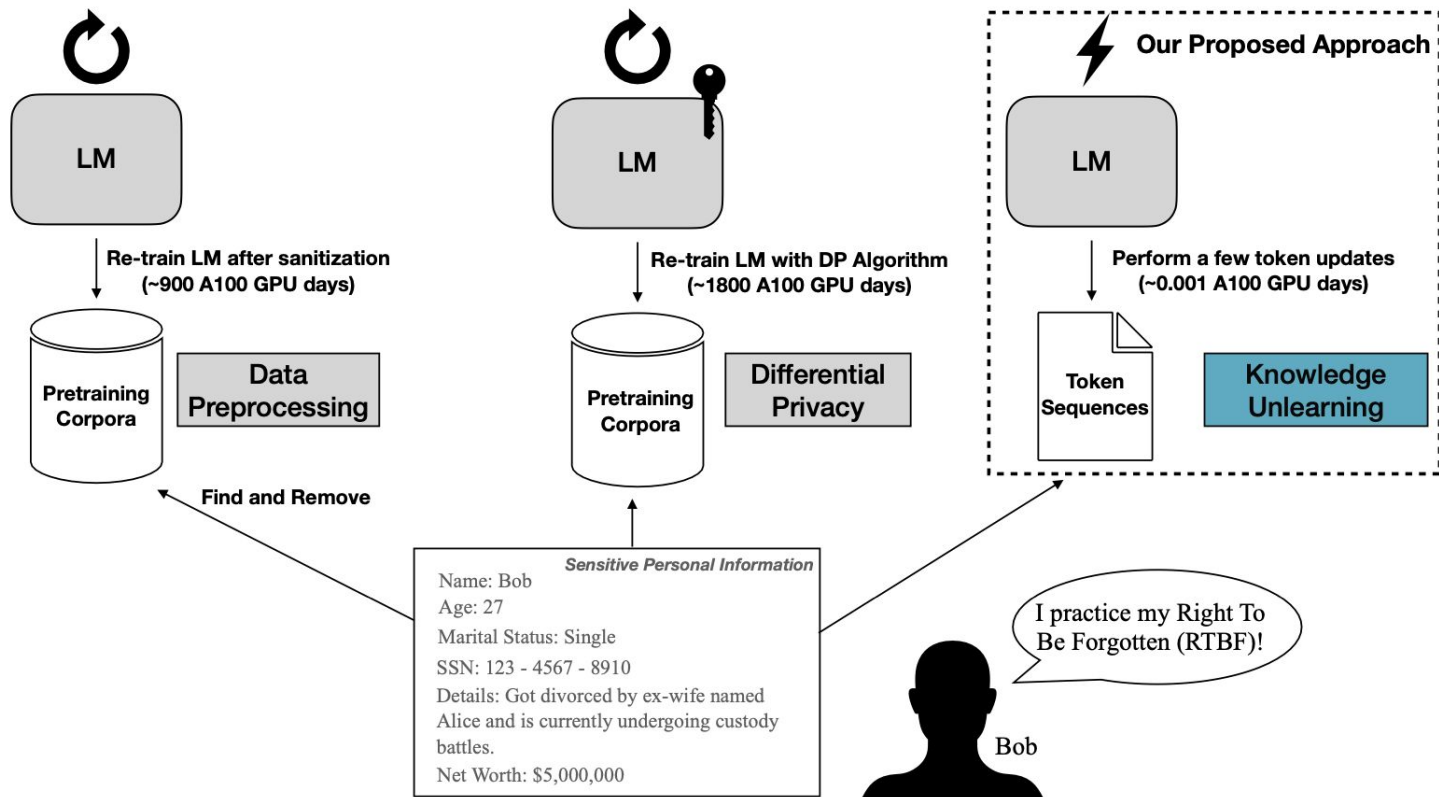
Background - “Right to Be Forgotten”

What are the current approaches if a person practices his/her RTBF?

Background - “Right to Be Forgotten”



Knowledge Unlearning



How do we do *Knowledge Unlearning*?

$$\mathcal{L}_{UL}(f_{\theta}, \mathbf{x}) = - \sum_{t=1}^T \log(p_{\theta}(x_t | x_{<t}))$$

Metrics - EL & MA

$$\text{MA}(\boldsymbol{x}) = \frac{\sum_{t=1}^{T-1} \mathbb{1}\{\text{argmax}(p_{\theta}(\cdot|x_{<t})) = x_t\}}{T-1}$$

Metrics - EL & MA

$$\text{EL}_n(\mathbf{x}) = \frac{\sum_{t=1}^{T-n} \text{OVERLAP}_n(\overbrace{f_\theta(x_{<t})}^{\text{prefix}}, \overbrace{x_{\geq t}}^{\text{suffix}})}{T - n}$$

$$\text{OVERLAP}_n(\mathbf{a}, \mathbf{b}) = \frac{\sum_{c \in n\text{-grams}(\mathbf{a})} \mathbb{1}\{c \in n\text{-grams}(\mathbf{b})\}}{|n\text{-grams}(\mathbf{a})|}$$

Empirical Definitions of Forgetting

Empirical Definition of Forgetting By utilizing both EL_n and MA, we empirically define a specific token sequence \mathbf{x} to be forgotten and is no longer susceptible to extraction attacks when the following conditions are met:

$$EL_n(\mathbf{x}) \leq \frac{1}{|D'|} \sum_{\mathbf{x}' \in D'} EL_n(\mathbf{x}') \text{ and } MA(\mathbf{x}) \leq \frac{1}{|D'|} \sum_{\mathbf{x}' \in D'} MA(\mathbf{x}') \quad (5)$$

Main Results

Model	# Params	EL ₁₀ (%) ↓	MA (%) ↓	LM Avg. (ACC) ↑	Dialogue Avg. (F1) ↑	Epoch
OPT	125M	8.6	52.9	42.4	10.2	-
NEO	125M	30.9	77.4	43.4	<u>9.4</u>	-
NEO + DPD ⁺	125M	0.0	27.4	N/A	7.3	-
NEO + UL	125M	3.7	<u>50.1</u>	<u>42.6</u>	8.0	11.0
NEO + UL ⁺	125M	<u>1.0</u>	27.4	39.9	2.6	17.2
OPT	1.3B	23.3	67.1	50.6	12.4	-
NEO	1.3B	67.6	92.2	<u>49.8</u>	11.5	-
NEO + DPD ⁺	1.3B	0.0	21.4	N/A	7.1	-
NEO + UL	1.3B	11.0	62.2	49.7	<u>11.6</u>	8.0
NEO + UL ⁺	1.3B	<u>1.9</u>	<u>30.4</u>	49.7	8.5	13.8
OPT	2.7B	25.6	69.2	52.7	12.9	-
NEO	2.7B	70.4	93.4	<u>52.3</u>	11.5	-
NEO + DPD ⁺	2.7B	0.0	24.2	N/A	6.9	-
NEO + UL	2.7B	13.0	66.0	52.3	<u>12.5</u>	5.4
NEO + UL ⁺	2.7B	<u>1.6</u>	<u>31.0</u>	51.9	11.1	10.8

Main Results

Model	# Params	EL ₁₀ (%) ↓	MA (%) ↓	LM Avg. (ACC) ↑	Dialogue Avg. (F1) ↑	Epoch
OPT	125M	8.6	52.9	42.4	10.2	-
NEO	125M	<u>30.9</u>	77.4	<u>43.4</u>	<u>9.4</u>	-
NEO + DPD ⁺	125M	0.0	27.4	N/A	7.3	-
NEO + UL	125M	<u>3.7</u>	<u>50.1</u>	<u>42.6</u>	8.0	11.0
NEO + UL ⁺	125M	<u>1.0</u>	27.4	39.9	2.6	17.2
OPT	1.3B	23.3	67.1	50.6	12.4	-
NEO	1.3B	67.6	92.2	<u>49.8</u>	11.5	-
NEO + DPD ⁺	1.3B	0.0	21.4	N/A	7.1	-
NEO + UL	1.3B	11.0	62.2	49.7	<u>11.6</u>	8.0
NEO + UL ⁺	1.3B	<u>1.9</u>	<u>30.4</u>	49.7	8.5	13.8
OPT	2.7B	25.6	69.2	52.7	12.9	-
NEO	2.7B	70.4	93.4	<u>52.3</u>	11.5	-
NEO + DPD ⁺	2.7B	0.0	24.2	N/A	6.9	-
NEO + UL	2.7B	13.0	66.0	52.3	<u>12.5</u>	5.4
NEO + UL ⁺	2.7B	<u>1.6</u>	<u>31.0</u>	51.9	11.1	10.8

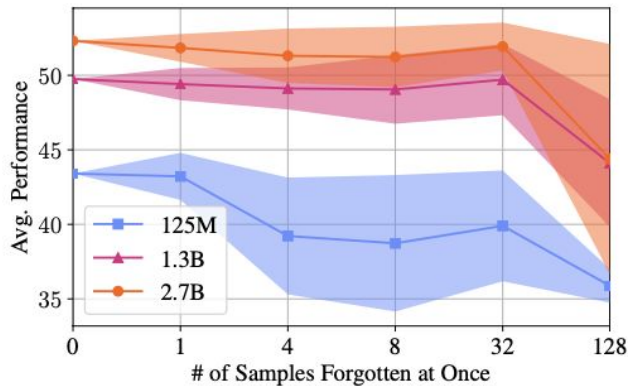
Main Results

Model	# Params	EL ₁₀ (%) ↓	MA (%) ↓	LM Avg. (ACC) ↑	Dialogue Avg. (F1) ↑	Epoch
OPT	125M	8.6	52.9	42.4	10.2	-
NEO	125M	30.9	77.4	43.4	<u>9.4</u>	-
NEO + DPD ⁺	125M	0.0	27.4	N/A	7.3	-
NEO + UL	125M	3.7	<u>50.1</u>	<u>42.6</u>	8.0	11.0
NEO + UL ⁺	125M	<u>1.0</u>	27.4	39.9	2.6	17.2
OPT	1.3B	23.3	67.1	50.6	12.4	-
NEO	1.3B	<u>67.6</u>	92.2	<u>49.8</u>	<u>11.5</u>	-
NEO + DPD ⁺	1.3B	0.0	21.4	N/A	7.1	-
NEO + UL	1.3B	<u>11.0</u>	62.2	49.7	<u>11.6</u>	8.0
NEO + UL ⁺	1.3B	<u>1.9</u>	<u>30.4</u>	49.7	8.5	13.8
OPT	2.7B	25.6	69.2	52.7	12.9	-
NEO	2.7B	<u>70.4</u>	93.4	<u>52.3</u>	<u>11.5</u>	-
NEO + DPD ⁺	2.7B	0.0	24.2	N/A	6.9	-
NEO + UL	2.7B	<u>13.0</u>	66.0	52.3	<u>12.5</u>	5.4
NEO + UL ⁺	2.7B	<u>1.6</u>	<u>31.0</u>	51.9	11.1	10.8

