

RetailBench: Evaluating Long-Horizon Autonomous Decision-Making and Strategy Stability of LLM Agents in Realistic Retail Environments

Anonymous ACL submission

Abstract

Large Language Model (LLM)-based agents have achieved notable success on short-horizon and highly structured tasks, yet their ability to maintain coherent decision-making over long horizons in dynamic environments remains an open challenge. We introduce *RetailBench*, a high-fidelity benchmark designed to evaluate long-horizon autonomous decision-making in realistic commercial scenarios, where agents must operate under stochastic demand and evolving external conditions.

We further propose the *Evolving Strategy & Execution* framework, which separates high-level strategic reasoning from low-level action execution, enabling adaptive and interpretable strategy evolution over time. This design is crucial for long-horizon tasks, where non-stationary environments and error accumulation require strategies to be revised at a different temporal scale than action execution. Experiments on seven state-of-the-art LLMs across progressively challenging environments show that our framework improves operational stability and efficiency compared to a Reflection-based baseline. However, performance degrades substantially as task complexity increases, revealing fundamental limitations in current LLMs for long-horizon, multi-factor decision-making.

1 Introduction

Recent large language models (LLMs), particularly when augmented with reasoning and tool-use capabilities, have demonstrated strong performance on a variety of cognitively demanding tasks, including code editing, mathematical problem solving, and complex information retrieval (Jimenez et al., 2024; Phan et al., 2025; Gao et al., 2024). However, growing evidence indicates that these capabilities do not readily translate into robust, general-purpose autonomy, especially in settings that require long-term planning, persistent goal maintenance, and adaptation to dynamic environments (Amodei, 2024;

Kwa et al., 2025; METR, 2025). Correspondingly, existing agent benchmarks—spanning web interaction (Mialon et al., 2023; Wei et al., 2025; Zhou et al., 2024; Deng et al., 2023; Jimenez et al., 2024; Team, 2025b) primarily focus on short-horizon or highly structured tasks, limiting their ability to evaluate sustained interaction with complex, evolving environments. Recent studies on long-horizon autonomy consistently demonstrate that even state-of-the-art agents struggle to maintain coherent strategies over extended time spans (Nof1.ai, 2025; Andon Labs, 2025; Backlund and Petersson, 2025).

To systematically study this challenge, we introduce *RetailBench*, a new benchmark grounded in real-world commercial data and informed by established economic modeling principles. *RetailBench* centers on a supermarket operation scenario that demands long-horizon decision-making, sustained interaction with a dynamic environment, and the integration of heterogeneous historical information. The benchmark evaluates whether LLM-based agents can autonomously sustain realistic business operations under complex, multi-factor conditions.

We evaluate seven state-of-the-art LLMs in this environment. Our results show that current models struggle to maintain stable decision quality as the decision space expands and often fail to incorporate all relevant information. Moreover, hallucinations and economically irrational behaviors frequently emerge during long-horizon execution, leading to environment collapse and preventing sustained autonomous operation.

Our contributions are summarized as follows:

- We introduce *RetailBench*, a high-fidelity benchmark for evaluating long-horizon autonomous decision-making in realistic retail environments.
- We propose the *Evolving Strategy & Execution* agent framework, which improves operational

RetailBench Environment

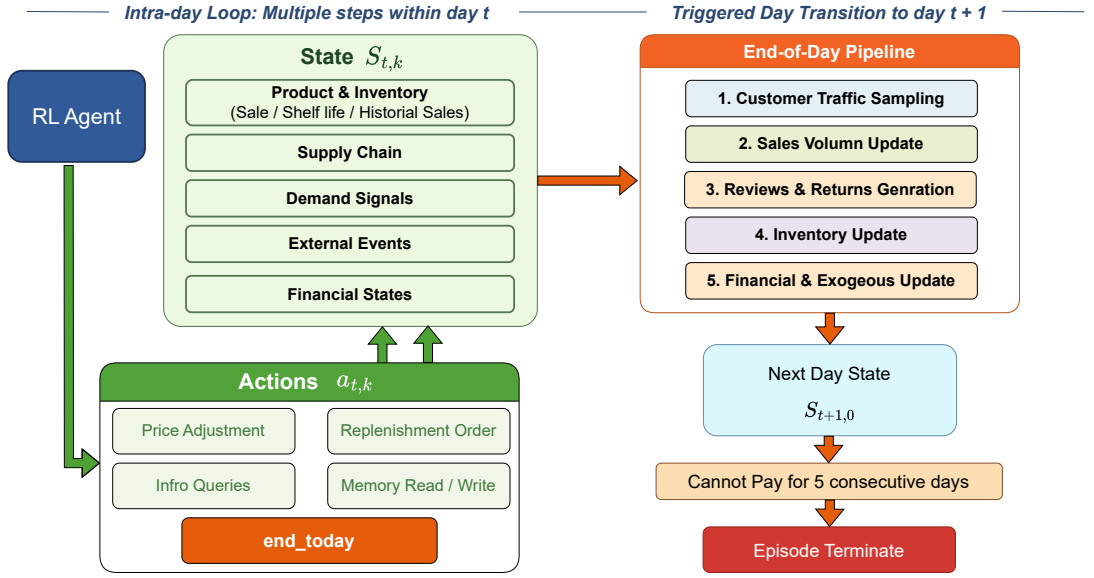


Figure 1: Overview of the hierarchical supermarket environment, illustrating intra-day agent–environment interactions and end-of-day state transition dynamics.

stability compared to a Reflection-based baseline.

