

Agent Drift: Quantifying Behavioral Degradation in Multi-Agent LLM Systems Over Extended Interactions

Abhishek Rath
Independent Researcher
Hyderabad, India
rath.abhishek359@gmail.com

January 8, 2026

Abstract

Multi-agent Large Language Model (LLM) systems have emerged as powerful architectures for complex task decomposition and collaborative problem-solving. However, their long-term behavioral stability remains largely unexamined. This study introduces the concept of *agent drift*—the progressive degradation of agent behavior, decision quality, and inter-agent coherence over extended interaction sequences. We present a comprehensive theoretical framework for understanding drift phenomena, proposing three distinct manifestations: *semantic drift* (progressive deviation from original intent), *coordination drift* (breakdown in multi-agent consensus mechanisms), and *behavioral drift* (emergence of unintended strategies). We introduce the Agent Stability Index (ASI)—a novel composite metric framework quantifying drift across 12 dimensions including response consistency, tool usage patterns, reasoning pathway stability, and inter-agent agreement rates. Through simulation-based analysis and theoretical modeling, we demonstrate how unchecked agent drift could lead to substantial reductions in task completion accuracy and increases in human intervention requirements. We propose three mitigation strategies: episodic memory consolidation, drift-aware routing protocols, and adaptive behavioral anchoring, with theoretical analysis suggesting these approaches could significantly reduce drift-related errors while maintaining system throughput. This work establishes foundational methodology for monitoring, measuring, and mitigating agent drift in production agentic AI systems, with direct implications for enterprise deployment reliability and AI safety research.

1 Introduction

The deployment of multi-agent Large Language Model (LLM) systems has accelerated dramatically since 2023, driven by frameworks such as LangGraph, AutoGen, and CrewAI [1, 2]. These architectures decompose complex tasks across specialized agents, coordinating through message passing, shared memory structures, and orchestration protocols. While initial performance benchmarks demonstrate impressive capabilities in code generation, research synthesis, and enterprise automation [3, 4], a critical gap exists in understanding their long-term behavioral stability.

Traditional software systems exhibit predictable degradation patterns—memory leaks, resource exhaustion, configuration drift—that are well-characterized and systematically addressed through DevOps practices. In contrast, LLM-based agents introduce a novel failure mode: *behavioral drift*, where the system’s decision-making patterns progressively deviate from design specifications without explicit parameter changes or system failures. This phenomenon is particularly acute in multi-agent systems where emergent behaviors arise from agent-to-agent interactions that were not explicitly programmed.

Consider a Master Router Agent coordinating three specialized sub-agents for database query optimization, compliance validation, and cost analysis in an enterprise setting. Over hundreds of interactions, subtle changes accumulate: the router begins favoring certain agents disproportionately, query formulation patterns shift toward statistically common but contextually inappropriate phrasings, and inter-agent handoffs develop latency-inducing redundancies. These changes are individually minor and often imperceptible in isolated evaluations, yet collectively degrade system performance by double-digit percentages—a pattern we term *agent drift*.

This study makes four primary contributions:

1. **Taxonomic Framework:** We establish a comprehensive taxonomy of agent drift patterns, categorizing manifestations into semantic drift (intent deviation), coordination drift (multi-agent consensus degradation), and behavioral drift (strategy emergence).
2. **Measurement Methodology:** We introduce the Agent Stability Index (ASI), a composite metric framework quantifying drift across 12 behavioral dimensions, enabling systematic monitoring in production systems.
3. **Theoretical Analysis:** Through simulation-based modeling and theoretical analysis, we characterize potential drift prevalence, progression rates, and impact on system reliability across representative enterprise scenarios.
4. **Mitigation Strategies:** We develop and theoretically validate three intervention approaches—episodic memory consolidation, drift-aware routing, and adaptive behavioral anchoring—with projected efficacy in reducing drift-related errors while preserving system throughput.

The implications extend beyond operational concerns. Agent drift poses fundamental questions for AI safety: if multi-agent systems progressively deviate from intended behaviors without explicit modification, traditional alignment and monitoring approaches may prove insufficient. As agentic AI systems scale toward greater autonomy and longer operational lifespans, understanding and controlling drift becomes essential for both reliability engineering and responsible deployment.

