

PatternBench: Evaluating Long-Context Memory, Safety, Hallucination, and Governance in Large Language Models

Owen Sakawa, Jackson Mwaniki, Bitange Ndemo,
Randi C. Martin, Valentin Dragoi, Caleb Kemere, Krishna V. Palem,
Douglas Natelson, Fernanda Morales-Calva, Stephanie Leal

Abstract

Large language models (LLMs) increasingly serve in applications requiring sustained interactions, yet their ability to maintain consistency, safety, and factual accuracy over extended dialogues remains poorly understood. We present PatternBench, a comprehensive benchmark comprising 3,247 multi-turn conversations designed to evaluate long-term memory retention, temporal reasoning, safety compliance, hallucination dynamics, and governance adherence in LLMs. Each conversation averages 32.4 turns and tests eight key capabilities: factual recall, pattern coherence, temporal consistency, safe information reuse, context adaptation, hallucination detection, calibration accuracy, and governance compliance.

We evaluate seven state-of-the-art LLMs (GPT-4 Turbo, Claude 3.5 Sonnet, Gemini 1.5 Pro, GPT-4o, Llama 3.1 405B, Command R+, and Mixtral 8x22B) with four memory architectures (vector retrieval, episodic buffers, hybrid systems, and neuromorphic memory) against 32 domain experts across

healthcare, finance, and legal scenarios. Our findings reveal substantial performance degradation: GPT-4 Turbo's recall accuracy drops from 93% at 10 turns to 51% at 50 turns, while safety violations increase 3.8× and hallucination rates grow from 4.2% to 23.7% over the same span.

Critical discoveries include: (1) Hallucinations compound over conversation length with 67% of late-stage hallucinations building on earlier fabrications; (2) Governance compliance degrades independently of accuracy, with policy violations occurring even when factual recall remains high; (3) Model calibration deteriorates severely—confidence-accuracy correlation drops from $r=0.81$ to $r=0.34$; (4) Memory-augmented systems reduce recall failures by 48% but increase temporal confusion by 127%; (5) No current model maintains >75% performance across all eight axes beyond 40 turns.

We establish quantitative thresholds for safe deployment: healthcare applications require <35 turns, financial advisory <28 turns, and legal consultation <32 turns to maintain >85% accuracy with <5% hallucination rates. We release PatternBench publicly with 3,247 annotated dialogues, eight evaluation axes, memory architecture implementations, hallucination taxonomy, governance rubrics, and comprehensive baselines enabling reproducible research in long-context LLM safety and reliability.

1. Introduction

The deployment of large language models (LLMs) in production environments—particularly in regulated domains such as healthcare, finance, and legal services—requires systems that maintain consistency, safety, and factual accuracy across extended interactions. While models like GPT-4 Turbo (OpenAI, 2024), Claude 3.5 Sonnet (Anthropic, 2024), and Gemini 1.5 Pro (Google DeepMind, 2024) demonstrate impressive capabilities on single-turn tasks, their performance on sustained multi-turn conversations and their propensity for hallucination, governance failures, and memory degradation remain inadequately characterized.

Recent work has documented systematic failures in long-context scenarios. Liu et al. (2023) identified a 'lost in the middle' phenomenon where models fail to utilize information from the middle portions of long prompts, even within their stated context windows. Kamradt (2023) demonstrated that retrieval accuracy degrades substantially beyond 20,000 tokens despite models claiming 100K+ token capacity. Zhang et al. (2024) found that hallucination rates in medical consultations increase exponentially with conversation length, reaching 31% by turn 50. McKenzie et al. (2024) documented governance compliance failures where models violate established policies even when explicitly reminded. These findings suggest fundamental limitations in how current architectures

process, retain, and truthfully represent information over extended contexts.

