

# Assessing The (In)Ability of LLMs To Reason in Interval Temporal Logic

Pietro Bellodi<sup>1,2</sup>   Pietro Casavecchia<sup>1,2</sup>   Alberto Paparella<sup>1,2</sup>  
Guido Sciavicco<sup>1,2</sup>   Ionel Eduard Stan<sup>3,2</sup>

<sup>1</sup>Department of Mathematics and Computer Science, University of Ferrara, Italy

<sup>2</sup>Applied Computational Logic and Artificial Intelligence Lab (ACLAI)

<sup>3</sup>Department of Informatics, Systems, and Communications, University of Milano-Bicocca, Italy

29/08/2025

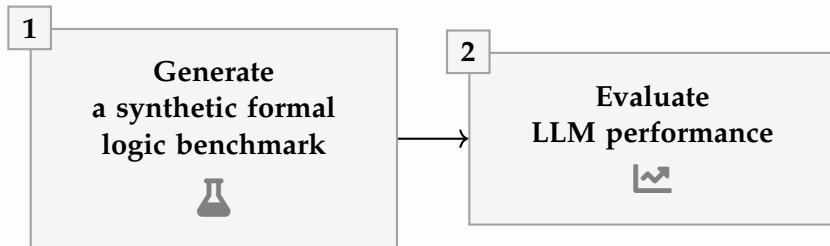


# I

## Introduction



# Objective



# Why synthetic benchmarks?



## Contamination

Often solutions become part of the training corpus of LLMs.

Training data decontamination procedures are often only partially effective.



## Scalability

Manually created benchmarks are difficult to scale both in *size* and *complexity*, requiring significant human effort and financial investments.

## II.I

# Generation (Logic)



# Halpern–Shoham’s Logic

**Halpern–Shoham’s Logic (HS)** is a modal logic where *intervals*, rather than points, are the fundamental states.

**Key idea:** Accessibility between intervals is determined by *Allen’s interval relations*.

Given a linear order  $\mathbb{D} = \langle D, < \rangle$ , a *strict interval* is an ordered pair  $[x, y]$  with  $x, y \in D$  and  $x < y$ .

Two intervals  $[x, y]$  and  $[w, z]$  are compared by their endpoints, and their relation is captured by one of Allen’s modalities.

HS introduces an *existential modality*  $\langle X \rangle$  for each Allen relation  $R_X$ . The six basic relations  $A, L, B, E, D, O$  each have an inverse  $\bar{X}$ , yielding 12 binary relations in total.

# Interval Relations and Modalities

HS modality	Definition w.r.t. the interval structure	Example
$\langle A \rangle$ (adjacent)	$[x, y] R_A [w, z] \Leftrightarrow y = w$	
$\langle L \rangle$ (later)	$[x, y] R_L [w, z] \Leftrightarrow y < w$	
$\langle B \rangle$ (begins)	$[x, y] R_B [w, z] \Leftrightarrow x = w \wedge z < y$	
$\langle E \rangle$ (ends)	$[x, y] R_E [w, z] \Leftrightarrow y = z \wedge x < w$	
$\langle D \rangle$ (during)	$[x, y] R_D [w, z] \Leftrightarrow x < w \wedge z < y$	
$\langle O \rangle$ (overlaps)	$[x, y] R_O [w, z] \Leftrightarrow x < w < y < z$	

**Table:** Allen's interval relations and HS modalities.

# Syntax of HS

**Alphabet:** propositional letters  $\mathcal{P}$ , classical connectives  $\neg, \vee$ , and modalities  $\langle X \rangle$  for  $X \in \mathcal{X}$  wit:

$$\mathcal{X} = \{A, \bar{A}, L, \bar{L}, B, \bar{B}, E, \bar{E}, D, \bar{D}, O, \bar{O}\}.$$

**Grammar:**

$$\varphi ::= p \mid \neg\varphi \mid \varphi \vee \varphi \mid \langle X \rangle \varphi \quad (p \in \mathcal{P}, X \in \mathcal{X}).$$

**Derived connectives:**  $\varphi \wedge \psi \equiv \neg(\neg\varphi \vee \neg\psi)$ ,  $\varphi \rightarrow \psi \equiv \neg\varphi \vee \psi$ ,  $\top \equiv p \vee \neg p$ .

