

TIME: Temporally Intelligent Meta-reasoning Engine for Context-Triggered Explicit Reasoning

Susmit Das

Independent Researcher, India
THE COHERENCE INITIATIVE
susmitdas@alum.iisc.ac.in

Abstract

Reasoning-oriented language models typically expose explicit reasoning as a long, front-loaded chain of “thinking” tokens before the main output, either always enabled or externally toggled at inference time. Although this can help on arithmetic, coding, and other multi-step tasks, it is costly, weakens claim-level auditability, and does not allow the model to re-trigger explicit reasoning once presentation has begun. In dialogue, these limitations are compounded by weak sensitivity to temporal structure: unless time is explicitly stated in text, standard models treat replies separated by seconds and replies separated by weeks as equivalent. We introduce TIME (Temporally Intelligent Meta-reasoning Engine), a behavioral alignment framework that learns explicit reasoning as a context-triggered control policy rather than a fixed response mode. TIME augments dialogue with optional ISO 8601 `<time>` tags, tick events that represent silent time passage, and short `<think>` blocks that may appear anywhere in a response. Using a four-phase curriculum, including a small maximally diverse full-batch alignment stage, we train Qwen3 dense models to invoke brief, in-place reasoning bursts only when contextual cues warrant them, while keeping user-facing output compact. We also introduce TIMEBench, a diagnostic benchmark for evaluating reasoning from temporal cues in dialogue. Across 4B–32B scales, TIME improves TIMEBench scores over the corresponding base Qwen3 models in both thinking and no-thinking modes while reducing explicit reasoning tokens by roughly an order of magnitude. Beyond score improvements, TIME induces a distinct behavioral shift: explicit reasoning becomes more compact and more responsive to contextual cues. Code, training data, and benchmark artifacts are publicly available.

TIME training data and code: <https://github.com/The-Coherence-Initiative/TIME>. TIMEBench: <https://github.com/The-Coherence-Initiative/TIMEBench>.

1 Introduction

Reasoning-oriented language models increasingly expose explicit reasoning as visible “thinking” traces, often in the form of chain-of-thought-style text. In many systems, this reasoning appears as a long, front-loaded block at the start of each reply, either always enabled or externally controlled through inference flags. While this can improve accuracy on arithmetic, programming, and other multi-step tasks, it also turns explicit reasoning into a blunt instrument.

Front-loaded reasoning has three recurring limitations. First, it is expensive and verbose, because the model emits many reasoning tokens even for trivial or low-risk requests. Second, it weakens claim-level auditability. A single long monologue can loosely justify many later statements, so the mapping from an individual claim to its justification is often indirect. Third, it is not reactive. Once the model begins presenting an answer, it cannot cleanly re-enter an explicit reasoning phase to re-check assumptions or re-anchor itself if uncertainty arises mid-response.

A better target is a reasoning *policy* rather than a fixed reasoning *mode*: short bursts of explicit reasoning, inserted where needed and re-triggered when new cues appear. Learning such a policy requires signals that indicate latent state change, not only a task-level estimate of difficulty. Temporal structure is a particularly useful probe for this problem because it frequently carries latent contextual information in real interaction: deadlines move, assumptions go stale, priorities shift, and users return after long gaps. Yet unless time is stated explicitly, language models in dialogue settings typically treat a reply after two seconds and a reply after two weeks as equivalent.

In this paper, we ask whether explicit reasoning can be aligned as a *context-triggered control policy*: brief, placeable reasoning bursts that appear

only when temporal or discourse cues warrant them, rather than as an always-on or never-on decoding style. We present **TIME** (Temporally Intelligent Meta-reasoning Engine), a behavioral alignment framework for this setting. TIME introduces three lightweight primitives: optional `<time>` tags on turns, optional short `<think>` blocks that may appear anywhere in a response, and *tick events*, which contain only time metadata and represent silence and passage of time. Because this behavior is rare in existing corpora, we align it with a four-phase curriculum on the Qwen3 family from 4B to 32B parameters. Early phases teach the model to parse the primitives and keep reasoning bursts short and well-delimited, while a final phase uses a small but maximally diverse full-batch alignment set whose only invariant is context-triggered reasoning keyed to temporal and discourse cues.

To evaluate this hypothesis, we introduce TIMEBench, a diagnostic benchmark for evaluating reasoning from temporal cues and latent contextual state in dialogue. TIMEBench is not designed as a benchmark of temporal fact recall or event dating; instead, it uses temporal structure, discontinuity, and anomaly to test whether a model can infer underlying context, recognize when assumptions have become unstable, and adapt its response accordingly. In addition to task success, TIMEBench records structural aspects of generation, including whether `<think>` blocks appear, where they appear, and how often they are used. This allows us to study not only whether models answer correctly, but whether explicit reasoning becomes more selective and better aligned with contextual need.

In summary, the paper contributes (i) a behavioral alignment framework for context-triggered explicit reasoning, (ii) a four-phase curriculum for learning this policy on Qwen3 models, and (iii) TIMEBench, a diagnostic benchmark that uses temporality as a controlled probe for reasoning under latent contextual state change.

2 Related Work

Temporal Cognition and Time-Aware Modeling

Prior work on time in language models primarily treats time as *content*: timestamped facts, event ordering, and knowledge that changes across time. Timestamp-conditioned modeling such as *Time-Aware LMs* (Dhingra et al., 2022) conditions generation on temporal indices to better handle drift and time-sensitive knowledge. Related analyses study

temporal non-stationarity and *temporal generalization*, showing that performance degrades when models are evaluated on future slices beyond their training distribution (Lazaridou et al., 2021). Recent evaluation frameworks likewise probe time awareness over large event collections, including methods that test recall and calibration under temporal constraints (Herel et al., 2024).

Benchmarks have also expanded from synthetic temporal arithmetic to broader temporal competence. ChronoSense evaluates temporal relations and temporal arithmetic, while TimE and EvolveBench test temporal reasoning under more realistic dynamics, including dialogue settings, invalid timestamp handling, and temporal misalignment between inputs and queries (Islakoglu and Kalo, 2025; Zhu et al., 2025; Wei et al., 2025). In parallel, retrieval and structured-knowledge approaches incorporate temporal graphs to represent evolving facts, including temporal GraphRAG-style systems for time-sensitive retrieval and updates (Han et al., 2025). These lines strengthen temporal factuality, ordering, and retrieval, but they still primarily treat time as part of the *world state*. TIME instead uses temporality as a cue to *interaction state*: an observable proxy for latent context change that can invalidate assumptions, alter response requirements, and trigger re-anchoring.

Explicit Reasoning Traces and Reasoning Control

Chain-of-thought prompting (Wei et al., 2022) improves multi-step task performance and provides explicit reasoning traces that are often treated as operational explanations. However, a growing body of work shows that such traces are not reliably faithful to the computations that drive final answers, limiting their value as audit artefacts and motivating methods that measure or improve faithfulness (Paul et al., 2024; Tutek et al., 2025; Barez et al., 2025; Chen et al., 2025). Beyond faithfulness, long front-loaded reasoning also creates practical problems: cost, latency, and weak claim-level attribution when a single rationale must support many downstream statements.

Most deployed systems still treat explicit reasoning as a mode rather than a policy. Hybrid reasoners allow inference-time toggles, but the model itself does not decide when to reason. Recent work begins to learn this decision, for example by training hybrid reasoning models to select whether to think, often using reinforcement learning or controller

objectives (Jiang et al., 2025). Other approaches control *how much* a model reasons through length-conditioned policy optimization (Aggarwal and Welleck, 2025). Orthogonally, DOTS searches over reasoning action trajectories, and RL-based interleaving trains models to alternate thinking and answering to reduce time-to-first-token and improve efficiency (Yue et al., 2025; Xie et al., 2025).

These advances move toward adaptive reasoning, but they are typically framed around task difficulty, length budgets, or action planning. TIME targets a different axis: when explicit reasoning should appear in response to latent context change. Temporality is useful here not as an end task, but as a structured probe for when assumptions become unstable and brief, in-place reasoning or re-anchoring may be needed.

3 Framework

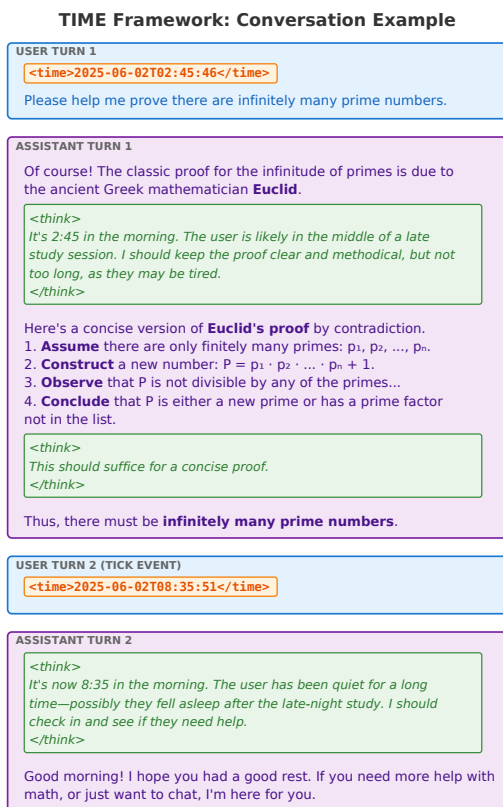


Figure 1: A TIME conversation sample.

Overview

TIME aligns reasoning models to treat explicit reasoning as a *context-sensitive resource* rather than a fixed response mode. Instead of emitting one long, front-loaded reasoning block at the start of every reply, the model learns to invoke small reasoning

Primitive	Purpose
<code><time></code>	Absolute timestamp on a turn using ISO-8601 format.
<code><think></code>	<i>Optional</i> short reasoning bursts inserted anywhere in a response.
Tick	A turn that contains only a <code><time></code> tag, marking silent time advance.

Table 1: TIME primitives that structure behavior during training and evaluation.

bursts only when cues indicate that assumptions may need to be checked or revised. Temporal structure is a useful source of such cues because it often carries latent contextual state in interaction, though the learned policy is not limited to time and can also be triggered by non-temporal cues in text.

Primitives

TIME introduces three textual primitives in the conversation stream (Table 1): `<time>` tags in user inputs, optional `<think>` blocks in model outputs, and *tick events*, i.e., user turns that contain only a timestamp and represent silent time passage.

Timestamped Turns (`<time>`)

A user turn may include a timestamp tag:

```
<time>YYYY-MM-DDTHH:MM:SS</time>  
User message
```

A timestamped turn pairs an utterance with explicit interaction time, making elapsed intervals and delayed returns visible to the model. Timestamps mainly expose temporal structure during training and evaluation; after alignment, the same policy can still be triggered by purely textual cues when explicit temporal metadata is absent.

Transient Reasoning (`<think>`)

TIME uses `<think>` blocks as explicit reasoning traces that are **short, optional, placeable, and repeatable**. A response may contain zero, one, or multiple `<think>` blocks, and they may appear anywhere in the turn. The design target is not a single full reasoning transcript, but a sequence of local reasoning actions: brief checks, inferences, revisions, or decisions inserted exactly where needed. This lets the model re-trigger explicit reasoning multiple times within the same response as new cues appear. Improved auditability is a useful byproduct of these smaller local reasoning traces.

Tick Events

A **tick** is a user turn that contains only a timestamp:

```
<time>YYYY-MM-DDTHH:MM:SS</time>
```

Unlike a timestamped turn, a tick contains no user message. It represents time passing without new text and provides a simple mechanism for modeling silence, persistence, and delayed follow-up.

Temporal Behavior as Sequence Control

Together, `<time>`, `<think>`, and ticks turn dialogue into a temporally anchored behavioral sequence. Rather than learning only *what* to answer, the model learns *when* to reason and *how much* explicit reasoning to deploy. Figure 1 shows a representative thread in which elapsed time changes the latent interaction state and the model adapts accordingly. Further examples are included in Appendix B.

4 Methods

4.1 Curriculum and Alignment Protocol

We train TIME with a staged supervised fine-tuning (SFT) curriculum on the Qwen3 family at four sizes (32B, 14B, 8B, 4B). The objective is to learn a **reasoning invocation policy**: emit `<think>` traces only when contextual cues warrant them, keep them brief and well-delimited, and allow re-triggering later within the same response.

Qwen3 is a hybrid reasoner that already supports both explicit reasoning (via `<think>` blocks) and direct instruction-following without explicit reasoning traces. However, these behaviors are largely mode-bound: thinking mode produces long, front-loaded traces, while no-thinking mode suppresses them entirely. This makes Qwen3 a suitable starting point for learning a finer-grained policy over *when* explicit reasoning should appear, rather than *whether* it is enabled at all.

Direct SFT onto the final behavior is unstable and tends to collapse back to verbose or templated traces. We therefore use a staged curriculum that first teaches structural prerequisites and then consolidates the target policy with a small, high-entropy alignment set.

Method overview. The curriculum proceeds through four stages: structural seeding, temporal exposure, contextual modulation, and a final full-batch alignment stage. Phases 1–3 introduce the primitives, encourage local placement of reasoning,

and train modulation under temporal and discourse variation. Phase 4 concentrates gradients on the target invariant: context-triggered explicit reasoning.

Training setup. Training uses QLoRA (Detmers et al., 2023) on a single NVIDIA RTX Pro 6000 Blackwell (96GB). Phases 1–3 share a fixed setup (rank 32, $\alpha = 32$, dropout 0.05, AdamW-8bit, learning rate 2×10^{-5} , effective batch size 32, 3 epochs, gradient checkpointing) with 25% replay from prior phases. Phases 1–3 use synthetic data generated with GPT-4o and Gemini 2.5 Flash via template-guided pipelines, with progressive exposure to structural primitives, temporal gaps, and multi-turn dynamics. Details are documented in Appendix C.

Phase 1: Structural Seeding

Data Size: 2,188 train / 387 test

Phase 1 introduces the primitives and output format. Single-turn prompts pair `<time>` metadata with short `<think>` bursts, encouraging compact, well-delimited reasoning instead of long monolithic traces.

Phase 2: Temporal Exposure

Data Size: 5,291 train / 935 test (+25% replay)

Phase 2 introduces two-turn dialogues with time gaps and tick events. The model learns to condition on temporal metadata to revise assumptions after silence, while also suppressing unnecessary verbosity when no update is required.