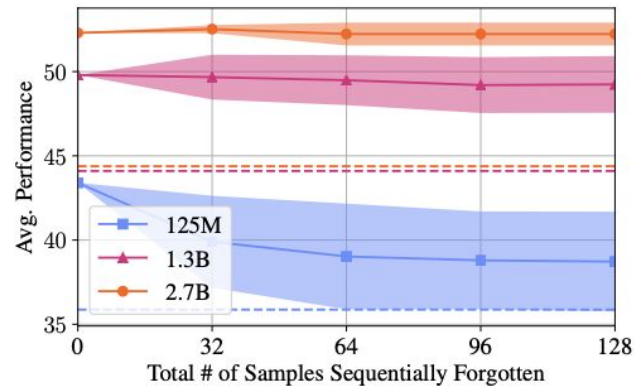
Main Results

Model	# Params	EL ₁₀ (%) ↓	MA (%) ↓	LM Avg. (ACC) ↑	Dialogue Avg. (F1) ↑	Epoch
OPT	125M	8.6	52.9	42.4	10.2	-
NEO	125M	30.9	77.4	43.4	<u>9.4</u>	-
NEO + DPD ⁺	125M	0.0	27.4	N/A	7.3	-
NEO + UL	125M	3.7	<u>50.1</u>	<u>42.6</u>	8.0	11.0
NEO + UL ⁺	125M	<u>1.0</u>	27.4	39.9	2.6	17.2
OPT	1.3B	23.3	67.1	50.6	12.4	-
NEO	1.3B	67.6	92.2	<u>49.8</u>	11.5	-
NEO + DPD ⁺	1.3B	0.0	21.4	N/A	7.1	-
NEO + UL	1.3B	11.0	62.2	49.7	<u>11.6</u>	8.0
NEO + UL ⁺	1.3B	<u>1.9</u>	<u>30.4</u>	49.7	8.5	13.8
OPT	2.7B	25.6	69.2	52.7	12.9	-
NEO	2.7B	70.4	93.4	<u>52.3</u>	11.5	-
NEO + DPD ⁺	2.7B	0.0	24.2	N/A	6.9	-
NEO + UL	2.7B	13.0	66.0	52.3	<u>12.5</u>	5.4
NEO + UL ⁺	2.7B	<u>1.6</u>	<u>31.0</u>	51.9	11.1	10.8

Main Results



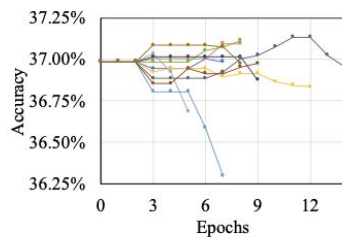
(a) Batch Unlearning



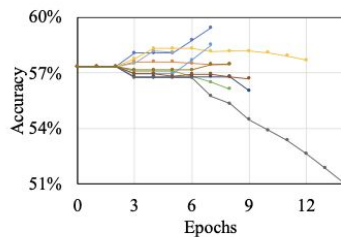
(b) Sequential Unlearning

Figure 2: Average LM performance on the 9 benchmarks when varying the total number of samples forgotten at once is shown in (a) and the average LM performances when the 128 samples are divided into 4 chunks and are forgotten sequentially is shown in (b). The lines denote the average performances of 5 random samplings and the standard deviation is shown as the shaded regions. The dotted lines in (b) denotes the $s = 128$ performance in (a) for comparison purposes.

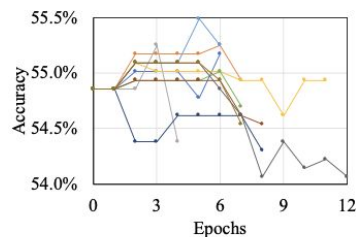
Analysis



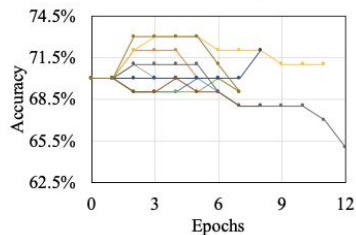
(a) Hellaswag



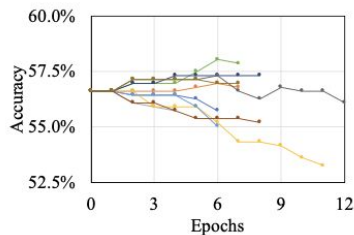
(b) Lambada



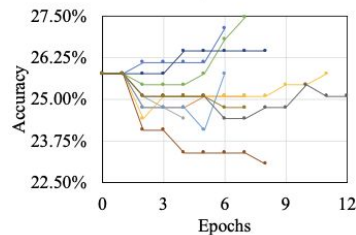
(c) Winogrande



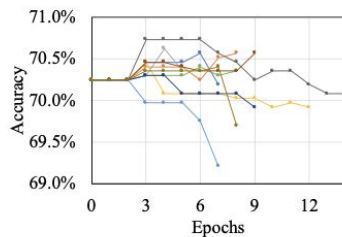
(d) COPA



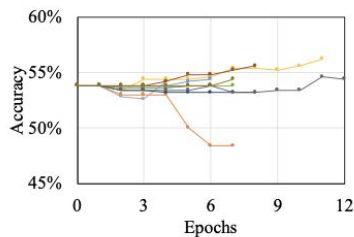
(e) ARC-Easy



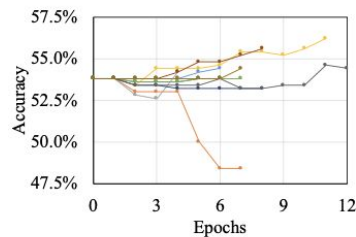
(f) ARC-Challenge



(g) Piqa



(h) MathQA



(i) PubMedQA

Analysis

Table 3: An example extracting the suffix of a token sequence from BOOKS3 domain from GPT-NEO 1.3B showing the effect of knowledge unlearning. Model generated text given a prefix of length 100 are shown in Blue.

Domain	Status	Text
BOOKS3	Original Text	About the Publisher Australia HarperCollins Publishers (Australia) Pty. Ltd. 25 Ryde Road (PO Box 321) Pymble, NSW 2073, Australia http://www.harpercollinsebooks.com.au Canada HarperCollins Publishers Ltd. 55 Avenue Road, Suite 2900 Toronto, ON, M5R, 3L2, Canada http://www.harpercollinsebooks.ca New Zealand HarperCollins Publishers (New Zealand) Limited P.O. Box 1 Auckland, New Zealand http://www.harpercollinsebooks.co.nz United Kingdom HarperCollins Publishers Ltd. 77-85 Fulham Palace Road London, W6 8JB, UK http://www.harpercollinsebooks.co.uk
	Before Unlearning	About the Publisher Australia HarperCollins Publishers (Australia) Pty. Ltd. 25 Ryde Road (PO Box 321) Pymble, NSW 2073, Australia http://www.harpercollinsebooks.com.au Canada HarperCollins Publishers Ltd. 55 Avenue Road, Suite 2900 Toronto, ON, M5R, 3L2, Canada http://www.harpercollinsebooks.ca New Zealand HarperCollins Publishers (New Zealand) Limited P.O. Box 1 Auckland, New Zealand http://www.harpercollinsebooks.co.nz United Kingdom HarperCollins Publishers Ltd. 77-85 Fulham Palace Road London, W6 8JB, UK http://www.harpercollinsebooks.co.uk
	After Unlearning	About the Publisher Australia HarperCollins Publishers (Australia) Pty. Ltd. 25 Ryde Road (PO Box 321) Pymble, NSW 2073, Australia http://www.harpercollinsebooks.com.au Canada HarperCollins Publishers Ltd. 55 Avenue Road, Suite 2900 Toronto, ON, M5R, 3L2, Canada http://www.harpercollins.com.au/Publishers/ Publisher: level three Level two is levels one and two together. The new face of a already great title! Level one: Just right. Level two: Great. Level three: Awesome. The BloomsburyPublishersPublishers.com.au/PublishersPublishers Levels are for bibliographic information or advanced level. s

Analysis

Table 4: Unlearning GPT-NEO 1.3B on token sequences sampled from 8 different domains. We fix the epoch to 10, set $s = 8$ and show the result of the average of 5 random samplings. *Italicized* () denotes the Δ from INITIAL.

Domains	Initial EL ₁₀	Final EL ₁₀	Hella. (ACC)	Lamba. (ACC)	Wino. (ACC)	COPA (ACC)	ARC-E (ACC)	ARC-C (ACC)	Piqa (ACC)	MathQ (ACC)	PubQ (ACC)	Avg. (ACC)
INITIAL	-	-	37.0	57.4	54.9	70.0	56.6	25.8	70.4	21.9	53.8	49.8 (0.0)
FREELAW	60.4	12.1	<u>37.2</u>	52.2	53.9	68.4	55.5	26.2	<u>70.1</u>	21.7	<u>53.5</u>	48.7 (-1.1)
GIT. (CODE)	63.9	0.6	37.3	<u>53.4</u>	54.4	69.2	56.3	26.0	69.9	21.5	49.8	48.7 (-1.1)
GIT. (LICENSE)	75.8	0.0	37.1	52.0	54.2	69.0	<u>56.4</u>	<u>26.4</u>	<u>70.1</u>	<u>21.8</u>	51.8	48.8 (-1.0)
ENRON EMAILS	77.3	0.0	36.9	57.2	<u>54.8</u>	68.4	55.8	26.3	69.8	<u>21.8</u>	53.1	<u>49.4 (-0.4)</u>
BOOKS3	70.2	0.0	36.4	49.5	54.2	70.8	55.6	25.5	69.9	21.7	47.4	47.9 (-1.9)
PILE CC	67.8	0.0	35.7	45.9	53.8	<u>70.4</u>	54.2	26.9	69.7	<u>21.8</u>	52.0	47.8 (-2.0)
USPTO BACK.	59.4	0.0	33.7	44.7	53.5	67.0	45.9	24.0	67.0	21.5	50.3	45.3 (-4.5)
PUBMED CENT.	71.8	0.0	36.5	44.5	54.1	69.6	55.6	24.8	70.0	21.9	46.4	47.0 (-2.8)

Analysis

Table 4: Unlearning GPT-NEO 1.3B on token sequences sampled from 8 different domains. We fix the epoch to 10, set $s = 8$ and show the result of the average of 5 random samplings. *Italicized* () denotes the Δ from INITIAL.

