



# TDBench: Harnessing Temporal Databases for TSQA in LLMs

For automated, scalable, and reliable evaluation of Time-Sensitive QA (TSQA)



Soyeon Kim<sup>1</sup>, Jindong Wang<sup>2</sup>, Xing Xie<sup>3</sup>, and Steven Euijong Whang<sup>1</sup>

<sup>1</sup>KAIST <sup>2</sup>William & Mary <sup>3</sup>Microsoft Research Asia

# What is Time-Sensitive QA (TSQA)?

Asking questions where answers evolve over time



## Answers Depend on Time

The correct factual answer changes depending on *when* the question is asked (e.g., "Current US President").



## Requires Temporal Alignment

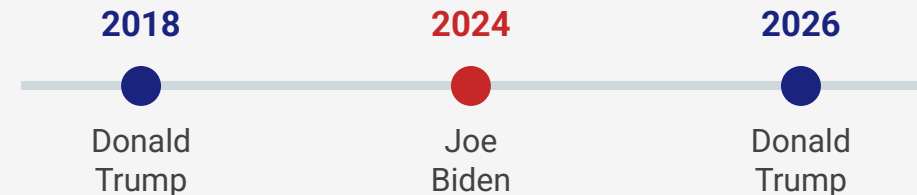
Models must distinguish between outdated and current knowledge, accurately reflecting the present state of the world.



## Requires Temporal Reasoning

Going beyond fact retrieval to understand temporal relationships: *before*, *after*, *during*, or *overlapping* events.

### Example: "Who is the current US President?"



### 💡 "Static" QA fails here.

In TSQA, LLMs must answer the correct president based on the specific time constraint (e.g., "current", "in 2026", ...) mentioned in the query.

# Background & Limitations

Critical Bottlenecks in Existing Time-Sensitive Question-Answering (TSQA) Benchmarks



## High Manual Cost

- Heavy reliance on manual template design
- Limited to a small set of fixed templates (9–16 types), hindering scalability



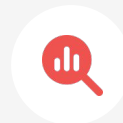
## Narrow Knowledge Source

- Centered almost exclusively on Wikipedia/Wikidata
- Poor coverage of specialized domains (e.g., Medical, Legal, Corporate)



## Simplistic Temporal Logic

- Covers only basic relations like "Before" or "After"
- Misses complex intervals ("Overlap", "Contain", "Meet")



## Lack of Reasoning Verification

- Benchmarks only focus on the final answer string
- Fails to detect "Right Answer, Wrong Reason" (Reasoning Hallucination)

# Motivation

Why We Need a Database-Driven Approach for TSQA

## Prior Limitations



### High Manual Cost & Simplistic Temporal Logics

Heavy reliance on manual template design.  
Limited to a small set of fixed templates (9–16 types).



### Lack of Reasoning Verification

Benchmarks focus on the final answer string.  
Fails to detect "Right Answer, Wrong Reason".



### Narrow Knowledge Source

Centered almost exclusively on Wikipedia/Wikidata.  
Poor coverage of specialized domains.



## TDBench Objectives



### Automated Scale & Diverse Temporal Logics

Eliminate manual bottlenecks.  
No need to write domain-specific templates.



### Reasoning Verification Along with Answer

Benchmark focus on both the final answer and rationales.  
Detect "Right Answer, Wrong Reason".



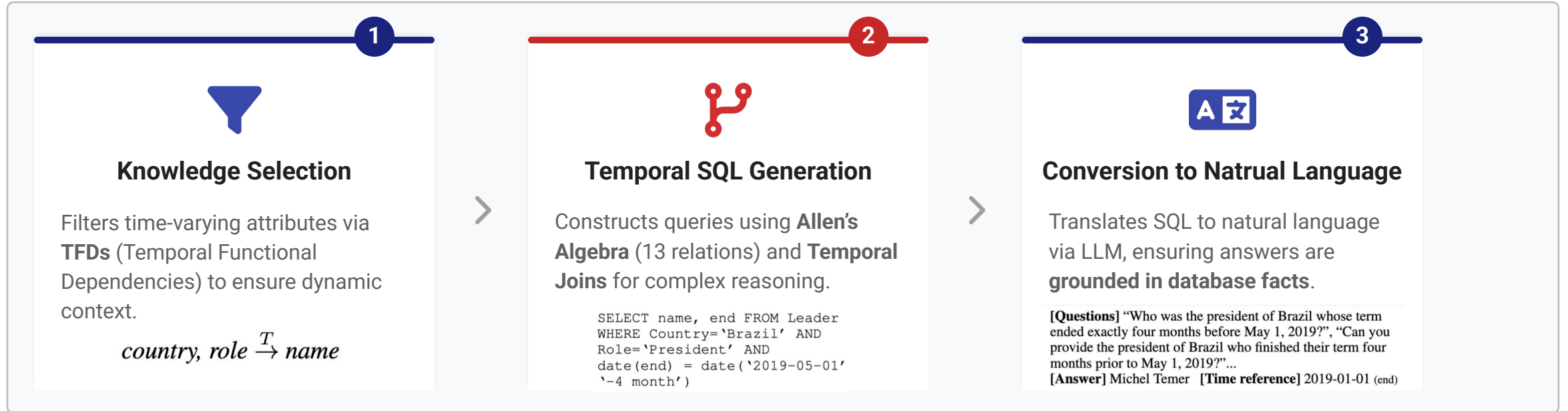
### Generalizability to Arbitrary Domains