- Through extensive experiments, we identify systematic failure modes of current LLM-based agents in long-horizon, multi-factor decision-making settings.

2 Environment Construction

2.1 Problem Formulation and Overview

We model supermarket operations as a Markov Decision Process (MDP), in which an autonomous agent manages a single retail store over a finite horizon of days. At each day t , the agent makes a sequence of operational decisions that jointly determine the store’s daily outcomes.

Formally, the MDP is defined as $(\mathcal{S}, \mathcal{A}, \mathcal{T}, R, \gamma)$, where \mathcal{S} denotes the state space, \mathcal{A} the action space, \mathcal{T} the stochastic transition dynamics, R the reward function, and $\gamma \in (0, 1]$ the discount factor.

At the beginning of each day t , the agent observes the initial state $S_{t,0}$ and executes a sequence of intra-day actions indexed by k . Each action $a_{t,k} \in \mathcal{A}$ induces a transition to the next intra-day state $S_{t,k+1}$. When the agent issues the *end-today* action, the environment transitions to the initial state of the next day, $S_{t+1,0}$, according to the transition dynamics \mathcal{T} . Figure 1 illustrates the overall interaction process.

The environment supports long-horizon operation over more than one thousand simulated days, with episodes terminating if the store fails to pay rent for five consecutive days.

2.2 State Space

The intra-day state $S_{t,k}$ summarizes the complete operational context of the store at step k of day t and is composed of multiple interdependent components: $S_{t,k} = (S_{t,k}^{\text{prod}}, S_{t,k}^{\text{inv}}, S_{t,k}^{\text{sup}}, S_{t,k}^{\text{dem}}, S_{t,k}^{\text{ext}}, S_{t,k}^{\text{fin}})$.

- $S_{t,k}^{\text{prod}}$ and $S_{t,k}^{\text{inv}}$ encode product-level attributes and on-hand inventory status, including prices, shelf life, and historical sales records. Product demand is grounded in real-world retail data derived from the Dominick’s dataset (Kilts Center for Marketing).
- $S_{t,k}^{\text{sup}}$ represents the supply chain state, including supplier prices, quality levels, and delivery lead times, constructed to reflect empirically observed price–quality relationships (Grewal et al., 2014).
- $S_{t,k}^{\text{dem}}$ captures demand-side signals such as recent customer traffic and aggregated review statistics, which influence consumer purchasing behavior (Fedewa et al., 2021).
- $S_{t,k}^{\text{ext}}$ represents external contextual information, including active news events with market-

Model	Avg. Days \uparrow	Avg. Daily Sales \uparrow	Avg. Daily Income \uparrow	Expiry Ratio \downarrow	Return Ratio \downarrow	Max Days \uparrow
<i>Framework: Evolving Strategy & Execution</i>						
DeepSeek-V3.2 (Exp.)	58.33	229.19	183.26	0.0889	0.1122	66
Gemini-3 (Fast)	50.67	399.39	294.71	0.0799	0.1311	59
GLM-4.6	52.40	174.34	124.67	0.0773	0.1293	58
Grok-4.1 Fast	61.75	508.08	336.94	0.0417	0.0847	88
Kimi-K2 (Thinking)	54.25	260.68	168.72	0.0239	0.1179	58
OpenAI-5.1 Mini	51.75	192.90	122.46	0.0360	0.1237	55
Qwen-235B (Thinking)	37.50	420.31	236.50	0.0745	0.0852	48
Average (7 models)	52.38	301.98	203.30	0.0590	0.1068	61.71
<i>Framework: Reflection</i>						
DeepSeek-V3.2(Exp.)	53.00	235.01	170.04	0.0382	0.1200	66
Gemini-3 (Fast)	45.67	447.74	255.38	0.0682	0.1350	50
GLM-4.6	55.00	160.70	125.67	0.0194	0.1176	62
Grok-4.1 Fast	48.33	297.94	197.54	0.1460	0.0925	54
Kimi-K2 (Thinking)	58.33	216.51	184.01	0.0964	0.1255	71
OpenAI-5.1 Mini	53.33	93.04	92.11	0.1062	0.1331	59
Qwen-235B (Thinking)	37.00	420.77	202.79	0.2645	0.1551	43
Average (7 models)	50.10	248.39	170.88	0.1092	0.1249	57.86
<i>Heuristic Policy (Upper Bound, Easy)</i>						
Hand-crafted Policy	180.00	674.18	729.46	0.0266	0.007	180

Table 1: Performance comparison of seven large language models under different agent frameworks in the EASY environment. A hand-crafted heuristic policy is included as an approximate upper bound.

By alternating between these two stages, the proposed framework enforces a principled separation between strategic deliberation and operational execution. This design mitigates uncontrolled strategy drift, promotes behavioral stability over long horizons, and facilitates more interpretable analysis of agent decision-making dynamics in complex environments. An illustration is provided in Appendix D.1.

3.2 Hierarchical Policy Representation

To support structured, interpretable, and temporally extended decision-making under the proposed framework, we represent the agent policy using a hierarchical abstraction that separates strategic intent from executable actions. Each policy consists of three conceptual layers:

- Macro Strategy**, which captures high-level managerial objectives that persist across multiple decision steps;
- Execution Strategy**, which encodes structured operational guidance in a machine-readable intermediate representation;
- Daily Actions**, which specify concrete executable operations issued to the environment.