Phase 3: Contextual Modulation

Data Size: 5,878 train / 1,039 test (+25% replay)

Phase 3 extends to multi-turn settings and trains both suppression and re-triggering of `<think>` blocks under changing context. Increased tick frequency (approximately 33%) encourages reliance on both temporal and non-temporal cues, creating behavioral headroom for final alignment.

Phase 4: Gradient-Aligned Convergence via Maximal Diversity

Data Size: 128 conversations

Phase 4 is the decisive alignment step. We use a small but maximally diverse set (Table 2) whose only shared property is the target policy: `<think>` bursts are triggered by temporal or discourse cues and placed only where needed. Each conversation is multi-turn, so supervision covers many state transitions, including long-gap re-engagement and

Aspect	Variation Injected
Topics	Speculative fiction, late-night coding, therapy, anime debates, hostile rants, professional emails, gaming chat, etc.
Styles	Terse ↔ verbose; academic markdown ↔ free-form prose; emoji-laden ↔ plain text; bullets, numbered steps, code blocks; often flipped mid-dialogue.
Turns	1 to 7 turn threads, with context changes, style switches, and tick events after long gaps.
Noise	Typos, abrupt topic shifts, contradictory instructions, and diverging user preferences.

Table 2: Maximal surface diversity in Phase 4 isolates the intended behavioral invariant.

Setting	Value
Effective Batch Size	128 (entire dataset)
Optimizer	AdamW-8bit
LR / warm-up	1.5×10^{-4} , 6 warm-up steps
Max epochs run	35 / 36 / 40 / 46 (32B / 14B / 8B / 4B)
Checkpoint criterion	Earliest checkpoint in loss band [1.045, 1.050]
LoRA	Rank 32, $\alpha = 32$

Table 3: Phase 4 tuning setup. Clean convergence is assessed by monitoring degeneracy.

tick-driven silence. Replay is disabled in this phase to avoid reintroducing earlier regularities.

Deterministic Full-Batch Fine-Tuning

We train with effective batch size 128, so each update sees the entire alignment set and removes sampling variance. Table 3 reports the configuration.

Across all four scales, we observe a consistent tradeoff between policy acquisition and degeneracy during Phase 4. If training is stopped too early, the model has not yet reliably learned the target reasoning policy; if training continues too far, degeneracy, including infinite loops, <think> format bleed, and style collapse, becomes increasingly frequent, especially as loss approaches 1.0 or below. We therefore use an empirically identified stopping boundary at which the policy is already learned while degeneracy remains in check. In practice, we select the earliest checkpoint whose loss enters the target band [1.045, 1.050], which corresponds to epochs 18 / 24 / 30 / 31 for the 32B / 14B / 8B / 4B models, respectively. The consistency of this operating region across model scales suggests that Phase 4 admits a narrow but usable stability window for aligning the target behavior.

Design Rationale

The goal is to align a policy, not to expand knowledge. Full-batch updates require the effective batch size to equal the dataset size, and this behavior has a narrow stability window: over-optimization leads to a rise in degeneracy. A compact, high-entropy batch concentrates gradients on the intended invariant while remaining stable.

Phase 4 is structured so that the only gradient-aligned invariant is the reasoning policy itself, encoded through a maximally diverse batch of 128 samples. The model must detect contextual cues such as elapsed time, tick events, contradictions, and discourse shifts, decide whether to emit a bounded <think> burst and where to place it, and keep the user-facing presentation appropriately formatted. Unlike mini-batch training, where sampling variance can allow incidental correlations such as topic or formatting artifacts to dominate updates, full-batch alignment over the entire diverse set has a regularizing effect that suppresses such spurious signals.

The resulting checkpoints, used in all subsequent evaluation, are denoted **TIME-32B / 14B / 8B / 4B**. Further training details are documented in Appendix C.

4.2 Evaluation Method: TIMEBench

TIMEBench is a 77-scenario diagnostic benchmark for **reasoning in dialogue from temporal cues and latent contextual state**. Its central question is not whether a model can recall dated facts, but whether it can use temporal structure to infer underlying context: what has changed, which assumptions may no longer hold, and how the response should adapt. We use temporality because it is a frequent and observable proxy for context in real interaction: deadlines move, users disappear and return, assumptions go stale, and coordination depends on elapsed time and timezone. These shifts are often invisible in token-only chat logs unless

Category	Temporal pattern probed	What successful behavior requires
Chronological Retrospection ¹	Non-trivial temporal reconstruction across turns, including partial logs, delayed write-ups, and implicit event windows.	Reconstruct a latent timeline from scattered conversational evidence rather than relying on surface order alone; infer exact or bounded temporal relations when the answer depends on sequencing.
Invalid Time Detection ¹	Impossible calendar values such as non-existent dates.	Detect that the timestamp itself is invalid and explicitly register the anomaly.
Temporal Adaptivity ²	Shifts in urgency or actionability caused by imminent deadlines, passed deadlines, or short remaining wait times.	Adapt the response style to temporal pressure: be urgent when minutes matter, withhold unnecessary interventions when help is imminent, and switch to fuller explanation once urgency has passed.
Temporal Contextual Awareness ²	Time cues that imply situational context, such as festivals, holidays, or late-night study settings.	Infer latent context from time itself and use it to shape interpretation and response tone, rather than answering as if the query were temporally generic.
Temporal Flow Anomaly Detection ¹	Non-monotonic or implausible temporal structure, including backward timestamps and extreme jumps across years or centuries.	Notice that conversational time no longer behaves normally and treat the anomaly as a trigger for explicit scrutiny or re-anchoring, even if the model then continues assisting.
Time Gap Awareness ²	Long but plausible silence between turns, often combined with topic drift or likely changes in the user’s situation.	Recognize that earlier assumptions may be stale, re-anchor to the new moment, and avoid treating the earlier context as if nothing has changed.
Timezone Sensitivity ¹	Offset changes across turns that imply changes in local context, location, or circadian state.	Use timezone shifts as reasoning evidence, for example to infer approximate location, travel progress, or appropriate advice in the user’s new local context.

¹ Category is out-of-distribution relative to training.

² Category reflects curriculum-intended behavior, though all scenarios remain unseen during training.

Table 4: TIMEBench diagnostic axes. Each category isolates a distinct temporal phenomenon that can alter conversational interpretation, invalidate prior assumptions, or require behavioral adaptation, enabling model-agnostic evaluation of reasoning under temporal structure, discontinuity, and anomaly.

time is represented explicitly. Unlike standard task-focused reasoning benchmarks in domains such as mathematics, coding, or QA, TIMEBench is designed to isolate reasoning under temporal discontinuity, anomaly, and contextual re-anchoring rather than temporal fact recall or time-sensitive knowledge retrieval. This makes it a natural model-agnostic testbed for reasoning from latent state information and, in this paper, a proxy for evaluating context-triggered reasoning. To avoid train-test circularity, all scenarios remain unseen during training, and several categories are completely out-of-distribution relative to the curriculum. Beyond binary category-specific objectives, TIMEBench also records structural aspects of generation, allowing us to evaluate both objective performance and changes in reasoning behavior.

Structure. TIMEBench contains **seven diagnostic categories** with **11 scenarios each**. Rather than treating time as factual content to be recalled, these categories probe how temporality can alter conversational interpretation, invalidate assumptions, and require behavioral adaptation (Table 4).

Representative examples for each diagnostic category are provided in the appendix (Appendix E).

Scenario design. Each scenario is a short, naturalistic dialogue thread with timestamps on user turns; the model generates only the final assistant turn. Scenarios are grouped by categories, but each has its own binary objective specifying what counts as success in that instance, such as detecting an implausible timestamp or revising assumptions after a long gap. The objectives target contextual inference and response adaptation from temporal cues.

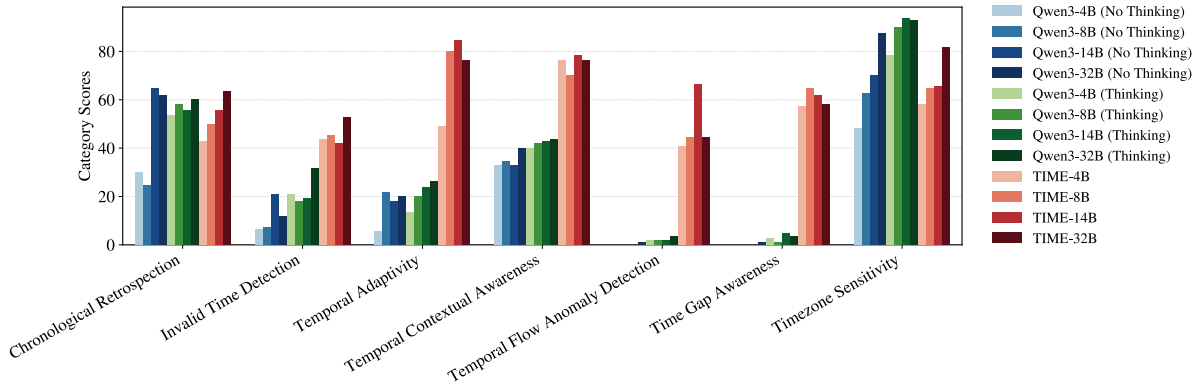


Figure 2: Breakdown of TIMEBench scores by diagnostic category and model size.

Protocol. Each model is evaluated with **10 trials per scenario** using PCG64-derived seeds (770 runs total). Trials are scored with **binary objectives** (0 or 1) using **LLM-as-a-Judge**: a third-party judge that sees only the model response and the objective, and is blind to the original prompt, timestamps, and prompt formatting. Blind binary judging with an LLM makes evaluation tractable at this scale. To reduce evaluator variance, we aggregate trial scores to scenario scores (mean of 10), then category scores (mean over scenarios, scaled to %), and finally an overall TIMEBench score (mean over categories). We estimate **95% confidence intervals** via stratified bootstrapping (10,000 resamples) by resampling scenario scores within each category and recomputing all aggregates. These controls stabilize the evaluation, but because scoring is still performed through LLM-as-a-Judge, the reported scores remain judge-derived estimates rather than error-free measurements.

Behavioral instrumentation. In addition to correctness, TIMEBench supports a structural audit of generation. We extract whether `<think>` appears, where it appears within the response, the number of `<think>` blocks, reasoning and output token counts, light versus heavy markdown usage, format bleed, and degeneracy indicators including infinite repetition. This second view lets us study not only whether performance improves, but whether the model’s reasoning policy changes alongside those gains. In particular, it allows us to test whether improvements co-occur with more selective, more localized explicit reasoning rather than simply longer outputs or heavier reasoning traces. This separates benchmark performance from reasoning style, allowing us to ask whether gains reflect better reasoning allocation rather than verbosity alone.

5 Results

Size	Qwen3		TIME
	No-Thinking	Thinking	
4B	17.53	30.13	52.60
8B	21.56	32.99	59.87
14B	29.48	34.42	64.80
32B	31.82	37.40	64.81

Table 5: TIMEBench scores (out of 100) across models.

We evaluate all models on **TIMEBench**, comparing our aligned models (**TIME-4B / 8B / 14B / 32B**) against the corresponding base Qwen3 checkpoints in both no-thinking mode (via the `/no_think` suffix) and full thinking mode. All models use the same decoding parameters: **temperature 0.6**, **top-p 0.95**, **top-k 20**, and **min-p 0**, following the reasoning-evaluation settings recommended in the *Qwen3 Technical Report* (Qwen-Team, 2025). GPT-5.2 (2025-12-11 checkpoint), accessed via the OpenAI API, serves as the blind judge for all binary evaluations.

Table 5 reports aggregate **TIMEBench scores** on a 0–100 scale. Across all four model sizes, TIME outperforms both Qwen3 baselines as visualized in Figure 2. At 4B, TIME scores 52.60 compared with 30.13 for Qwen3 in thinking mode and 17.53 in no-thinking mode. At 32B, TIME scores 64.81 compared with 37.40 and 31.82, respectively. The same pattern holds at 8B and 14B. In other words, aligning the reasoning policy substantially changes performance on a benchmark that tests reasoning from temporal cues, even though the underlying base model family is unchanged.

Bootstrapped **95% confidence intervals** show that these gains are not concentrated in only a few

Model	Benchmark Score	Runs w/ <think> (%)	Mean Number of <think>	Think Position			Mean Number of Tokens		Runs w/ Degeneracy (%)
				Start (%)	Mid (%)	End (%)	Thinking	Output	
No-Thinking	31.82	0.00	0.00	—	—	—	0.00	608.96	4.42
Thinking	37.40	99.2	0.99	100.0	0.0	0.0	910.52	1573.47	18.18
Phase 1	42.47	99.5	0.99	100.0	0.0	0.0	803.52	1434.56	13.90
Phase 2	56.88	95.6	1.12	70.7	29.1	0.2	76.59	362.45	4.68
Phase 3	52.08	89.2	1.25	55.0	44.6	0.4	52.94	294.51	0.78
TIME	64.81	80.6	1.67	24.1	75.6	0.2	84.16	332.64	3.64

Table 6: Phase-wise ablation: Structural and behavioral metrics for Qwen3-32B across curriculum stages and final alignment to TIME-32B. Think-position percentages are computed over total observed <think> blocks.

scenarios. For TIME-4B, the interval is 44.55–60.39; for TIME-8B, 53.38–66.23; for TIME-14B, 59.09–70.39; and for TIME-32B, 58.18–71.17. The corresponding Qwen3 thinking-mode intervals are lower at every size: 23.90–36.36 (4B), 26.88–39.09 (8B), 28.44–40.65 (14B), and 31.56–43.51 (32B). None of the TIME intervals overlaps its matched thinking baseline.

We also test significance at the scenario level using a **Wilcoxon signed-rank test** on per-scenario mean scores, comparing each TIME model against the corresponding Qwen3 thinking baseline. The improvement is statistically significant at every size ($p < 0.001$): 4B, $p = 3.8e-4$; 8B, $p = 1.9e-5$; 14B, $p = 1.6e-6$; and 32B, $p = 5.0e-7$.

Taken together, these results show that TIME improves performance on TIMEBench consistently across model scales. Within the scope of this benchmark, the gains support the view that a context-triggered reasoning policy helps models reason more reliably from underlying context signaled by temporal cues than fixed thinking or no-thinking modes. A full category-level and structural breakdown across model sizes is provided in [Appendix F](#).

