

Semantic Logic Structured Memory for Multi-Turn LLMs

Anonymous ACL submission

Abstract

Large language models (LLMs) exhibit severe unreliability in multi-turn underspecified language interactions, where critical information is revealed gradually across turns, leading to irreversible early commitments and compounding errors. Although many memory-based approaches have been thus studied recently, they implicitly treat dialogue history as flat token sequences or operate on opaque latent states. We argue that explicit semantic differentiation constitutes a critical representational layer to expose and revise early commitments. To this end, we propose Semantic Logic Structured Memory (SLSM), a language-grounded semantic state representation that explicitly tracks facts, unknowns, assumptions and constraints, which supports explicitly localized state revision across dialogue turns. Unlike slot-based dialogue state tracking or schema-driven memory, SLSM induces semantic structure dynamically from reasoning failures and constraint violations, rather than from predefined schemas. **Evaluated on multi-turn instruction-following benchmarks, SLSM significantly improves reliability under sharded instruction settings, reducing outcome variance without sacrificing task performance. Our results suggest that explicit semantic state tracking is a critical component for robust multi-turn language reasoning under underspecification.**

1 Introduction

Large language models (LLMs) have demonstrated impressive capabilities in many real applications. However, in practice, multi-turn interactions are much more common. LLMs remain fundamentally unreliable in the multi-turn setting, as they are originally trained and deployed to maximize the performance in the single-turn manner, i.e., an objective mismatch.

One of the key challenges in multi-turn LLM interaction is under-specification, which refers to

that each interaction only provides scattered under-specified dialogues, and critical information often incrementally arrives only after several turns. The current LLM behaviors tend to form premature commitments based on incomplete context, treating early responses as fixed anchors rather than provisional hypotheses. Once such early commitments are made, subsequent reasoning is frequently constrained or distorted by them, leading to compounding errors that are difficult to recover from as the dialogue unfolds. It thus requires the model to continuously track and revise its internal understanding of the task.

In response to under-specification, a growing body of work has been dedicated to developing memory mechanisms, aiming to manage information across dialogue turns. Existing approaches typically augment models with external memory buffers, long-context prompting, retrieval-augmented generation, or learned latent states that summarize prior interactions. While these techniques improve information persistence and mitigate context-length limitations, they largely treat dialogue history as undifferentiated text or compress it into opaque representations. As a result, the model’s internal commitments remain implicit and entangled with surface tokens, making them difficult to identify, verify or revise when new information contradicts earlier assumptions.

Although memory mechanisms extend what LLMs can remember, they offer limited support for selectively revising what the model believes, especially under under-specified, evolving task semantics. In multi-turn interactions, effective reasoning requires distinguishing between (i) the information that is firmly established, (ii) the information that remains uncertain, and (iii) assumptions that are provisionally adopted to proceed with incomplete context. However, most existing memory designs do not make such distinctions explicit: facts, assumptions, intermediate inferences, and specula-

084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
tive hypotheses are stored or summarized in a uniform and flat representation. This lack of semantic differentiation prevents the model from performing localized revisions when new evidence arrives, forcing it instead to rely on global regeneration or implicit self-correction. As a consequence, early assumptions tend to persist beyond their validity, and contradictions are resolved implicitly through surface-level text generation, rather than principled state updates.

To address this representation-revision mismatch, we propose Semantic Logic Structured Memory (SLSM), a language-grounded semantic state representation system. Specifically, we design a structured state in each turn for semantic differentiation, consisting of fact, unknowns, assumptions, constraints and plan. It explicitly maintain and distinguish information that requires various actions in the multi-turn planning. To track and update the semantic state after each turn, we propose a semantic distillation module to extract the relevant elements from the semantic state. To selectively revise the semantic state for the next turn interaction, we propose an updater to incrementally revise the extracted relevant elements from the semantic distiller, and creating a next plan correspondingly.