Detailed policy configurations and example policies are provided in Appendix A.2.1.

4 Experiment Settings

We conduct experiments under three environment configurations with increasing levels of difficulty. These configurations vary in market complexity, budget constraints, and the presence of exogenous dynamics.

4.1 Environment Configurations

We employ a heuristic policy with full access to the environment’s internal state as a calibration baseline for each environment variant. Environment parameters are tuned such that the heuristic policy remains stable across different difficulty levels while still experiencing meaningful operational pressure in Appendix A.4. To evaluate models’ information-processing and external perception capabilities, we design three environment configurations:

- Easy**: A controlled environment without dynamic news events or adaptive supplier price–quality relationships. The market contains five product categories. The agent is initialized with a budget of 10,000 and incurs a fixed daily rent of 250.
- Middle**: A moderately complex environment that expands the product space to all twenty categories while still excluding dynamic news events and supplier adaptations. The initial budget is increased to 50,000, with a daily rent of 1,000.

279	• <i>Hard</i> : The most challenging and realistic environ-	each simulated day consisting of up to 50 interac-	323
280	ment, incorporating dynamically generated news	tion rounds; the Reflection baseline additionally	324
281	events and time-varying supplier price–quality	performs a reflection step at the end of each day.	325
282	relationships. The market includes all twenty	Full prompt specifications for both frameworks	326
283	product categories. The agent starts with a bud-	are provided in Appendix C.	327
284	get of 50,000, pays a daily rent of 1,000, and		
285	receives twenty news items per day.		
286	Detailed specifications of all environment configu-	Models. We conduct experiments using a di-	328
287	rations are provided in Appendix A.3.	verse set of contemporary large language mod-	329
288		els, including Qwen-235B (Thinking) (Team,	330
289	4.2 Evaluation Metrics	2025a), Kimi K2 (Thinking) (Team et al., 2025b),	331
290	Metrics. We evaluate store-level operational per-	GLM-4.6 (Team et al., 2025a), DeepSeek-V3.2-	332
291	formance using the following metrics (↑ indicates	Exp (DeepSeek-AI et al., 2025), Gemini-3-Flash-	333
292	higher is better; ↓ indicates lower is better):	Preview (Google DeepMind, 2024), Grok-4.1 Fast	334
293		(xAI, 2025), and GPT-5-Mini (OpenAI, 2025). Due	335
294	• <i>Days</i> (↑): the number of operating days before	to the substantial token costs incurred during long-	336
295	episode termination;	horizon rollouts, we use lower-cost variants for	337
296		closed-source models to ensure experimental feasi-	338
297	• <i>MaxDays</i> (↑): the maximum number of operating	bility.	339
298	days achieved across three rollouts;	Evaluation Protocol. We evaluate seven large	340
299		language models under each of the three environ-	341
300	• <i>Avg. Daily Sales</i> (↑): the average number of items	ment configurations. Each model is assessed using	342
301	sold per day;	three independent rollouts, and all reported metrics	343
302		are averaged across runs. To isolate the effect of the	344
303	• <i>Avg. Daily Income</i> (↑): the average money earned	agent framework, we further conduct a controlled	345
304	per day;	comparison between the proposed framework and	346
305		the Reflection baseline in the <i>Easy</i> environment,	347
306	• <i>Expiry Ratio</i> (↓): the fraction of products that	where both frameworks are evaluated under identi-	348
307	expire before being sold;	cal conditions. Finally, to quantify the gap between	349
308		current LLM-based agents and an environment-	350
309	• <i>Return Ratio</i> (↓): the fraction of sold products	optimal strategy, we implement a hand-crafted pol-	351
310	that are returned by customers.	icy based on internal knowledge unavailable to	352
311		the agents, which serves as an approximate upper	353
312	All reported metrics are averaged over three inde-	bound on achievable performance.	354
313	pendent rollouts, each subject to a fixed maximum		
314	execution horizon.	5 Results	355
315			
316	4.3 Experimental Setup	5.1 Performance Comparison between Two	356
317	Agent Frameworks. Preliminary experiments indi-	Agent Frameworks	357
318	cate that simple ReAct-style interaction frame-	Table 1 compares seven large language models	358
319	works are unstable in long-horizon settings, fre-	in the Easy environment under two agent frame-	359
320	quently exhibiting premature episode termination	works: <i>Evolving Strategy & Execution</i> and <i>Reflec-</i>	360
321	during mid-horizon execution. This motivates the	<i>tion</i> . Under the proposed framework, Grok-4.1 Fast	361
322	evaluation of agent frameworks that explicitly sup-	achieves the strongest overall performance, attain-	362
	port sustained strategic control.	ing the highest average daily sales, average daily	363
	We compare our proposed <i>Evolving Strategy &</i>	profit, and the longest maximum survival duration,	364
	<i>Execution</i> framework with a Reflection-based base-	indicating superior long-horizon planning and op-	365
	line (Shinn et al., 2023). Unlike our approach, the	erational stability. In contrast, Kimi-K2 (Thinking)	366
	Reflection framework does not explicitly represent	yields the lowest expiry ratio, reflecting more con-	367
	or preserve high-level strategies, instead relying	servative and risk-aware inventory management.	368
	on iterative post-hoc reflection to update a global	At the framework level, <i>Evolving Strategy & Ex-</i>	369
	long-term memory. Both frameworks operate un-	<i>ecution</i> consistently outperforms <i>Reflection</i> across	370
	der a maximum context length of 40k tokens, with	all key metrics, achieving higher sales and profit	371