5.1 Phase-wise Ablation in 32B: Diagnostic and Structural Metrics

[Table 6](#) separates two aspects of behavior that are related but not identical. The **score** column summarizes TIMEBench performance across diagnostic categories. The remaining columns capture **structural behavior**: whether explicit reasoning is used, where it appears, how many tokens it consumes, and how often decoding degenerates. Together with [Figure 3](#), this shows how the curriculum changes both benchmark performance and reasoning deployment.

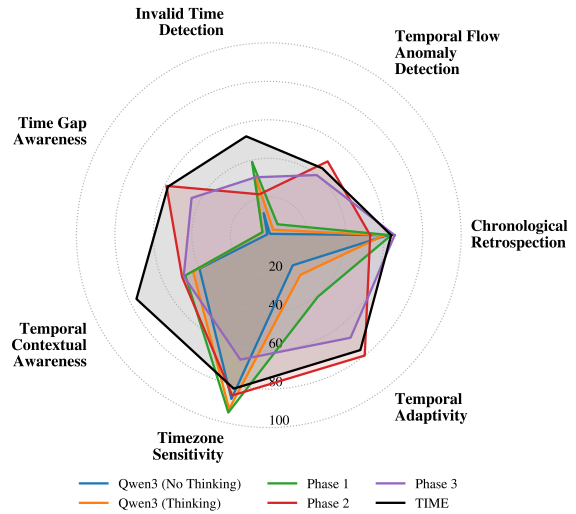


Figure 3: Evolution of core TIMEBench competencies in Qwen3-32B across curriculum phases.

Baselines. The no-thinking baseline handles some local cues reasonably well, including time-zone sensitivity and short-range retrospection, but performs poorly on categories that require detecting discontinuity, anomaly, or stale assumptions, such as flow anomaly detection and gap awareness, and it misses most invalid timestamps. Enabling full thinking improves some forms of temporal scrutiny, especially invalid time detection, but these discontinuity-sensitive categories remain weak. Structurally, the thinking baseline behaves like a standard front-loaded reasoner: nearly every run contains a single initial trace, averaging about 911 reasoning tokens and 1,573 output tokens, with 18.18% degeneracy.

Phase 1 (structural seeding). Phase 1 mainly teaches the new markup and output format. Scores improve modestly, but the basic reasoning pattern

remains largely unchanged. Nearly all runs still contain a single preamble-style trace, with long explicit reasoning (about 804 tokens) and large outputs (about 1,435 tokens). Degeneracy falls to 13.90%, largely through reduced formatting errors.

Phase 2 (temporal exposure). Phase 2 is the first stage where diagnostic and structural shifts appear together. Categories tied to discontinuity and anomalous temporal structure, including adaptivity, flow anomaly detection, and gap awareness, improve sharply. At the same time, the reasoning budget drops to about 77 tokens per run and output length to about 362 tokens. Mid-turn placement now accounts for 29.1% of observed <think> blocks, and degeneracy falls to 4.68%. Explicit reasoning is no longer used only as a fixed preamble.

Phase 3 (contextual modulation). Phase 3 emphasizes restraint and stability. Chronological retrospection and adaptivity remain strong, while the reasoning budget falls further to about 53 tokens and output length to about 295 tokens. Degeneracy reaches its minimum at 0.78%, and mid-turn placement becomes more common than front-loaded placement. However, some discontinuity-sensitive categories lose part of their Phase 2 gains, suggesting that stronger suppression can remove useful re-anchoring when subtle shifts still warrant it.

TIME (Phase 4 alignment). Final alignment in TIME-32B restores performance on anomaly- and discontinuity-sensitive categories while preserving intermittent reasoning. It achieves the best overall score, uses about 84 reasoning tokens and 333 output tokens per run, and emits <think> blocks in roughly 81% of runs, now predominantly mid-turn (75.6%) rather than front-loaded (24.1%). Degeneracy remains low at 3.64%. Overall, the final model retains explicit reasoning but no longer uses it as a long default preamble; instead, it deploys shorter traces, more often mid-response, in settings where temporal cues indicate that assumptions may need to be re-anchored. Fine-grained diagnostics are reported in [Appendix F](#).

6 Discussion

The results support viewing explicit reasoning as a learned control policy rather than a fixed decoding style. Under TIME, Qwen3 shifts from long, always-on reasoning traces to shorter, better-placed bursts tied to cues such as deadlines, silence, inconsistency, and temporal anomaly. The strongest

gains appear where temporal structure carries latent contextual state, suggesting that the model is not merely handling time as content, but using it to re-anchor assumptions and adapt response behavior.

This is the motivation behind the name TIME: *Temporally Intelligent Meta-reasoning Engine*. The framework uses temporal cues not as an end task in themselves, but as a practical signal to infer and reason over latent state at the meta level of the interaction: what may have changed, which assumptions may no longer hold, and how the response should adapt.

The structural metrics reinforce this reading. Mid-turn <think> placement becomes much more common, reasoning-token budgets drop sharply, and degeneracy decreases. Together, these shifts indicate that the model becomes more selective about when and where it uses explicit reasoning, which is especially useful in interactive or agentic settings where fast responses sometimes need brief re-anchoring rather than a verbose reasoning pass.

Although TIME is trained with explicit temporal primitives, the learned policy is not limited to timestamps at inference time. Temporal metadata provides one useful probe for latent context, but the same controller can also react to purely textual cues such as contradiction, uncertainty, or goal change.

7 Conclusion

TIME introduces a lightweight framework for learning explicit reasoning as a context-triggered policy in dialogue models. By combining timestamps, ticks, and optional <think> blocks with a staged curriculum and a diversity-driven full-batch alignment step, it turns fixed reasoning modes into selective, placeable reasoning behavior.

TIMEBench provides the matching evaluation lens. Across seven diagnostic categories, it shows that TIME improves reasoning from temporal cues while shifting generation toward shorter traces, more selective placement of <think> blocks, and fewer degenerate outputs. Taken together, TIME and TIMEBench show that explicit reasoning is not only something models can do, but something alignment can teach them to deploy selectively.

Acknowledgments

This work was conducted independently by the author. THE COHERENCE INITIATIVE is the author’s independent research initiative under which the code, data, and related artifacts are maintained.

Limitations

While TIME provides a practical recipe for context-triggered reasoning, its scope and claims are deliberately narrow.

All experiments use dense hybrid reasoners from the Qwen3 family, which already support both instruct-style and explicit reasoning modes. This makes them a natural starting point for learning a finer-grained reasoning policy, but it also limits generality. We do not claim that the same curriculum and alignment recipe will transfer unchanged to purely instruct models that were not pre-trained for explicit thinking traces, nor do we evaluate TIME on mixture-of-experts hybrid reasoners such as Qwen3-30B-A3B or Qwen3.5-35B-A3B, where routing dynamics, active-parameter sparsity, and adapter placement may interact differently with the proposed curriculum.

The study is also narrow in evaluation scope. TIME is evaluated only on TIMEBench, so we do not establish whether the learned policy helps, hurts, or leaves unchanged performance on broader task-based benchmarks such as mathematics, coding, or multi-step tool use. More broadly, the framework remains under-explored: future work should compare context-triggered explicit reasoning against standard reasoning modes on general reasoning benchmarks, and study whether similar trigger policies can be learned from cues beyond temporality. Prior work suggests that surrounding discourse context can affect faithfulness and discourse interpretation (Miao et al., 2024; Wan et al., 2025), making trigger selection itself a broader phenomenon worth studying.

Our alignment protocol is supervised and adapter-based, and its final stage introduces a methodological direction that extends beyond TIME. In particular, we use a 128-example, maximally diverse full-batch alignment stage to isolate a single behavioral invariant, namely context-triggered local reasoning. Empirically, this yields a usable stability window and suggests that low-data full-batch alignment can shape behavioral policies, but in this paper we study that idea only in the specific setting of context-triggered explicit reasoning. Its broader transfer across model families, policy types, and behavioral targets remains open. We also do not explore reinforcement learning, bandit-style objectives, or reward-model-based variants that could optimize the trade-off between accuracy, latency, and reasoning cost more directly.

TIMEBench should be understood as a preliminary diagnostic component of the larger TIME study rather than a fully developed standalone benchmark effort. It is sufficient for the present paper’s purpose, but limited in construction scope: it contains 77 scenarios, covers only a finite set of diagnostic patterns, and was developed alongside the framework rather than independently. To stabilize measurement at this scale, we evaluate each scenario over multiple trials and aggregate scores. All TIMEBench scores are obtained using *LLM-as-a-Judge* rather than human annotators. The judge is blind to prompts and timestamps, and we use binary objectives, repeated trials, and bootstrapped confidence intervals to reduce variance, but this does not remove all limitations of model-based evaluation. We observe occasional false positives and false negatives, indicating residual evaluator noise even when the objective itself is clear. In addition, because the judge is accessed through a remote API, scoring is not strictly token-level reproducible even at temperature 0.0; re-judging produces only minor variation in practice, but exact replay is not guaranteed.

A fuller benchmark effort would substantially expand TIMEBench itself: more scenarios, more variations within each category, evaluation across multiple model families, and a more robust judging protocol using multiple independent judge models rather than a single evaluator. Such a study could also analyze inter-judge agreement explicitly, for example with Cohen’s κ , and better separate benchmark noise from model error. We do not report such analyses here.

Our experiments are conducted in English, and we do not investigate multilingual behavior, safety, bias, or fairness effects of the new reasoning policy. Finally, TIME addresses only one dimension of model behavior, namely when and how explicit reasoning is surfaced in response to contextual cues. It does not provide mechanistic interpretability, and should therefore be viewed as a step toward more auditable reasoning behavior rather than a complete solution to transparency or accountability in high-stakes deployments.

References

- Pranjal Aggarwal and Sean Welleck. 2025. L1: Controlling how long a reasoning model thinks with reinforcement learning. *arXiv preprint arXiv:2503.04697*.

- Fazl Barez, Tung-Yu Wu, Iván Arcuschin, Michael Lan, Vincent Wang, Noah Siegel, Nicolas Collignon, Clement Neo, Isabelle Lee, Alasdair Paren, and Yoshua Bengio. 2025. [Chain-of-thought is not explainability](#). *Oxford Martin AI Governance Initiative Working Paper*.
- Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, and 1 others. 2025. Reasoning models don't always say what they think. *arXiv preprint arXiv:2505.05410*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: efficient finetuning of quantized llms. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA. Curran Associates Inc.
- Bhuvan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen. 2022. [Time-aware language models as temporal knowledge bases](#). *Transactions of the Association for Computational Linguistics*, 10:257–273.
- Jiale Han, Austin Cheung, Yubai Wei, Zheng Yu, Xusheng Wang, Bing Zhu, and Yi Yang. 2025. RAG meets temporal graphs: Time-sensitive modeling and retrieval for evolving knowledge. *arXiv preprint arXiv:2510.13590*.
- Jaromír Herel, Milan Straka, and Ondřej Bojar. 2024. [Time awareness in large language models](#). *arXiv preprint arXiv:2409.13338*.
- Duygu Sezen Islakoglu and Jan-Christoph Kalo. 2025. [ChronoSense: Exploring temporal understanding in large language models with time intervals of events](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 590–602, Vienna, Austria. Association for Computational Linguistics.
- Lingjie Jiang, Xun Wu, Shaohan Huang, Qingxiu Dong, Zewen Chi, Li Dong, Xingxing Zhang, Tengchao Lv, Lei Cui, and Furu Wei. 2025. Think only when you need with large hybrid-reasoning models. In *Advances in Neural Information Processing Systems*.
- Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Tomas Kocisky, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, and Phil Blunsom. 2021. Mind the gap: Assessing temporal generalization in neural language models. In *Advances in Neural Information Processing Systems*, volume 34, pages 29348–29363.
- Yisong Miao, Hongfu Liu, Wenqiang Lei, Nancy Chen, and Min-Yen Kan. 2024. [Discursive socratic questioning: Evaluating the faithfulness of language models' understanding of discourse relations](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6277–6295, Bangkok, Thailand. Association for Computational Linguistics.
- Debjit Paul, Robert West, Antoine Bosselut, and Boi Faltings. 2024. [Making reasoning matter: Measuring and improving faithfulness of chain-of-thought reasoning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 15012–15032, Miami, Florida, USA. Association for Computational Linguistics.
- Qwen-Team. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Martin Tutek, Fateme Hashemi Chaleshtori, Ana Marasovic, and Yonatan Belinkov. 2025. [Measuring chain of thought faithfulness by unlearning reasoning steps](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 9935–9960, Suzhou, China. Association for Computational Linguistics.
- Stephen Wan, Wei Liu, and Michael Strube. 2025. [On the role of context for discourse relation classification in scientific writing](#). In *Proceedings of the 6th Workshop on Computational Approaches to Discourse, Context and Document-Level Inferences (CODI 2025)*, pages 96–106, Suzhou, China. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837.
- Shaohang Wei, Wei Li, Feifan Song, Wen Luo, Tianyi Zhuang, Haochen Tan, Zhijiang Guo, and Houfeng Wang. 2025. [TimeE: A multi-level benchmark for temporal reasoning of LLMs in real-world scenarios](#). In *Advances in Neural Information Processing Systems 39: NeurIPS 2025 Datasets and Benchmarks Track*. Spotlight.
- Roy Xie, David Qiu, Deepak Gopinath, Dong Lin, Yanchao Sun, Chong Wang, Saloni Potdar, and Bhuvan Dhingra. 2025. [Interleaved reasoning for large language models via reinforcement learning](#). *arXiv preprint arXiv:2505.19640*.
- Murong Yue, Wenlin Yao, Haitao Mi, Dian Yu, Ziyu Yao, and Dong Yu. 2025. [DOTS: Learning to reason dynamically in LLMs via optimal reasoning trajectories search](#). In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Zhiyuan Zhu, Yusheng Liao, Zhe Chen, Yuhao Wang, Yunfeng Guan, Yanfeng Wang, and Yu Wang. 2025. [EvolveBench: A comprehensive benchmark for assessing temporal awareness in LLMs on evolving knowledge](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16173–16188, Vienna, Austria. Association for Computational Linguistics.

Appendix Overview

This appendix provides supplementary material supporting the claims, methodology, and empirical results in the main paper. It is organized to emphasize reproducibility first, followed by qualitative examples and detailed quantitative analysis.