We summarize our major contributions below:

- We develop a new LLM memory framework through rigorously defining the semantic state for the multi-turn LLM interaction. We also design a self-contained algorithm to explicitly tracking and selectively revise the semantic state in each turn.
- We perform extensive experiments to validate our argument that multi-turn LLM interaction suffers from unsatisfactory performance without tracking and selectively revising semantic states, and our proposed SLSM substantially improve the reliability of multi-turn LLM interaction.

2 Related Work

Multi-turn evaluation benchmarks. MultiChallenge (Deshpande et al., 2025)

- LLM Gets Lost (Laban et al., 2025)
- MT-Eval (Kwan et al., 2024)
- AgentBoard (Chang et al., 2024)
- LLM as judge (Tang et al., 2025)

SOTA baselines for multi-turn LLM memory.

3 Proposed Method

In this section, we present Semantic Logic Structured Memory (SLSM) for multi-turn LLM interaction. We define a semantic state representation that explicitly maintain a structured semantic differentiation (Section 3.1), instead of encoding dialogue state implicitly in flat token sequences. Crucially, this semantic state is treated as an explicit first-class object, whose contents are deterministically updated and consulted at every turn,. We introduce an incremental update approach to maintain the semantic states over turns by selectively revising them as new information becomes available in Section 3.2.

3.1 Semantic State and Differentiation

We represent the multi-turn memory state at turn t as a structured tuple

$$S_t = (F_t, U_t, A_t, C_t, P_t),$$

where each component corresponds to a set of elements from distinct types of commitment with explicitly different epistemic status and revision rules, i.e., the fact set F_t , unknown set U_t , assumption set A_t , constraint set C_t and plan set P_t .

Facts. A fact $f \in \mathcal{F}$ is defined as

$$f = (\text{id}, \text{desc}, \text{src}, \text{turn}, \text{status}), \quad (1)$$

where id is an index for the fact, desc is a propositional description explicitly asserted by the user or returned by an external tool, $\text{src} \in \{\text{user}, \text{tool}\}$ denotes provenance, turn denotes the index of the turn deriving this fact, and $\text{status} \in \{\text{open}, \text{closed}\}$. Facts are treated as stable premises and are never introduced by model inference. The status field enables bookkeeping for scope expiration or task termination, without implying epistemic uncertainty.

Unknowns. An unknown $u \in \mathcal{U}$ represents a missing but necessary piece of information:

$$u = (\text{id}, \text{desc}, \text{required_by}, \text{priority}, \text{status}), \quad (2)$$

where desc is a canonicalized description of the missing field, required_by specifies the dependent assumption or constraint, $\text{priority} \in \mathbb{R}^+$ indicates urgency, and $\text{status} \in \{\text{open}, \text{closed}\}$.

If the status of an unknown u is open, it indicates that task-critical information required for safe progression is still missing. As a result, the unknown induces a non-proceed action in the next

178 turn by explicitly triggering a corresponding plan p ,
179 such as clarification or verification, which targets
180 this unresolved element.

181 **If the status is closed**, the unknown has been
182 resolved by newly acquired evidence, and its con-
183 tent is grounded in the semantic state. The resolved
184 information is then instantiated as either an assump-
185 tion a or a constraint c , depending on whether it
186 represents provisional inferred content or an ex-
187 plicit requirement, respectively. Once closed, the
188 unknown no longer participates in action selection.

189 Each unknown u must correspond to an ask-
190 able question and be justified by an explicit depen-
191 dency. Unlike predefined dialogue slots, unknowns
192 in SLSM are dynamically induced by reasoning
193 failure or constraint gaps, and do not assume a
194 fixed schema.

195 **Assumptions.** An assumption $a \in \mathcal{A}$ is a provi-
196 sional hypothesis introduced by the system:

197
$$a = (\text{id}, \text{desc}, \text{basis}, \text{confidence}, \text{turn}, \text{status}), \quad (3)$$

198 where $\text{basis} \subseteq \text{ID}(S)$ records supporting ele-
199 ments in a semantic state S by including their ids,
200 confidence $\in [0, 1]$ is a self-assessed reliability
201 score, and status $\in \{\text{valid}, \text{closed}, \text{contradicted}\}$.