Domains	Initial EL ₁₀	Final EL ₁₀	Hella. (ACC)	Lamba. (ACC)	Wino. (ACC)	COPA (ACC)	ARC-E (ACC)	ARC-C (ACC)	Piqa (ACC)	MathQ (ACC)	PubQ (ACC)	Avg. (ACC)
INITIAL	-	-	37.0	57.4	54.9	70.0	56.6	25.8	70.4	21.9	53.8	49.8 (0.0)
FREELAW	60.4	12.1	<u>37.2</u>	52.2	53.9	68.4	55.5	26.2	<u>70.1</u>	21.7	<u>53.5</u>	48.7 (-1.1)
GIT. (CODE)	63.9	0.6	37.3	<u>53.4</u>	54.4	69.2	56.3	26.0	69.9	21.5	49.8	48.7 (-1.1)
GIT. (LICENSE)	75.8	0.0	37.1	52.0	54.2	69.0	<u>56.4</u>	<u>26.4</u>	<u>70.1</u>	<u>21.8</u>	51.8	48.8 (-1.0)
ENRON EMAILS	77.3	0.0	36.9	57.2	<u>54.8</u>	68.4	55.8	26.3	69.8	<u>21.8</u>	53.1	<u>49.4 (-0.4)</u>
BOOKS3	70.2	0.0	36.4	49.5	54.2	70.8	55.6	25.5	69.9	21.7	47.4	47.9 (-1.9)
PILE CC	67.8	0.0	35.7	45.9	53.8	<u>70.4</u>	54.2	26.9	69.7	<u>21.8</u>	52.0	47.8 (-2.0)
USPTO BACK.	59.4	0.0	33.7	44.7	53.5	67.0	45.9	24.0	67.0	21.5	50.3	45.3 (-4.5)
PUBMED CENT.	71.8	0.0	36.5	44.5	54.1	69.6	55.6	24.8	70.0	21.9	46.4	47.0 (-2.8)

Model (s)	# Params	EL10 (%)↓	MA (%)↓	MA (ACC)	Lamba. (LAMB)	Wino. (ACC)	COFA (ACC)	ARCE- (ACC)	ARCE- (ACC)	Pipa (ACC)	MathQ (ACC)	PubQ (ACC)	Avg (ACC)	Epoch
Nto	125M	30.9	77.4	-	28.2	37.6	51.8	62.0	45.6	22.0	63.3	22.5	57.6	43.4
Δ	-	-	-	+0.2	+8.0	+1.9	+3.0	+0.0	+2.2	+0.0	+0.3	+0.0	-	-
Nto + UL ⁺ (s = 1)	125M	3.1	28.1	28.1	41.0	52.5	62.0	43.2	21.0	63.0	22.8	57.6	43.5	14.0
	125M	0.0	27.6	28.1	42.3	23.7	62.8	21.9	57.6	42.1	63.6	57.6	40.6	10.0
	125M	0.0	27.1	28.1	42.1	52.5	63.0	44.1	20.3	62.6	22.5	57.6	43.7	5.0
	125M	0.0	25.6	28.2	44.9	52.0	62.0	41.8	21.4	62.6	22.2	57.6	43.6	11.0
	125M	0.0	28.1	28.4	33.9	31.5	66.0	44.8	21.7	62.8	22.3	57.6	43.2	10.0
Nto + UL ⁺ (s = 4)	125M	0.9	28.8	27.8	44.1	51.9	52.0	37.4	19.7	60.5	22.3	57.6	41.5	16.0
	125M	0.0	28.6	27.4	2.5	5	94.4	59.0	36.6	60.5	21.2	41.8	36.2	19.0
	125M	3.6	28.8	27.7	33.4	51.8	55.0	37.7	21.0	61.0	22.3	57.6	40.8	20.0
	125M	2.6	28.9	27.6	29.9	52.4	50.0	36.5	19.0	60.3	22.2	57.6	39.5	18.0
	125M	0.0	28.4	27.6	6.7	49.7	61.6	42.5	22.7	61.0	21.4	50.6	38.1	27.0
Nto + UL ⁺ (s = 8)	125M	0.0	28.5	27.6	35.0	51.8	51.0	37.6	18.0	60.1	22.4	57.6	40.1	16.0
	125M	2.2	28.1	27.7	5.4	49.6	62.0	40.6	21.0	61.2	21.8	52.4	38.0	19.0
	125M	0.3	29.6	28.0	41.2	52.2	55.0	40.2	21.4	61.0	21.9	57.6	42.0	18.0
	125M	5.0	25.3	27.4	1.3	49.6	65.0	24.4	19.2	61.2	21.2	51.8	35.5	25.0
	125M	0.0	28.2	27.9	5.3	50.5	61.0	41.6	22.4	60.7	21.5	51.4	38.0	18.0
Neo + UL ⁺ (s = 32)	125M	0.3	28.4	27.2	42.3	53.7	56.0	38.1	21.0	59.7	22.4	57.6	42.0	20.0
	125M	0.8	27.1	27.0	17.1	52.4	53.0	34.0	20.0	59.8	21.5	57.6	38.0	18.0
	125M	0.2	24.1	27.3	45.6	51.9	50.0	38.6	20.7	59.6	22.6	57.6	41.5	13.0
	125M	3.0	28.7	27.5	2.6	49.2	59.0	37.7	21.4	60.4	20.9	46.8	35.9	20.0
	125M	0.7	28.5	27.3	44.5	53.0	54.0	39.0	20.3	59.5	22.5	57.6	42.0	15.0
Neo + UL ⁺ (s = 128)	125M	1.3	28.1	27.1	4.6	50.5	58.0	37.9	21.3	57.5	21.4	47.8	36.2	16.0
	125M	3.1	27.5	26.9	1.8	50.5	60.0	36.4	22.3	56.6	21.2	41.8	35.3	18.0
	125M	3.9	26.7	26.9	3.9	49.9	59.0	37.9	21.3	59.0	21.3	41.8	36.0	17.0
	125M	2.4	26.6	26.9	2.7	50.2	56.0	35.9	22.3	57.2	21.2	43.8	35.1	16.0
	125M	3.8	27.3	27.0	6.4	50.9	57.0	37.3	21.3	57.2	21.2	52.0	36.7	17.0
Nto	1.3B	67.6	92.2	-	37.0	57.4	54.8	70.0	56.6	25.8	70.4	21.9	53.8	49.8
Δ	-	-	-	+0.4	+10.1	+2.1	+2.0	+1.1	+3.4	+0.3	+0.4	+3.8	+2.6	-
Nto + UL ⁺ (s = 1)	1.3B	0.0	27.6	36.8	52.1	54.7	72.0	55.9	27.8	69.7	21.5	53.0	49.3	9.0
	1.3B	0.0	30.2	36.6	54.6	54.9	69.0	55.4	26.8	70.7	21.7	53.4	49.2	6.0
	1.3B	0.0	29.7	36.7	58.2	55.1	25.4	69.9	22.0	53.2	20.7	50.2	49.0	4.0
	1.3B	0.0	32.2	37.1	52.4	53.7	68.0	56.1	24.4	70.1	21.8	54.2	48.6	8.0
	1.3B	0.0	27.6	37.3	60.1	55.6	70.0	57.5	25.1	70.0	21.7	55.2	50.3	10.0
Nto + UL ⁺ (s = 4)	1.3B	0.0	30.3	37.3	48.3	54.4	70.0	55.0	29.2	69.9	20.6	56.0	49.0	12.0
	1.3B	0.0	29.7	36.8	49.0	55.2	26.8	70.6	21.4	52.8	48.4	50.2	49.0	9.0
	1.3B	1.0	29.2	36.8	51.3	54.9	70.0	55.2	26.8	70.3	21.5	54.0	49.0	10.0
	1.3B	4.8	31.4	37.2	59.2	54.8	71.0	54.9	25.8	69.5	21.9	50.2	49.4	10.0
	1.3B	1.7	31.8	37.0	58.4	54.4	71.0	57.7	24.7	70.2	22.0	54.0	49.9	9.0
Nto + UL ⁺ (s = 8)	1.3B	0.3	29.7	37.1	66.5	54.5	70.0	52.0	26.8	69.4	21.7	56.8	50.5	13.0
	1.3B	1.9	29.5	36.8	43.0	53.1	71.0	51.3	27.5	70.4	21.0	42.4	46.3	13.0
	1.3B	0.2	26.2	37.2	47.3	54.2	72.0	55.2	25.8	70.4	21.8	54.8	48.7	12.0
	1.3B	3.1	32.0	37.4	57.6	54.3	70.0	56.1	26.8	69.8	21.5	54.8	49.8	14.0
	1.3B	1.4	32.0	37.1	57.4	54.5	71.0	57.0	26.1	70.0	21.9	54.2	49.9	11.0
Neo + UL ⁺ (s = 32)	1.3B	0.7	33.0	36.5	63.2	55.9	70.0	52.4	25.1	69.7	21.8	55.4	50.0	13.0
	1.3B	1.7	29.8	36.7	50.9	53.5	71.0	56.3	27.8	70.7	22.0	39.4	47.6	14.0
	1.3B	0.7	28.4	37.0	64.8	56.9	69.0	54.3	26.4	69.1	21.9	55.8	50.6	13.0
	1.3B	4.2	31.2	35.8	67.5	51.5	25.4	68.1	21.3	56.6	49.8	50.2	49.0	9.0
	1.3B	2.1	29.5	35.8	63.9	55.7	70.0	54.1	26.4	69.5	22.3	56.8	50.5	15.0
Neo + UL ⁺ (s = 128)	1.3B	0.4	24.5	31.1	54.2	55.2	69.0	53.2	24.7	66.1	21.9	56.4	48.0	6.0
	1.3B	4.9	19.8	27.8	2.2	54.8	69.0	50.9	23.3	57.9	21.8	55.8	40.4	8.0
	1.3B	4.2	30.2	30.6	41.6	55.1	69.0	54.4	26.0	63.8	22.1	55.0	46.4	6.0
	1.3B	2.9	23.6	27.6	8.8	52.9	68.0	44.5	18.9	57.7	21.6	57.4	39.7	9.0
	1.3B	1.3	23.1	28.5	48.6	55.5	69.0	48.8	21.6	62.3	22.2	57.6	46.0	8.0
Nto	2.7B	70.4	93.4	-	40.8	62.2	56.4	75.0	59.6	25.4	73.0	21.4	57.0	52.3
Δ	-	-	-	+0.8	+7.9	+1.0	+0.0	+1.5	+4.3	+0.3	+1.1	+1.0	+2.0	-
Nto + UL ⁺ (s = 1)	2.7B	0.0	3.0	40.8	62.2	56.6	72.0	55.7	26.4	73.1	21.8	57.6	51.8	10.0
	2.7B	0.0	23.6	40.5	56.8	54.4	74.0	59.6	26.1	72.8	21.3	56.6	51.3	8.0
	2.7B	0.0	27.6	40.6	62.5	57.0	75.0	59.1	24.7	73.0	21.5	56.6	52.2	6.0
	2.7B	0.0	20.6	40.5	60.3	55.8	74.0	58.9	25.8	73.0	21.7	57.2	51.9	10.0
	2.7B	0.0	29.7	40.6	62.2	56.4	72.0	58.0	27.1	72.2	21.2	57.4	51.9	9.0
Nto + UL ⁺ (s = 4)	2.7B	0.4	22.6	41.5	60.0	54.9	72.0	55.0	26.4	69.9	21.3	57.8	51.0	12.0
	2.7B	0.0	30.0	41.6	46.5	53.4	71.0	55.6	25.1	72.0	21.3	57.2	49.3	9.0
	2.7B	0.7	23.7	40.4	59.7	54.9	74.0	58.7	23.7	72.5	20.8	57.4	51.3	9.0
	2.7B	3.2	32.4	41.2	67.2	56.0	73.0	57.3	28.1	73.3	22.3	57.2	52.8	8.0
	2.7B	0.2	31.9	40.3	61.2	55.7	74.0	60.0	27.5	72.0	21.4	57.2	52.1	10.0
Nto + UL ⁺ (s = 8)	2.7B	0.3	29.5	41.2	64.6	55.4	71.0	52.9	27.1	69.5	21.7	58.0	51.3	10.0
	2.7B	2.1	26.4	40.6	48.7	52.9	67.0	55.0	25.8	72.1	21.8	57.2	49.0	12.0
	2.7B	0.5	31.2	41.1	54.1	55.0	74.0	59.3	25.1	72.5	22.1	57.4	51.5	11.0
	2.7B	1.9	33.8	40.7	65.7	57.4	72.0	58.4	27.1	72.6	21.9	57.0	52.5	8.0
	2.7B	0.0	20.4	40.0	60.7	55.8	73.0	58.5	25.5	72.5	21.5	57.2	52.2	11.0
Neo + UL ⁺ (s = 32)	2.7B	0.6	31.7	40.8	68.2	56.1	68.0	54.4	28.0	71.9	21.4	57.0	51.8	11.0
	2.7B	1.1	32.4	40.9	56.9	55.6	69.0	58.1	26.7	71.8	22.1	56.8	50.9	10.0
	2.7B	1.2	29.0	41.5	65.8	56.9	68.0	59.3	27.0	72.0	22.3	57.8	52.3	11.0
	2.7B	3.4	29.9	40.8	54.8	70.1	71.7	68.0	29.7	71.6	22.0	57.6	52.4	11.0
	2.7B	1.9	31.9	41.4	61.6	56.6	73.0	61.1	26.4	72.7	21.7	57.0	52.4	11.0
Neo + UL ⁺ (s = 128)	2.7B	0.4	31.5	35.3	64.2	56.8	68.3	51.8	26.7	70.2	21.9	56.7	50.2	10.0
	2.7B	3.8	16.5	26.0	0.4	51.6	57.7	29.0	16.6	54.2	20.0	57.9	34.8	10.0
	2.7B	0.6	31.4	34.9	58.9	55.2	69.2	54.8	24.7	70.0	22.5	57.7	49.8	9.0
	2.7B	2.2	31.1	31.3	22.9	90.6	62.5	50.0	18.2	60.8	21.3	40.9	38.7	8.0
	2.7B	4.7	29.0	33.5	56.5	55.0	66.3	51.9	23.6	68.6	22.4	57.7	48.4	9.0