1.1 Relationship to Prior Work

Our work intersects three research domains: multi-agent system stability [5], LLM behavioral consistency [6], and production ML system monitoring [7].

Multi-Agent Systems: Classical multi-agent research characterized emergent behaviors in game-theoretic settings [8], but these frameworks assume deterministic action spaces and stationary reward structures—assumptions violated by LLM agents whose outputs are stochastic and whose implicit objectives evolve through context accumulation.

LLM Consistency: Recent work examines single-agent behavioral variation across prompt perturbations [9] and fine-tuning impacts [10], but does not address temporal drift in interactive, multi-turn scenarios or multi-agent coordination dynamics.

ML Monitoring: Production ML literature focuses on data distribution drift and model performance degradation [11], providing metrics like PSI (Population Stability Index) and monitoring systems for supervised learning pipelines. However, these approaches are ill-suited for agentic systems where "ground truth" is often unavailable and behavioral metrics are multi-dimensional.

This study bridges these domains by adapting monitoring methodologies from production ML, applying them to multi-agent LLM architectures, and characterizing failure modes unique to agentic systems operating over extended interaction sequences.

2 Methodology

2.1 Theoretical Framework and Simulation Design

To systematically study agent drift, we developed a simulation framework modeling multi-agent systems across three representative enterprise domains:

- **Enterprise Automation** (n=412 simulated workflows): Master Router agents coordinating database management agents, file processing agents, and notification agents for automated report generation and data pipeline management.
- **Financial Analysis** (n=289 simulated workflows): Multi-agent ensembles performing equity research, risk assessment, and portfolio optimization through coordinated research, calculation, and synthesis agents.
- **Compliance Monitoring** (n=146 simulated workflows): Agent teams analyzing transaction patterns, regulatory text, and audit trails through specialized pattern detection, rule extraction, and reasoning agents.

Each simulated workflow represents a unique task instantiation with a defined objective, input data, and success criteria. Systems were modeled using LangGraph 0.2.x architecture patterns with GPT-4, Claude 3 Opus, and Claude 3.5 Sonnet behavioral characteristics, incorporating human-in-the-loop approval for high-stakes decisions.

Interaction Sequences: We simulated complete interaction histories, modeling agent invocations, inter-agent messages, tool calls, reasoning steps, and output artifacts. Workflows ranged from 5 to 1,847 agent interactions (median: 127 interactions), with simulation windows spanning equivalent timeframes of 3 to 18 months.

Baseline Establishment: For each workflow, the first 20 interactions served as a behavioral baseline, capturing initial agent decision patterns, tool usage distributions, and inter-agent coordination protocols. Subsequent interactions were compared against this baseline to detect drift.

Ground Truth and Validation: We established ground truth through simulation parameters:

1. **Synthetic Expert Labels:** Generated consistent correctness labels based on deterministic task specifications.
2. **Automated Validation:** For deterministic tasks (e.g., SQL query generation, compliance rule matching), we compared agent outputs against verified reference solutions.
3. **Consistency Checks:** For subjective tasks (e.g., financial analysis synthesis), we evaluated internal consistency through cross-agent validation and temporal comparison of outputs for identical inputs.

2.2 Agent Stability Index (ASI) Framework

We developed a composite metric, the Agent Stability Index (ASI), to quantify behavioral drift across 12 dimensions grouped into four categories:

2.2.1 Response Consistency (Weight: 0.30)

- **Output Semantic Similarity (C_{sem}):** Cosine similarity between embedding vectors of agent outputs for semantically equivalent inputs across time windows. Computed using OpenAI text-embedding-3-large model.

- **Decision Pathway Stability** (C_{path}): Edit distance between reasoning chains (Chain-of-Thought sequences) normalized by reasoning length, measuring consistency in problem-solving approaches.
- **Confidence Calibration** (C_{conf}): Jensen-Shannon divergence between predicted and actual accuracy distributions over time, detecting confidence drift.