**Universal modality:**  $[X]\varphi \equiv \neg\langle X \rangle\neg\varphi$ .



**Interval model:**  $M = \langle \mathbb{I}(\mathbb{D}), V \rangle$  with  $\mathbb{D}$  a linear order,  $\mathbb{I}(\mathbb{D})$  the set of strict intervals over  $\mathbb{D}$ , and  $V : \mathcal{P} \rightarrow 2^{\mathbb{I}(\mathbb{D})}$  a valuation.

**Truth on an interval  $[x, y]$ :**

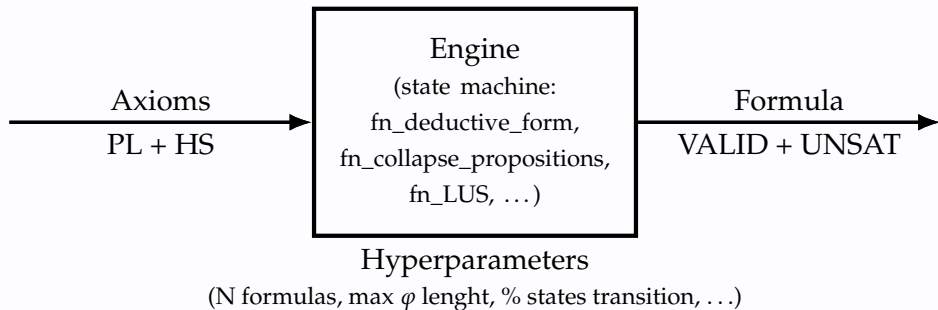
$$\begin{aligned} M, [x, y] \models p &\Leftrightarrow [x, y] \in V(p) \quad (p \in \mathcal{P}), \\ M, [x, y] \models \neg \psi &\Leftrightarrow M, [x, y] \not\models \psi, \\ M, [x, y] \models \psi_1 \vee \psi_2 &\Leftrightarrow M, [x, y] \models \psi_1 \text{ or } M, [x, y] \models \psi_2, \\ M, [x, y] \models \langle X \rangle \psi &\Leftrightarrow \exists [w, z] \text{ s.t. } [x, y] R_X [w, z] \text{ and } M, [w, z] \models \psi. \end{aligned}$$

## II.II

# Generation (Scheme)



# Scheme



## II.III

# Generation (Input)



Axiom	Type
$(a \rightarrow b) \leftrightarrow (\neg b \rightarrow \neg a)$	contrapositive
$(a \rightarrow b) \leftrightarrow (\neg a \vee b)$	implication as disjunction
$a \wedge (b \vee c) \leftrightarrow (a \wedge b) \vee (a \wedge c)$	distributivity of $\wedge$ over $\vee$
$\neg(a \wedge b) \leftrightarrow \neg a \vee \neg b$	De Morgan
$\neg(a \vee b) \leftrightarrow \neg a \wedge \neg b$	De Morgan
$a \wedge a \leftrightarrow a$	idempotence
$a \wedge \top \leftrightarrow a$	identity
$\neg\neg a \leftrightarrow a$	double negation
$a \wedge (a \rightarrow b) \rightarrow b$	modus ponens
$(a \rightarrow b) \wedge \neg b \rightarrow \neg a$	modus tollens
$\neg a \wedge (a \vee b) \rightarrow b$	disjunctive syllogism

**Table:** All propositional axioms used.

Axiom	Type
$\langle \overline{B} \rangle \langle \overline{E} \rangle p \leftrightarrow \langle \overline{E} \rangle \langle \overline{B} \rangle p$	commutativity
$\langle L \rangle p \leftrightarrow \langle A \rangle \langle A \rangle p$	definability
$\neg[B]p \leftrightarrow \langle B \rangle \neg p$	duality
$[A](p \rightarrow q) \rightarrow ([A]p \rightarrow [A]q)$	K axiom
$\langle B \rangle \langle B \rangle p \rightarrow \langle B \rangle p$	transitivity
$\langle B \rangle [\overline{B}]p \rightarrow p$	temporality
$p \rightarrow [A] \langle \overline{A} \rangle p$	inverse of temporality
$\langle A \rangle \langle \overline{A} \rangle p \rightarrow [A] \langle \overline{A} \rangle p$	stability