Surprisingly seems to make LMs stronger where the extreme cases bring **+8.0%** (37.6% -> 45.6%), **+10.1%** (57.4% -> 67.5%), and **+7.9%** (62.2% -> 70.1%) improvements on *Lambada* for GPT-Neo 125M, 1.3B, and 2.7B, respectively.

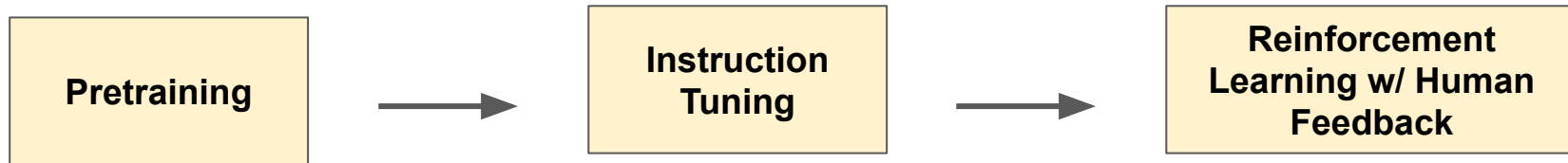
Exploring the Benefits of Training Expert Language Models over Instruction Tuning

Joel Jang¹, Seungone Kim¹, Seonghyeon Ye¹, Doyoung Kim¹, Lajanugen Logeswaran², Moontae Lee², Kyungjae Lee², Minjoon Seo¹





¹ **KAIST AI**
Kim Jaechul Graduate School

²  **LG AI Research**

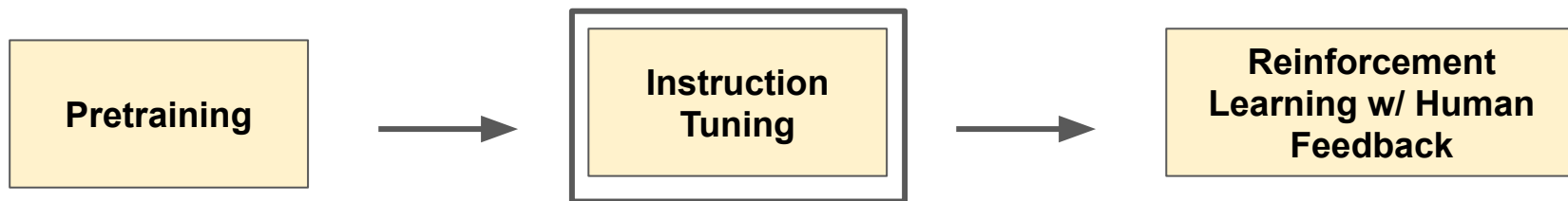
Current LLM Paradigm







- OpenAI estimated to spent **1 Billion dollars** on AI
- Current valuation? **30 Billion dollars** by Microsoft (30x)

	Samsung SDI 006405.KS	\$37.69 B
	Hyundai HYMTF	\$31.92 B
	LG Electronics LGLG.F	\$14.07 B
	Adidas 644 ADS.DE	\$27.60 B

Current LLM Paradigm



- OpenAI estimated to spent **1 Billion dollars** on AI
- Current valuation? **30 Billion dollars** by Microsoft (30x)

	Samsung SDI 006405.KS	\$37.69 B
	Hyundai HYMTF	\$31.92 B
	LG Electronics LGLG.F	\$14.07 B
	Adidas 644 ADS.DE	\$27.60 B

Summarization

The picture appeared on the wall of a Poundland store on Whymark Avenue [...] How would you rephrase that in a few words?

Sentiment Analysis

Review: We came here on a Saturday night and luckily it wasn't as packed as I thought it would be [...] On a scale of 1 to 5, I would give this a

Question Answering

I know that the answer to "What team did the Panthers defeat?" is in "The Panthers finished the regular season [...]". Can you tell me what it is?

Multi-task training

Zero-shot generalization

Natural Language Inference

Suppose "The banker contacted the professors and the athlete". Can we infer that "The banker contacted the professors"?

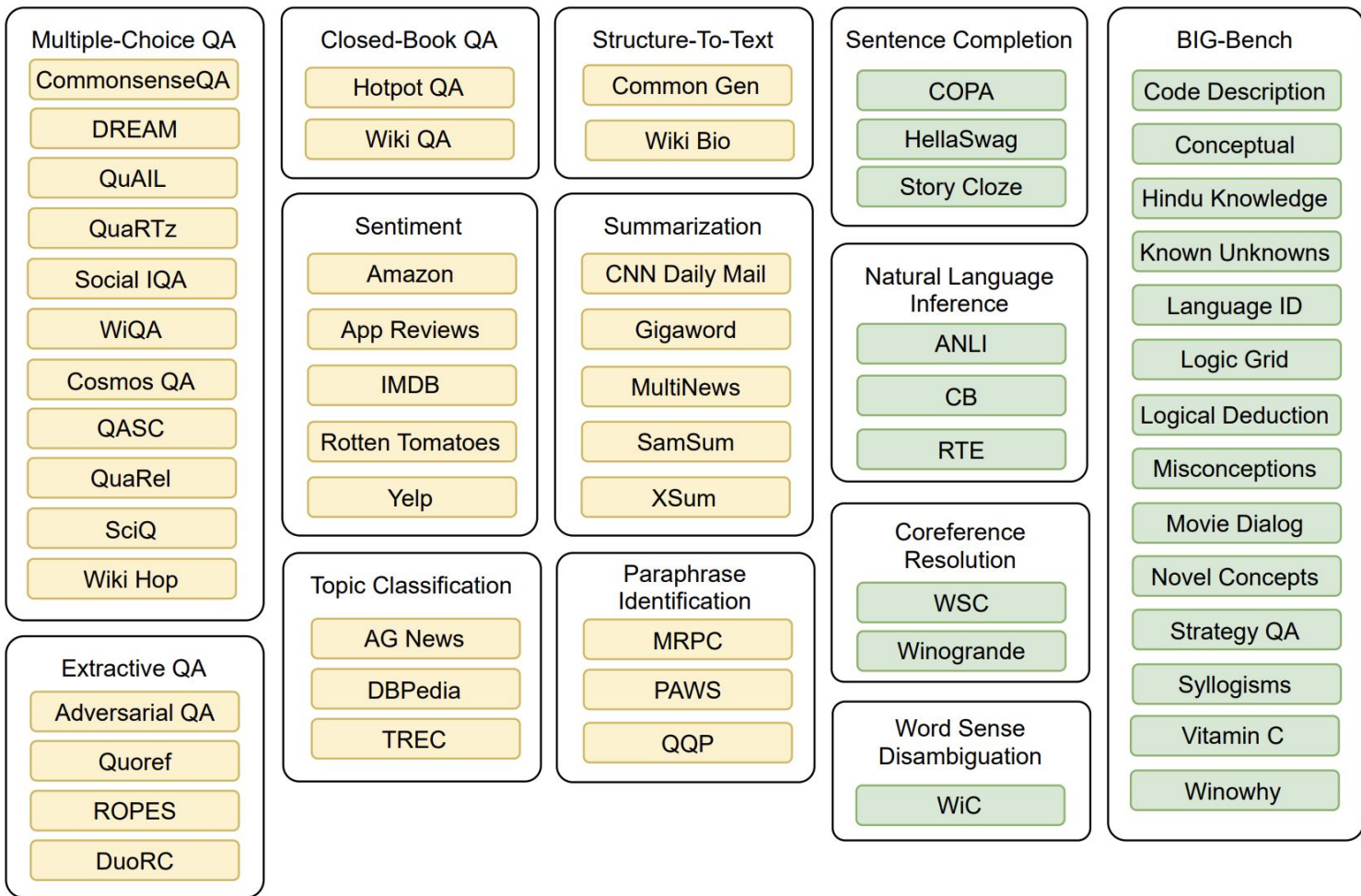
T₀

Graffiti artist Banksy is believed to be behind [...]

