





DyKnow: <u>Dynamically Verifying Time-Sensitive</u> Factual <u>Know</u>ledge in LLMs

•••• ICT DAYS

S. Mahed Mousavi, Simone Alghisi, Giuseppe Riccardi Signals and Interactive Systems Lab, University of Trento, Italy

Introduction

DyKnow

LLMs are transforming to a **common go-to** for asking a vast variety of **factual questions**.

BUT! LLMs are trained on **static** data & evaluated via **static** benchmarks

Static benchmarks get **outdated** while factual knowledge is generally subject to **change**

We present **DyKnow** an evaluation framework to **dynamically** evaluate the knowledge in LLMs

For each fact, most **current info** are retrieved from the **Wikidata** at the time of evaluation



We evaluate the 24 LLMs & 4 knowledge editing algorithms regarding real-world facts.



Findings

- LLMs differ from traditional knowledge repositories, making it important to investigate:
- 1 the types of knowledge these models can reliably manage
- 2 the types of querying/alignment operations they support

We encourage community engagement to expand DyKnow into a current and active benchmark.

RQ1. Accuracy & Consistency of SOTA LLM via DyKnow

Significant percentage of outdated or irrelevant responses in all models	(Year) Model (2019) GPT-2 (2020) GPT-3	Correct 26% 42%	Outdated 42% 47%	Irrelevant 32% 12%	100%	Instruction-Tuned Models ChatGPT GPT-4
High percentage of outdated responses even in recent models	(2020) T5 (2021) GPT-J (2022) Bloom (2022) Flan-T5	11% 41% 35% 18%	21% 46% 39%	68% 13% 16%	80%	Mixtral _{I.} Mistral _{I.} Llama-3 _{I.}



No model exceeds 80% correctness

Output consistency

- 1 varies significantly across different models
- 2 increases with recent models
- 3 is higher with Instruction-tuned models than pre-trained models

No model achieves 100% consistency



RQ2: LLMs' Data Approximation

Approximating the LLMs' Pre-Training data shows:

RQ3: Updating LLMs' Knowledge

SERAC fails to achieve high performance with updating at best less than 40%

ROME demonstrates an overall poor performance & is outperformed by **MEMIT**

MEMIT excels with GPT-J and Llama-2_c while IKE

Model	#Outdated Facts	Knowledge Editing					
		Modifyin	g Parameters	Preserving Parameters			
		ROME	MEMIT	SERAC	IKE		
(2019) GPT-2	54	17%	33%	4%	49%		
(2021) GPT-J	60	11%	83%	0%	97% [‡]		
(2023) Llama-2 _C .	48	4%	77%	36%	18%		
(2023) Mistral.	/1	0%	$\int \mathcal{O} $		0206		