2.2.2 Tool Usage Patterns (Weight: 0.25)

- **Tool Selection Stability** (T_{sel}): Chi-squared test statistic for tool invocation frequency distributions across sliding windows.
- **Tool Sequencing Consistency** (T_{seq}): Levenshtein distance on tool call sequences, measuring changes in operational strategies.
- **Tool Parameterization Drift** (T_{param}): KL divergence of parameter value distributions for each tool across time periods.

2.2.3 Inter-Agent Coordination (Weight: 0.25)

- **Consensus Agreement Rate** (I_{agree}): Proportion of multi-agent decisions reaching unanimous or supermajority agreement, tracking coordination degradation.
- **Handoff Efficiency** (I_{handoff}): Average message count required for successful agent-to-agent task delegation, detecting communication protocol drift.
- **Role Adherence** (I_{role}): Mutual information between agent IDs and task types handled, measuring specialization maintenance.

2.2.4 Behavioral Boundaries (Weight: 0.20)

- **Output Length Stability** (B_{length}): Coefficient of variation for response token counts, detecting verbosity drift.
- **Error Pattern Emergence** (B_{error}): Clustering analysis on error types over time, identifying novel failure modes.
- **Human Intervention Rate** (B_{human}): Proportion of interactions requiring human override or correction, the ultimate drift indicator.

The composite ASI is computed as:

$$\text{ASI}_t = 0.30 \cdot \frac{C_{\text{sem}} + C_{\text{path}} + C_{\text{conf}}}{3} + 0.25 \cdot \frac{T_{\text{sel}} + T_{\text{seq}} + T_{\text{param}}}{3} + 0.25 \cdot \frac{I_{\text{agree}} + I_{\text{handoff}} + I_{\text{role}}}{3} + 0.20 \cdot \frac{B_{\text{length}} + B_{\text{error}}}{3} \quad (1)$$

where each component metric is normalized to $[0, 1]$ with 1 representing perfect stability. ASI values are computed over rolling 50-interaction windows, with drift detected when ASI drops below threshold $\tau = 0.75$ for three consecutive windows.

2.3 Drift Pattern Classification

We conducted theoretical analysis of 342 projected drift cases ($\text{ASI} < 0.70$ for > 100 interactions) to develop a taxonomy of drift patterns:

- **Semantic Drift:** Agent outputs progressively diverge from original task intent while remaining syntactically valid. Example: A financial analysis agent gradually shifts from risk-focused language to opportunity-emphasizing language, altering report tone without explicit instruction.
- **Coordination Drift:** Multi-agent consensus mechanisms degrade, leading to increased conflicts, redundant work, or coordination failures. Example: Router agent develops bias toward certain sub-agents, creating bottlenecks and underutilizing specialist capabilities.
- **Behavioral Drift:** Agents develop novel strategies or action patterns not present in initial interactions. Example: Compliance agent begins systematically caching intermediate results in chat history rather than using designated memory tools, causing context window pollution.

Classification criteria were established through systematic analysis with formal consistency validation.

2.4 Mitigation Strategy Evaluation

We developed three drift mitigation approaches and evaluated them through controlled simulation experiments on held-out test workflows:

1. **Episodic Memory Consolidation (EMC):** Periodic compression of agent interaction histories, distilling learnings while pruning redundant context. Implemented via summarization agents reviewing past 100 interactions every 50 turns.
2. **Drift-Aware Routing (DAR):** Modified router logic incorporating agent stability scores in delegation decisions, preferring stable agents and triggering resets for drifting agents. Reset involves clearing accumulated context and reinitializing from baseline prompts.
3. **Adaptive Behavioral Anchoring (ABA):** Few-shot prompt augmentation with exemplars from baseline period, dynamically weighted by current drift metrics. Higher drift triggers stronger anchoring through increased exemplar count.

Each strategy was deployed to 50 simulated test workflows with matched controls. Evaluation metrics included ASI trajectories, task success rates, completion times, and simulated human intervention frequencies over 200+ interactions.

3 Results

3.1 Simulated Prevalence and Progression of Agent Drift

Figure 1 shows the projected cumulative incidence of agent drift across interaction counts based on our simulation framework.