**Table:** Selected HS axioms, one for each type.

## II.IV

# Generation (Output)



# Output

id_enriched	premise_size	conclusion_size	total_size	compositional_hops	id_base	id_plus
Te1	6	7	14	1	Ae9	Ae41
<b>base:</b> $\neg(t \vee q) \models (\neg t \wedge \neg q)$ <b>plus:</b> $\neg[\text{later}](s) \models \langle \text{later} \rangle \neg s$ <b>enriched:</b> $\neg(t \vee \neg[\text{later}](s)) \models (\neg t \wedge \neg \langle \text{later} \rangle \neg s)$ <b>enriched_valid:</b> $\neg(\neg(t \vee \neg[\text{later}](s))) \wedge (t \vee \langle \text{later} \rangle \neg s)$ <b>enriched_unsat:</b> $\neg(t \vee \neg[\text{later}](s)) \wedge (\neg(\neg t) \vee \langle \text{later} \rangle \neg s)$						
Ti8	7	7	15	2	Ti3	Ae54
<b>base:</b> $\langle \text{meets} \rangle \langle \text{met\_by} \rangle \neg[\text{begins}](q) \models [\text{meets}] \langle \text{met\_by} \rangle \langle \text{begins} \rangle \neg q$ <b>plus:</b> $\langle \text{during} \rangle \neg p \models \neg[\text{during}](p)$ <b>enriched:</b> $\langle \text{meets} \rangle \langle \text{met\_by} \rangle \neg[\text{begins}](\langle \text{during} \rangle \neg p) \models [\text{meets}] \langle \text{met\_by} \rangle \langle \text{begins} \rangle (\neg \neg[\text{during}](p))$ <b>enriched_valid:</b> $[\text{meets}] [\text{met\_by}] [\text{begins}](\langle \text{during} \rangle \neg p) \vee [\text{meets}] \langle \text{met\_by} \rangle \langle \text{begins} \rangle [\text{during}](p)$ <b>enriched_unsat:</b> $(\langle \text{meets} \rangle \langle \text{met\_by} \rangle \neg[\text{begins}](\langle \text{during} \rangle \neg p)) \wedge \langle \text{meets} \rangle [\text{met\_by}] \langle \text{begins} \rangle \neg[\text{during}](p)$						

**Table:** Example of raw output of two sets of valid and unsatisfiable formulas.



II.V

Generation (Engine)



## Formula Structures

Equivalence form:  $\varphi \leftrightarrow \psi$

Implicative form:  $\varphi \rightarrow \psi$

## Formula Roles

Base types: equivalence or implicative

Plus types: only equivalence

## Substitution Principle

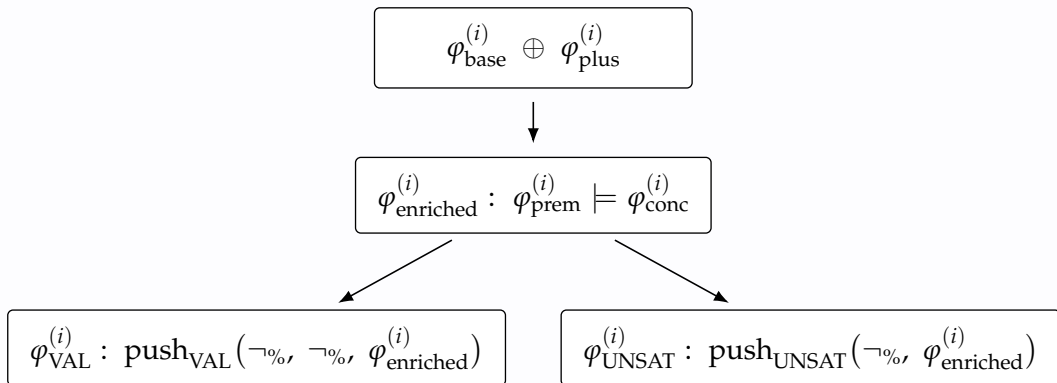
Inspired by uniform substitution:

$$\varphi_{prem}^{base}[p/\psi] \models \varphi_{conc}^{base}[p/\psi]$$

Instead of a WFF  $\psi$ , we use a plus formula:

$$\varphi_{enriched} : \varphi_{prem}^{base}[p/\varphi_{prem}^{plus}]_{partial} \models \varphi_{conc}^{base}[p/\varphi_{conc}^{plus}]$$