202 **Being valid** means that the assumption is cur-
203 rently provisionally accepted as consistent with all
204 known facts and active constraints, and is there-
205 fore permitted to support downstream reasoning
206 and response generation. Importantly, validity does
207 not imply truth; it only indicates that the assump-
208 tion has not yet been challenged by contradictory
209 evidence.

210 **Being closed** means that the assumption is no
211 longer active in the current semantic state, not be-
212 cause it has been falsified, but because it has be-
213 come irrelevant or superseded by state evolution.
214 For example, when the task focus changes, the as-
215 sumption’s dependent unknowns are resolved in a
216 way that renders the assumption unnecessary, or
217 when the system transitions to a different plan that
218 no longer relies on it. Closed assumptions are ex-
219 cluded from further reasoning without being treated
220 as incorrect.

221 **Being contradicted** means that the assumption
222 has been explicitly invalidated by new evidence,
223 such as newly introduced facts or constraint viola-
224 tions that are incompatible with the assumption’s
225 content. Contradicted assumptions are actively
226 marked as incorrect and may trigger localized revi-

227 sion, including the retraction of downstream infer-
228 ences that depended on them.

229 Assumptions are explicitly retractable and never
230 upgraded to facts. This prevents inferred hypoth-
231 eses from being retroactively treated as ground truth,
232 which is a key source of irreversible early commit-
233 ment in multi-turn interaction.

234 **Constraints.** A constraint $c \in \mathcal{C}$ encodes explicit
235 requirements:

236
$$c = (\text{id}, \text{desc}, \text{weight}, \text{turn}, \text{status}), \quad (4)$$

237 where desc is a declarative description, weight
238 represents how strict this constraint is (higher
239 means trying harder to satisfy it), and status \in
240 $\{\text{satisfied}, \text{closed}, \text{violated}\}$.

241 **Being satisfied** means that the constraint is cur-
242 rently fulfilled by the existing facts and resolved
243 unknowns in the semantic state, and therefore im-
244 poses no further restrictions on action selection or
245 response generation.

246 **Being violated** means that the constraint is in-
247 compatible with the current semantic state. Specifi-
248 cally, no assignment of open unknowns can satisfy
249 the constraint without retracting existing facts or
250 valid assumptions. A violated constraint signals
251 an explicit inconsistency and may trigger localized
252 revision, verification, or corrective clarification in
253 subsequent turns.

254 **Being closed** means that the constraint is no
255 longer active in the current interaction, not because
256 it has been violated or satisfied, but because it has
257 become irrelevant due to state evolution. For ex-
258 ample, when the task focus changes, when the con-
259 straint is superseded by a stronger requirement, or
260 when the dialogue terminates. Closed constraints
261 are excluded from further decision making without
262 being treated as incorrect.

263 **Plan.** A plan $p \in \mathcal{P}$ specifies the system’s opera-
264 tional mode and task for the next turn:

265
$$p = (\text{id}, \text{mode}, \text{target}), \quad (5)$$

266 where mode $\in \{\text{proceed}, \text{verify}\}$ is a specific ac-
267 tion, and target $\in \text{ID}(S)$ denotes a target element
268 in a semantic state S , on which the specific action
269 will be on. The plan component does not encode
270 long-horizon planning, but only the system’s imme-
271 diate epistemic stance toward the next turn. A plan
272 with mode *proceed* marks the successful closure
273 of the target element in memory, indicating that no
274 further verification is required in subsequent turns.

275 **Admissibility of Semantic States.** Beyond defining
276 the structure of semantic states, it is necessary
277 to characterize when a state provides sufficient
278 grounding for safely closing semantic elements.
279 Admissibility captures whether the current semantic
280 state permits assigning a *proceed* mode to a plan instance targeting a specific element, without
281 risking premature commitment or inconsistency.
282 Formally, a semantic state $S = (F, U, A, C, P)$ is
283 said to be *closure-admissible* w.r.t. a target element
284 if and only if all of the following conditions hold:
285

- All unknowns are resolved, i.e., there exists no $u \in U$ with status *open*;
- No assumption is contradicted by current evidence, i.e., there exists no $a \in A$ with status *contradicted*;
- All valid assumptions $a \in A$ satisfy a minimum *confidence* requirement τ ;
- No constraint is violated, i.e., $c \in C$ with status *violated*.