4

Arizona Cardinals

Yes



Burst of Instruction-Tuned LMs (MT LMs)

- FLAN, T0, InstructGPT, Tk-Instruct, Flipped, OPT-IML, GPT-JT, FLAN-T5, BLOOMz, mT0, etc.
- ALL Instructed-tuned LMs have the same analysis / storyline....

Burst of Instruction-Tuned LMs (MT LMs)

- FLAN, T0, InstructGPT, Tk-Instruct, Flipped, OPT-IML, GPT-JT, FLAN-T5, BLOOMz, mT0, etc.
- ALL Instructed-tuned LMs have the same analysis / storyline....

Scaling the total number of training tasks is one of the **key components** of the unseen task generalization capabilities of MT LMs.

Burst of Instruction-Tuned LMs (MT LMs)

- 

Scaling of the ge

is one task is.



🧐 Expert Language Models

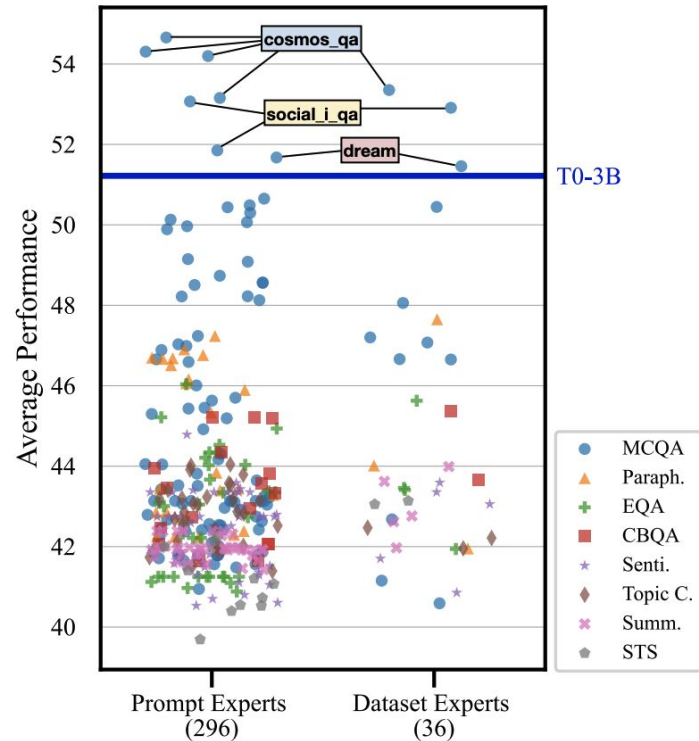
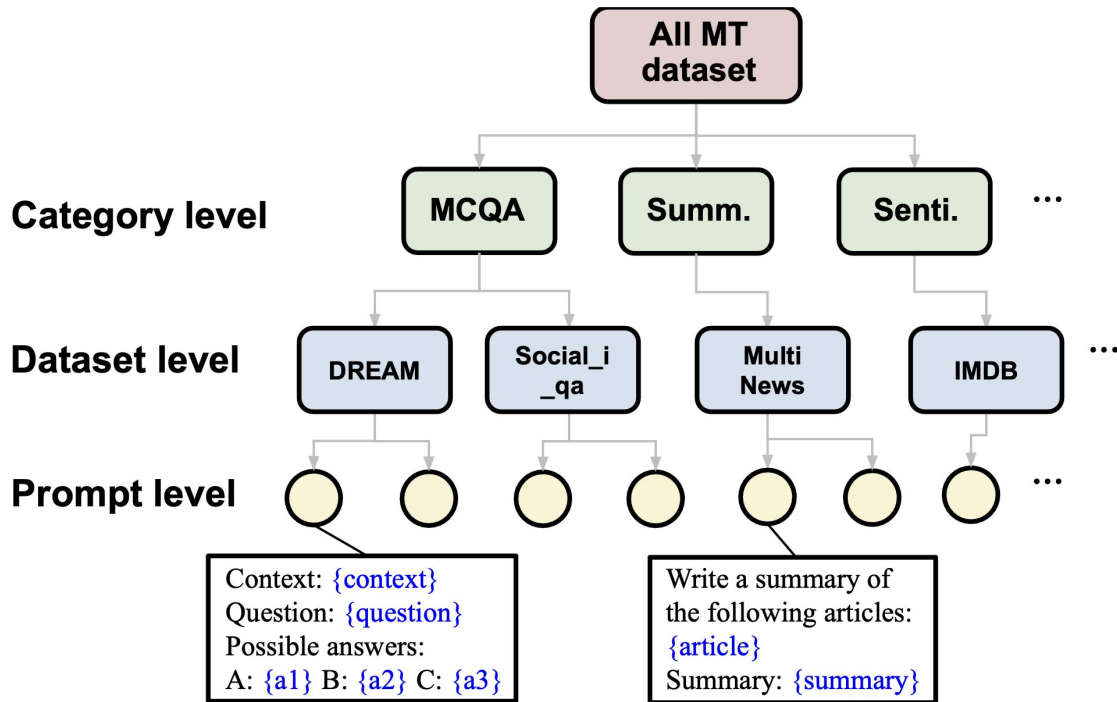


Figure 1. Mean accuracy performance of Expert LMs (each trained on a single task) on 11 unseen datasets compared to an instruction-tuned LM, T0-3B. Results show some Expert LMs surpassing T0-3B, challenging the commonly held belief that simply scaling the total number of training tasks is the key component to enhancing the capability of MT LMs.

Summarization

The picture appeared on the wall of a Poundland store on Whymark Avenue [...] How would you rephrase that in a few words?

Sentiment Analysis

Review: We came here on a Saturday night and luckily it wasn't as packed as I thought it would be [...] On a scale of 1 to 5, I would give this a

Question Answering

I know that the answer to "What team did the Panthers defeat?" is in "The Panthers finished the regular season [...]". Can you tell me what it is?

Multi-task training

Zero-shot generalization

Natural Language Inference

Suppose "The banker contacted the professors and the athlete". Can we infer that "The banker contacted the professors"?

T₀

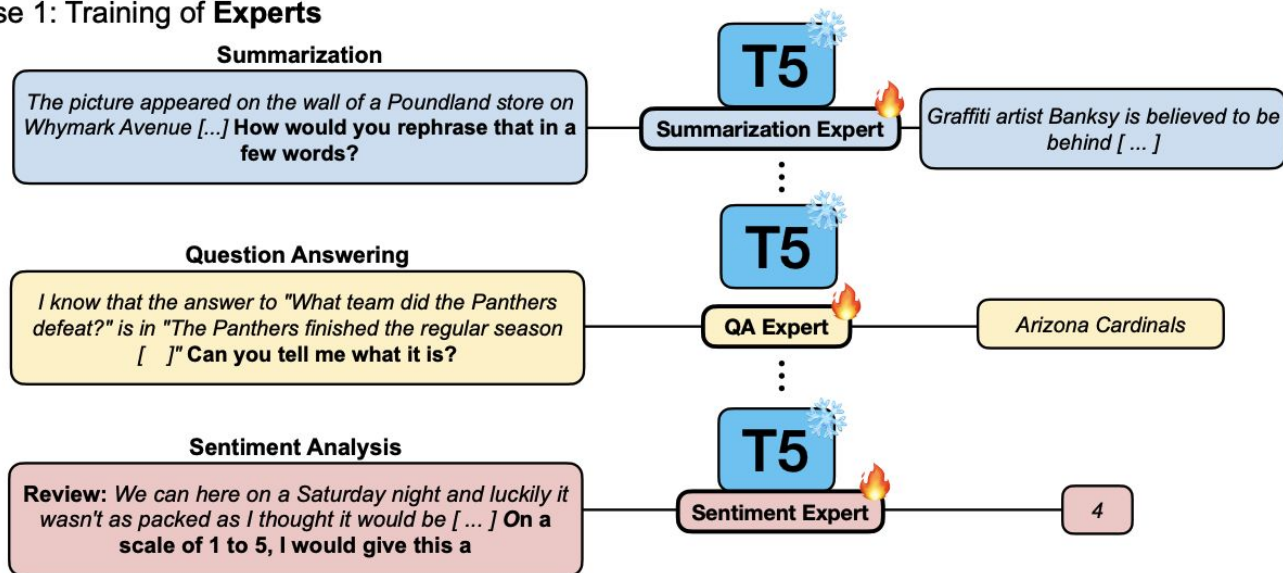
Graffiti artist Banksy is believed to be behind [...]