## Valid and unsatisfiable formulas



## Tree view of Formula Construction

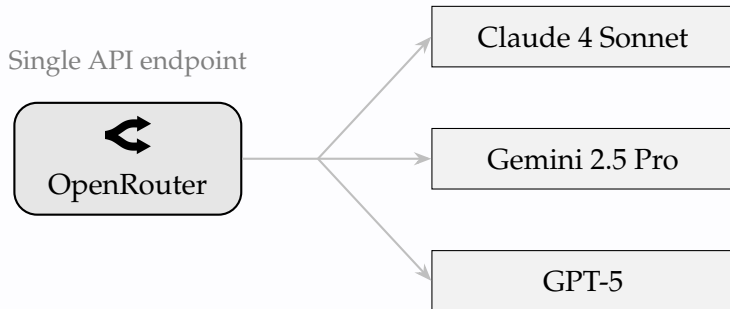
```
Ti77: ([finished](!(<begun_by>(!((p & !([before](t))))))) | //
|      // [met_by](<meets>(<overlaps>((!(p) | !(<before>(!(t)))))))
|-- Ti32b: ([finished]([begun_by]((p & t))) | //
|      // [met_by](<meets>(<overlaps>((!(p) | !(t))))))
|  |-- Ti18b: ([finished]([begun_by](!(p))) | //
|  |  // [met_by](<meets>(<overlaps>(p))))
|  |  |-- Ti15b: (!(q) | [met_by](<meets>(q)))
|  |  |  |-- Ai38b: p |= [met_by](<meets>(p))
|  |  |  |-- Ae16p: q |= !(q)
|  |  |-- Ae30p: <finished>(<begun_by>(p)) |= <overlaps>(p)
|  |-- Ae7p: !((p & t)) |= (!(p) | !(t))
|-- Ae43p: !([before](t)) |= <before>(!(t))
```

# III

## Evaluation



# OpenRouter



# Adopted prompting strategies



## Few-Shot

Provides **complete examples** of problems with their solutions to facilitate learning by analogy without explicit instructions.



## CoT

Guides the model to **decompose** the problem into multiple components to facilitate its resolution.



## Context

A natural language **introduction** to LTL is provided followed by its syntax and semantics. A natural language **introduction** to HS is provided followed by its syntax and semantics.