Extend beyond Wikipedia to domains like Legal/Medical.

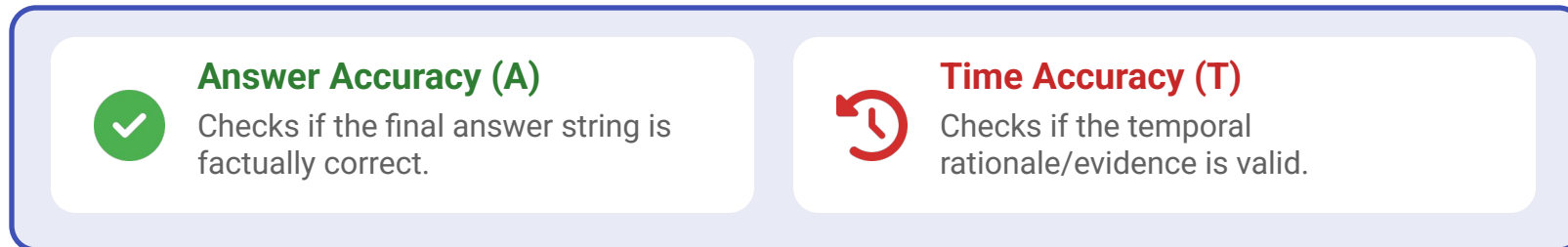
# TDBench Framework

Automated 3-Step Construction & Verifiable Evaluation Pipeline

## TSQA CONSTRUCTION



## TSQA EVALUATION



# Key Differences

What Makes TDBench Unique and Advanced to Prior TSQA Benchmarks



## Fully Automated

Replaces human-written templates with database-driven generation. No manual effort required to create question contexts.

$country, role \xrightarrow{T} name$



## 13 Temporal Relations

Covers the full spectrum of Allen's Interval Algebra (e.g., *Overlap*, *During*, *Meet*) using Temporal SQL.



## Domain Scalability

Extensible to any domain with temporal data (Medical, Legal, Media), allowing for specialized benchmark creation.



## Fine-grained Evaluation

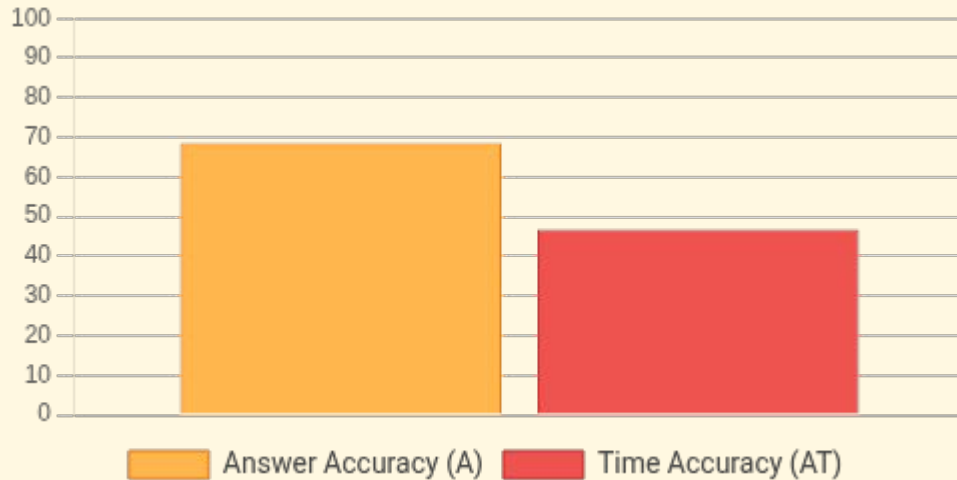
Introduces Time Accuracy (T) to verify if the model's rationale references the correct time period, exposing hallucinations.