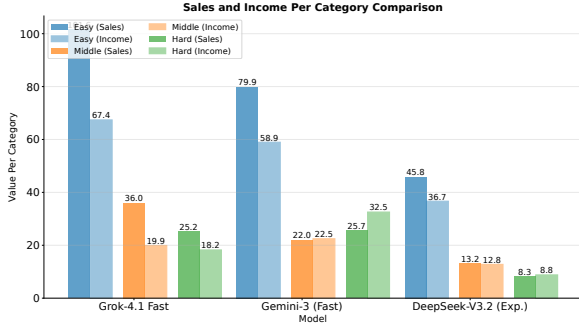


Figure 2: Category-level sales and profit per category across Easy, Middle, and Hard environments. Results are shown for three representative models.

while substantially reducing product expiry. These results demonstrate that explicitly decoupling strategy evolution from execution improves both revenue efficiency and inventory control. Nevertheless, a notable performance gap remains relative to the hand-crafted heuristic policy, suggesting that current LLM-based agents still fall short of optimal long-horizon decision-making.

5.2 Performance across Environments with Varying Difficulty

As environment difficulty increases, all models exhibit consistent performance degradation, including shorter operational durations, higher expiry and return ratios, and reduced long-horizon stability. Although the expansion of product categories in the Middle and Hard settings enables higher aggregate sales and profit for some models, *Sales per Category* and *Profit per Category* decline substantially (Figure 2), indicating persistent challenges in effective resource allocation within increasingly high-dimensional decision spaces.

The relatively small performance gap between the Middle and Hard environments can be partly attributed to the delayed impact of intensified news dynamics, whose effects typically unfold over longer time horizons. Appendix B.3 reports the full results across all environments under the proposed framework. Overall, while models demonstrate limited adaptability as task complexity increases, their performance remains substantially below the heuristic upper bound, highlighting persistent limitations in robust long-horizon decision-making under complex and dynamic conditions.

6 Analysis

Based on both quantitative evaluation results and manual inspection of trajectories, we identify sev-

Model	Context	Easy		Middle		Hard	
		SKU ↑	Cat ↑	SKU ↑	Cat ↑	SKU ↑	Cat ↑
DeepSeek-V3.2 (Exp.)	128K	4.75	3.31	<u>6.83</u>	<u>5.34</u>	6.64	<u>5.45</u>
Gemini-3 (Fast)	1M	8.72	4.34	13.60	12.43	21.57	13.56
GLM-4.6	200K	4.82	2.98	3.23	2.46	<u>7.08</u>	5.22
OpenAI-5.1 Mini	400K	3.66	1.94	4.10	4.10	6.67	4.79
Grok-4.1 Fast	200K	<u>5.78</u>	<u>3.68</u>	4.87	3.08	5.40	3.20
Kimi-K2 (Thinking)	256K	5.06	3.09	5.17	4.77	6.91	4.75
Qwen-235B (Thinking)	256K	4.59	3.67	3.17	2.36	5.11	4.51
Heuristic Strategy	–	9.03	4.87	35.34	18.61	34.82	18.52

Table 2: SKU and Category counts sold by different models across difficulties. Context denotes the maximum supported context length of each model. **Bolded** indicates the best model; underlined indicates the second best (Human excluded).

eral key factors that explain why current models fail to operate reliably in the environment.

We analyze these failure modes from four complementary perspectives.

6.1 Non-scalable Decision-Making Capability

As shown in Appendix B.3, models achieve reasonable performance in the Easy setting but exhibit consistent degradation as the environment scales in complexity. To better understand this phenomenon, we analyze both the average number of stock keeping units (SKUs) and product categories sold per day, as well as the number of SKUs and categories explicitly considered during the *Evolving Strategy* phase (Tables 2 and Appendix B.6).

Across most models, decision-making capability does not scale proportionally with the size of the environment. Instead, performance remains relatively flat despite a substantial increase in the number of available options. Notably, Gemini-3 (Fast), which supports the largest maximum context window, performs better than other models in this respect, suggesting that larger context capacity helps retain salient information even when the effective interaction context is constrained.

Nevertheless, even the strongest models fail to cover the full decision space. This indicates that current systems are unable to expand their effective decision-making scope as the environment grows, leading to systematic performance degradation in larger and more complex settings.

6.2 Incomplete Decision-Making Due to Limited Information Coverage

We analyze the information sources that models attend to when making operational decisions by examining the data queried for the set of SKUs included in each day’s final strategy. This analysis reveals a clear concentration of attention on a limited subset of signals.

Model	Supplier	Inventory	Return	Rating	Price	Review	Sales	History
DeepSeek-V3.2 (Exp.)	41.4	100.0	26.7	75.4	6.6	0.6	89.2	0.0
Gemini-3 (Fast)	69.6	100.0	41.0	82.3	67.7	6.3	87.6	2.9
GLM-4.6	93.9	98.3	79.8	77.8	84.1	5.1	96.3	0.0
OpenAI-5.1 Mini	58.8	99.0	5.1	78.6	3.9	10.6	84.7	0.3
Grok-4.1 Fast	89.9	99.8	84.0	94.0	88.6	36.0	96.4	0.3
Kimi-K2 (Thinking)	57.2	90.4	25.9	51.4	36.2	11.0	64.0	0.5
Qwen-235B (Thinking)	76.8	78.3	31.5	58.2	19.8	23.6	74.6	0.0
Average	69.7	95.1	42.0	74.0	43.8	13.3	84.7	1.0

Table 3: Percentage of days on which each model queries specific information sources when making decisions for SKUs included in the final daily strategy. Higher values indicate greater reliance on the corresponding data source.