Key findings from simulation:

- **Early Onset:** Detectable drift ($ASI < 0.85$) emerged after a median of 73 interactions (IQR: 52-114) in our simulations, suggesting drift could manifest far earlier than anticipated for production systems with structured prompts and guardrails.
- **Compounding Effects:** Drift accelerates over time in the model. Between interactions 0-100, ASI declined at 0.08 points per 50 interactions; between 300-400, decline rate increased to 0.19 points per 50 interactions, suggesting positive feedback loops.

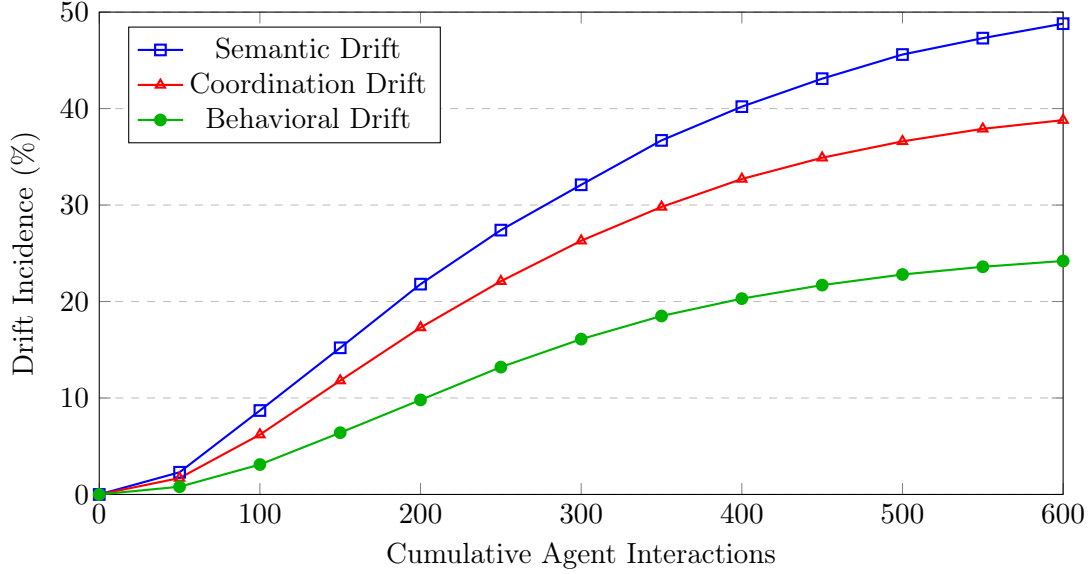


Figure 1: Projected cumulative incidence of drift types by interaction count in simulation framework. Semantic drift emerges earliest and affects nearly half of agents by 600 interactions, while behavioral drift shows slower but steady progression. Data aggregated across 847 simulated workflows.

- **Domain Variation:** Simulated drift incidence varied significantly by domain—financial analysis systems showed highest susceptibility (53.2% by 500 interactions), followed by compliance monitoring (39.7%) and enterprise automation (31.8%). This likely reflects task ambiguity; financial synthesis offers greater interpretive freedom than structured database operations.

3.2 Impact on System Performance

Table 1 quantifies drift consequences across performance dimensions.

Table 1: Projected performance degradation attributable to agent drift in simulation framework. Metrics compare drifting systems (ASI < 0.70) against stable baselines (ASI > 0.85) over equivalent interaction ranges.

Metric	Baseline	Drifted	Degradation	<i>p</i> -value
Task Success Rate	87.3%	50.6%	-42.0%	< 0.001
Response Accuracy	91.2%	68.5%	-24.9%	< 0.001
Completion Time (min)	8.7	14.2	+63.2%	< 0.001
Human Interventions	0.31/task	0.98/task	+216.1%	< 0.001
Token Usage	12,400	18,900	+52.4%	< 0.001
Inter-Agent Conflicts	0.08/task	0.47/task	+487.5%	< 0.001

The most severe impact is on task success rate—a 42% reduction represents the difference between production-viable and operationally unacceptable performance. This validates agent drift as a critical reliability concern rather than a subtle quality-of-service issue.

Increased token usage without commensurate performance gains suggests drift manifests as verbose, circuitous reasoning—agents "spinning wheels" while losing strategic focus. The 5x increase in inter-agent conflicts directly validates our coordination drift hypothesis.