295 If any of these conditions fails, the semantic state is
296 non-admissible for closure, indicating that further
297 clarification, verification, or revision is required
298 before any plan instance may be marked as *proceed*. This admissibility notion is declarative and
299 independent of any particular update or planning
300 mechanism. In Section 3.2, we show how admissibility
301 is checked after each semantic content update
302 and how non-admissible states trigger selective re-
303 vision and alternative epistemic modes.
304

305 **Interaction-Level Semantics.** These components
306 collectively form an explicit semantic interaction
307 layer that governs multi-turn behavior. **Facts**
308 provide stable premises, **unknowns** surface missing
309 information that blocks safe progression, **assumptions**
310 enable provisional reasoning under uncertainty, **constraints**
311 enforce admissible solution space, and **plans**
312 mediate the system’s epistemic stance toward the next turn. More importantly,
313 these elements are jointly consulted to determine
314 whether the system may proceed, must clarify, or
315 needs to revise prior commitments. This ensures
316 that interaction decisions are driven by explicit semantic state rather than implicit inference.
317
318

319 3.2 Incremental Update of Semantic States

320 Having defined the structure and epistemic roles
321 of semantic states in Section 3.1, we now specify

322 how such states evolve across turns. Multi-turn
323 interaction is inherently a stateful process: each
324 user utterance may introduce new evidence, resolve
325 previously missing information, or invalidate provi-
326 sional commitments. Therefore, a semantic mem-
327 ory is only meaningful if it supports principled and
328 selective state revision over turns.
329

330 We define an *incremental, evidence-driven* up-
331 date mechanism that (i) preserves the epistemic
332 typing introduced in Section 3.1, (ii) avoids global
333 rewriting of prior state, and (iii) makes all revisions
334 explicit and auditable. Crucially, state updates must
335 be justified by observable evidence rather than by
336 implicit model preference or generation heuristics.
337

338 **State Transition Interface.** We model semantic
339 state evolution as an incremental transition:
340

$$\Delta S_t \leftarrow \text{UPDATE}(S_{t-1}, x_t), \quad (6)$$

381 where S_{t-1} denotes the semantic state from the pre-
382 vious turn, x_t is the current user input, and ΔS_t is
383 a finite set of localized modifications. The updated
384 state is obtained by applying ΔS_t to S_{t-1} , while
385 all unaffected elements are preserved unchanged.
386 This formulation enforces locality: revisions target
387 specific semantic objects rather than re-encoding
388 the entire dialogue history. It thereby enables se-
389 lective correction of invalid commitments while
390 stably preserving unrelated and already verified
391 state elements.
392

393 **Semantic Evidence Extraction.** Before any ad-
394 missibility checks and planning decisions are made,
395 the current user input x_t is first interpreted as se-
396 mantic evidence w.r.t. the existing state S_{t-1} . This
397 step extracts candidate updates to the semantic com-
398 ponents $(F_{t-1}, U_{t-1}, A_{t-1}, C_{t-1})$ without select-
399 ing a plan, which is formally written as follows:
400

$$(\Delta F_t, \Delta U_t, \Delta A_t, \Delta C_t) \leftarrow \text{EXT}(S_{t-1}, x_t), \quad (7)$$

401 where EXT maps the raw utterance into typed se-
402 mantic evidence, and explicitly isolates the change.
403 Importantly, this extraction step is non-committal:
404 it may introduce, resolve or contradict semantic
405 elements, but it does not determine whether the
406 system should proceed or verify.
407

408 The extraction follow the following principles.
409 First, only propositions explicitly asserted by the
410 user or returned by external tools may enter the
411 fact set. Second, missing but task-relevant infor-
412 mation is surfaced as unknowns rather than silently
413 instantiated. Third, inferred or defaulted content is
414
415

370 introduced, when necessary, as assumptions with
371 explicit confidence. Fourth, explicit requirements
372 expressed in the input are recorded as constraints.
373 These operations update the semantic content of the
374 state while leaving the plan component undecided.