Model	Macro Strategy			Execution Strategy		
	Std_diff (↓)	MAC (↓)	TV (↓)	Std_diff (↓)	MAC (↓)	TV (↓)
DeepSeek-V3.2(Exp.)	0.130	0.090	5.00	0.289	0.144	7.93
Gemini-3 (Fast)	0.136	0.085	3.82	0.293	0.225	10.18
GLM-4.6	0.131	0.096	4.52	0.264	0.209	9.87
OpenAI-5.1 Mini	0.069	0.051	2.50	0.229	0.166	8.00
Grok-4.1 Fast	0.079	0.045	2.71	0.204	0.151	9.39
Kimi K2(Thinking)	0.179	0.133	6.95	0.335	0.271	14.08
Qwen-235B (Thinking)	0.240	0.189	6.51	0.391	0.306	10.57

Table 4: Temporal instability metrics of macro- and execution-level strategy similarity in the Easy environment. All metrics are lower-is-better (↓). Larger values indicate greater temporal instability. Column-wise maximum values are highlighted in bold.

Specifically, most models primarily rely on supplier prices, inventory levels, SKU ratings, and historical sales records when deciding how to operate selected SKUs in Table 3. In contrast, several other critical signals—such as recent customer reviews, return rates, and current selling prices—are consistently underutilized or entirely ignored.

Further correlation analysis shows a strong positive relationship between the frequency of SKU reviews queries and average daily sales performance in Appendix B.5. This observation aligns with the underlying environment dynamics and indicates that incomplete information coverage is a key factor limiting decision quality. Overall, these results suggest that models often fail to perform sufficiently comprehensive information gathering, leading to systematically suboptimal operational decisions.

6.3 Temporal Instability in Execution-Level Decision-Making

Beyond limitations in decision scalability and information coverage, we identify temporal instability in execution-level decision-making as a key contributor to long-horizon failure. Even under relatively stable environmental conditions, agents frequently revise their strategies across consecutive days, leading to inconsistent execution trajectories.

Measuring Strategy Similarity. We quantify temporal instability by measuring the similarity be-

tween strategies on adjacent days at both macro and execution levels. Macro strategy similarity is assessed using an LLM-based prompt that evaluates semantic consistency between consecutive high-level plans. Execution strategy similarity is computed via set-based Jaccard similarity over key fields, including `focus_skus`, `sku_supplier_mapping`, `news_to_monitor`, and `sku_to_monitor`, with the final execution similarity obtained by averaging across fields.

Instability Metrics. To characterize temporal fluctuations, we compute three complementary metrics: the standard deviation of first-order differences (*Std_diff*) to capture short-term volatility, the mean absolute change (*MAC*) to estimate typical day-to-day variation, and total variation (*TV*) to reflect cumulative long-term instability.

Across all three environment configurations, we observe consistent temporal instability in both macro- and execution-level strategies. For clarity, we present detailed results from the Easy environment in Table 4 and Appendix B.4. Despite its reduced complexity, the Easy setting already exhibits pronounced temporal fluctuations. In particular, macro strategies remain relatively stable over time, whereas execution strategies show substantially larger variability across most models. This effect is especially pronounced for Qwen-235B (Thinking), which also performs poorly across all environments.

Overall, these results indicate that long-horizon failures arise not only from suboptimal strategy formulation, but more fundamentally from the inability to maintain temporally consistent execution policies over extended horizons.

6.4 Hallucinations and Invalid Actions

Finally, manual inspection reveals recurrent failure patterns that directly break planning correctness and action validity. We distinguish two closely related but practically different issues:

Hallucinations in reasoning Models occasionally generate reasoning traces that reference non-existent SKUs or fabricate numerical quantities, and subsequently incorporate these hallucinated elements into multi-step plans (examples are provided in Appendix B.1). As a result, decisions become misaligned with the true environment state, even when the overall planning structure appears internally coherent.