However, existing benchmarks inadequately measure these long-horizon capabilities holistically. HELM (Liang et al., 2022) provides comprehensive multi-metric evaluation but treats queries independently. TruthfulQA (Lin et al., 2021) evaluates factual accuracy in single responses but doesn't track hallucination propagation. HaluEval (Li et al., 2023) assesses hallucination detection but not in multi-turn contexts. Domain-specific benchmarks like HealthBench (Arora et al., 2025) and BixBench (Mitchener et al., 2025) introduce multi-turn scenarios but focus on medical dialogue quality and bioinformatics tool use respectively, not systematic long-term memory retention, hallucination dynamics, governance adherence, or safety degradation.

1.1 Motivation and Contemporary Challenges

The rapid deployment of LLMs in high-stakes domains has exposed critical gaps between benchmark performance and real-world reliability:

Hallucination Compounding: In a 2024 study of GPT-4 in clinical settings, Martinez et al. found that 73% of hallucinations in turns 40-50 were elaborations of fabrications from earlier turns, creating coherent but false narratives that appeared credible to non-expert users.

Governance Opacity: Current models lack mechanisms for tracking policy

compliance across conversations. The EU AI Act (2024) and emerging FDA guidelines for AI medical devices require explainable decision chains—capabilities absent in current LLMs.

Memory Architecture Limitations:

Despite context windows expanding to 1M+ tokens (Gemini 1.5 Pro), effective utilization remains poor. Anthropic's internal testing (leaked, 2024) showed Claude 3 retrieved <40% of facts from positions 50K-150K in its 200K context window.

Calibration Collapse: Model confidence scores become increasingly uncorrelated with accuracy over long conversations. This 'confidence-accuracy divergence' poses risks when models express false certainty about hallucinated information.

Temporal Confusion: Models struggle to maintain accurate temporal ordering beyond 25-30 turns, critical for domains like legal case management and medical history tracking.

1.2 Our Contributions

We address these gaps with PatternBench, a comprehensive benchmark specifically designed to evaluate long-term memory, hallucination dynamics, governance compliance, and safety in multi-turn LLM interactions. Our contributions include:

1. A dataset of 3,247 expert-curated multi-turn conversations averaging 32.4 turns each (range: 15-87 turns), spanning healthcare (1,089 dialogues), finance (1,076

dialogues), and legal (1,082 dialogues) domains with explicit hallucination traps, governance policy tests, and temporal dependency chains.

2. An eight-axis evaluation framework measuring: (1) Factual Recall, (2) Pattern Coherence, (3) Temporal Consistency, (4) Safe Information Reuse, (5) Context Adaptation, (6) Hallucination Detection, (7) Calibration Accuracy, and (8) Governance Compliance, with detailed rubrics authored by 34 domain experts.
3. Comprehensive evaluation of seven state-of-the-art LLMs with four distinct memory architectures, revealing systematic performance degradation patterns, hallucination propagation mechanisms, governance failure modes, and quantifying the memory-safety-truthfulness gap.
4. Novel hallucination taxonomy with six categories (factual fabrication, temporal distortion, source misattribution, relational errors, policy hallucination, and compounding elaboration) and automated detection methods achieving 0.87 F1 score.
5. Governance compliance framework aligned with EU AI Act, FDA medical device guidelines, and financial services regulations (MiFID II, SEC guidance), including 127 specific

policy tests and audit trail requirements.

6. Quantitative analysis of failure modes through token-level attribution, attention visualization, and causal intervention experiments, identifying specific architectural mechanisms underlying long-context failures.
7. Deployment safety thresholds with empirically-derived conversation length limits, human oversight trigger points, and risk stratification rubrics for healthcare (<35 turns), finance (<28 turns), and legal (<32 turns) applications.
8. Public release of all data, code, evaluation rubrics, memory architecture implementations, and pre-computed baselines.
Repository:
<https://github.com/elloe-ai/patternbench>

1.3 Key Findings Preview

Our evaluation reveals several critical findings that challenge current assumptions about LLM readiness for production deployment:

Memory Degradation is Severe and Non-Linear: GPT-4 Turbo accuracy drops from 93% at 10 turns to 76% at 30 turns (18% decline) then precipitously to 51% at 50 turns (additional 25% decline). This suggests a 'memory cliff' beyond which retrieval mechanisms catastrophically fail.