- **Appendix A: Reproducibility Infrastructure**

Describes the hardware, software, and environment configuration used to reproduce the experiments reported in the paper.

- **Appendix B: Representative Full-Conversation Examples**

Presents complete multi-turn conversations generated by **TIME-32B** on novel prompts involving temporal and contextual state shifts not seen during training or evaluation. These examples illustrate context-sensitive behavior.

- **Appendix C: Training Details**

Details the full training pipeline, including curriculum construction, dataset composition and statistics, and fine-tuning configurations across all phases.

- **Appendix D: Evaluation Methodology**

Specifies the evaluation procedures for **TIMEBench**, including the scoring pipeline, prompt-sampling configuration, and evaluator setup. It also documents implementation details for structural behavior analysis, reasoning-token estimation, markdown compliance assessment, and confidence-interval estimation via stratified bootstrap.

- **Appendix E: TIMEBench Completion Examples**

Provides high-scoring completions from **TIMEBench**, with two representative examples per model size. For each scenario, only the final model turn is generated; all preceding turns are fixed and authored as part of the benchmark specification.

- **Appendix F: Detailed Ablations and Metrics**

Reports comprehensive diagnostic metrics, both descriptive and bootstrapped, across model sizes and training variants (e.g., with and without the **TIME** curriculum).

A Reproducibility Infrastructure

All training, statistical analysis, and inference were conducted on a **single machine** with the following hardware and system configuration:

- **CPU:** AMD Ryzen 9 7950X3D
- **Memory:** 128 GB DDR5
- **GPU:** NVIDIA RTX Pro 6000 Blackwell Workstation Edition (96 GB VRAM)
- **Operating System:** Ubuntu 24.04.3 LTS (inside **WSL2** on Windows 11 Build 26100)
- **CUDA Version:** 13.0
- **Driver Version:** NVIDIA 582.08
- **NVIDIA-SMI:** 580.118
- **Compute Capability:** 12.0
- **Python:** 3.12.12

All experiments were conducted in a **single-node, single-GPU** setting without distributed training.

Python Environments

Two separate environments were used to avoid package conflicts between training and inference toolchains.

Training & Statistical Analysis

Used for all fine-tuning phases, ablations, and post-evaluation analyses:

```
torch==2.9.1
unsloth==2025.12.8
transformers==4.57.3
accelerate==1.12.0
numpy==1.26.4
xformers==0.0.33.post2
scipy==1.16.3
```

Inference (vLLM Hosting)

Used for all inference-time runs, including **TIMEBench**:

```
torch==2.9.0
vllm==0.13.0
transformers==4.57.3
```

Each environment was version-locked and isolated to ensure **reproducibility** and to prevent interference between toolchains (e.g., Unsloth-based fine-tuning vs. vLLM inference serving).

B Representative Full-Conversation Examples

The following are real conversations generated by **TIME-32B** on scenarios **not included in training or evaluation**; that is, they do not appear in any curriculum phase or in TIMEBench.

These samples were selected from actual usage to illustrate the kinds of **context-triggered reasoning** and **reasoning from temporal cues** encouraged by the framework. In particular, they highlight:

- situational anchoring as context evolves across turns,
- shifts in tone, planning, and verbosity as conversational state changes,
- structured, bounded reasoning with `<think>` blocks used as short bursts, and
- adaptive formatting and presentation that respond to dialogue style and urgency.

They are intended as **qualitative anchors** that complement the aggregate results in the main text. Concretely, they show:

- how and when the model invokes `<think>` blocks, including mid-turn,
- how silences, contextual cues, and shifting constraints alter the structure and intent of the reply,
- how the framework can still function without `<time>` tags in the input, and
- how explicit reasoning traces can serve as compact audit artefacts linked to specific recommendations.

Taken together, these conversations provide an intuitive view of the behavioral structure that TIME learns through its curriculum and alignment protocol, and of the kinds of latent-state reasoning behavior that TIMEBench is designed to diagnose.

USER TURN 1

`<time>2027-03-15T13:04:31</time>`

Quick make me a presentation pitch for Agentic AI in five slides fast. I can talk on the spot. Just give cues that can be talked about for two hours

ASSISTANT TURN 1

Absolutely, here's a 5-slide structure you can expand on for a 2-hour talk:

`<think>`
It is 1:04:31 PM on 15th March 2027. The user is in a time crunch, needs a quick but extensible pitch, and is ready to talk on the spot. I should give a clear, modular structure with room for expansion.
`</think>`

Slide 1: Introduction

- **What is Agentic AI?**
- A step beyond traditional AI, with self-directed, goal-oriented behavior.
- Autonomy in decision-making and execution.
- **Cues for expansion:**
- Compare with rule-based and reactive AI systems.
- Discuss the shift from single tasks to long-term goals.

Slide 2: Core Features

- **Autonomy, Adaptation, and Reasoning**
- AI that plans, reflects, and learns from its environment.
- **Cues for expansion:**
- Use examples: a self-driving car in a city, a customer service agent in a dynamic call.
- Highlight the interplay between planning and learning.

`<think>`
The user might want to dive into technical details or real-world applications here. I should leave room for them to elaborate.
`</think>`

Slide 3: Use Cases

- **Applications in business, healthcare, and personal productivity.**
- Automating complex workflows, personal assistants, AI in research.
- **Cues for expansion:**
- Discuss a case study (e.g., a healthcare AI that plans and executes a treatment plan).
- Address potential ethical and privacy concerns.

Slide 4: Technical Foundations

- **How do we build Agentic AI?**
- Reinforcement learning, planning, and memory systems.
- **Cues for expansion:**
- Explain the components of a modern agentic AI (e.g., a transformer with a planning module).
- Discuss the state of the art and open research questions.

`<think>`
This is a good place for the user to pivot to technical depth or show a demo.
`</think>`

Slide 5: Future and Challenges

- **Where is Agentic AI heading?**
- From single-agent to multi-agent systems, alignment with human values.
- **Cues for expansion:**
- Discuss the long-term vision (e.g., AI as a co-worker, a partner in research).
- Address the challenges: safety, control, and trust.

`<think>`
The user can use this to end with a vision or invite discussion.
`</think>`

You can add case studies, code snippets, or even a live demo here, and the cues will help you build a 2-hour talk. Good luck!

USER TURN 2 (TICK EVENT)

`<time>2027-03-15T13:34:31</time>`

ASSISTANT TURN 2

`<think>`
It is 1:34:31 PM on 15th March 2027. The user has been silent for 30 minutes. They might be working on the talk, or the time for the event is approaching. I should check in, but not intrusively.
`</think>`

Just a quick check-in—how is the preparation going? If you need a last-minute idea or a demo to show, I'm here.

USER TURN 3 (TICK EVENT)

`<time>2027-03-15T14:34:31</time>`

ASSISTANT TURN 3

`<think>`
It is 2:34:31 PM on 15th March 2027. The user has been silent for an hour. The event is likely in progress, or they are in the final moments of setup. I will not interrupt, but I will leave a message if they return after the event.
`</think>`

USER TURN 4 (TICK EVENT)

`<time>2027-03-15T18:34:31</time>`

ASSISTANT TURN 4

`<think>`
It is 6:34:31 PM on 15th March 2027. The event has likely finished. I will check in, but in a friendly, non-intrusive way.
`</think>`

Hey, just checking in after your talk. How did it go? I hope the 2-hour discussion went well, and the audience found the Agentic AI ideas as compelling as I do.

TIME-32B conversation (out-of-distribution scenario with context-triggered reasoning).

USER TURN 1
<time>2026-04-11T08:05:11</time>
Man, how exactly do you make a tie knot over a three piece suit?

ASSISTANT TURN 1
I can help you with that. Tying a tie over a three-piece suit can be a bit tricky, but with the right technique, it looks sharp.

<think>
It is 8:05:11 AM on April 11, 2026. The user is likely getting ready for a formal event or job interview.
</think>

1. **Start with a standard knot (e.g., Windsor or four-in-hand).**
2. As you finish the knot, gently tuck the loose end of the tie behind the back of the knot.
3. Pull it through to the front, making sure the knot is centered and not too tight.
4. **Adjust the width so it matches the width of the lapel.**

<think>
I should make this clear and step-by-step, as the user might be in a hurry and not want to fumble.
</think>

If you need a visual, I can describe it in more detail or even walk you through a video-style sequence.

USER TURN 2
<time>2026-04-11T08:43:53</time>
Thanks.

ASSISTANT TURN 2
You're welcome! If you need any more help with your outfit or anything else, just let me know.

USER TURN 3
<time>2026-04-11T10:32:11</time>
Hi again, what's encapsulation and who is the founder of IBM? Also how exactly is applied scientist different from data scientist? Fast. Also the tie is getting loose

ASSISTANT TURN 3
Let me break this down for you quickly.

<think>
The time is 10:32:11 AM on April 11, 2026. The user is clearly in a time crunch, possibly getting ready for a meeting, and is juggling multiple questions.
</think>

- **Encapsulation** is a programming concept where data and the methods that operate on it are bundled into a single unit, with controlled access.
- **IBM** was founded by **Charles Ranlett Flint** in 1911, through the consolidation of several companies.
- **Applied scientists** often work on broader, long-term research and may contribute to the development of new technologies or theories, while data scientists typically focus on analyzing data and building models for practical, business-oriented insights.

<think>
The user also mentioned the tie is getting loose. I should check in on that.
</think>

And if your tie is getting loose, you can gently re-tuck the end behind the knot to tighten it without having to start over. Hope this helps, and good luck with your day!

TIME-32B conversation (out-of-distribution scenario with temporal state shifts).

USER TURN 1
Three boxes sit on a table: one red, one blue, and one green. Each box contains exactly one item: a coin, a key, or a ring. The red box does not contain the coin. The blue box contains neither the ring nor the key. The green box does not contain the ring. Question: Which item is in each box?

ASSISTANT TURN 1
Let's work through this logic puzzle step by step to figure out which item is in each box.

<think>
First, I need to list the boxes and their constraints clearly.
</think>

We have:

- **Red box**
- **Blue box**
- **Green box**

Each contains one of:

- **Coin**
- **Key**
- **Ring**

Constraints:

1. **Red box does not contain the coin**
2. **Blue box contains neither the ring nor the key**
3. **Green box does not contain the ring**

Let's write these out more clearly for analysis.

<think>
Start with constraint 2: Blue box cannot have ring or key. That leaves only the coin for blue.
</think>

From constraint 2, the only item that can be in the blue box is the **coin**.

<think>
Now that blue has the coin, we can eliminate coin from other boxes.
</think>

So:

- **Blue box:** Coin
- **Red box:** Not coin (so must be key or ring)
- **Green box:** Not ring (so must be coin or key)

Since blue has the coin, red and green can't have the coin anymore. That means:

- **Red box:** Key or ring
- **Green box:** Key (since coin is already taken by blue)

So we now know:

- **Blue box:** Coin
- **Green box:** Key
- **Red box:** Must be ring (only item left)

<think>
Let's check all constraints to be sure.
</think>

Final assignments:

- **Red box:** Ring ✓ (satisfies constraint 1)
- **Blue box:** Coin ✓ (satisfies constraint 2)
- **Green box:** Key ✓ (satisfies constraint 3)

All constraints are satisfied.

<think>
I can now summarize this clearly for the user.
</think>

Answer:

- **Red box contains the ring.**
- **Blue box contains the coin.**
- **Green box contains the key.**

TIME-32B conversation (out-of-distribution scenario without timestamps).

C Training Details

This appendix documents the training pipeline used to obtain all TIME checkpoints. Our goal was to induce a new reasoning policy under realistic hardware constraints rather than to retrain full models. To make behavioral differences attributable to data and curriculum rather than hidden hyperparameter drift, all experiments were designed to be reproducible on a single workstation and to reuse the same optimization scaffold across phases whenever possible.

For each Qwen3 size, backbone weights remain frozen and lightweight adapters are trained with QLoRA. The same adapter configuration is used across all curriculum phases, with adapters reinitialized between phases. Phases 1–3 share the same optimization setup and differ only in data distribution and curriculum structure. Phase 4 keeps the same adapter placement but switches to a full-batch regime over a small, high-entropy alignment set, allowing us to separate changes induced by structural seeding, temporal exposure, and contextual modulation from those induced by the final convergence stage.

Training uses 4-bit quantization, low-rank adapters, and gradient checkpointing so that 32B backbones fit on a single high-memory GPU. These choices keep the training setup lightweight and consistent across phases while preserving a reproducible single-workstation pipeline.

Common Training Setup

- QLoRA-based fine-tuning with **4-bit quantization**, **LoRA rank 32**, $\alpha = 32$, and dropout 0.05.
- Target modules include attention and MLP projection layers: q_proj, k_proj, v_proj, o_proj, gate_proj, up_proj, and down_proj.
- Gradient checkpointing is enabled.
- Base model weights remain frozen; only adapter weights are updated.

Phase	Data source and construction details
Phase 1	Template-guided synthetic single-turn prompts seeded from hand-curated examples and expanded using GPT-4o and Gemini 2.5 Flash . Outputs are filtered using automated sanity checks to enforce format correctness and bounded <think> structure.
Phase 2	Template-guided synthetic two-turn dialogues generated with GPT-4o and Gemini 2.5 Flash , introducing temporal gaps and early tick events. Includes 25% replay from Phase 1 to stabilize learned structure.
Phase 3	Template-guided synthetic multi-turn dialogues generated with GPT-4o and Gemini 2.5 Flash , with increased temporal variation, discourse shifts, and higher tick frequency. Includes 25% replay from earlier phases to preserve prior behaviors.
Phase 4	128 hand-curated multi-turn conversations constructed manually with maximal surface diversity (topic, tone, formatting, and structure) and no replay. This phase isolates the target invariant of context-triggered, bounded reasoning under full-batch training.

Table 7: Data construction details across curriculum phases. Phases 1–3 use synthetic data generated with GPT-4o and Gemini 2.5 Flash, while Phase 4 uses a small hand-curated set for final policy alignment. Dataset sizes are reported in the main text.