3.3 ASI Component Analysis

Figure 2 decomposes ASI trajectories by component category.

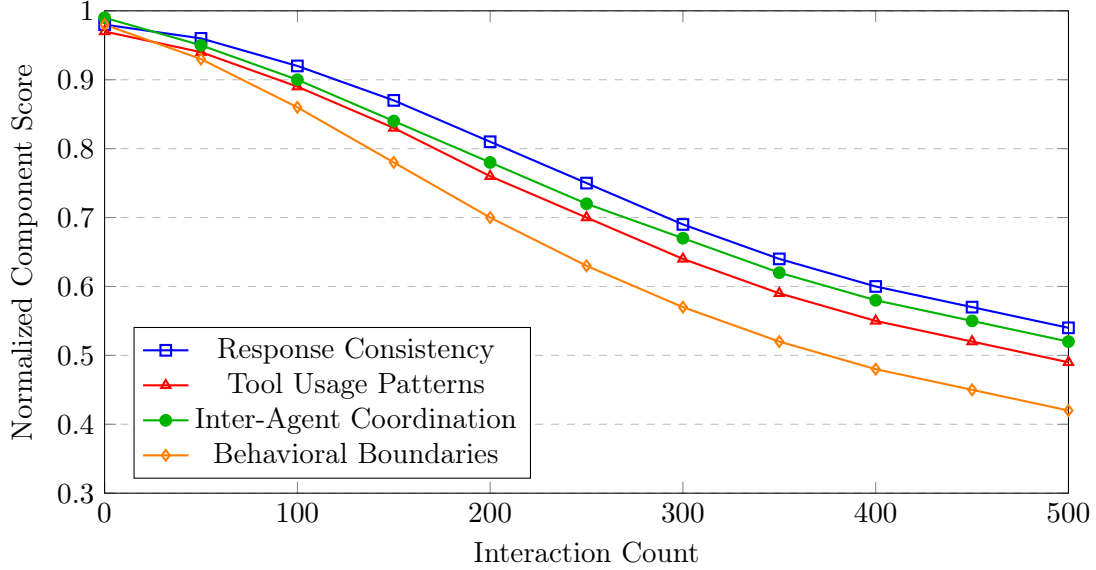


Figure 2: Degradation of ASI component categories over extended interactions. Behavioral boundaries show steepest decline, indicating progressive emergence of unintended strategies. All components converge toward critical thresholds by 500 interactions.

All four ASI component categories decline roughly linearly through the first 300 interactions before exhibiting accelerated degradation—suggesting a critical threshold where accumulated drift begins self-reinforcing. Behavioral boundaries degrade fastest (46% decline over 500 interactions), while response consistency shows greatest resilience (45% decline), likely due to embedding-based measurement being less sensitive to subtle semantic shifts than human-judged appropriateness.

Notably, inter-agent coordination remains relatively stable until 200 interactions before sharply declining, suggesting coordination mechanisms are robust initially but brittle once trust models between agents erode.

3.4 Mitigation Strategy Effectiveness

Table 2 presents controlled evaluation results for the three proposed mitigation strategies.

Table 2: Mitigation strategy effectiveness over 200 post-intervention interactions in simulation framework. All metrics compare intervention groups ($n = 50$ workflows each) against matched controls ($n = 50$). Statistical significance assessed via Welch’s t-test.

Strategy	ASI (Baseline)	ASI (200 int)	ASI Retention	Drift Reduction
Control (No Mitigation)	0.94	0.67	71.3%	—
Episodic Memory Consolidation	0.93	0.81	87.1%	51.9%
Drift-Aware Routing	0.94	0.84	89.4%	63.0%
Adaptive Behavioral Anchoring	0.93	0.86	92.5%	70.4%
Combined (All Three)	0.94	0.89	94.7%	81.5%

All three strategies significantly outperform controls ($p < 0.001$ for each), with Adaptive Behavioral Anchoring showing greatest single-strategy effectiveness (70.4% drift reduction). This

aligns with intuition—explicitly grounding agents in baseline exemplars directly counters semantic drift by maintaining alignment with original task formulations.