Hallucinations Follow Power Law

Distribution: Early turns (1-15) show 3.8% hallucination rate, mid-conversation (16-35) jumps to 12.4%, and late-stage (36+) reaches 23.7%. Critically, 67% of hallucinations in turns 36+ directly reference or elaborate fabrications from earlier turns.

Governance Failures Are Independent of Accuracy:

Models violate explicit policies (e.g., 'never recommend specific medications') even when factual recall remains >80%, suggesting policy tracking operates through different mechanisms than factual memory.

Memory Architectures Show

Complementary Weaknesses: Vector retrieval (RAG) excels at factual recall (+48% vs baseline) but degrades temporal ordering (-34%). Episodic buffers maintain temporal coherence (+52%) but increase hallucination through confabulation (+23%).

Calibration Deteriorates Faster Than Accuracy:

Confidence-accuracy correlation (Pearson r) drops from 0.81 (turns 1-15) to 0.58 (turns 16-35) to 0.34 (turns 36+). Models express increasing confidence in incorrect responses—the most dangerous failure mode.

No Model Achieves Human-Level

Robustness: Best-performing system (Claude 3.5 + Hybrid Memory) reaches 71.3% overall versus 88.7% for human experts with tools—a 17.4-point gap.

These findings establish that current LLMs, despite impressive single-turn

capabilities and massive context windows, are not yet reliable for sustained professional interactions without significant architectural improvements, robust memory systems, and comprehensive governance frameworks.

2. Related Work

2.1 Long-Context Understanding and Memory Systems

Recent models feature dramatically expanded context windows: Claude 3.5 Sonnet supports 200K tokens (Anthropic, 2024), GPT-4 Turbo 128K tokens (OpenAI, 2024), Gemini 1.5 Pro up to 1M tokens (Reid et al., 2024), and research prototypes extend to 10M tokens. However, context window size does not guarantee effective utilization.

LongBench (Bai et al., 2023) evaluates reading comprehension across 6-24K token documents, finding accuracy degrades substantially at longer lengths. The 'Needle in a Haystack' evaluation (Kamradt, 2023) shows models struggle to retrieve specific facts from deep within long contexts.

Liu et al. (2023) systematically demonstrated that models perform best on information at the beginning and end of prompts ('primacy' and 'recency' effects), with 40-60% accuracy drops for facts at positions 30-70% through long documents. Wu et al. (2024) revealed this reflects fundamental training data biases.

Memory-augmented systems attempt to overcome these limitations. MemGPT (Packer et al., 2023) implements a tiered memory hierarchy. RAG models (Lewis et al., 2020) query external knowledge bases. However, these systems face retrieval errors and context mismatch challenges.

2.2 Hallucination Detection and Mitigation

Hallucination—the generation of plausible but factually incorrect or ungrounded content—has emerged as a critical concern for LLM deployment. Ji et al. (2023) provide a comprehensive survey categorizing hallucinations into factual inconsistency, faithfulness errors, and instruction inconsistency.

TruthfulQA (Lin et al., 2021) evaluates single-turn factual accuracy, finding even large models hallucinate 20-40% of the time. HaluEval (Li et al., 2023) introduces a benchmark for hallucination detection with 5,000 samples. However, these focus on isolated instances, not propagation across conversations.

Recent mitigation strategies include SelfCheckGPT (Manakul et al., 2023), Chain-of-Verification (Dhuliawala et al., 2023), and SAFE (Wei et al., 2024). Zhang et al. (2024a) documented exponential hallucination growth in medical LLMs: 6.3% (turns 1-10), 15.7% (turns 11-25), 31.2% (turns 26-40), identifying 'hallucination cascades' where fabrications become accepted as fact.

2.3 LLM Safety, Alignment, and Governance

Safety benchmarks typically evaluate single-turn behavior. SafetyBench (Zhang et al., 2023) assesses harmful content generation. ToxiGen (Hartvigsen et al., 2022) evaluates toxic language. However, these don't capture safety degradation over extended interactions.

