

NAACL 2025

Transient Tables : Evaluating LLMs' Reasoning on Temporally Evolving Semi-structured Tables

Abhilash Shankarampeta*, Harsh Mahajan*, Tushar Kataria, Dan Roth, Vivek Gupta Affiliations: UC San Diego, University of Utah, University of Pennsylvania, Arizona state university







Motivation

- Information is inherently transient and constantly updated
 - Examples: company profits, political figures, sports rankings, etc.
- LLMs are typically trained on static datasets
- **Research question:** Can LLMs effectively reason over temporal changes in information through in-context learning?

India				India			India		
	۲			3				(
					10g	TO BE	Nickname(s)	N	Men in Blue
	Flag of India				1000		Association	B	Board of Control for Cricke
Personnel					вс	CI		Per	reconnel
Test captain	Virat Kohli			Nickname(s)	Men ir	Blue, Team India	Captain	F	Rohit Sharma
One Day captain	Virat Kohli			Association	Board	of Control for Cricket in	Coach	F	Rahul Dravid
[20] captain	n Virat Kohli			India		History			
Coach	Bavi Shaetri		-	O	Perse	onnel	Test status	1	931
boach	History			Captain	Bowi S	lhastri	acquired	national	Cricket Council
	HIStory		l	oouon	His	orv	ICC status	F	full Member (1926)
lest status	1932			Test status	1931		ICC region	A	Asia
latamat	land Orbitation	0		acquired			ICC Rankings	Curren	It ^[3] Best-ever
International Cricket Council				International Cricket Council			Test	1st	1st (1 April 1973)
ICC Rankings	Current ^[1]	Best-ever		ICC status	Full M	ember (1926)	ODI	1st	1st (January 2013)
Test	1st	1st		ICC region	Asia		T20I	1st	1st ^{[1][2]} (28 March 2014)
ODI	2nd	1st		ICC Rankings	Current	1 Best-ever		т	ests
T20I	2nd	1st	Ц.,	ODI	2nd	1st (1 December 1994)	First Test	V L	England at Lord's, ondon: 25–28 June 1932
	Tests			T20I	3rd	1st (28 March 2014)	Last Test	v	West Indies at
First Test v + England at Lord's,			Tests			C	Queen's Park Oval, Port of		
	London; 25-28 June 1932			First Test	v. + England at Lord's,			S	Spain; 20–24 July 2023
Last Test	v ICI Sri Lan	ka at Feroz Shah			Londo	n; 25–28 June 1932	Total ⁴	572	173/176
	Kotla Ground,	Delhi; 2–6		Last Test	V.	Australia at Melbourne	Total	UTL	(222 draws, 1 tie)
	December 20	17			26-29	December 2020	This year ⁽⁵⁾	7	3/2
Tests P	laved Won/I	_ost		Tests	Played	Won/Lost			(2 draws)
Total ^[2] 5	18 143/1	58 Iraws, 1 tie)	1	Total ^[2]	544	158/168 (217 draws, 1 tie)	World Test Championshi appearances	р Р	! (first in 2019–2021)
This	1 7/1 (3	draws)		This year ^[3]	4	1/3 (0 draws)	Best result	9	Runners-up (2019-21,

Sample question - How many Test matches did the Indian Cricket Team play between 2020 and 2023?

Key Contributions

- Transient Tables Dataset
 - A novel QA dataset with **3,971 questions** from **14,000+ tables** spanning **1,238 entities**
 - **Template-based question** generation pipeline using LLMs'
- Baseline results with **state-of-the-art models**
 - GPT-40, Llama3-70B, Gemini 1.5, GPT-40-mini, Llama3-8B, Mixtral
- Novel modeling strategies using task decomposition to enhance performance

Dataset Creation

- Entity timeline selection from Wikipedia infoboxes
- Timeline cleaning and filtering (8-12 tables per entity)
- Query-answer generation through templates
- Statistics: 3,971 questions, 1,238 entities, average 11.42 tables per entity



Question Categorization

- Time information:
 - 2,985 implicit vs. 986 explicit questions
- Reasoning types
 - extraction, counting, comparison, etc.
- Complexity:
 - 2,113 single key vs. 1,858 multiple key questions







Reasoning Operation

Modeling Techniques

- Information granularity variations: closed book, single table, full timeline, oracle timeline
- Task decomposition approaches:
 - Without decomposition
 - Information retrieval
 - Information extraction
 - Information retrieval-extraction



Task Decomposition: Why It Matters

Question: "Who was the coach of the Indian Cricket Team when Virat Kohli was captain and the team had its highest ICC Test ranking?"

Approach 1: Without Decomposition

Process full timeline of tables simultaneously





Task Decomposition: Why It Matters

Question: "Who was the coach of the Indian Cricket Team when Virat Kohli was captain and the team had its highest ICC Test ranking?"

Approach 2: With Decomposition

O1 Table Retrieval

Identify tables from periods when Kohli was captain 02 Key Extraction

> Extract "Captain," "Coach," and "ICC Test ranking" keys from these tables

O3 Answer Generation Identify when ranking was highest during Kohli's captaincy and determine coach during that specific period

Task Decomposition : Levels

Information Retrieval (IR)

Information Extraction (IE)

Stage 1: "Table Retrieval" (identify relevant tables from timeline)

Stage 2: "Answer Generation" (reason over retrieved tables) **Stage 1:** "Key Extraction" (extract relevant attributes from tables)

Stage 2: "Answer Generation" (reason over extracted keys) Information Retrieval Extraction (IRE)

Stage 1: "Table Retrieval" (identify relevant tables)

Stage 2: "Key Extraction" (extract relevant attributes from retrieved tables)

Stage 3: "Answer Generation" (reason over extracted keys)

Experimental Setup

- Models evaluated
 - GPT-40, Llama3-70B, Gemini-1.5-flash, GPT-40-mini, Llama3-8B, Mixtral-7x8B
- Prompting Techniques
 - Zero shot, Few Shot, Chain of Thought
- Evaluation metrics
 - F1, Exact Match (EM), Rouge-1, Rouge-L
- Human evaluation baseline for comparison

Results: Context Decomposition



Results: Different Models



Performance of Models

Results in different in-context variations and different intermediate task decompositions with various prompting methods.

Results : Temporal Models

Temporal Specific Models



CoT prompting with oracle tables and Key Extraction for task decomposition.

Conclusion & Future Work

- A novel task of question answering on temporally evolving tables.
- A new Transient Tables dataset
 - 3,971 question-answer pairs.
 - From over 14k tables and 1,238 entities across various time periods.
- First study on LLM reasoning over entity-centric temporal tables.

In future:

- Currently its confined to Wikipedia infoboxes. Extending it to diverse structures beyond tables.
- Neuro-symbolic learning for better interpretability.



Thanks

Scan the QR to checkout Transient Tables