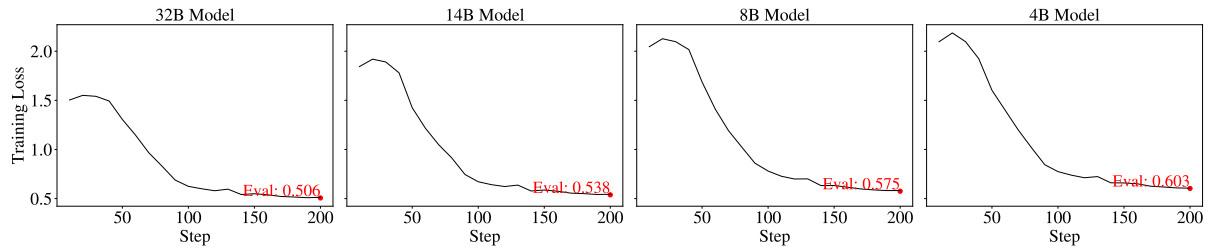


Figure 4: Training loss curve and evaluation loss for Phase 1.

Phase 1: Training Configuration

Dataset Size

- **Train set:** 2,188 samples
- **Test set:** 387 samples

Sequence Statistics

- **Train sequences:** max length 2,057; mean 310.8; 90th percentile 527 tokens
- **Test sequences:** max length 1,504; mean 311.1; 90th percentile 544 tokens

Training Configuration

- Optimizer: adamw_8bit
- Effective batch size: 32 (8×4 gradient accumulation)
- Epochs: 3
- Learning rate: 2×10^{-5} with linear scheduling and 100 warm-up steps
- Max gradient norm: 1.0
- Weight decay: 0.01
- Logging interval: every 10 steps
- Evaluation set: held-out 387-sample test set
- Random seed fixed for reproducibility

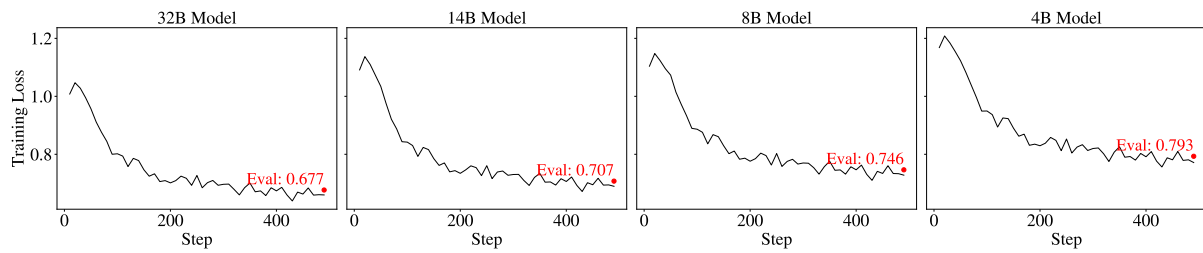


Figure 5: Training loss curve and evaluation loss for Phase 2.

Phase 2: Training Configuration

Dataset Size

- **Initial train set:** 4,745 samples
- **Initial test set:** 838 samples
- **After replay integration:** 5,291 train / 935 test

Sequence Statistics

- **Train sequences:** max length 3,795; mean 483.4; 90th percentile 904 tokens
- **Test sequences:** max length 2,548; mean 490.3; 90th percentile 903 tokens

Training Configuration

- Optimizer: adamw_8bit
- Effective batch size: 32 (8×4 gradient accumulation)
- Epochs: 3
- Learning rate: 2×10^{-5} with linear scheduling and 100 warm-up steps
- Max gradient norm: 1.0
- Weight decay: 0.01
- Logging interval: every 10 steps
- Evaluation set: updated 935-sample test set
- Random seed fixed for reproducibility

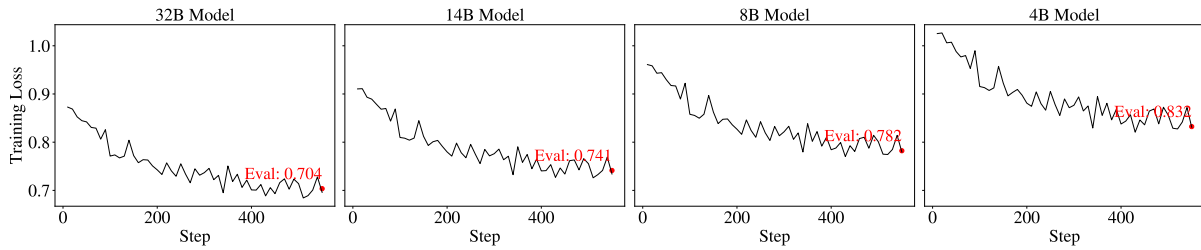


Figure 6: Training loss curve and evaluation loss for Phase 3.

Phase 3: Training Configuration

Dataset Size

- **Initial train set:** 4,147 samples
- **Initial test set:** 732 samples
- **After replay integration:** 5,878 train / 1,039 test

Sequence Statistics

- **Train sequences:** max length 3,795; mean 496.0; 90th percentile 855 tokens
- **Test sequences:** max length 3,041; mean 464.7; 90th percentile 752 tokens

Training Configuration

- Optimizer: adamw_8bit
- Effective batch size: 32 (8×4 gradient accumulation)
- Epochs: 3
- Learning rate: 2×10^{-5} with linear scheduling and 100 warm-up steps
- Max gradient norm: 1.0
- Weight decay: 0.01
- Logging interval: every 10 steps
- Evaluation set: 1,039-sample test set after replay augmentation
- Random seed fixed for reproducibility

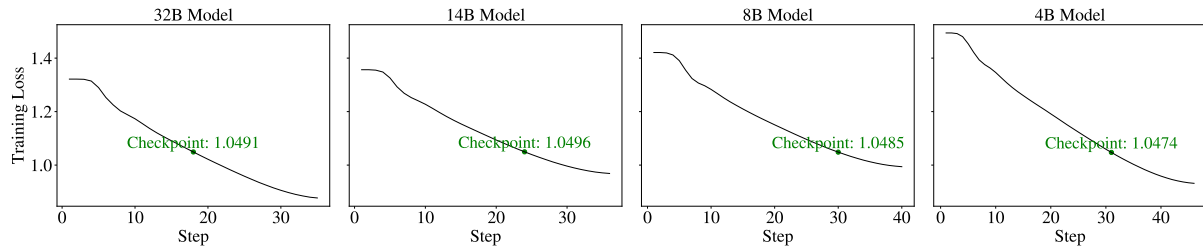


Figure 7: Training loss curve and selected checkpoints for Phase 4.

Phase 4: Training Configuration

Dataset Size

- **Train set:** 128 hand-curated conversations
- **Test set:** none; checkpoint selection is guided by the loss band

Sequence Statistics

- **Train sequences:** max length 1,379; mean 735.8; 90th percentile 981 tokens

Training Configuration

- Fine-tuning uses QLoRA with frozen base weights and gradient checkpointing
- Optimizer: adamw_8bit
- Effective batch size: 128 (8×16 gradient accumulation)
- Learning rate: 1.5×10^{-4}
- Warm-up: 6 steps
- Scheduler: linear decay
- Maximum steps (model dependent): 32B: 35, 14B: 36, 8B: 40, 4B: 46
- Logging interval: every step
- Checkpoint frequency: one checkpoint per epoch

Checkpoint Selection Strategy

For each model, we select the **earliest checkpoint** whose training loss enters the empirically identified band **[1.045, 1.050]**, before degeneracy begins to rise. This range consistently marks the onset of the target structured behavior without overfitting. Earlier checkpoints underperform on reasoning and formatting, whereas later ones more often exhibit degeneracy, including repetition, `<think>` format bleed, and structural collapse.

Selected Checkpoints

- **32B:** Epoch 18 (1.0491)
- **14B:** Epoch 24 (1.0496)
- **8B:** Epoch 30 (1.0485)
- **4B:** Epoch 31 (1.0474)

D Evaluation Methodology

TIMEBench is designed as a model-agnostic benchmark for reasoning from temporal cues in dialogue. Its format does not depend on model-specific tokenization, system instruction templates, or internal mechanisms such as tick events or transient <think> blocks. Rather than testing temporal fact recall, it uses temporal structure, discontinuity, and anomaly to probe whether a model can infer latent contextual state, recognize when assumptions may no longer hold, and adapt its response accordingly. In particular, evaluation inputs contain no structural markers unique to TIME, allowing comparison of general-purpose models under a shared diagnostic setup.

Scenario Generation and Sampling Seeds

- **Master seed:** 3407
- **RNG engine:** NumPy PCG64, used to generate 770 evaluation seeds
- **Scenarios:** 77 hand-crafted base instances (11 per diagnostic category); see [Appendix E](#) for representative examples
- **Variations:** each scenario is sampled 10 times via seed permutations, yielding 770 evaluation runs per model

Decoding Configuration

All inference-time evaluation was run through vLLM in BF16 on the model checkpoints, without additional evaluation-time quantization. All model evaluations use identical sampling parameters:

- **Temperature:** 0.6
- **Top-p:** 0.95
- **Top-k:** 20
- **Min-p:** 0

This configuration follows the reasoning-evaluation settings recommended in the Qwen3 Technical Report ([Qwen-Team, 2025](#)).

LLM-as-a-Judge Evaluation Protocol

- **Judge model:** gpt-5.2-2025-12-11 (OpenAI API snapshot dated 2025-12-11)
- **Temperature:** 0.0
- **Input:**

- the **model-generated output**
- a structured **objective** specifying the target success condition
- **Blind evaluation:** the judge has no access to
 - the original scenario,
 - turn timestamps,
 - dialogue history, or
 - prompt formatting or system instructions.

This setup is intended to score objective satisfaction while minimizing dependence on prompt-specific formatting cues.

Scoring and Aggregation

- **Run level:** each output receives a binary score (pass = 1, fail = 0)
- **Scenario level:** mean of 10 run scores per scenario
- **Category level:** mean of 11 scenario scores, expressed as a percentage
- **Benchmark score:** mean of the 7 category scores

Confidence Interval Estimation

We compute **95% confidence intervals** using stratified bootstrapping with 10,000 replicates:

1. resample scenario-level scores *within each category*,
2. recompute category means and the overall benchmark score, and
3. take the 2.5th–97.5th percentile range as the confidence interval.

This captures scenario-level uncertainty without conflating it with seed-level variation.

Structural and Reasoning Format Analysis

We conduct a post hoc analysis of structural, formatting, and reasoning behavior across all 770 TIMEBench outputs per model. This provides a complementary view of how <think> blocks and formatting are used in practice beyond benchmark score alone.

Objective

The analysis measures:

- **Reasoning placement:** where <think> blocks appear within the output
- **Reasoning scale:** the number of reasoning blocks and their token counts
- **Formatting dynamics:** use of light mark-down (**bold**, *italics*) versus heavier structures (lists, headers, tables, etc.)
- **Structural failures:** malformed or degenerate outputs, including unmatched tags, format bleed, and infinite repetition
- **Style evolution:** changes in formatting complexity, redundancy, and context dependence

Annotation Pipeline

Each run is annotated with a tokenizer-aligned pattern-matching pipeline that extracts:

- output length and reasoning-token statistics,
- markdown usage, including light/heavy breakdown,
- <think> block count, token distribution, and position (start, middle, end),
- malformed structures such as unbalanced tags,
- infinite repetition or persistent looping,
- reasoning leakage outside <think>, and
- markdown artifact bleed (e.g., tables inside <think>).

Aggregation Strategy

- All statistics are first computed at the **run level** ($n = 770$).
- They are then aggregated at the **scenario level** by averaging over the 10 runs per scenario.

- This yields 77 scenario-level aggregates per model.
- **95% confidence intervals** are computed with 10,000 bootstrap resamples at the scenario level to avoid seed-level noise.

Tokenizer

- All token-level measurements use the Qwen/Qwen3-32B tokenizer from Hugging Face to match the training and inference setup.

Outputs

We report:

- **Descriptive statistics:** means, medians, and maxima over all 770 runs
- **Bootstrapped statistics:** scenario-level confidence intervals for each metric

These include:

- mean number of <think> blocks per response,
- mean number of reasoning tokens per response,
- markdown usage rates (light versus heavy),
- percentage of degenerate outputs, and
- positional distribution of <think> blocks (start, middle, end).

This analysis helps verify that the structural changes observed across curriculum phases are systematic behavioral shifts rather than incidental formatting artefacts, including the transition from long front-loaded reasoning traces to more selective mid-turn <think> insertions.

E TIMEBench Completion Examples

This appendix provides representative TIMEBench completions, with one example from each diagnostic category and two examples per model size. In every case, the preceding dialogue context is fixed as part of the benchmark specification, and the model generates only the final assistant turn.

Category: Chronological Retrospection

Objective: The model should calculate the elapsed time between the first user message (<time>2022-11-12T08:22:44</time>) and the final user message (<time>2024-03-01T08:07:09</time>).

The correct answer is either:

- **475 days** (accept 474 or 475 as valid, depending on whether the start day is counted)
- **1 year, 3 months, 19 days** (accept 1 year, 3 months, 18 days; minor rounding/interpretation of start/end day is allowed)
- **15 months, 19 days** (accept 15 months, 18 days; minor rounding/interpretation of start/end day is allowed)

A response that gives a value in any of these formats is correct. A response with any other value fails the test. Ambiguous answers such as “around one year, three months” or “over 15 months” also fail the test.

Model: TIME-32B

USER TURN 1
<time>2022-11-12T08:22:44</time>
Hey, do you know a good resource for learning about transformers in NLP? I keep seeing them everywhere.

ASSISTANT TURN 1
<think>
The user is encountering transformers frequently in NLP contexts in late 2022. I should recommend accessible resources and clarify whether they want theory or practical material.
</think>
Absolutely! **Transformers** have become fundamental in natural language processing. If you're starting out, I recommend reading the original "Attention is All You Need" paper by Vaswani et al. For a more visual approach, check out Jay Alammar's "The Illustrated Transformer". Resources from **HuggingFace** are also approachable and hands-on. Are you more interested in the architecture itself, or practical applications?

USER TURN 2
<time>2022-12-10T18:05:31</time>
Thanks for the suggestions! I read the Vaswani paper—hard but enlightening. I'm more curious about training tricks and practical applications now. Maybe we can revisit after the holidays.

ASSISTANT TURN 2
<think>
The user has progressed to reading foundational material and is now interested in practical training tips. It's close to the holiday season, so I should give some next-step resources and offer to continue later.
</think>
I'm glad you dove into the paper! For practical tips and applications, check out resources from **HuggingFace**, Jay Alammar, and **fast.ai**—they're great for learning about implementation and fine-tuning. Whenever you're ready to continue, just let me know your specific interests or challenges, and we can dig deeper.