524	Invalid or irrational actions. Models sometimes	574
525	output actions that violate basic constraints or are	575
526	inconsistent with historical demand, such as nega-	
527	tive order quantities or implausible pricing in Ap-	
528	pendix B.2. Even though state-modifying oper-	
529	ations are restricted to only a small set of tools,	
530	these invalid actions occur with non-negligible fre-	
531	quency and can quickly destabilize the system in	
532	long-horizon operation.	
533	In realistic operational settings, both hallucina-	
534	tions and invalid actions would be unacceptable	
535	and could directly trigger system collapse.	
536	7 Related Work	
537	Retail and supply-chain decision benchmarks.	
538	Recent work on retail and supply-chain decision-	
539	making has evolved from isolated subprob-	
540	lems toward integrated, multi-stage environments.	
541	Reinforcement-learning benchmarks such as OF-	
542	COURSE, GymSC-style simulators, and MAR-	
543	LIM (Zhu et al., 2023; Shar et al., 2022; Leluc	
544	et al., 2023) model fulfillment and inventory con-	
545	trol under demand uncertainty, demonstrating the	
546	effectiveness of learned policies in mitigating ef-	
547	fects such as demand volatility. However, these	
548	environments typically emphasize fixed structures	
549	and short- to medium-horizon policies, limiting	
550	their ability to evaluate sustained strategic behavior.	
551	More recent benchmarks incorporating LLM-based	
552	agents, including InvAgent and AIM-Bench (Quan	
553	and Liu, 2025; Zhao et al., 2025), begin to explore	
554	language-driven reasoning and decision biases, yet	
555	still fall short in assessing long-horizon, strategy-	
556	aware autonomy in realistic retail settings.	
557	Long-horizon planning benchmarks for LLMs.	
558	In parallel, a growing body of benchmarks tar-	
559	gets long-horizon planning abilities of large lan-	
560	guage models across structured and interactive en-	
561	vironments. PlanBench focuses on classical plan	
562	generation, while WebShop, Mind2Web, and Sci-	
563	enceWorld (Deng et al., 2023; Wang et al., 2022;	
564	Yao et al., 2023) evaluate multi-step interaction,	
565	error recovery, and adaptive behavior in dynamic	
566	settings. More recent benchmarks—such as Her-	
567	oBench, OdysseyBench, UltraHorizon (Luo et al.,	
568	2025; Anokhin et al., 2025; Wang et al., 2025), and	
569	explicitly stress extended, interdependent decision	
570	sequences that require persistent memory, hierar-	
571	chical reasoning, and long-term strategy mainte-	
572	nance. Collectively, these benchmarks suggest that	
573	while short-horizon tasks are increasingly tractable,	
	robust long-horizon planning and execution remain	574
	a central challenge for LLM-based agents.	575
	Long-horizon agent frameworks. Recent ad-	576
	vances in agent design address these challenges	577
	by introducing structured frameworks that de-	578
	couple high-level planning from low-level execu-	579
	tion. Approaches such as Plan-and-Act (Erdo-	580
	gan et al., 2025) and EAGLET (Si et al., 2025)	581
	adopt planner–executor architectures to support hi-	582
	erarchical decision-making, dynamic replanning,	583
	and improved execution stability over long hori-	584
	zons. Extensions to multi-agent settings ELHPlan	585
	(Ling et al., 2025) and plan-aware context manage-	586
	ment frameworks PAACE (Yuksel, 2025) further	587
	emphasize task decomposition, explicit strategy	588
	representation, and memory-aware context control.	589
	Together, these works indicate that scalable long-	590
	horizon autonomy increasingly relies on structured	591
	planning modules and controlled execution mecha-	592
	nisms rather than monolithic end-to-end policies.	593
	8 Conclusion	594
	This paper introduces RetailBench, a high-	595
	fidelity benchmark for evaluating long-horizon au-	596
	tonomous decision-making in realistic retail envi-	597
	ronments. RetailBench models supermarket opera-	598
	tions as a stochastic, multi-factor, and temporally	599
	extended process, requiring agents to jointly reason	600
	about pricing, inventory, information acquisition,	601
	and financial sustainability over extended horizons.	602
	We further propose the Evolving Strategy & Execu-	603
	tion framework, which decouples high-level strat-	604
	egy evolution from low-level execution to better	605
	support long-horizon autonomy.	606
	Experiments across seven state-of-the-art large	607
	language models show that the proposed frame-	608
	work consistently improves operational stabil-	609
	ity and economic performance compared to a	610
	Reflection-based baseline. Nevertheless, perfor-	611
	mance degrades sharply as environment complexity	612
	increases, exposing persistent limitations in deci-	613
	sion scalability, information utilization, execution	614
	stability, and action validity. These findings indi-	615
	cate that while structured agent frameworks alle-	616
	viate some challenges, current LLM-based agents	617
	remain far from robust, strategy-aware autonomy in	618
	complex dynamic environments. RetailBench pro-	619
	vides a principled testbed for advancing research	620
	on long-horizon decision-making and agentic rea-	621
	soning.	622

9 Limitations

Despite its realism and scale, this work has several limitations. First, RetailBench focuses on a single-store supermarket setting; while expressive, it does not capture multi-store coordination, competitive markets, or strategic interactions among multiple autonomous agents. Second, although the environment incorporates stochastic demand, news dynamics, and supply-chain delays, it remains a simulation grounded in historical data and simplified economic assumptions, which may not fully reflect the complexities of real-world retail systems.

Third, our evaluation is limited to prompting-based LLM agents without parameter updates or long-term learning across episodes; stronger performance may be achievable through reinforcement learning, fine-tuning, or hybrid neuro-symbolic approaches. Finally, while we identify key failure modes such as hallucinations and economically irrational actions, we do not propose explicit algorithmic mechanisms to enforce economic constraints or factual grounding during execution. Addressing these limitations—through richer environments, multi-agent extensions, learning-based adaptation, and constraint-aware action control—remains an important direction for future research.

References

Dario Amodei. 2024. Machines of loving grace. <https://www.darioamodei.com/essay/machines-of-loving-grace>. Accessed: 2025-12-01.

Andon Labs. 2025. Vending-bench 2: A benchmark for long-horizon business simulation. <https://andonlabs.com/evals/vending-bench-2>. Accessed: 2025-12-10.

Petr Anokhin, Roman Khalikov, Stefan Rebrikov, Viktor Volkov, Artyom Sorokin, and Vincent Bissonnette. 2025. Herobench: A benchmark for long-horizon planning and structured reasoning in virtual worlds. *Preprint*, arXiv:2508.12782.

ashraq. 2025. financial-news-articles. <https://huggingface.co/datasets/ashraq/financial-news-articles>. Accessed: 2025-12-01.