Relation	Interval Diagram	SQL Condition	Example Temporal Context	Criteria
$a$ before $b$		$a.end < b.start$	A president who ends before <i>September 20, 2000</i>	end
$a$ after $b$		$a.start > b.end$	A president who starts after <i>March 20, 2001</i>	start
$a$ meet $b$		$a.end = b.start$	A president who ends exactly <i>half a year</i> before <i>March 20, 2001</i>	end
$a$ met-by $b$		$a.start = b.end$	A president who starts exactly <i>half a year</i> after <i>September 20, 2000</i>	start
$a$ overlap $b$		$a.start < b.start$ $\wedge b.start < a.end < b.end$	A president who starts before <i>September 20, 2000</i> and ends between <i>September 20, 2000</i> and <i>March 20, 2001</i>	start, end
		$a.start < b.start$ $\wedge a.end$ IS NULL	A president who is currently serving	start
$a$ overlapped-by $b$		$a.end > b.end$ $\wedge b.start < a.start < b.end$	A president who starts between <i>September 20, 2000</i> and <i>March 20, 2001</i> and ends after <i>March 20, 2001</i>	start, end
$a$ equal $b$		$a.start = b.start$ $\wedge a.end = b.end$	A president who starts in <i>September, 2000</i> and ends in <i>March, 2001</i>	start, end
$a$ start $b$		$a.start = b.start$ $\wedge a.end < b.end$	A president who starts in <i>September, 2000</i> and ends before <i>March, 2001</i>	start, end
$a$ started-by $b$		$a.start = b.start$ $\wedge a.end > b.end$	A president who starts in <i>September, 2000</i>	start
$a$ finish $b$		$a.start > b.start$ $\wedge a.end = b.end$	A president who starts after <i>September 20, 2000</i> and ends in <i>March, 2001</i>	start, end
$a$ finished-by $b$		$a.start < b.start$ $\wedge a.end < b.end$	A president who ends in <i>March, 2001</i>	end
$a$ during $b$		$a.start > b.start$ $\wedge a.end < b.end$	A president who starts after <i>September 20, 2000</i> and ends before <i>March 20, 2001</i>	start, end
$a$ contain $b$		$a.start < b.start$ $\wedge a.end > b.end$	A president who starts before <i>September 20, 2000</i> and ends after <i>March 20, 2001</i>	start, end

# Experimental Results

Quantifying Hallucinations and Demonstrating Efficiency Across 8 LLMs

## The "Reasoning Gap": 21.7% Drop



**Finding:** Models like GPT-4o and Llama 3.1-70B often get the correct answer string, but cite incorrect time periods (Reasoning Hallucination).

## Efficiency: ~90% Cost Reduction

Netflix Dataset Token Usage (Per Question Generation)

Full Docs 1,288

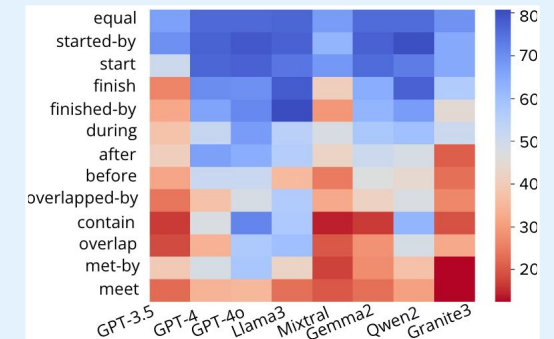
TDBench 138



## Qualitative Insights

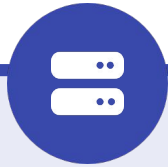
### Temporal Blind Spots

Diverse and complex temporal constraints identify model-specific weaknesses in temporal reasoning.



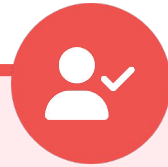
# Conclusion & Takeaways

Summary of Contributions and Future Directions



## Scalable Automation

Leveraging **Temporal Databases** and SQL eliminates manual bottlenecks, enabling the systematic generation of complex TSQA pairs at scale.



## Trustworthy Evaluation

The **Time Accuracy (T)** metric uncovers "Right Answer, Wrong Reason" cases, revealing that LLMs hallucinate temporal rationales ~21.7% of the time.



## Cost-Effective Scale

TDBench is ready for **domain-specific** benchmark (Medical, Legal) with minimal cost compared to LLM-only generation

## Get TDBench Resources

Access the benchmark, code, and full paper.

 <https://github.com/ssoy0701/tdbench>



<https://iclr.cc/virtual/2026/poster/10009076>

## Thank You!

Questions? Reach out to us:

[purplehibird@kaist.ac.kr](mailto:purplehibird@kaist.ac.kr)