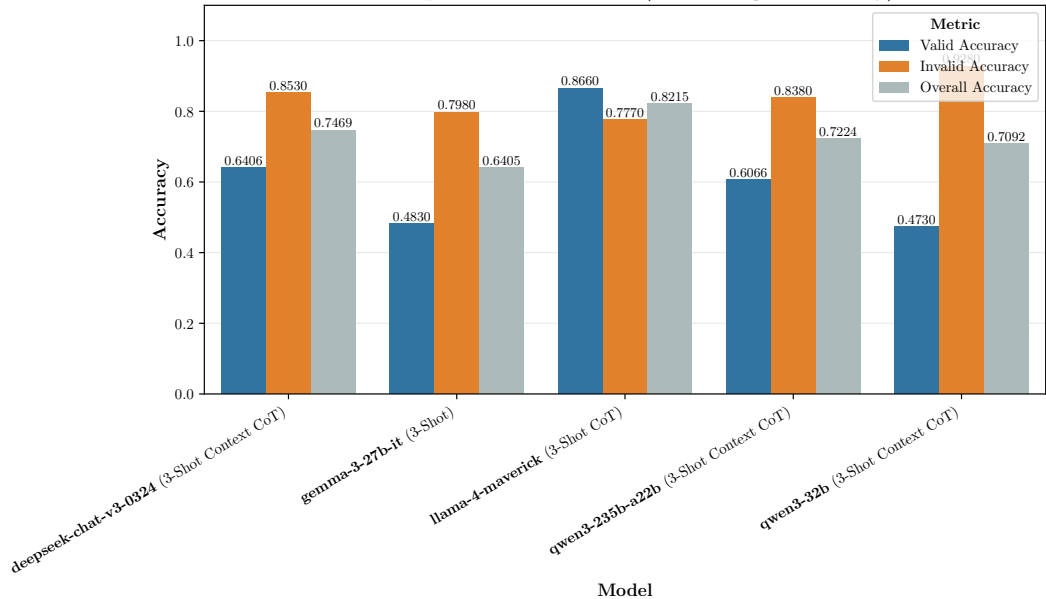
IV

Results



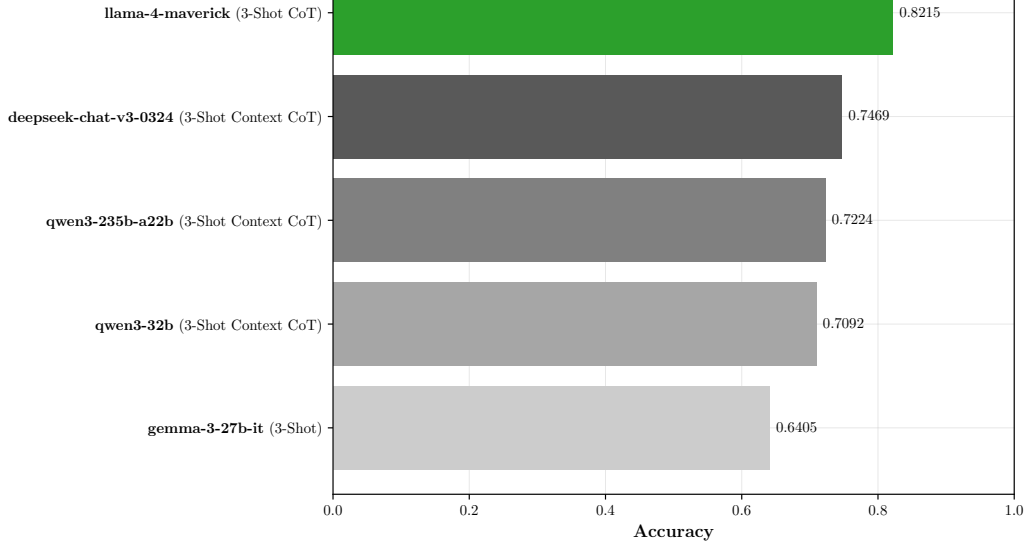


Model performance metrics (best configuration only)

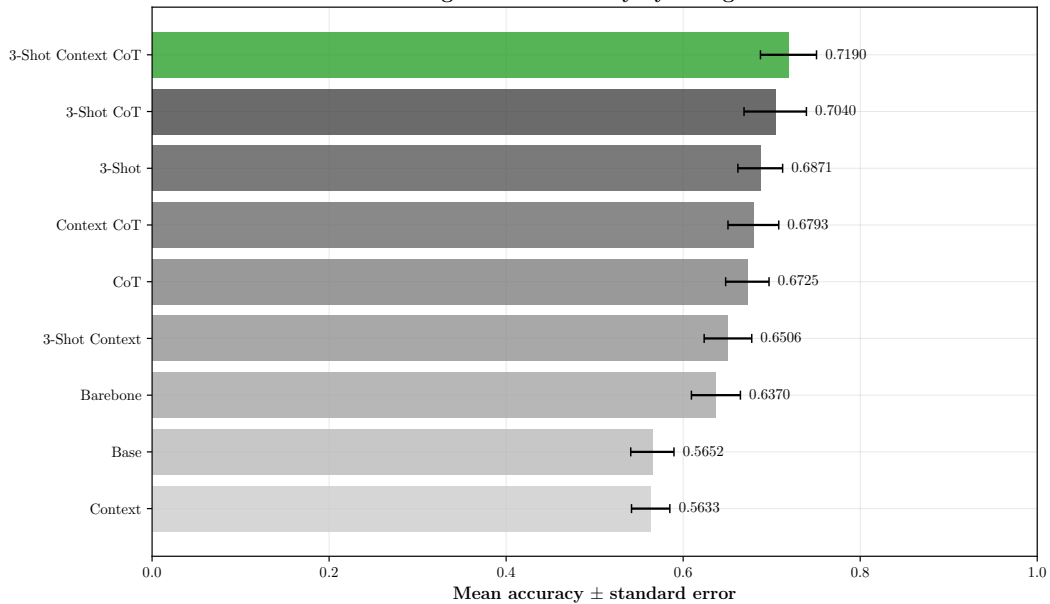


Accuracy (best configuration only)

Model



Average overall accuracy by configuration



V

Appendix

• • • • • • • • •

Overall accuracy for all model configurations