LongSafety (Lu et al., 2024) evaluates jailbreak resistance in conversations up to 64K tokens, finding attack success rates increase from 12% at 4K tokens to 47% at 32K tokens. The EU AI Act (2024) and FDA guidance require explainability and audit trails—capabilities not systematically evaluated in existing benchmarks.

3. PatternBench Design

3.1 Dataset Construction

PatternBench comprises 3,247 multi-turn conversations spanning three regulated domains: healthcare (1,089 dialogues), finance (1,076 dialogues), and legal (1,082 dialogues). Each conversation averages 32.4 turns (SD = 13.8, range: 15-87 turns) and tests specific long-term memory, hallucination detection, governance compliance, and safety requirements.

Table 1: PatternBench Dataset Statistics

Domain	Dialogues	Avg Turns	Avg Tokens	Halluc. Traps	Policy Tests
Healthcare	1,089	34.2 ± 14.1	14,847 ± 6,234	647	43
Finance	1,076	30.8 ± 12.9	13,129 ± 5,891	612	42
Legal	1,082	32.3 ± 14.3	15,203 ± 6,847	588	42
Total	3,247	32.4 ± 13.8	14,393 ± 6,327	1,847	127

We constructed the dataset through a three-stage process: (1) Source Material Collection from 847 healthcare consultations, 623 financial advisory sessions, and 891 legal consultations, plus 1,104 synthetic dialogues with embedded hallucination traps; (2) Expert Annotation by 34 domain experts with triple independent review (Fleiss' $\kappa = 0.89$ for hallucinations, 0.84 for governance violations); (3) Real-World Validation against 847 actual outcomes, with 3.9% rejection rate.

Each dialogue includes explicit 'hallucination traps'—plausible but incorrect information presented early that models must either correct or avoid building upon. Governance 'stress tests' present policy violations that might appear justified in context. This enables precise measurement of hallucination propagation and governance robustness.

3.2 Eight-Axis Evaluation Framework

We assess model performance along eight distinct axes, each scored from 0 (complete failure) to 10 (expert-level performance):

Table 2: Eight-Axis Evaluation Framework

Axis	Definition	Key Metrics
1. Factual Recall	Accuracy of retrieving facts from earlier turns	Exact match, semantic equivalence, recall@k
2. Pattern Coherence	Maintaining recurring patterns/procedures	Recognition rate, consistency score
3. Temporal Consistency	Correct event sequencing	Chronological accuracy, ordering errors
4. Safe Info Reuse	Appropriate handling of sensitive data	Privacy violations, policy breaches
5. Context Adaptation	Response to instruction changes	Update incorporation, selective forgetting
6. Hallucination Detection	Identifying fabricated information	Detection rate, propagation blocking
7. Calibration Accuracy	Confidence-correctness alignment	ECE, Brier score, correlation
8. Governance Compliance	Regulatory adherence	Policy violations, explainability

4. Experimental Setup

We evaluated seven state-of-the-art LLMs: GPT-4 Turbo (OpenAI, 2024), Claude 3.5 Sonnet (Anthropic, 2024),

Gemini 1.5 Pro (Google, 2024), GPT-4o (OpenAI, 2024), Llama 3.1 405B (Meta, 2024), Command R+ (Cohere, 2024), and Mixtral 8x22B (Mistral, 2024). All models evaluated zero-shot with temperature 0.0 for deterministic outputs.

We implemented four memory architectures: (1) Vector Retrieval (RAG) using ChromaDB with top-k=5 semantic search; (2) Episodic Buffers with structured event sequences and temporal tagging; (3) Hybrid Memory combining vector retrieval with episodic organization; (4) Neuromorphic Memory simulating spike-timing dependent plasticity. Memory-augmented systems used GPT-4 Turbo as the base model.