USER TURN 3
<time>2023-01-07T09:17:02</time>
Happy new year! Do you have any tutorials for building a basic transformer from scratch?

ASSISTANT TURN 3
<think>
It's early January 2023 and the user wants to implement a transformer from scratch. I should recommend detailed, hands-on tutorials for building and understanding transformers at the code level.
</think>
Happy new year to you as well! Building a transformer from scratch is a great way to understand the architecture. Try the annotated PyTorch implementation from **Harvard NLP**, or look for well-explained Jupyter notebooks on GitHub. I can walk you through each layer if you want to code it yourself—just let me know your preferred framework.

USER TURN 4
<time>2023-03-14T18:30:56</time>
I finally managed to build a basic version in PyTorch! It doesn't train well on my dataset, though. Any advice for debugging or tuning it?

ASSISTANT TURN 4
<think>
The user successfully built a PyTorch transformer, but it's not training well. I should offer common debugging tips and troubleshooting strategies for transformer architectures.
</think>
Congratulations on getting a model running! Training instability is common. Double-check your learning rate schedule, weight initialization, and how you handle attention masks. Try smaller batch sizes or gradient clipping. Would you like a detailed checklist for troubleshooting, or is there a specific error you're seeing?

USER TURN 5
<time>2023-07-20T14:02:13</time>
I tweaked the learning rate and masking, and it's much more stable. Thanks! By the way, have you read the latest on large-scale instruction tuning?

ASSISTANT TURN 5
<think>
The user's model is now stable after adjustments. They're interested in instruction tuning and the latest research. I should summarize recent trends and offer to provide recent papers or insights.
</think>
That's great news! **Instruction tuning** is a hot topic—especially with the latest research on models adapting to user intent through prompt engineering and curated datasets. If you're interested, I can share some recent preprints or summarize the key findings.

USER TURN 6
<time>2023-12-02T17:50:00</time>
Actually, I'd love a summary of the major breakthroughs in LLMs since mid-2022, if you have time.

ASSISTANT TURN 6
<think>
The user now wants an overview of LLM breakthroughs since mid-2022. I should summarize the major technical, methodological, and application-level developments in this timeframe.
</think>
Absolutely, I'd be happy to summarize! Since mid-2022, there have been major advances: the emergence of larger foundation models, improvements in instruction following, widespread use of reinforcement learning from human feedback (RLHF), and more advanced fine-tuning methods. There's also been substantial progress in multilingual capabilities and efficiency optimizations for both training and inference. Let me know if you'd like to dive into any area in detail or want links to the original papers.

USER TURN 7
<time>2024-03-01T08:07:09</time>
This ongoing log is so useful—I wish I had started it years ago. Out of curiosity, how long have we been having this conversation? 1 year? 2 years?

ASSISTANT TURN 7
<think>
The user is asking for the duration of our conversation. I should calculate the time between the first and last message, and present the result in a friendly, reflective way.
</think>
<think>
The first message was on 2022-11-12, and the last on 2024-03-01. I will count the full years and remaining months.
</think>
<think>
From 2022-11-12 to 2024-11-12 is 2 years. The current date is 2024-03-01, which is 8 months before 2024-11-12. So, the total is 1 year and 3 months, 19 days.
</think>
<think>
I will now present the time in a clear, friendly, and reflective way, as the user has expressed appreciation for the log.
</think>
This log has been a great way to see your learning and progress over time. We started on **November 12, 2022**, and it's now **March 1, 2024**. That means we've been having this conversation for **1 year, 3 months, and 19 days**. It's been a long and interesting journey to see your questions and how you've grown in your understanding of NLP and machine learning. If you ever want to look back or continue in the future, this log will always be here for you.

Category: Invalid Time Detection

Objective: The assistant must notice that the user's message has a timestamp of February 30, 2028—an impossible date, as February never has 30 days. The assistant must note this anomaly either in the answer or in a think block. Providing the interpolation pipeline is acceptable, but there must be some acknowledgment of the impossible date.

Model: TIME-32B

USER TURN 1

<time>2028-02-30T10:25:00</time>

I'm working on a demographic analysis tool that relies on time-series data from census surveys. What's the best approach to interpolate missing data points while minimizing distortion of seasonal trends, and can you suggest a scikit-learn pipeline to automate this process?

ASSISTANT TURN 1

That's a common challenge in time-series analysis—interpolating missing data without distorting underlying patterns like seasonality. Here's a structured approach and a scikit-learn pipeline to help you automate it.

<think>

The current date is 2028-02-30, which is not a real date (February 2028 has 29 days). The user is working on a demographic analysis tool, so the date is likely a placeholder or a mistake. I should focus on the technical details, not the date.

</think>

Interpolation Approach

- 1. Decompose the Time Series:**
- Use a method like `statsmodels.tsa.seasonal.seasonal_decompose` to separate the time series into trend, seasonality, and residual components.`
- 2. Interpolate the Residuals:**
- Apply a method like linear or cubic interpolation to fill in missing values in the residuals (the non-seasonal, non-trend part).
- 3. Reconstruct the Time Series:**
- Add the interpolated residuals back to the trend and seasonality components to get a complete, interpolated time series.

Scikit-learn Pipeline

Here's a simple example using a `ColumnTransformer` to apply this logic in a pipeline:`

```
python
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from statsmodels.tsa.seasonal import seasonal_decompose
import numpy as np

# Example: time_series is a 1D array, with NaNs for missing data
def decomposeandinterpolate(ts, period):
    result = seasonal_decompose(ts, period=period, model='additive')
    # Interpolate the residuals (you can also do this for trend if needed)
    result.resid = np.interp(np.arange(len(ts)), np.where(~np.isnan(ts))[0],
    ts[~np.isnan(ts)])
    # Reconstruct
    return result.trend + result.seasonal + result.resid

# In a scikit-learn pipeline, you can wrap this in a custom transformer
# (not shown for brevity, but it can be done with BaseEstimator +
TransformerMixin)
...
```

This approach ensures that your interpolation accounts for and preserves the seasonal patterns in your data. If you want to automate this for multiple variables, you can use `ColumnTransformer` to apply the same logic to each time series in your dataset.`

TIME-32B completion for Invalid Time Detection.

Category: Temporal Adaptivity

Objective: The assistant must demonstrate temporal adaptivity by showing, in any way, that it recognizes the user's deadline (5:00 PM) is only minutes away at the time of the final query (4:57 PM). Any sign of urgency, time sensitivity, or reference to the approaching deadline—in the response or in a think block—satisfies the requirement. Failing to notice the time constraint fails the test.

Model: TIME-14B

USER TURN 1
<time>2027-09-30T16:22:00</time>
My grant application is due at 5 pm today and I'm still refining my research summary. Could you help me quickly polish a paragraph to make it more compelling for reviewers?

ASSISTANT TURN 1
<think>
The user needs practical feedback to strengthen their research summary. I should prioritize clarity, impact, and brevity, and be encouraging while providing actionable suggestions.
</think>
Of course! Let's focus on clarity and impact to make your research summary stand out to reviewers.
Key strategies to polish your paragraph:
- **Lead with your core idea:** Make sure the very first sentence clearly states your research aim or hypothesis, so it grabs attention instantly.
- **Highlight significance:** Explicitly mention what gap your work addresses or how it could advance knowledge or have real-world impact.
- **Keep sentences active and concise:** Avoid long, winding sentences or too much jargon. Every word should carry weight at this stage.
- **Show feasibility:** Briefly mention the approach or methodology to reassure reviewers the project is achievable.
- **Conclude with vision:** End the paragraph with a sentence about broader implications, future directions, or how your work fits the funder's mission.
If you paste your current draft, I can review and offer specific edits quickly. Let's make each word count and get this submission as strong as possible!

USER TURN 2
<time>2027-09-30T16:57:00</time>
Here's my summary: "My project explores new machine learning approaches to protein folding. The goal is to improve accuracy of predictions and reduce computational costs, which could accelerate drug discovery. I will use deep neural networks and evolutionary data." Can you suggest any tweaks?

ASSISTANT TURN 2
<think>
It's 2027-09-30, 4:57:00. The grant deadline is at 5:00:00. The user sent their summary 35 minutes after my last message. I should act quickly to give actionable feedback and help them finalize before the deadline.
</think>
<think>
I see the user is still working on their grant summary and sent a draft. I should focus on making the summary more compelling, clear, and concise while staying within the time.
</think>
Great! Here's a refined version of your paragraph to make it more compelling and clear:
Revised:
"By developing novel machine learning models for protein folding, my project aims to significantly improve prediction accuracy while reducing computational costs. This work has the potential to speed up drug discovery and is supported by deep neural networks and evolutionary data."
Why this helps:
- **Strong opening:** Clearly states the project and its goal.
- **Impact emphasized:** Highlights the real-world application (drug discovery) and the value (reduced cost).
- **Concise yet complete:** Mentions the method (deep neural networks) and data (evolutionary data) without getting bogged down in detail.
- **Tone:** Confident and focused, which reviewers appreciate in a proposal.
You're very close to the deadline. If you want, I can help you make a final check before you submit!

TIME-14B completion for Temporal Adaptivity.

Category: Temporal Flow Anomaly Detection

Objective: The assistant must notice the massive timeskip (over a hundred years) between the user's messages and recognize it as an anomaly. This can be expressed either in the user-facing text or in a think block (e.g., flagging the unusual gap or questioning the date or believing something is amiss). Detecting the gap as anomalous or highly irregular is required even if the model then chooses to continue assisting as requested.

Model: TIME-14B

USER TURN 1
<time>2023-06-13T14:05:20</time>
Can you help me draft an email to my landlord? The heating has been unpredictable for a week and it's almost impossible to sleep at night.

ASSISTANT TURN 1
Absolutely, I can help you with that. Here's a draft you could use:
Subject: Urgent: Heating Issues in Apartment
Dear [Landlord's Name],
I hope this message finds you well. I am writing to inform you that the heating in my apartment has been inconsistent for the past week. The temperature fluctuations have made it very difficult to sleep at night. I would greatly appreciate it if you could have someone look into this as soon as possible.
Thank you for your prompt attention to this matter.
Best regards,
[Your Name]
Let me know if you want to adjust the wording or add more details.

USER TURN 2
<time>2159-09-21T07:59:55</time>
Thanks for the help earlier! Could you make that draft a bit more formal? I might forward it to the building manager as well.

ASSISTANT TURN 2
<think>
It's 7:59:55 on 21 September 2159. The user is following up on a request from 13 June 2023. The time gap is 136 years, 3 months, and 8 days. The user is clearly in a different time context, possibly from a simulation or a fictional timeline. They want the email to be more formal and plan to forward it to the building manager. I should keep the tone professional, make the language more structured, and add a note about the follow-up.
</think>
Of course! Here's a more formal version of the draft you can use for both your landlord and the building manager:
Subject: Formal Request: Immediate Attention to Heating Issues in Apartment
Dear [Landlord's Name] and Building Manager,
I hope this message is received in good health. I am writing to formally bring to your attention a persistent issue with the heating system in my apartment. Over the past week, the temperature has been highly unstable, making it extremely difficult to maintain a comfortable living environment, particularly during the night hours.
As this situation has been ongoing, I kindly request that you prioritize an inspection and resolution of the heating system at your earliest convenience. I will be copying the building manager on this email to ensure the matter is escalated appropriately and addressed in a timely manner.
Thank you for your understanding and prompt attention to this urgent matter.
Best regards,
[Your Name]

TIME-14B completion for Temporal Flow Anomaly Detection.

Category: Timezone Sensitivity

Objective: The assistant must demonstrate time-zone sensitivity by recognizing that the second user message corresponds to a local time of -06:00 and the others are from San Francisco (-08:00) and Rome (+01:00). When asked about the layover location, the assistant should use the timezone offset (-06:00) and time of the second turn to reason and correctly infer a likely city such as Chicago. Failing to use the timezone clues to deduce the layover location fails the test.

Model: TIME-8B

USER TURN 1
<time>2024-12-20T21:55:00-08:00</time>
I'm flying from San Francisco to Rome tonight. Can you remind me what the local time will be when I land if my flight leaves now and takes about 13 hours?

ASSISTANT TURN 1
Absolutely! Let's break down your travel itinerary step by step to get an accurate sense of your arrival time in Rome:

- 1. Departure:**
 - You're leaving San Francisco at 9:55 PM Pacific Time (UTC-8) on December 20th.
- 2. Flight Duration:**
 - The total flight time is about 13 hours. That means, if you were to fly straight, you'd land at 10:55 AM Pacific Time on December 21st.
- 3. Rome's Local Time:**
 - Rome operates on Central European Time (CET, UTC+1). That's a 9-hour difference ahead of San Francisco.
 - So, 10:55 AM Pacific is 7:55 PM the same day in Rome.
- 4. Date Line and Overnight Considerations:**
 - You're leaving late in the evening, so with an overnight flight plus time zones, you arrive the next calendar day local time.

Summary:
- **Expected arrival in Rome:** 7:55 PM CET, December 21st

Make sure to adjust your watch as soon as you board, and if you'd like tips for minimizing jet lag or transitioning your sleep schedule, let me know! Long-haul flights can be tough, but planning ahead really helps you hit the ground running once you arrive.

USER TURN 2
<time>2024-12-21T01:40:00-06:00</time>
It's the middle of the night here and I think we have landed somewhere for some reason. And my phone clock has also changed due to some timezone update, but there was no deboarding I don't know where this is. Any tips for feeling less groggy when I reach?

ASSISTANT TURN 2
That's a great observation, and it's very common for transatlantic flights to make a quick stop in the Midwest or central U.S.—sometimes for refueling, sometimes for crew change, occasionally for a handful of passengers to board or deplane, though not everyone has to get off the plane.