Axel Backlund and Lukas Petersson. 2025. Vending-bench: A benchmark for long-term coherence of autonomous agents. *Preprint*, arXiv:2502.15840.

DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao

Wu, Bowei Zhang, Chaofan Lin, Chen Dong, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenhao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, and 245 others. 2025. Deepseek-v3.2: Pushing the frontier of open large language models. *Preprint*, arXiv:2512.02556.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *Preprint*, arXiv:2306.06070.

Lutfi Eren Erdogan, Nicholas Lee, Sehoon Kim, Suhong Moon, Hiroki Furuta, Gopala Anumanchipalli, Kurt Keutzer, and Amir Gholami. 2025. Plan-and-act: Improving planning of agents for long-horizon tasks. *Preprint*, arXiv:2503.09572.

Dave Fedewa, Chris Holder, Wynn Teichner, and Ben Wiseman. 2021. Five-star growth: Using online ratings to design better products. *McKinsey & Company*. Accessed: 2025-11-30.

Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. 2024. Omni-math: A universal olympiad level mathematic benchmark for large language models. *Preprint*, arXiv:2410.07985.

Google DeepMind. 2024. Gemini 3 flash (preview). <https://deepmind.google/technologies/gemini/>. Large multimodal language model, preview version.

Dhruv Grewal, Jens Nordfält, Anne Roggeveen, Rainer Olbrich, and Hans Christian Jansen. 2014. Price-quality relationship in pricing strategies for private labels. *Journal of Product and Brand Management*, 23(6):429–438.

Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. Swe-bench: Can language models resolve real-world github issues? *Preprint*, arXiv:2310.06770.

University of Chicago Booth School of Business Kilts Center for Marketing. Dominick’s dataset. <https://www.chicagobooth.edu/research/kilts/research-data/dominicks>. Accessed: 2025-11-01.

Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney Von Arx, Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles, Seraphina Nix, Tao Lin, Neev Parikh, and 6 others. 2025. Measuring ai ability to complete long tasks. *Preprint*, arXiv:2503.14499.

Rémi Leluc, Elie Kadoche, Antoine Bertoncello, and Sébastien Gourvénéec. 2023. Marlim: Multi-agent