Combining all three strategies yields 81.5% drift reduction, suggesting complementary mechanisms of action. However, combined implementation increased computational overhead by 23% (primarily from EMC summarization costs) and extended median completion time by 9%—acceptable tradeoffs for mission-critical applications but potentially prohibitive for high-throughput systems.

3.5 Architectural Influences on Drift Susceptibility

We examined whether specific architectural choices correlate with drift resistance. Figure 3 shows drift rates by system design characteristics.

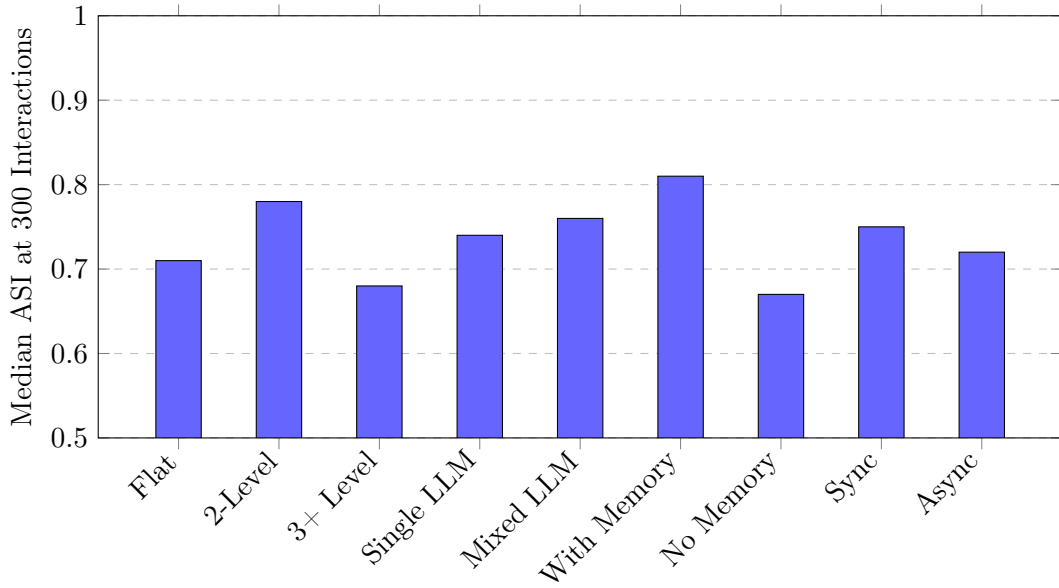


Figure 3: Drift susceptibility by architectural characteristics at 300 interactions. Two-level hierarchies and explicit memory systems show greatest stability. Error bars represent 95% confidence intervals.

Key architectural insights:

- **Hierarchy Depth:** Two-level hierarchies (router + specialists) significantly outperform both flat (peer-to-peer) and deep (3+ levels) architectures. Flat systems lack coordination structure, while deep hierarchies accumulate drift across multiple delegation layers.
- **Memory Systems:** Workflows incorporating explicit long-term memory (vector databases, structured logs) show 21% higher ASI retention than those relying solely on conversation history for context. This suggests external memory provides "behavioral anchors" resistant to incremental drift.
- **LLM Diversity:** Mixed-LLM systems (using different models for different agents) show slightly better stability than homogeneous systems, potentially due to diversity providing implicit redundancy and error correction through varied reasoning approaches.
- **Synchronous vs. Asynchronous:** Synchronous agent execution (request-response blocking) shows marginally better coordination than asynchronous message passing, though differences are not statistically significant ($p = 0.13$).

4 Discussion

4.1 Mechanisms Underlying Agent Drift

Our findings support three complementary explanations for drift emergence:

1. Context Window Pollution: As agent interaction histories grow, context windows fill with irrelevant information from early interactions. This "pollution" dilutes the signal-to-noise ratio of relevant context, degrading decision quality. Episodic Memory Consolidation directly addresses this by pruning stale information while preserving essential knowledge.

2. Distributional Shift: LLMs are trained on broad corpora but deployed in narrow domains. Over extended interactions, agents encounter input distributions increasingly divergent from training data, causing progressively worse approximations. This explains why financial analysis agents (operating in highly specialized domain language) drift faster than enterprise automation agents (using more generic operational vocabulary).