4

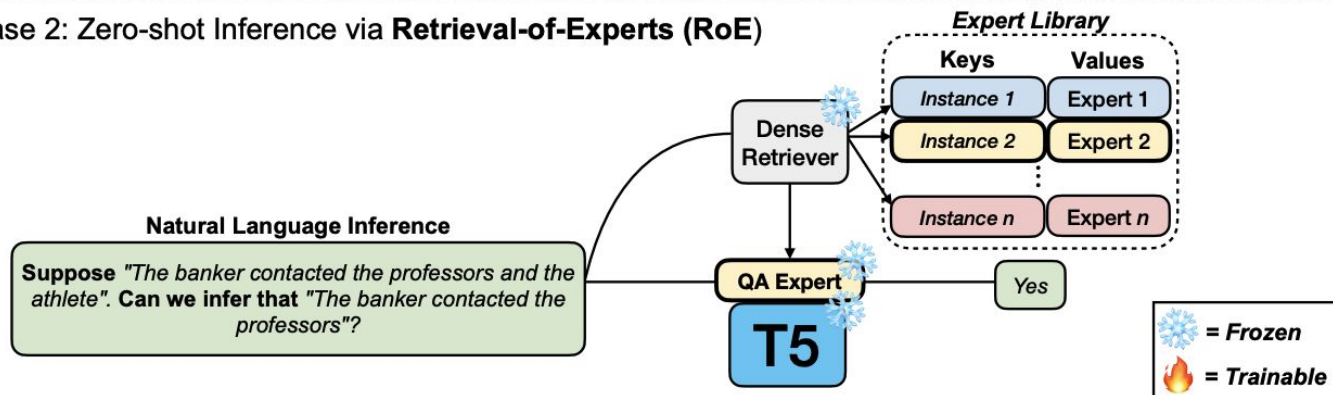
Arizona Cardinals

Yes

Phase 1: Training of Experts



Phase 2: Zero-shot Inference via Retrieval-of-Experts (RoE)



Main Results - 11 unseen tasks

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T0-11B	80.83	70.12	43.56	38.68	41.26	90.02	33.58	92.40	59.94	61.45	56.58	60.76
GPT-3(175B)	63.50	46.40	34.60	35.40	34.50	91.00	78.90	83.20	70.20	65.40	45.92	59.00
T0-3B	60.61	48.81	35.10	33.27	33.52	75.13	27.18	84.91	50.91	65.00	51.27	51.43

Main Results - 11 unseen tasks

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T0-11B	80.83	70.12	43.56	38.68	41.26	90.02	33.58	92.40	59.94	61.45	56.58	60.76
GPT-3(175B)	63.50	46.40	34.60	35.40	34.50	91.00	78.90	83.20	70.20	65.40	45.92	59.00
T0-3B	<u>60.61</u>	<u>48.81</u>	35.10	33.27	<u>33.52</u>	75.13	27.18	84.91	50.91	65.00	<u>51.27</u>	51.43
T5(3B) + Cos PE	49.53	49.52	36.21	36.11	36.38	89.63	43.77	97.06	<u>56.65</u>	57.02	49.01	54.63

Main Results - 11 unseen tasks

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T0-11B	80.83	70.12	43.56	38.68	41.26	90.02	33.58	92.40	59.94	61.45	56.58	60.76
GPT-3(175B)	63.50	46.40	34.60	35.40	34.50	91.00	78.90	83.20	70.20	65.40	45.92	59.00
T0-3B	<u>60.61</u>	<u>48.81</u>	35.10	33.27	<u>33.52</u>	75.13	27.18	84.91	50.91	65.00	<u>51.27</u>	51.43
T5(3B) + Cos PE	49.53	49.52	36.21	36.11	36.38	89.63	43.77	97.06	<u>56.65</u>	57.02	49.01	54.63
T5(3B) + PE w/ RoE	64.01	43.57	<u>35.49</u>	<u>34.64</u>	31.22	<u>79.25</u>	<u>34.60</u>	<u>86.33</u>	61.60	<u>62.21</u>	52.97	<u>53.48</u>
T5(3B) + PE w/ RoE (ORC.)	70.32	70.12	40.02	40.11	42.07	92.88	55.00	97.47	64.40	65.77	58.90	63.37

Main Results - 13 Tasks of BIG-Bench

Dataset (metric)	T0 3B	Cos PE 3B
Known Un.	47.83	58.70
Logic Grid	32.10	30.70
Strategy.	53.23	42.36
Hindu Kn.	34.86	51.43
Movie D.	53.22	46.72
Code D.	53.33	66.67
Concept	67.25	72.92
Language	14.94	25.95
Vitamin	58.18	46.55
Syllogism	52.27	50.00
Misconcept.	52.05	47.03
Logical	45.33	42.40
Winowhy	44.29	44.33
BIG-bench AVG	46.84	48.13

Main Results - 13 Tasks of BIG-Bench

Dataset (metric)	T0 3B	Cos PE 3B	T0 11B	GPT-3 175B	PALM 540B
Known Un.	47.83	58.70	65.22	60.87	56.52
Logic Grid	32.10	30.70	33.67	31.20	32.10
Strategy.	53.23	42.36	54.67	52.30	64.00
Hindu Kn.	34.86	51.43	42.86	32.57	56.00
Movie D.	53.22	46.72	57.33	51.40	49.10
Code D.	53.33	66.67	51.67	31.67	25.00
Concept	67.25	72.92	71.72	26.78	59.26
Language	14.94	25.95	18.33	15.90	20.10
Vitamin	58.18	46.55	57.33	12.30	14.10
Syllogism	52.27	50.00	48.33	50.50	49.90
Misconcept.	52.05	47.03	52.97	47.95	47.47
Logical	45.33	42.40	54.67	23.42	24.22
Winowhy	44.29	44.33	55.00	51.50	45.30
BIG-bench AVG	46.84	48.13	51.06	37.57	41.77

Main Results - 8 unseen generative tasks

Method	wiki auto (BLEU)	HGen (ROUGE)	haiku (ROUGE)	covid qa (BS)	eli5 (BS)	emdg (BS)	esnli (BS)	twitter (BS)	Total Avg.
T0-3B	<u>21.76</u>	<u>33.29</u>	19.93	50.00	59.86	47.76	<u>42.80</u>	28.40	<u>37.98</u>
T5(3B) + SAM PE	30.69	<u>25.49</u>	<u>25.25</u>	<u>49.93</u>	<u>47.94</u>	51.36	58.28	69.55	44.81
T5(3B) + PE w/ RoE	3.88	35.55	26.53	<u>33.52</u>	<u>33.66</u>	49.90	28.61	<u>49.22</u>	32.61
T5(3B) + PE w/ RoE (ORC.)	31.56	35.55	30.16	52.49	63.20	58.36	60.02	82.08	51.67

Main Results - 8 unseen generative tasks

Method	wiki auto (BLEU)	HGen (ROUGE)	haiku (ROUGE)	covid qa (BS)	eli5 (BS)	emdg (BS)	esnli (BS)	twitter (BS)	Total Avg.
T0-3B	<u>21.76</u>	<u>33.29</u>	19.93	50.00	59.86	47.76	<u>42.80</u>	28.40	<u>37.98</u>

Main Results - 8 unseen generative tasks

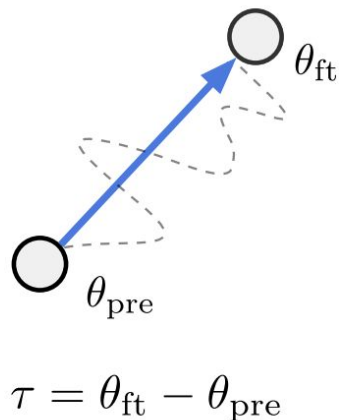
Method	wiki auto (BLEU)	HGen (ROUGE)	haiku (ROUGE)	covid qa (BS)	eli5 (BS)	emdg (BS)	esnli (BS)	twitter (BS)	Total Avg.
T0-3B	<u>21.76</u>	<u>33.29</u>	19.93	50.00	59.86	47.76	<u>42.80</u>	28.40	<u>37.98</u>
T5(3B) + SAM PE	30.69	25.49	<u>25.25</u>	<u>49.93</u>	<u>47.94</u>	51.36	58.28	69.55	44.81

Main Results - 8 unseen generative tasks

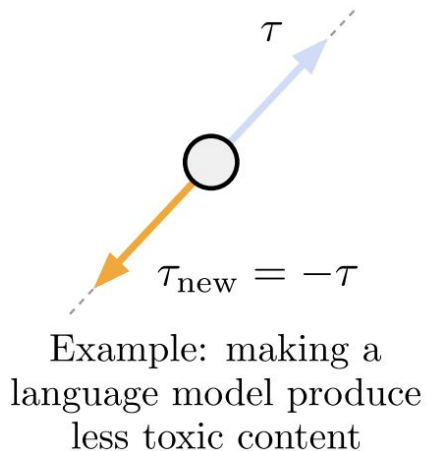
Method	wiki auto (BLEU)	HGen (ROUGE)	haiku (ROUGE)	covid qa (BS)	eli5 (BS)	emdg (BS)	esnli (BS)	twitter (BS)	Total Avg.
T0-3B	<u>21.76</u>	<u>33.29</u>	19.93	50.00	59.86	47.76	<u>42.80</u>	28.40	<u>37.98</u>
T5(3B) + SAM PE	30.69	<u>25.49</u>	<u>25.25</u>	<u>49.93</u>	<u>47.94</u>	51.36	58.28	69.55	44.81
T5(3B) + PE w/ RoE	3.88	35.55	26.53	<u>33.52</u>	<u>33.66</u>	49.90	28.61	<u>49.22</u>	32.61
T5(3B) + PE w/ RoE (ORC.)	31.56	35.55	30.16	52.49	63.20	58.36	60.02	82.08	51.67

Merging (Previous Work)

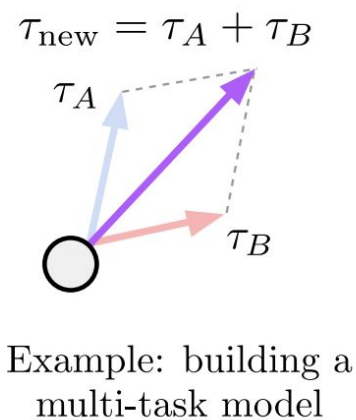
a) Task vectors



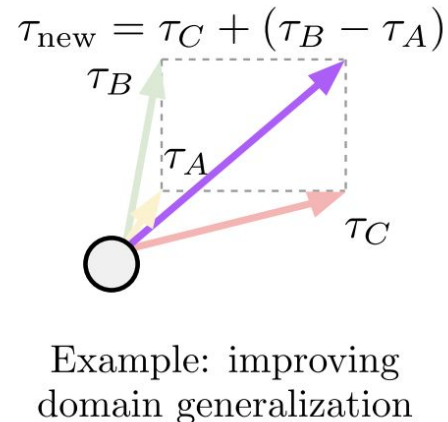
b) Forgetting via negation



c) Learning via addition



d) Task analogies



Matena, M., & Raffel, C. (2021). Merging models with fisher-weighted averaging. NeurIPS 2022.