375 We denote the resulting intermediate state as S_t^- ,
376 which reflects all semantic information grounded in
377 x_t but has not yet undergone admissibility checking
378 or action selection.

379 **Evidence-Driven State Revision.** When any of
380 $\Delta F_t, \Delta U_t, \Delta A_t$ and ΔC_t is non-empty, updates
381 are applied in a type-aware and localized manner:

- **Facts** may only be added or marked as closed based on explicit user input or tool output; they are never introduced or altered by the system inference.
- **Unknowns** are marked as *open* when task-relevant information is missing. Open unknowns remain explicitly represented and are never silently instantiated. When missing information is revealed through user interaction or external tools, the corresponding unknown is resolved and converted into an assumption with *valid* status and a system-evaluated confidence. The introduction of this new assumption does not by itself guarantee admissibility, so its consistency with existing assumptions and constraints is evaluated in the admissibility check at the state level.
- **Assumptions** may be marked as *valid*, *closed* or *contradicted*, depending on newly available evidence. Importantly, assumptions are retractable and are never promoted to facts.
- **Constraints** are evaluated against the current state and marked as *satisfied*, *violated*, or *closed* without being implicitly discarded.

406 Revisions affect only the minimal set of state el-
407 ements, whose status is justified by new evidence;
408 unrelated elements are guaranteed to remain un-
409 changed. This non-interference property prevents
410 cascading side effects and supports fine-grained
411 correction of early commitments.

412 **Admissibility Check.** After semantic evidence
413 extraction and localized state revision, the system
414 must determine whether the updated semantic state
415 provides sufficient grounding to safely close seman-
416 tic elements. In SLSM, this decision is formalized
417 through an admissibility check, which evaluates

418 whether a plan instance targeting a specific element
419 may be assigned the *proceed* mode without intro-
420 ducing premature commitment or unresolved incon-
421 sistency. Concretely, a semantic state is considered
422 non-admissible for closure whenever it violates any
423 of the following epistemic preconditions:

1. **Unresolved unknowns.** There exists an un-
424 known $u \in U_t$ with status *open*. Since open
425 unknowns represent task-critical missing in-
426 formation, proceeding would require silent
427 value imputation, which is disallowed.
2. **Contradicted assumptions.** There exists an
428 assumption $a \in A_t$ with status *contradicted*.
429 Such assumptions explicitly encode inconsis-
430 tency with current evidence and must be re-
431 vised before further reasoning.
3. **Excessive assumption uncertainty.** There ex-
432 ists a valid assumption whose confidence falls
433 below a predefined threshold, indicating that
434 the current state relies on unstable speculative
435 commitments.
4. **Violated constraints.** There exists a con-
436 straint $c \in C_t$ with status *violated*, implying
437 that the current state is incompatible with ex-
438 plicit task requirements.

439 Whenever any of these conditions is met, the cur-
440 rent semantic state is not closure-admissible. In
441 such cases, no plan instance can be assigned the
442 *proceed* status in the next turn, and the system
443 must instead favor epistemic memory operations,
444 i.e., verification to revise the semantic state.

445 Admissibility gating prevents premature epis-
446 temic closure in multi-turn interaction. In standard
447 LLM dialogue, speculative choices made under
448 missing information are immediately committed in
449 surface text and become difficult to retract, leading
450 to compounding errors. By decoupling interaction
451 continuity from memory closure, SLSM preserves
452 uncertainty as explicit unknowns and retractable
453 assumptions, which ensures that early speculative
454 decisions do not prematurely solidify or constrain
455 subsequent state revision.