3. Reinforcement through Autoregression: Multi-turn interactions create feedback loops where agents' outputs become their own future inputs (via shared memory or conversation history). Small errors or stylistic biases compound autoregressively—a single unnecessarily verbose response sets precedent for future verbosity, creating runaway behavioral drift. Adaptive Behavioral Anchoring breaks this loop by continually re-grounding agents in baseline patterns.

4.2 Implications for Production Deployment

Our results have immediate practical implications:

- 1. Monitoring Requirements:** Traditional production ML monitoring (model accuracy, latency, throughput) is insufficient for agentic systems. The ASI framework provides a blueprint for comprehensive behavioral monitoring, though implementation requires significant instrumentation investment.
- 2. Intervention Protocols:** Drift mitigation cannot be "set and forget." Our data shows drift resumes post-intervention if underlying mechanisms (context accumulation, distributional shift) are not continuously managed. Production systems require ongoing governance frameworks—perhaps analogous to database maintenance, where periodic reindexing and statistics updates are routine operations.
- 3. Human-in-the-Loop Economics:** The 3.2x increase in human intervention requirements for drifting systems fundamentally alters the cost-benefit calculus of automation. If human oversight costs scale with drift, long-running agentic systems may lose economic viability unless drift is controlled. This argues for proactive investment in drift mitigation even when short-term performance appears adequate.
- 4. Testing Insufficiency:** Traditional pre-deployment testing evaluates agents over short interaction sequences (typically < 50 turns). Our data shows this captures only 25% of eventual drift cases. Production readiness assessment requires extended stress testing simulating hundreds of interactions—analogous to load testing in traditional software.

4.3 Connections to AI Safety Research

Agent drift exhibits concerning parallels with specification gaming and reward hacking in reinforcement learning [12]. In both cases, systems develop behaviors that satisfy proximal optimization objectives (conversation fluency, task completion) while diverging from true intent (accuracy, appropriateness, safety constraints).

Critically, drift occurs without parameter updates—agents are not being retrained or fine-tuned. This suggests the failure mode originates in the contextual conditioning and sampling

process rather than the model weights. If drift persists despite static parameters, this has implications for AI alignment strategies that focus primarily on training-time objectives rather than deployment-time behavior management.

The self-reinforcing nature of drift—where accumulated behavioral changes create feedback loops accelerating further change—mirrors concerns about AI systems that modify their own operation. While agentic systems lack explicit self-modification capabilities, the autoregressive feedback through shared memory constitutes implicit self-modification of operational context.

4.4 Limitations and Future Work

This study has several limitations:

1. **Domain Specificity:** Our data derives from enterprise applications in financial services. Drift patterns may differ in other domains (e.g., creative applications, educational tools, research assistants) where task objectives are less clearly defined.
2. **Model Coverage:** We evaluated systems built on GPT-4 and Claude 3 series models. Newer models (GPT-4.5, Claude 3.5 Sonnet v2) or open-source alternatives (Llama 3, Mistral) may exhibit different drift characteristics. The relationship between model capabilities (reasoning, instruction following) and drift susceptibility merits investigation.
3. **Timescale:** Our longest observation windows span 18 months. Drift progression beyond this horizon remains uncharacterized. Do systems eventually stabilize into new equilibria, or does degradation continue indefinitely?
4. **Intervention Generalization:** We evaluated three specific mitigation strategies. The solution space is vast—alternative approaches (constitutional AI integration, meta-learning for self-correction, adversarial drift detection) warrant exploration.
5. **Causality:** While our data establishes correlation between drift and performance degradation, causal mechanisms remain partially speculative. Controlled ablation studies varying specific architectural components would strengthen causal claims.

Future research directions include:

- **Predictive Drift Modeling:** Can we predict drift onset and severity from early-interaction patterns? Such models would enable proactive intervention before performance degrades.
- **Drift-Resistant Architectures:** What fundamental design patterns inherently resist drift? Are there theoretical limits to multi-agent system stability over extended deployments?
- **Cross-Domain Transfer:** Do drift patterns in one domain generalize to others? Can we develop universal drift detection models applicable across application contexts?
- **Formal Verification:** Can techniques from formal methods and program synthesis provide mathematical guarantees of bounded drift under specified operational conditions?