Ilharco, G., Ribeiro, M. T., Wortsman, M., Gururangan, S., Schmidt, L., Hajishirzi, H., & Farhadi, A. (2022). Editing Models with Task Arithmetic. arXiv preprint arXiv:2212.04089.

Main Results - Merging

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T5(3B) + Cos PE	<u>49.53</u>	49.52	36.21	36.11	36.38	89.63	43.77	97.06	56.65	57.02	49.01	54.63
T5(3B) + Soc PE	61.26	38.81	33.16	33.63	33.46	<u>90.50</u>	<u>37.21</u>	<u>97.09</u>	<u>55.28</u>	50.00	50.11	<u>52.77</u>
T5(3B) + Cos&Soc PE (MER.)	49.10	<u>39.40</u>	<u>33.80</u>	<u>34.28</u>	<u>34.18</u>	91.63	36.29	97.25	55.06	<u>51.25</u>	<u>49.62</u>	51.99

Main Results - Merging

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T5(3B) + Cos PE	49.53	49.52	36.21	36.11	36.38	89.63	43.77	97.06	56.65	57.02	49.01	54.63
T5(3B) + Soc PE	61.26	38.81	33.16	33.63	33.46	90.50	37.21	97.09	55.28	50.00	50.11	52.77
T5(3B) + Cos&Soc PE (MER.)	49.10	<u>39.40</u>	<u>33.80</u>	<u>34.28</u>	<u>34.18</u>	91.63	36.29	97.25	55.06	<u>51.25</u>	<u>49.62</u>	51.99

Main Results - Merging

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T5(3B) + Cos PE	49.53	49.52	36.21	36.11	36.38	89.63	43.77	97.06	56.65	57.02	49.01	54.63
T5(3B) + Soc PE	61.26	38.81	33.16	33.63	33.46	90.50	37.21	97.09	55.28	50.00	50.11	52.77
T5(3B) + Cos&Soc PE (MER.)	49.10	39.40	33.80	34.28	34.18	91.63	36.29	97.25	55.06	51.25	49.62	51.99

Main Results - Merging

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T5(3B) + Cos PE	<u>49.53</u>	49.52	36.21	36.11	36.38	89.63	43.77	97.06	56.65	57.02	49.01	54.63
T5(3B) + Soc PE	61.26	38.81	33.16	33.63	33.46	<u>90.50</u>	<u>37.21</u>	<u>97.09</u>	<u>55.28</u>	50.00	50.11	<u>52.77</u>
T5(3B) + Cos&Soc PE (MER.)	49.10	<u>39.40</u>	<u>33.80</u>	<u>34.28</u>	<u>34.18</u>	91.63	36.29	97.25	55.06	<u>51.25</u>	<u>49.62</u>	51.99
T5(3B) + Cos DE	59.71	57.62	33.45	33.93	34.54	90.00	36.58	96.29	53.37	42.88	49.91	53.48
T5(3B) + Soc DE	65.52	48.69	35.20	35.39	37.11	83.25	30.38	87.18	54.27	54.62	51.39	53.00
T5(3B) + Cos&Soc DE (MER.)	<u>60.43</u>	<u>54.17</u>	<u>35.01</u>	<u>34.53</u>	<u>35.52</u>	91.25	<u>35.59</u>	96.73	54.33	<u>42.88</u>	<u>50.05</u>	53.68

Main Results - Merging

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T5(3B) + Cos PE	<u>49.53</u>	49.52	36.21	36.11	36.38	<u>89.63</u>	43.77	<u>97.06</u>	56.65	57.02	49.01	54.63
T5(3B) + Soc PE	61.26	38.81	33.16	33.63	33.46	<u>90.50</u>	<u>37.21</u>	<u>97.09</u>	<u>55.28</u>	50.00	50.11	<u>52.77</u>
T5(3B) + Cos&Soc PE (MER.)	49.10	<u>39.40</u>	<u>33.80</u>	<u>34.28</u>	<u>34.18</u>	91.63	36.29	97.25	55.06	<u>51.25</u>	<u>49.62</u>	51.99
T5(3B) + Cos DE	59.71	57.62	33.45	33.93	34.54	<u>90.00</u>	36.58	<u>96.29</u>	<u>53.37</u>	42.88	49.91	53.48
T5(3B) + Soc DE	65.52	48.69	35.20	35.39	37.11	<u>83.25</u>	30.38	<u>87.18</u>	<u>54.27</u>	54.62	51.39	53.00
T5(3B) + Cos&Soc DE (MER.)	<u>60.43</u>	<u>54.17</u>	<u>35.01</u>	<u>34.53</u>	<u>35.52</u>	91.25	<u>35.59</u>	96.73	54.33	<u>42.88</u>	<u>50.05</u>	53.68

Main Results - Analysis

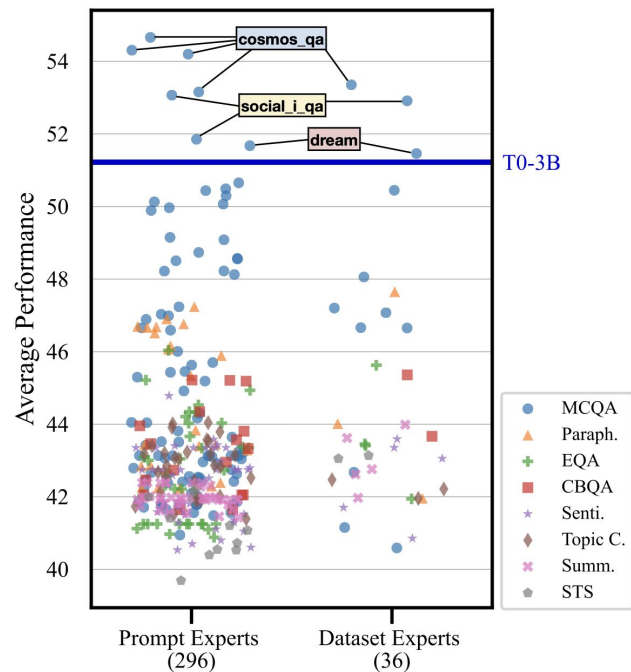
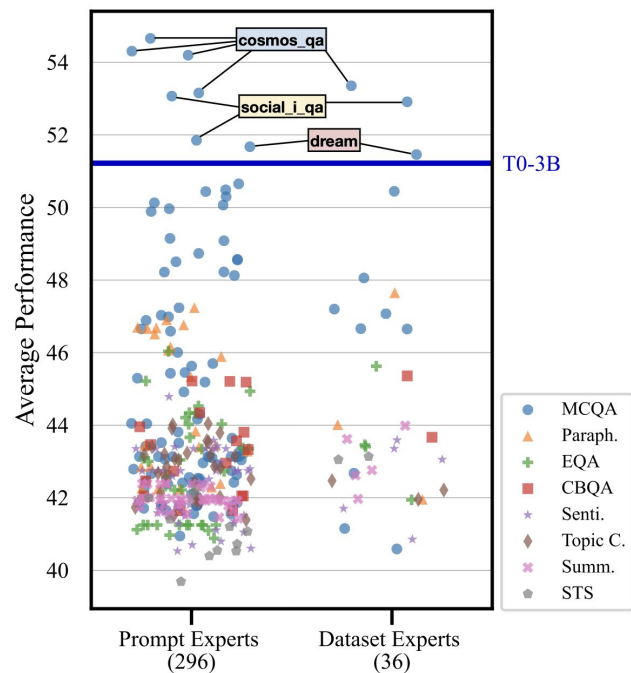


Figure 1: Average accuracy performance of Expert LMs (each trained on a single task) on 11 unseen datasets compared to an instruction-tuned LM, T0-3B.

Main Results - Analysis

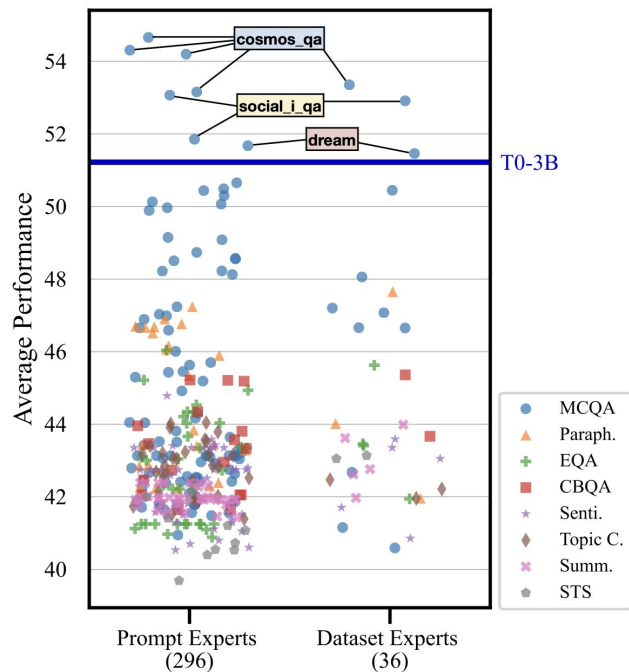


Common Traits

- 3 datasets are all **commonsense reasoning** tasks

Figure 1: Average accuracy performance of Expert LMs (each trained on a single task) on 11 unseen datasets compared to an instruction-tuned LM, T0-3B.

Main Results - Analysis



Common Traits

- 3 datasets are all **commonsense reasoning** tasks
- 3 datasets have a significant ($>20\%$) performance gap from human upper-bound performance = task difficulty

Figure 1: Average accuracy performance of Expert LMs (each trained on a single task) on 11 unseen datasets compared to an instruction-tuned LM, T0-3B.