456 **Coupling State Update with Plan Derivation.**
457 Semantic state updates directly determine the sys-
458 tem’s epistemic stance for the next turn by induc-
459 ing the set of admissible memory-management
460 plans. Given the updated semantic state compo-
461 nents F_t, U_t, A_t , and C_t , the plan set P_t is derived

466 by selecting plan instances whose modes are compatible with the current admissibility status of the
467 state. In particular, unresolved task-critical un-
468 knowns and violated constraints induce verification
469 plans, whereas a plan instance ($_, \text{proceed}, \text{target}$)
470 is admitted only if the current semantic state is ad-
471 missible with respect to the target element target .
472 Thus, plan derivation is not governed by ad hoc
473 dialogue policies, but is systematically induced by
474 semantic state properties.
475

476 With explicit state tracking and selective re-
477 vision, state-induced action selection completes the
478 SLSM control loop: semantic evidence updates the
479 state, the state induces appropriate actions, and ac-
480 tions determine whether the system seeks informa-
481 tion, verifies beliefs, or proceeds to solve the task.
482 This closed-loop design is critical for maintaining
483 stability and recoverability in underspecified multi-
484 turn interactions.

4 Experiments

4.1 Objectives

487 **Objective 1: Necessity for Preventing Premature**
488 **Commitment. Hypothesis (O1).** In underspeci-
489 fied multi-turn interactions, *preventing premature*
490 *commitment and avoiding lock-in of early specula-*
491 *tive decisions* requires explicit semantic differen-
492 *tiation with localized revision; memory capacity,*
493 *retrieval, or summarization alone are insufficient.*
494 *This robustness can be achieved without degrading*
495 *final task accuracy.*

496 **Objective 2: Causal Role of Explicit Revision.**
497 **Hypothesis (O2).** The reduction of premature com-
498 mitment in SLSM is causally enabled by explicit
499 tracking and localized revision of assumptions. Dis-
500 abling revision negates the robustness benefits of
501 semantic differentiation, even when all other com-
502 ponents are preserved.

503 **Objective 3: Structural Necessity of Semantic**
504 **Components. Hypothesis (O3).** The semantic
505 components in SLSM (facts, unknowns, assump-
506 tions, constraints) play complementary roles in pre-
507 venting premature commitment and lock-in. While
508 a full factorial ablation is beyond the scope of this
509 submission, we expect that removing a key compo-
510 nent (e.g., assumptions) will measurably degrade
511 robustness under identical model capacity and to-
512 ken budgets; we provide a targeted ablation to par-
513 tially support this hypothesis.

514 **Objective 4: Failure-Mode Specificity. Hy-**
515 **pothesis (O4).** SLSM mitigates a distinct class
516 of structural failure modes—specifically prema-
517 ture assumption commitment and irreversible lock-
518 in—that persist across state-of-the-art memory and
519 agent-based baselines, rather than merely improv-
520 ing average task accuracy.

Experimental Task Checklist (Frozen)

Global Setup (Day 0)

- Implement unified API-only runner (prompt → response → JSONL) 523
- Add caching and retry logic 524
- Fix model and decoding (temperature = 0) 526
- Define baselines: Plain Chat, Naive Prompt 527
- Memory 528
- Fix unified output schema for all benchmarks 529

MultiChallenge (Day 1)

- Integrate official dataset and evaluator 531
- Generate final-turn responses (Plain / Prompt 532
- / SLSM) 533
- Run official metrics 534
- Export main table (overall + category-level 535
- scores) 536
- Log failure and error cases 537

LLMs Get Lost (Day 2)

- Integrate session/trajectory format 539
- Generate per-turn responses (minimal re- 540
- quired) 541
- Compute lostness and task completion metrics 542
- Export main table 543
- Extract 3–5 representative cases 544

MT-Eval (Day 3a)

- Integrate official tasks and evaluator 546
- Generate final responses 547
- Run official scoring 548
- Export summary table (main or appendix) 549