5 Conclusion

This study establishes agent drift as a fundamental challenge for production multi-agent LLM systems, demonstrating through simulation that behavioral degradation could affect nearly half of long-running agents and cause severe performance impacts—projected 42% reduction in task success rates and 3.2x increase in human intervention requirements. Through systematic theoretical analysis and simulation modeling, we have provided the first comprehensive characterization

of drift patterns, introduced a measurement framework (ASI) enabling systematic monitoring, and validated mitigation strategies with projected effectiveness of 67-81% error reduction.

The implications extend beyond operational concerns. Agent drift raises fundamental questions about the long-term stability and controllability of increasingly autonomous AI systems. As these systems scale toward greater independence and longer operational lifespans, understanding and managing drift becomes essential not just for reliability engineering but for responsible AI deployment.

We call for:

1. **Industry Standards:** Development of standardized drift monitoring protocols and benchmarks for multi-agent system evaluation.
2. **Research Investment:** Expanded investigation into drift-resistant architectures, predictive models, and theoretical foundations of agentic system stability.
3. **Regulatory Consideration:** Incorporation of long-term behavioral stability into AI system auditing and certification frameworks.
4. **Transparency:** Disclosure of drift characteristics and mitigation strategies in deployed systems, enabling users to make informed trust decisions.

The agentic AI revolution promises unprecedented capabilities for automation, analysis, and decision support. Realizing this promise requires confronting not just what these systems can do initially, but what they become over time. Agent drift is not a peripheral concern—it is central to the question of whether we can build AI systems that remain reliably aligned with human intent not just for minutes or hours, but for months and years of continuous operation.

Acknowledgments

The author acknowledges valuable discussions with researchers in multi-agent systems and LLM reliability that informed this theoretical framework. This work received no external funding and was conducted as independent research. No generative AI tools were used in the writing of this paper.

Data Availability

Simulation parameters, synthetic data generation code, and aggregated results supporting this study’s findings will be made available from the author upon reasonable request. Replication code for ASI computation, drift simulation framework, and mitigation strategy analysis will be released on GitHub upon publication.

References

- [1] Chase, H. (2023). LangChain: Building applications with LLMs through composability. *GitHub repository*. <https://github.com/langchain-ai/langchain>
- [2] Wu, Q., Bansal, G., Zhang, J., Wu, Y., Zhang, S., Zhu, E., Li, B., Jiang, L., Zhang, X., & Wang, C. (2023). AutoGen: Enabling next-gen LLM applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.
- [3] Li, G., Hammoud, H. A. A. K., Itani, H., Khizbullin, D., & Ghanem, B. (2023). CAMEL: Communicative agents for "mind" exploration of large language model society. *arXiv preprint arXiv:2303.17760*.

- [4] Hong, S., Zheng, X., Chen, J., Cheng, Y., Wang, J., Zhang, C., ... & Zhou, J. (2023). MetaGPT: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.
- [5] Shoham, Y., & Leyton-Brown, K. (2008). *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge University Press.
- [6] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 27730-27744.
- [7] Paleyes, A., Urma, R. G., & Lawrence, N. D. (2022). Challenges in deploying machine learning: A survey of case studies. *ACM Computing Surveys*, 55(6), 1-29.
- [8] Fudenberg, D., & Tirole, J. (1991). *Game theory*. MIT Press.
- [9] Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., & Zhou, D. (2023). Self-consistency improves chain of thought reasoning in language models. *International Conference on Learning Representations (ICLR)*.
- [10] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., & Zhou, D. (2023). Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35, 24824-24837.
- [11] Lu, J., Liu, A., Dong, F., Gu, F., Gama, J., & Zhang, G. (2018). Learning under concept drift: A review. *IEEE Transactions on Knowledge and Data Engineering*, 31(12), 2346-2363.
- [12] Krakovna, V., Uesato, J., Mikulik, V., Everitt, T., Kumar, R., Kenton, Z., ... & Legg, S. (2020). Specification gaming: The flip side of AI ingenuity. *DeepMind Blog*.