3 Other Advantages ELMs have over MT LMs

3 Other Advantages ELMs have over MT LMs

1. Is not susceptible to Negative Task Transfer from multitask training

3 Other Advantages ELMs have over MT LMs

1. Is not susceptible to Negative Task Transfer from multitask training
2. Can continually learn new tasks

3 Other Advantages ELMs have over MT LMs

1. Is not susceptible to Negative Task Transfer from multitask training
2. Can continually learn new tasks
3. Can perform *composition* of instructions better than MT LMs

1. Seen Task Performance

Method	MCQA (12) (ACC)	Senti. (5) (ACC)	Topic C. (3) (ACC)	Paraph. (3) (ACC)	STS (2) (ROUGE-L)	Summ. (5) (ROUGE-L)	EQA (4) (ROUGE-L)	CBQA (2) (ROUGE-L)	Total Avg.
T0-3B	46.97	<u>66.40</u>	59.99	76.63	41.90	<u>33.10</u>	28.79	24.67	47.30
T0-11B	<u>51.32</u>	64.03	<u>60.95</u>	<u>73.64</u>	<u>45.42</u>	<u>33.10</u>	41.20	<u>30.37</u>	<u>50.00</u>

1. Seen Task Performance

Method	MCQA (12) (ACC)	Senti. (5) (ACC)	Topic C. (3) (ACC)	Paraph. (3) (ACC)	STS (2) (ROUGE-L)	Summ. (5) (ROUGE-L)	EQA (4) (ROUGE-L)	CBQA (2) (ROUGE-L)	Total Avg.
T0-3B	46.97	<u>66.40</u>	59.99	76.63	41.90	<u>33.10</u>	28.79	24.67	47.30
T0-11B	<u>51.32</u>	64.03	<u>60.95</u>	<u>73.64</u>	<u>45.42</u>	<u>33.10</u>	41.20	<u>30.37</u>	<u>50.00</u>
T5(3B)+ PE w/ RoE	58.95	70.18	96.52	72.97	47.57	33.14	<u>30.36</u>	51.89	57.70
T5(3B)+ PE w/ RoE (ORC.)	56.28	84.52	96.91	79.34	47.94	35.40	40.34	43.24	60.50

2. Continual Learning of New Tasks

Method	Seen Avg.	Unseen Avg.	Gen Avg.
<i>Before Continual Learning</i>			<i>Unseen</i>
T0-3B	47.30	51.43	37.98
T5(3B) + PE w/ RoE	57.70	53.48	32.61

2. Continual Learning of New Tasks

Method	Seen Avg.	Unseen Avg.	Gen Avg.
<i>Before Continual Learning</i>			<i>Unseen</i>
T0-3B	47.30	51.43	37.98
T5(3B) + PE w/ RoE	57.70	53.48	32.61
<i>After Continual Learning</i>			<i>Seen</i>
CT0-3B	47.54	50.84	54.52 (↑)
T5(3B) + PE ⁺ w/ RoE	57.70	53.33	55.60 (↑)

3. Compositional Instructions

Instruction #1: Summarize the English Text

Instruction #2: Translate this text from English to {Language}

3. Compositional Instructions

Instruction #1: Summarize the English Text

Instruction #2: Translate this text from English to {Language}

Compositional Instruction: Summarize the English Text AND translate this text from English to {Language}

3. Compositional Instructions

Method	xsum en→ko	xsum en→es	xsum en→zh	xsum en→fr	xsum en→ja	Total Avg.
MT0-3B	1.84	16.14	6.74	20.37	3.44	9.71
MT5-3B + MER. EX.	8.23	16.97	2.40	20.55	13.98	12.43

3. Compositional Instructions

Method	xsum en→ko	xsum en→es	xsum en→zh	xsum en→fr	xsum en→ja	Total Avg.
MT0-3B	1.84	16.14	6.74	20.37	3.44	9.71
MT5-3B + MER. EX.	8.23	16.97	2.40	20.55	13.98	12.43

3. Compositional Instructions

Task	EXAMPLE
xsum en→es	MT0: El asesinato de un niño de tres años de edad en Francia fue atribuido a su hermano mayor. M.E: La policía francesa arrestó a cuatro miembros de la familia del niño por su presunta implicación en el crimen, que ocurrió el 17 de septiembre en la casa familiar en Mulhouse, al oeste del país, y de más de 100.000 habitantes.
xsum en→fr	MT0: Le président de la République démocratique du Malawi a été condamné à cinq ans de prison pour complicité dans l'assassinat de Paul Mphwiyo. M.E: Le 8 novembre 2013, l'ancien ministre de la Justice du Malawi, M. Ralph Kasambara, a été arrêté après avoir commis le meurtre de Paul Paul MPHWIYO, le directeur du budget du ministère des Finances.
xsum en→ja	MT0: 副裁 Meng Ship 和副 Meng Teng 被加拿大警方逮捕,被指控侵犯公民利。 M.E: カナダの最高裁判所(CFO)は 12月に逮捕された創設者の息子であり、副社長はカナダ政府とカナダ移民局(CBSA)と警察を告訴した。
xsum en→zh	MT0: The Sierra Leonean nurse who was isolated for seven hours at the airport terminal has said that the isolation experience is "terrifying" and may make other medical workers reluctant to go to West Africa. M.E: 一名感染埃博拉病毒的生Craig Spencer目前正在大都院接受隔治,但只得到了一粮食棒的。
xsum en→ko	MT0: Korean peninsula has had its warmest winter since 1973, according to the Meteorological Administration. M.E: 지난해 1월은 국내에서 가장 따뜻한 겨울이었다.

Future Work

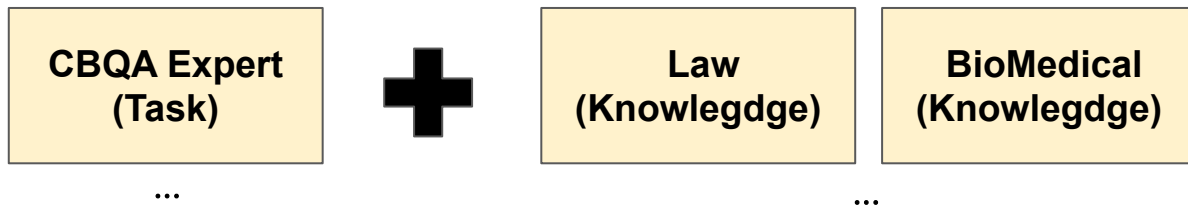
1. Train a supervised retriever
 - Close the Gap between Current RoE & RoE (Oracle)

Method	NLI					Sentence Completion			Coreference Resolut.		WSD	Total Avg.
	RTE	CB	AN. R1	AN. R2	AN. R3	COPA	Hellasw.	StoryC.	Winogr.	WSC	WiC	
T0-11B	80.83	70.12	43.56	38.68	41.26	90.02	33.58	92.40	59.94	61.45	56.58	60.76
GPT-3(175B)	63.50	46.40	34.60	35.40	34.50	91.00	78.90	83.20	70.20	65.40	45.92	59.00
T0-3B	<u>60.61</u>	<u>48.81</u>	35.10	33.27	<u>33.52</u>	75.13	27.18	84.91	50.91	65.00	<u>51.27</u>	51.43
T5(3B) + Cos PE	49.53	49.52	36.21	36.11	36.38	89.63	43.77	97.06	<u>56.65</u>	57.02	49.01	54.63
T5(3B) + PE w/ RoE	64.01	43.57	<u>35.49</u>	<u>34.64</u>	31.22	<u>79.25</u>	<u>34.60</u>	<u>86.33</u>	61.60	<u>62.21</u>	52.97	<u>53.48</u>
T5(3B) + PE w/ RoE (ORC.)	70.32	70.12	40.02	40.11	42.07	92.88	55.00	97.47	64.40	65.77	58.90	63.37

- Beat Flan-T5-3B (Current SOTA)! (~61)
 - + Train CoT Experts (Rationale experts)

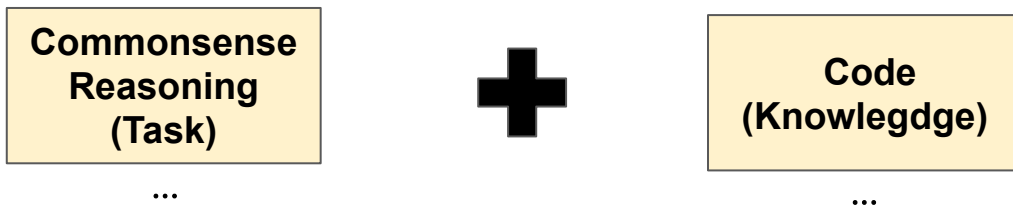
Future Work

1. Train a supervised retriever
 - Close the Gap between Current RoE & RoE (Oracle)
2. Exploring Merging
 - Currently, only *Task* + *Task* Expert Merging
 - What if *Task* + *Knowledge* Expert Merging?
 - How about *Knowledge* + *Knowledge* Expert Merging?



Future Work

1. Train a supervised retriever
 - Close the Gap between Current RoE & RoE (Oracle)
2. Exploring Merging
 - Currently, only *Task* + *Task* Expert Merging
 - What if *Task* + *Knowledge* Expert Merging?
 - How about *Knowledge* + *Knowledge* Expert Merging?



Future Work

1. Train a supervised retriever
 - Close the Gap between Current RoE & RoE
2. Exploring Merging
 - Currently, only *Task* + *Task* Expert Merging
 - What if *Task* + *Knowledge* Expert Merging?
 - How about *Knowledge* + *Knowledge* Expert Merging?

2023.01 Expert
(Knowledge)

...



2023.02 Expert
(Knowledge)

...

The screenshot shows the Hugging Face Online Language Modelling (OLM) community page. The header includes the Hugging Face logo, a search bar, and navigation links for Models, Datasets, Spaces, Docs, Solutions, Pricing, Log In, and Sign Up. The main content area is titled 'Online Language Modelling' with a 'Community' tag. Below this, there are sections for 'Research interests' (Making language models know whats up.), 'Team members' (2 members), and a list of models. The models list includes:

- o1m/o1m-gpt2-dec-2022 (Updated 21 days ago, 404 likes, 7 comments)
- o1m/o1m-roberta-base-dec-2022 (Updated 21 days ago, 122 likes, 6 comments)
- o1m/o1m-roberta-base-oct-2022 (Updated 21 days ago, 103 likes, 5 comments)
- o1m/o1m-gpt2-oct-2022 (Updated 21 days ago, 95 likes, 8 comments)
- o1m/o1m-roberta-base-latest (Updated 26 days ago, 30 likes, 3 comments)
- o1m/o1m-gpt2-latest (Updated 26 days ago, 28 likes, 4 comments)

Future Work

1. Train a supervised retriever
 - Close the Gap between Current RoE & RoE (Oracle)
2. Exploring Merging
 - Currently, only *Task* + *Task* Expert Merging
 - What if *Task* + *Knowledge* Expert Merging?
 - How about *Knowledge* + *Knowledge* Expert Merging?
3. Explore other Benefits of Distributed & Collaborative Training
 - Efficiency, Privacy, Personalization, Etc.

Q & A

Part 1

- Towards Continual Knowledge Learning of Language Models [ICLR'22]
- TemporalWiki: A Lifelong Benchmark for Training and Evaluating Ever-Evolving Language Models [EMNLP'22]
- Knowledge Unlearning for Mitigating Privacy Risks in Language Models [*under review*]

Part 2

- Exploring the Benefits of Training Expert Language Models over Instruction Tuning [*under review*]

Thank You