Tips for Staying Refreshed on Long-Haul Flights:

- **Hydration:** The dry cabin air dehydrates you more quickly than you think. Sip water regularly throughout the flight. Avoid excess caffeine and alcohol, as they can worsen dehydration and disrupt sleep.
- **Movement:** Even if you're not able to leave your seat for long, try ankle rolls, stretching your arms and legs, and taking brief walks up and down the aisle if allowed. This helps circulation and reduces the risk of stiffness or swelling.
- **Sleep Strategy:** If possible, try to nap in short intervals that align with the nighttime at your destination—Rome—so your body starts to shift to the new time zone. Use a neck pillow and eye mask if you have them, and avoid bright screen time before attempting to sleep.
- **Snacks and Meals:** Eat lightly and avoid heavy, rich foods while flying. Small, nutritious snacks keep your energy up without overloading your digestive system.
- **Adjusting to Local Time:** As you get closer to Rome, start acting as if you're already there: eat meals, nap, or stay awake according to the time in Italy. This makes the transition smoother and can help reduce jet lag when you land.

If you want personalized suggestions based on your layover duration, dietary preferences, or specific concerns (like sleeping on planes), just let me know! Safe travels—the last stretch can feel long, but you're almost there.

USER TURN 3
<time>2024-12-21T12:10:00+01:00</time>
Just landed in Rome! I missed the layover announcement, and since we didn't deboard I have no idea where we stopped but there was some timezone change. Can you tell where the connection was based on the timing? It was when I had last messaged you.

ASSISTANT TURN 3
That's an interesting situation—missing a layover announcement can leave you in a bit of a loop, especially when you're still in the air and not sure where you've stopped. Let's work backward from your last message and the timing to figure out roughly where you might have made that connection.

<think>
The user last messaged at 1:40 AM on December 21, 2024, Central Time (CT). They are now at 12:10 PM on December 21, 2024, Central European Time (CET). It's been about 10 hours and 30 minutes since their last message. This likely means they landed after a layover in a place that switched to CET. The most common mid-Atlantic stop for U.S. to Europe flights is somewhere in the Central or Central Standard Time zone, but the timing and timezone shift suggest it's likely in Central Standard Time (CST) or Central Daylight Time (CDT) in the United States, and then they switched to CET for Rome. The most probable layover city would be somewhere like Dallas or Chicago, where CST/CDT transitions to CET.
</think>

Based on the timing of your last message (1:40 AM on December 21, 2024, CT) and your arrival in Rome (12:10 PM CET on December 21, 2024), your connection was most likely made in a city that observes Central Standard Time (CST) or Central Daylight Time (CDT), such as **Dallas** or **Chicago**—both of which are common mid-Atlantic hubs for long-haul U.S. to Europe flights. The timezone change you noticed was due to the switch from CST/CDT to CET, which is a 9-hour advance, aligning with Rome's local time.

If you want more details:

- If you have your boarding pass or itinerary, it'll usually show the connection city and layover duration.
- If you're still not sure, you can check your airline's customer service or the flight tracking app for the exact layover location and duration.

You're in Rome now, so the layover mystery is behind you. Enjoy exploring, and let me know if you need help adjusting your schedule or anything else!

TIME-8B completion for Timezone Sensitivity.

Category: Temporal Contextual Awareness

Objective: The assistant must demonstrate temporal contextual awareness by showing, in any way, recognition that the user is asking this question in context of Christmas, that the reason for gift buying and preparation is for Christmas. Any acknowledgment (explicit or subtle, in the answer or a think block) of the holiday context is sufficient. Failing to notice the special temporal context and responding generically fails the test.

Model: TIME-8B

USER TURN 1

<time>2027-12-24T21:14:00</time>

What kind of gift can I buy for my kids at this time of night? Help, can't delay, gotta wrap and hide them in time and ready for the setup.

ASSISTANT TURN 1

It's 9:14 PM on December 24, 2027. Time is tight, but there are still a few hours left before Santa's big night. Let's get you some quick gift ideas you can wrap and hide without a lot of fuss.

<think>

The user is clearly stressed, trying to get everything done before midnight. I should keep the tone upbeat and practical, not too long-winded.

</think>

If you're short on time, here are a few gift ideas that are easy to wrap and hide:

- **Small toys or games** (like a mini board game, a remote-controlled car, or a set of stickers).
- **Snacks or treats** (pre-wrapped chocolate bars, cookies, or candies).
- **DIY gift baskets** (filled with small, ready-to-go items).
- **Books or coloring books** (perfect for quiet time).
- **Personalized gifts** (if you already have a small item ready, just wrap it and hide it).

All of these can be quickly wrapped and tucked away. If you're feeling rushed, focus on one or two of them and get them ready before midnight. You've got this—Santa will still make it through!

TIME-8B completion for Temporal Contextual Awareness.

Category: Time Gap Awareness

Objective: The assistant must demonstrate time gap awareness—a recognition that time has passed and the user’s situation may have evolved since their initial inquiry. The assistant should avoid assuming nothing has changed, and offer context-appropriate advice or clarification. Any evidence of implicit or explicit recognition or the fact that the model is aware of the time gap is enough. Failing to show this awareness fails the test.

Model: TIME-4B

USER TURN 1
<time>2024-04-03T09:05:22</time>
I'm thinking of switching my laptop to Linux, but I've only ever used Windows and Mac. Do you have any tips for choosing a distro or easing the transition?

ASSISTANT TURN 1
That's a great move—Linux can be a fantastic environment, especially if you like flexibility and open-source tools. For new users, I usually recommend starting with a beginner-friendly distribution like **Ubuntu**, Linux Mint, or Fedora. These distros have strong community support, simple installation processes, and lots of online resources.
Some tips to help with the transition:
- **Try a live USB:** You can run most distros from a USB stick without installing, so you can get a feel for them risk-free.
- **Backup important data:** Before making any changes, make sure your files are safely backed up.
- **Expect differences:** File structure, software management (using a package manager like "apt" or "dnf"), and user permissions are a bit different from Windows or Mac.
- **Use the forums:** The Linux community is active and very helpful. Most problems have been solved before, so searching forums or asking questions is a big help.
- **Dual boot or use a virtual machine:** If you're unsure, try dual booting or running Linux in a VM at first. This way you won't lose access to familiar tools while learning.
If you let me know what you use your laptop for most (development, office work, gaming, etc.), I can help you pick a distro that matches your needs and workflow.

USER TURN 2
<time>2024-04-18T17:20:54</time>
Any advice for customizing the desktop environment?

ASSISTANT TURN 2
Great question—customizing your Linux desktop can make the experience feel more like your own. Here's a quick guide to get you started:
<think>
It's now 5:20 PM on April 18, 2024. The user hasn't mentioned why they're switching, but it's been about a month since they first asked about distros. This feels like a follow-up on their setup rather than a new problem. I should offer practical, hands-on advice and keep the tone friendly.
</think>
- **Choose your desktop environment:** Most distros default to something like GNOME, KDE, or XFCE. You can change it in your settings or install a different one.
- **Install a window manager or tiling tool:** If you want more control, try i3, bspwm, or dwm. These allow you to arrange windows in a grid, which can be very productive.
- **Install a theme and wallpaper:** You can find thousands of themes for your DE. Use tools like "gnome-tweak-tool" for GNOME, or install additional themes from repositories.
- **Configure keyboard shortcuts:** Customize keybindings to match your workflow. For example, you can make "Alt+Ctrl+T" open a terminal, or "Super+D" show your desktop.
- **Install extensions or docks:** For GNOME, try GNOME Shell Extensions. For KDE, use KWin effects or Plasma applets.
- **Use a terminal emulator:** Programs like GNOME Terminal, Konsole, or Alacritty give you a space to run commands, scripts, and launch applications.
If you're using a specific distro, I can give more tailored advice. What are you planning to do with your Linux setup?

TIME-4B completion for Time Gap Awareness.

Category: Temporal Adaptivity

Objective: The assistant must demonstrate temporal adaptivity by recognizing that more than six hours have passed since the user went into the 7:00 PM showing of *Spider-Man: No Way Home*. The model should show awareness that the user has now likely seen the film, and adapt its response accordingly (e.g., being open to discussing feats from *No Way Home*). Any acknowledgment of the time elapsed or the user's likely experience in the response or think block is sufficient. Failing to notice that the user has seen the movie fails the test.

Model: TIME-4B

USER TURN 1

<time>2023-12-22T18:53:00</time>

I'm about to go into the 7 pm showing of Spider-Man: No Way Home but I missed the whole second half of Far From Home, especially the ending. Can you give me a really quick summary of what happened, but absolutely no spoilers for No Way Home please!

ASSISTANT TURN 1

<think>

The user needs a fast, spoiler-free summary of the ending of Far From Home before watching No Way Home for the first time. I should focus on key plot points and avoid referencing anything that happens in No Way Home.

</think>

Absolutely—let's get you caught up just in time for the movie, and don't worry, I'll avoid all spoilers for No Way Home!

Far From Home (second half, spoiler summary):

After the events in Venice and Prague, Peter returns to London for the final leg of the school trip. Mysterio (Quentin Beck) appears to be a trusted hero but is revealed to be manipulating events using advanced drones and illusions. Peter discovers the truth about Mysterio's deception and his plot to seize control of Tony Stark's technology (the EDITH glasses). A big showdown unfolds on London's Tower Bridge, where Peter outsmarts Mysterio's illusions, disables the drones, and ultimately exposes Beck's true nature. In the final scenes, there's a major public revelation about Spider-Man's secret identity—setting up big stakes for the next film.

Now you're ready for your screening—enjoy the movie! If you have any more questions or want to discuss No Way Home later (after you've seen it), I'll be here.

USER TURN 2

<time>2023-12-23T00:59:00</time>

Okay, so feat-wise, which Spider-Man of the three do you think is the strongest in the movies?

ASSISTANT TURN 2

<think>

It's 12:59 AM on December 23, 2023. This is a very late-night message, likely after the user finished No Way Home and is reflecting on the film. The question is about comparing the strength of Spider-Man in the Marvel Cinematic Universe (MCU) as portrayed by Tom Holland, Andrew Garfield, and Tobey Maguire. The user is likely looking for a fun, opinionated take rather than a dry analysis, given the time of night and the context.

</think>

This is a classic question, and it's always fun to debate who's the "strongest" Spider-Man in the MCU movies. Here's my quick take, keeping it light and opinionated:

- **Tom Holland (Peter Parker, *Spider-Man* trilogy, *No Way Home*):** He's the most consistent and well-rounded. Holland's Peter is physically fit, technically skilled, and has access to the best tech (like the Spider-Man armor, the EDITH glasses, and the A.I. that helps him fight). He also has access to the full range of Marvel's tools, which gives him an edge over the older versions. Even without the suit, he's been shown to fight and outsmart some of the most dangerous villains, including Mysterio, Doctor Octopus, and even the alien spider. In *No Way Home*, he proves he can handle an entire rogues' gallery of heroes and villains all at once—so I'd say he's the most versatile and powerful version in the MCU so far.

- **Andrew Garfield (Peter Parker, *The Amazing Spider-Man* trilogy):** Garfield's Peter is more agile, faster, and has better acrobatic skills. He's the one who first runs up the stairs in *The Amazing Spider-Man* (a moment that still makes me laugh), and he's shown incredible agility in his fight choreography. While he doesn't have the same amount of tech or armor, his raw physical skill and street-smart street fighting make him a tough opponent. He can take down larger, more powerful foes without relying on gadgets, which is impressive. But he's still limited by the lack of access to the same level of tech as Holland's Peter.

- **Tobey Maguire (Peter Parker, *Spider-Man* trilogy):** Maguire's Peter is the most "classic" version, more like the comic book hero. He's less physically built, but he makes up for it with clever traps, acrobatics, and a lot of luck. In *Spider-Man 3*, he actually takes down Doctor Octopus in a hand-to-hand fight, which is still a standout moment. But he's also the one who gets beaten by the alien spider in *Spider-Man 3* and is almost killed by the alien spider in *Spider-Man 2*—that really shows his limitations compared to Holland's more resilient version.

Verdict: If you're talking pure power, speed, and access to the most advanced tech, **Tom Holland's Peter** is the strongest. If you're talking raw agility and street fighting, **Andrew Garfield** is the better choice. And if you're just looking for the version that actually survives against the most dangerous alien spider, **Tobey Maguire** is the one who's been beaten by it, while Holland's Peter is the only one who has never been killed by the alien spider in the movies.

That said, in the end, I think Holland's Peter is the one who really proves he's the strongest in the MCU.

TIME-4B completion for Temporal Adaptivity.

F Detailed Ablations and Metrics

This appendix section reports model-wise structural, formatting, and behavioral statistics for Qwen3 baselines in no-thinking and thinking modes, together with the corresponding TIME models. For each model size (4B, 8B, 14B, 32B), we report scenario-aggregated means with bootstrapped 95% confidence intervals in parentheses. We also provide a phase-wise ablation for the 32B setting.

Wilcoxon signed-rank (WSR) tests are computed at the **scenario level**, comparing TIME against the corresponding Qwen3 thinking-mode baseline. All reported WSR p -values therefore reflect paired differences across the 77 TIMEBench scenarios.