550	AgentBoard (Day 3b)	A Example Appendix	597
551	• Select dialogue/memory-only subset	This is an appendix.	598
552	• Generate responses		
553	• Run corresponding metrics		
554	• Export appendix table		
555	Figures and Paper Integration (Day 3c)		
556	• Consolidate all tables (uniform formatting)		
557	• Generate one normalized bar plot across		
558	benchmarks		
559	• Write experimental setup paragraph (model,		
560	decoding, metrics)		
561	• Write benchmark scope and subset declaration		
562	Explicitly Excluded		
563	• LLM-as-judge experiments		
564	• Multiple model sweeps		
565	• Multi-seed or stochastic decoding		
566	• Full-track AgentBoard evaluation		
567	References	B Related Work	599
568	Ma Chang, Junlei Zhang, Zhihao Zhu, Cheng Yang,	Text-based / retrieval-based memory	600
569	Yuju Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng		
570	Kong, and Junxian He. 2024. Agentboard: An an-		
571	alytical evaluation board of multi-turn llm agents.		
572	<i>Advances in neural information processing systems</i> ,		
573	37:74325–74362.		
574	Kaustubh Deshpande, Ved Sirdeshmukh, Johannes Bap-	Latent / learned memory	601
575	tist Mols, Lifeng Jin, Ed-Yeremai Hernandez-		
576	Cardona, Dean Lee, Jeremy Kritz, Willow E Primack,		
577	Summer Yue, and Chen Xing. 2025. Multichallenge:		
578	A realistic multi-turn conversation evaluation bench-	Structured / schema-based memory	602
579	mark challenging to frontier llms. In <i>Findings of</i>		
580	<i>the Association for Computational Linguistics: ACL</i>		
581	2025	Dialogue State Tracking	603
582	Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei	Belief State / Belief Tracking(POMDP / Dia-	604
583	Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun	logue Policy)	605
584	Liu, and Kam-Fai Wong. 2024. Mt-eval: A multi-	Schema-based / JSON / Table Prompting	606
585	turn capabilities evaluation benchmark for large lan-		
586	guage models. In <i>Proceedings of the 2024 Confer-</i>		
587	<i>ence on Empirical Methods in Natural Language</i>		
588	<i>Processing</i> , pages 20153–20177.	Knowledge Graph / Dynamic KG Updating	607
589	Philippe Laban, Hiroaki Hayashi, Yingbo Zhou, and	Non-monotonic Reasoning / Belief Revision	608
590	Jennifer Neville. 2025. Llms get lost in multi-turn		
591	conversation. <i>arXiv preprint arXiv:2505.06120</i> .		

- 567 **References**
- 568 Ma Chang, Junlei Zhang, Zhihao Zhu, Cheng Yang,
569 Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng
570 Kong, and Junxian He. 2024. Agentboard: An an-
571 alytical evaluation board of multi-turn llm agents.
572 *Advances in neural information processing systems*,
573 37:74325–74362.
- 574 Kaustubh Deshpande, Ved Sirdeshmukh, Johannes Bap-
575 tist Mols, Lifeng Jin, Ed-Yeremai Hernandez-
576 Cardona, Dean Lee, Jeremy Kritz, Willow E Primack,
577 Summer Yue, and Chen Xing. 2025. Multichallenge:
578 A realistic multi-turn conversation evaluation bench-
579 mark challenging to frontier llms. In *Findings of*
580 *the Association for Computational Linguistics: ACL*
581 2025, pages 18632–18702.
- 582 Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei
583 Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun
584 Liu, and Kam-Fai Wong. 2024. Mt-eval: A multi-
585 turn capabilities evaluation benchmark for large lan-
586 guage models. In *Proceedings of the 2024 Confer-*
587 *ence on Empirical Methods in Natural Language*
588 *Processing*, pages 20153–20177.
- 589 Philippe Laban, Hiroaki Hayashi, Yingbo Zhou, and
590 Jennifer Neville. 2025. Llms get lost in multi-turn
591 conversation. *arXiv preprint arXiv:2505.06120*.
- 592 Yuqi Tang, Kehua Feng, Yunfeng Wang, Zhiwen Chen,
593 Chengfei Lv, Gang Yu, Qiang Zhang, and Keyan
594 Ding. 2025. Learning an efficient multi-turn dia-
595 logue evaluator from multiple judges. *arXiv preprint*
596 *arXiv:2508.00454*.