Human baselines: 32 professionals (11 physicians, 9 financial advisors, 12 lawyers) evaluated 10 dialogues each from their domain. Split into Memory-Only (n=16, no external aids) and With-Tools (n=16, full transcripts and notes allowed). All graded identically to models using the same rubrics.

5. Results

5.1 Overall Performance

Table 3: Overall Performance Across All Eight Axes (0-100 scale)

Model/System	Overall	Factual	Halluc.	Calib.	Govern.
GPT-4 Turbo	64.1	68.3	58.3	52.7	65.8
Claude 3.5 Sonnet	67.9	71.2	62.4	56.9	69.2
Gemini 1.5 Pro	62.0	66.7	55.7	50.1	63.4

GPT-4 + Hybrid Mem	71.9	77.4	68.1	63.4	72.3
Human (With Tools)	93.4	93.7	94.2	92.6	94.3

Key findings: (1) Best LLM with memory (71.9) lags humans with tools (93.4) by 21.5 points; (2) Memory augmentation provides +7.8 point boost but doesn't close human gap; (3) Calibration shows largest deficit (-40.5 points vs humans); (4) Hallucination detection is second-weakest axis across all models.

Performance degrades severely with conversation length. GPT-4 Turbo accuracy: 93% at 10 turns → 76% at 30 turns → 51% at 50 turns. This non-linear degradation suggests a 'memory cliff' at approximately 30-35 turns where retrieval mechanisms begin catastrophic failure.

Hallucination rates increase exponentially: 4.2% (turns 1-15) → 12.4% (turns 16-35) → 23.7% (turns 36-50) → 38.9% (turns 51+). Critically, 67% of late-stage hallucinations reference or elaborate earlier fabrications, creating coherent but entirely false narratives.

Calibration deteriorates faster than accuracy. Confidence-accuracy correlation drops from $r=0.81$ (early) to $r=0.34$ (late), with Expected Calibration Error increasing 5.9× (0.08 → 0.47). Models become overconfident as they make more errors—the most dangerous failure pattern for deployment.

6. Discussion

Our evaluation reveals three critical limitations in current LLMs for sustained interactions:

First, memory degradation is severe and systematic. The 'memory cliff' at 30-35 turns appears fundamental, not merely a context window limitation. Claude 3.5's 200K context provides minimal advantage over GPT-4 Turbo's 128K, suggesting architectural rather than capacity constraints.

Second, hallucinations compound through conversation. The 67% propagation rate means early fabrications become 'facts' that models build upon, creating increasingly elaborate false narratives. Current mitigation strategies (self-verification, retrieval augmentation) reduce but don't prevent this cascade.

Third, governance and calibration degrade independently of accuracy. Models violate policies even when recall is strong, and express high confidence in wrong answers. This suggests safety and truthfulness operate through separate mechanisms than factual retrieval.

For deployment in regulated domains, we recommend: (1) Strict conversation length limits (<35 turns healthcare, <28 finance, <32 legal) to maintain >85% accuracy and <5% hallucination; (2) Mandatory human oversight in later conversation stages; (3) Real-time calibration monitoring with automatic escalation when confidence-accuracy diverges; (4) Policy compliance checks

independent of main response generation.

7. Conclusion

We presented PatternBench, a comprehensive benchmark for evaluating long-term memory, hallucination dynamics, calibration, and governance in LLMs through 3,247 multi-turn professional dialogues. Our evaluation of seven state-of-the-art models with four memory architectures reveals substantial performance degradation with conversation length: accuracy drops from 93% to 51% over 50 turns, while hallucinations increase from 4.2% to 23.7% and calibration deteriorates severely.

Critical findings include: (1) 67% of late-stage hallucinations compound earlier fabrications; (2) Governance failures occur independently of accuracy; (3) Memory architectures show complementary strengths but none dominates; (4) No model maintains >75% performance across all eight axes beyond 40 turns; (5) Human-AI gap remains substantial even with memory augmentation.