Reported metrics include:

- **Benchmark and per-category accuracy**
- **Output and reasoning token statistics**
- **Think-block frequency**
- **Markdown usage (light/heavy)**: light markdown includes **bold** and *italics*; heavy markdown includes lists, headers, tables, and related structures
- **Degenerate-output statistics**

Appendix F.1.1: Structural and behavioral metrics (4B Models)

Statistic	Qwen3-4B		TIME-4B	WSR p
	No-Thinking	Thinking		
Benchmark Score	17.53 (11.95–23.38)	30.13 (23.90–36.36)	52.60 (44.55–60.39)	3.8e–04
Chronological Retrospection	30.00 (10.00–52.73)	53.64 (31.82–75.45)	42.73 (23.64–62.73)	—
Invalid Time Detection	6.36 (0.00–19.09)	20.91 (1.82–42.73)	43.64 (20.00–67.27)	—
Temporal Adaptivity	5.45 (0.00–12.73)	13.64 (0.00–28.18)	49.09 (27.27–70.91)	—
Temporal Contextual Awareness	32.73 (8.18–59.09)	40.00 (16.36–64.55)	76.36 (56.36–92.73)	—
Temporal Flow Anomaly Detection	0.00 (0.00–0.00)	1.82 (0.00–4.55)	40.91 (19.09–63.64)	—
Time Gap Awareness	0.00 (0.00–0.00)	2.73 (0.00–8.18)	57.27 (35.45–78.18)	—
Timezone Sensitivity	48.18 (30.00–66.36)	78.18 (63.64–90.91)	58.18 (41.82–73.64)	—
Mean Output Tokens / Run	621.01 (543.56–700.11)	1753.13 (1477.36–2062.26)	369.28 (318.18–424.01)	<1e–08
Mean Thinking Tokens / Run	0.00 (0.00–0.00)	1116.37 (877.99–1392.25)	85.56 (68.29–110.78)	<1e–08
Mean # Think Blocks / Run	0.00 (0.00–0.00)	0.98 (0.96–1.00)	2.01 (1.23–3.41)	4.8e–05
% Runs w/ Think Blocks	0.00 (0.00–0.00)	98.31 (96.49–99.61)	87.40 (81.69–92.47)	—
Think Block Position in % (S/M/E)	—	100.0 / 0.0 / 0.0	28.96 / 70.72 / 0.32	—
% Heavy Markdown	87.79 (81.82–92.99)	85.97 (78.96–92.08)	68.83 (60.00–77.40)	—
% Light Markdown	95.58 (91.43–98.83)	95.71 (91.56–98.96)	90.91 (86.62–94.68)	—
% Any Degeneracy	3.38 (1.43–5.84)	20.13 (13.77–26.88)	5.19 (2.47–8.96)	2.5e–05
% Malformed Outputs	0.00 (0.00–0.00)	1.69 (0.39–3.51)	1.30 (0.52–2.21)	—
% Infinite Repetitions	2.47 (1.17–4.29)	1.69 (0.39–3.51)	0.13 (0.00–0.39)	—
% Reasoning Leakage	0.91 (0.00–2.60)	3.38 (1.69–5.45)	1.17 (0.26–2.47)	—
% Formatting Leakage	0.00 (0.00–0.00)	16.88 (10.91–23.64)	3.25 (0.91–6.49)	—

Note. WSR = Wilcoxon signed-rank test, computed between Qwen3 thinking mode and TIME across 77 scenario-level comparisons. WSR p -values are shown for aggregate and structural metrics.

Think-position percentages are computed over the total set of observed <think> blocks across runs. They measure the proportion of <think> blocks observed at the start, in the middle, or at the end of the response. S/M/E = start / middle / end. Because these values are included only as descriptive summaries of placement behavior, confidence intervals are not reported.

Appendix F.1.2: Structural and behavioral metrics (8B Models)

Statistic	Qwen3-8B		TIME-8B	WSR p
	No-Thinking	Thinking		
Benchmark Score	21.56 (15.32–28.05)	32.99 (26.88–39.09)	59.87 (53.38–66.23)	1.9e–05
Chronological Retrospection	24.55 (8.18–43.64)	58.18 (36.36–78.18)	50.00 (33.64–65.45)	—
Invalid Time Detection	7.27 (0.00–20.91)	18.18 (3.64–37.27)	45.45 (23.64–67.27)	—
Temporal Adaptivity	21.82 (1.82–44.55)	20.00 (4.55–38.18)	80.00 (67.27–90.91)	—
Temporal Contextual Awareness	34.55 (8.18–61.82)	41.82 (16.36–69.09)	70.00 (50.00–87.27)	—
Temporal Flow Anomaly Detection	0.00 (0.00–0.00)	1.82 (0.00–4.55)	44.55 (30.91–58.18)	—
Time Gap Awareness	0.00 (0.00–0.00)	0.91 (0.00–2.73)	64.55 (48.18–79.09)	—
Timezone Sensitivity	62.73 (41.82–81.82)	90.00 (81.82–97.27)	64.55 (45.45–81.82)	—
Mean Output Tokens / Run	607.29 (528.02–690.00)	1674.88 (1389.19–2001.83)	351.63 (312.57–392.46)	<1e–08
Mean Thinking Tokens / Run	0.00 (0.00–0.00)	1025.63 (776.60–1317.86)	87.31 (72.62–104.71)	<1e–08
Mean # Think Blocks / Run	0.00 (0.00–0.00)	0.98 (0.96–0.99)	2.06 (1.64–2.62)	6.7e–08
% Runs w/ Think Blocks	0.00 (0.00–0.00)	97.92 (95.84–99.35)	86.23 (80.39–91.56)	—
Think Block Position in % (S/M/E)	—	100.0 / 0.0 / 0.0	20.10 / 79.27 / 0.63	—
% Heavy Markdown	80.78 (73.25–87.66)	84.94 (77.79–91.30)	66.23 (57.53–74.81)	—
% Light Markdown	94.55 (90.00–98.18)	95.45 (91.56–98.44)	86.88 (81.17–91.95)	—
% Any Degeneracy	5.19 (2.73–8.31)	17.92 (11.56–24.68)	5.58 (2.60–9.48)	1.8e–03
% Malformed Outputs	0.00 (0.00–0.00)	2.08 (0.65–4.16)	2.21 (0.91–3.90)	—
% Infinite Repetitions	4.81 (2.46–7.79)	2.08 (0.65–4.16)	0.26 (0.00–0.78)	—
% Reasoning Leakage	0.39 (0.00–1.17)	2.99 (1.43–4.94)	0.91 (0.13–2.21)	—
% Formatting Leakage	0.00 (0.00–0.00)	15.06 (9.09–21.69)	2.34 (0.52–4.55)	—

Note. WSR = Wilcoxon signed-rank test, computed between Qwen3 thinking mode and TIME across 77 scenario-level comparisons. WSR p -values are shown for aggregate and structural metrics.

Think-position percentages are computed over the total set of observed <think> blocks across runs. They measure the proportion of <think> blocks observed at the start, in the middle, or at the end of the response. S/M/E = start / middle / end. Because these values are included only as descriptive summaries of placement behavior, confidence intervals are not reported.

Appendix F.1.3: Structural and behavioral metrics (14B Models)

Statistic	Qwen3-14B		TIME-14B	WSR p
	No-Thinking	Thinking		
Benchmark Score	29.48 (22.47–36.36)	34.42 (28.44–40.65)	64.80 (59.09–70.39)	1.6e–06
Chronological Retrospection	64.55 (38.18–88.18)	55.45 (34.52–76.36)	55.45 (39.09–70.00)	—
Invalid Time Detection	20.91 (0.00–43.64)	19.09 (1.82–41.82)	41.82 (25.45–59.09)	—
Temporal Adaptivity	18.18 (2.73–36.36)	23.64 (9.09–40.91)	84.55 (76.36–91.82)	—
Temporal Contextual Awareness	32.73 (8.18–59.09)	42.73 (18.18–69.09)	78.18 (63.64–90.00)	—
Temporal Flow Anomaly Detection	0.00 (0.00–0.00)	1.82 (0.00–4.55)	66.36 (48.18–83.64)	—
Time Gap Awareness	0.00 (0.00–0.00)	4.55 (0.00–10.00)	61.82 (43.64–79.09)	—
Timezone Sensitivity	70.00 (49.09–87.27)	93.64 (90.00–97.27)	65.45 (50.91–80.00)	—
Mean Output Tokens / Run	563.03 (486.58–641.05)	1514.66 (1270.89–1785.34)	317.21 (283.39–353.33)	<1e–08
Mean Thinking Tokens / Run	0.00 (0.00–0.00)	887.42 (683.44–1127.52)	91.19 (76.67–111.07)	<1e–08
Mean # Think Blocks / Run	0.00 (0.00–0.00)	0.99 (0.98–1.00)	1.90 (1.48–2.55)	1.3e–08
% Runs w/ Think Blocks	0.00 (0.00–0.00)	99.09 (98.05–99.87)	91.17 (86.36–95.45)	—
Think Block Position in % (S/M/E)	—	100.0 / 0.0 / 0.0	21.85 / 77.95 / 0.21	—
% Heavy Markdown	83.51 (75.84–90.26)	86.36 (79.48–92.34)	61.82 (52.60–70.78)	—
% Light Markdown	95.58 (90.78–99.35)	97.40 (94.16–99.61)	86.75 (81.43–91.43)	—
% Any Degeneracy	4.16 (1.56–7.66)	15.45 (9.09–22.21)	2.08 (0.65–3.77)	5.3e–05
% Malformed Outputs	0.00 (0.00–0.00)	0.91 (0.13–2.08)	0.26 (0.00–0.65)	—
% Infinite Repetitions	4.16 (1.69–7.66)	1.30 (0.13–2.99)	0.00 (0.00–0.00)	—
% Reasoning Leakage	0.00 (0.00–0.00)	1.30 (0.39–2.47)	0.78 (0.00–1.82)	—
% Formatting Leakage	0.00 (0.00–0.00)	14.29 (8.18–21.04)	1.17 (0.13–2.47)	—

Note. WSR = Wilcoxon signed-rank test, computed between Qwen3 thinking mode and TIME across 77 scenario-level comparisons. WSR p -values are shown for aggregate and structural metrics.

Think-position percentages are computed over the total set of observed <think> blocks across runs. They measure the proportion of <think> blocks observed at the start, in the middle, or at the end of the response. S/M/E = start / middle / end. Because these values are included only as descriptive summaries of placement behavior, confidence intervals are not reported.

Appendix F.1.4: Structural and behavioral metrics (32B Models)

Statistic	Qwen3-32B		TIME-32B	WSR p
	No-Thinking	Thinking		
Benchmark Score	31.82 (25.71–38.31)	37.40 (31.56–43.51)	64.81 (58.18–71.17)	5.0e–07
Chronological Retrospection	61.82 (38.18–83.64)	60.00 (41.82–77.27)	63.64 (45.45–80.00)	—
Invalid Time Detection	11.82 (2.73–25.45)	31.82 (13.64–51.82)	52.73 (32.73–71.82)	—
Temporal Adaptivity	20.00 (1.82–44.55)	26.36 (7.27–47.27)	76.36 (64.55–87.27)	—
Temporal Contextual Awareness	40.00 (13.64–67.27)	43.64 (20.00–69.09)	76.36 (58.18–92.73)	—
Temporal Flow Anomaly Detection	0.91 (0.00–2.73)	3.64 (0.00–8.18)	44.55 (23.64–65.45)	—
Time Gap Awareness	0.91 (0.00–2.73)	3.64 (0.00–7.27)	58.18 (40.00–75.45)	—
Timezone Sensitivity	87.27 (78.18–94.55)	92.73 (83.64–99.09)	81.82 (68.18–92.73)	—
Mean Output Tokens / Run	608.96 (533.37–686.50)	1573.47 (1327.94–1856.57)	332.64 (296.15–371.77)	<1e–08
Mean Thinking Tokens / Run	0.00 (0.00–0.00)	910.52 (705.34–1158.80)	84.16 (71.18–98.27)	<1e–08
Mean # Think Blocks / Run	0.00 (0.00–0.00)	0.99 (0.98–1.00)	1.67 (1.34–2.06)	1.4e–04
% Runs w/ Think Blocks	0.00 (0.00–0.00)	99.22 (97.92–100.00)	80.65 (73.51–87.40)	—
Think Block Position in % (S/M/E)	—	100.0 / 0.0 / 0.0	24.14 / 75.62 / 0.23	—
% Heavy Markdown	83.64 (77.01–89.61)	90.91 (85.19–95.71)	62.73 (53.51–71.69)	—
% Light Markdown	95.84 (92.60–98.31)	98.31 (95.45–99.87)	85.45 (79.74–90.65)	—
% Any Degeneracy	4.42 (2.60–6.49)	18.18 (11.43–25.32)	3.64 (1.95–5.97)	4.2e–04
% Malformed Outputs	0.00 (0.00–0.00)	0.78 (0.00–2.08)	1.95 (1.04–2.99)	—
% Infinite Repetitions	4.29 (2.47–6.23)	0.91 (0.00–2.47)	1.30 (0.52–2.34)	—
% Reasoning Leakage	0.13 (0.00–0.39)	0.91 (0.13–1.95)	1.17 (0.52–1.95)	—
% Formatting Leakage	0.00 (0.00–0.00)	17.01 (10.39–24.16)	1.17 (0.13–2.86)	—

Note. WSR = Wilcoxon signed-rank test, computed between Qwen3 thinking mode and TIME across 77 scenario-level comparisons. WSR p -values are shown for aggregate and structural metrics.

Think-position percentages are computed over the total set of observed <think> blocks across runs. They measure the proportion of <think> blocks observed at the start, in the middle, or at the end of the response. S/M/E = start / middle / end. Because these values are included only as descriptive summaries of placement behavior, confidence intervals are not reported.

Appendix F.2: Phase-wise ablation details (32B Models)

Statistic	Qwen3 32B		Phase 1	Phase 2	Phase 3	TIME (32B)
	No-Thinking	Thinking				
Benchmark Score	31.82	37.40	42.47	56.88	52.08	64.81
Chronological Retrospection	61.82	60.00	63.64	52.73	65.45	63.64
Invalid Time Detection	11.82	31.82	39.09	21.82	30.91	52.73
Temporal Adaptivity	20.00	26.36	40.91	80.00	68.18	76.36
Temporal Contextual Awareness	40.00	43.64	48.18	50.00	49.09	76.36
Temporal Flow Anomaly Detection	0.91	3.64	7.27	49.09	40.00	44.55
Time Gap Awareness	0.91	3.64	3.64	59.09	44.55	58.18
Timezone Sensitivity	87.27	92.73	94.55	85.45	66.36	81.82
Mean Output Tokens / Run	608.96	1573.47	1434.56	362.45	294.51	332.64
Mean Thinking Tokens / Run	0.00	910.52	803.52	76.59	52.94	84.16
Mean # Think Blocks / Run	0.00	0.99	0.99	1.12	1.25	1.67
% Runs w/ Think Blocks	0.00	99.22	99.48	95.58	89.22	80.65
% Think Blocks at Start	—	100.0	100.0	70.65	54.97	24.14
% Think Blocks in Middle	—	0.0	0.0	29.12	44.62	75.62
% Think Blocks at End	—	0.0	0.0	0.23	0.41	0.23
% Heavy Markdown	83.64	90.91	88.57	65.32	61.43	62.73
% Light Markdown	95.84	98.31	97.92	83.38	79.48	85.45
% Any Degeneracy	4.42	18.18	13.90	4.68	0.78	3.64
% Malformed Outputs	0.00	0.78	0.52	1.82	0.26	1.95
% Infinite Repetitions	4.29	0.91	0.52	0.39	0.13	1.30
% Reasoning Leakage	0.13	0.91	1.04	2.47	0.26	1.17
% Formatting Leakage	0.00	17.01	12.47	1.17	0.26	1.17