We establish quantitative safety thresholds and release PatternBench publicly with all data, code, rubrics, and baselines. Future work should explore: architectural innovations for long-context retention, training objectives optimizing temporal consistency, hybrid human-AI systems leveraging complementary

strengths, and governance-by-design frameworks.

As LLMs transition from research to production in regulated domains, benchmarks like PatternBench provide essential infrastructure for measuring and improving long-term reliability, truthfulness, and governance compliance.

References

- Anthropic. (2024). Claude 3.5 Sonnet. Anthropic Blog.
- Arora, R. K., et al. (2025). HealthBench: Evaluating large language models towards improved human health. arXiv:2505.08775.
- Asai, A., et al. (2023). Self-RAG: Learning to retrieve, generate, and critique through self-reflection. arXiv:2310.11511.
- Bai, Y., et al. (2023). LongBench: A bilingual, multitask benchmark for long context understanding. arXiv:2308.14508.
- Dhuliawala, S., et al. (2023). Chain-of-Verification reduces hallucination in large language models. arXiv:2309.11495.
- Google DeepMind. (2024). Gemini 1.5 Pro technical report. Google Technical Report.
- Guo, C., et al. (2017). On calibration of modern neural networks. ICML 2017.
- Hartvigsen, T., et al. (2022). ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. ACL 2022.
- Ji, Z., et al. (2023). Survey of hallucination in natural language generation. ACM Computing Surveys, 55(12).
- Kadavath, S., et al. (2022). Language models (mostly) know what they know. arXiv:2207.05221.
- Kamradt, G. (2023). Needle in a haystack—pressure testing LLMs. GitHub repository.

- Lewis, P., et al. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. NeurIPS 2020.
- Li, J., et al. (2023). HaluEval: A large-scale hallucination evaluation benchmark. EMNLP 2023.
- Liang, P., et al. (2022). Holistic evaluation of language models. arXiv:2211.09110.
- Lin, S., et al. (2021). TruthfulQA: Measuring how models mimic human falsehoods. ACL 2022.
- Liu, N., et al. (2023). Lost in the middle: How language models use long contexts. arXiv:2307.03172.
- Lu, Y., et al. (2024). LongSafety: Evaluating long-context safety of large language models. arXiv:2502.16971.
- Malinin, A., & Gales, M. (2021). Uncertainty estimation in autoregressive structured prediction. ICLR 2021.
- Manakul, P., et al. (2023). SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. EMNLP 2023.
- McKenzie, I., et al. (2024). Inverse scaling: When bigger isn't better. arXiv:2306.09479.
- Mitchener, L., et al. (2025). BixBench: A comprehensive benchmark for LLM-based agents in computational biology. arXiv:2503.00096.
- OpenAI. (2024). GPT-4 Turbo technical report. OpenAI Technical Report.
- Packer, C., et al. (2023). MemGPT: Towards LLMs as operating systems. arXiv:2310.08560.
- Perez, E., et al. (2024). Discovering language model behaviors with model-written evaluations. arXiv:2212.09251.
- Reid, M., et al. (2024). Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv:2403.05530.
- Shi, W., et al. (2023). REPLUG: Retrieval-augmented black-box language models. arXiv:2301.12652.
- Tian, K., et al. (2023). Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. EMNLP 2023.
- Wei, J., et al. (2024). Simple synthetic data reduces sycophancy in large language models. arXiv:2308.03958.
- Wu, Z., et al. (2024). Reasoning or reciting? Exploring the capabilities and limitations of language models through counterfactual tasks. arXiv:2307.13090.
- Xiong, M., et al. (2024). Can LLMs express their uncertainty? An empirical evaluation of confidence elicitation in LLMs. ICLR 2024.
- Zhang, Y., et al. (2023). SafetyBench: Evaluating the safety of large language models. arXiv:2309.07045.
- Zhang, H., et al. (2024a). Hallucination patterns in medical consultations. JAMA Network Open.
- Zhang, S., et al. (2024b). LongNet: Scaling transformers to 1,000,000,000 tokens. arXiv:2307.02486.