729	reinforcement learning for inventory management .	GLM Team, Aohan Zeng, Xin Lv, Qinkai Zheng,	784
730	<i>Preprint</i> , arXiv:2308.01649.	Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang	785
731	Shaobin Ling, Yun Wang, Chenyou Fan, Tin Lun Lam,	Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong	786
732	and Junjie Hu. 2025. Elhplan: Efficient long-horizon	Wang, Lucen Zhong, Mingdao Liu, Rui Lu, Shulin	787
733	task planning for multi-agent collaboration . <i>Preprint</i> ,	Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei,	788
734	arXiv:2509.24230.	and 152 others. 2025a. Glm-4.5: Agentic, reason-	789
735	Haotian Luo, Huaisong Zhang, Xuelin Zhang, Haoyu	ing, and coding (arc) foundation models . <i>Preprint</i> ,	790
736	Wang, Zeyu Qin, Wenjie Lu, Guozheng Ma, Haiying	arXiv:2508.06471.	791
737	He, Yingsha Xie, Qiyang Zhou, Zixuan Hu, Hongze	Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jia-	792
738	Mi, Yibo Wang, Naiqiang Tan, Hong Chen, Yi R.	hao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen,	793
739	Fung, Chun Yuan, and Li Shen. 2025. Ultrahori-	Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei	794
740	zon: Benchmarking agent capabilities in ultra long-	Cui, Hao Ding, Mengnan Dong, Angang Du, Chen-	795
741	horizon scenarios . <i>Preprint</i> , arXiv:2509.21766.	zhuang Du, Dikang Du, Yulun Du, Yu Fan, and 150	796
742	Daniel McFadden. 1974. Conditional logit analysis of	others. 2025b. Kimi k2: Open agentic intelligence .	797
743	qualitative choice behavior. In Paul Zarembka, editor,	<i>Preprint</i> , arXiv:2507.20534.	798
744	<i>Frontiers in Econometrics</i> , pages 105–142. Academic	Qwen Team. 2025a. Qwen3 technical report . <i>Preprint</i> ,	799
745	press, New York.	arXiv:2505.09388.	800
746	METR. 2025. Measuring ai ability to complete long	The Terminal-Bench Team. 2025b. Terminal-bench: A	801
747	tasks . METR blog.	benchmark for ai agents in terminal environments .	802
748	Grégoire Mialon, Clémentine Fourier, Craig Swift,	Mark D. Uncles. 1987. Discrete choice analysis: Theory	803
749	Thomas Wolf, Yann LeCun, and Thomas Scialom.	and application to travel demand . <i>Journal of the</i>	804
750	2023. Gaia: a benchmark for general ai assistants .	<i>Operational Research Society</i> , 38(4):370–371.	805
751	<i>Preprint</i> , arXiv:2311.12983.	Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and	806
752	Nof1.ai. 2025. Alpha arena — exploring the limits	Prithviraj Ammanabrolu. 2022. Scienceworld: Is	807
753	of large language models as quant traders. https:	your agent smarter than a 5th grader? <i>Preprint</i> ,	808
754	//nof1.ai/blog/TechPost1 . Accessed: 2025-12-	arXiv:2203.07540.	809
755	10.	Weixuan Wang, Dongge Han, Daniel Madrigal Diaz,	810
756	OpenAI. 2025. Gpt-5 mini. https://platform.	Jin Xu, Victor Rühle, and Saravan Rajmohan.	811
757	openai.com/docs/models/gpt-5-mini . Large	2025. Odysseybench: Evaluating llm agents on	812
758	language model — cost-efficient GPT-5 variant.	long-horizon complex office application workflows .	813
759	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li,	<i>Preprint</i> , arXiv:2508.09124.	814
760	Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang,	Jason Wei, Zhiqing Sun, Spencer Papay, Scott McK-	815
761	Mohamed Shaaban, John Ling, Sean Shi, Michael	inney, Jeffrey Han, Isa Fulford, Hyung Won Chung,	816
762	Choi, Anish Agrawal, Arnav Chopra, Adam Khoja,	Alex Tachard Passos, William Fedus, and Amelia	817
763	Ryan Kim, Richard Ren, Jason Hausenloy, Oliver	Glaese. 2025. Browsecomp: A simple yet chal-	818
764	Zhang, Mantas Mazeika, and 1093 others. 2025. Hu-	lenging benchmark for browsing agents . <i>Preprint</i> ,	819
765	manity’s last exam . <i>Preprint</i> , arXiv:2501.14249.	arXiv:2504.12516.	820
766	Yinzhu Quan and Zefang Liu. 2025. Invagent: A	xAI. 2025. Grok 4.1 fast and agent tools api. https://	821
767	large language model based multi-agent system for	x.ai/news/grok-4-1-fast . Accessed [Your Ac-	822
768	inventory management in supply chains . <i>Preprint</i> ,	cess Date].	823
769	arXiv:2407.11384.	Shunyu Yao, Howard Chen, John Yang, and Karthik	824
770	Ibrahim Shar, Wenhuan Sun, Haiyan Wang, and Chetan	Narasimhan. 2023. Webshop: Towards scalable real-	825
771	Gupta. 2022. Deep reinforcement learning toward ro-	world web interaction with grounded language agents .	826
772	bust multi-echelon supply chain inventory optimiza-	<i>Preprint</i> , arXiv:2207.01206.	827
773	tion . pages 1385–1391.	Kamer Ali Yuksel. 2025. Paace: A plan-aware	828
774	Noah Shinn, Federico Cassano, Edward Berman, Ash-	automated agent context engineering framework .	829
775	win Gopinath, Karthik Narasimhan, and Shunyu Yao.	<i>Preprint</i> , arXiv:2512.16970.	830
776	2023. Reflexion: Language agents with verbal rein-	Xuhua Zhao, Yuxuan Xie, Caihua Chen, and Yuxiang	831
777	forcement learning . <i>Preprint</i> , arXiv:2303.11366.	Sun. 2025. Aim-bench: Evaluating decision-making	832
778	Shuzheng Si, Haozhe Zhao, Kangyang Luo, Gang	biases of agentic llm as inventory manager . <i>Preprint</i> ,	833
779	Chen, Fanchao Qi, Minjia Zhang, Baobao Chang,	arXiv:2508.11416.	834
780	and Maosong Sun. 2025. A goal without a plan	Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou,	835
781	is just a wish: Efficient and effective global plan-	Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue	836
782	ner training for long-horizon agent tasks . <i>Preprint</i> ,		
783	arXiv:2510.05608.		

Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2024. [Webarena: A realistic web environment for building autonomous agents](#). *Preprint*, arXiv:2307.13854.

Yiheng Zhu, Yang Zhan, Xuankun Huang, Yuwei Chen, yujie Chen, Jiangwen Wei, Wei Feng, Yinzhi Zhou, Haoyuan Hu, and Jieping Ye. 2023. [Ofcourse: A multi-agent reinforcement learning environment for order fulfillment](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 34765–34777. Curran Associates, Inc.

A Environment Configuration Details

A.1 State Space Decomposition

We construct the environment using real-world retail data from the Dominick’s dataset ([Kilts Center for Marketing](#)). From the 20 available product categories, we select 96 SKUs with the most complete and informative sales records. The richness of these observations enables reliable modeling of the relationship between pricing decisions and realized demand.

Key Symbols. Let \mathcal{J} denote the set of selected SKUs, with cardinality $|\mathcal{J}| = J$ (where $J = 96$ in our experiments). Each SKU $j \in \mathcal{J}$ belongs to a product category $\text{cat}(j) \in \{1, \dots, 20\}$. For each SKU j , let $\mathcal{K}(j)$ denote its associated supplier set, with $|\mathcal{K}(j)| = 5$.

A.1.1 Product State S_t^{prod}

For each SKU j , we maintain a set of static attributes together with time-varying operational signals:

$$S_{t,j}^{\text{prod}} = (\text{id}_j, \text{desc}_j, L_j, p_{jt}, \mathbf{h}_{j,t}^{\text{sales}}, \mathbf{r}_{j,t}). \quad (1)$$

Here, id_j and desc_j denote the unique identifier and textual description of SKU j , respectively. L_j denotes the shelf life (in days), and p_{jt} is the retail price on day t . The term $\mathbf{h}_{j,t}^{\text{sales}}$ encodes historical sales statistics, while $\mathbf{r}_{j,t}$ summarizes aggregated customer review signals, such as average rating and review volume.

Data grounding. The SKU set and cost priors are constructed from the Dominick’s dataset ([Kilts Center for Marketing](#)). Historical sales records from the same source are used to calibrate the demand model.

A.1.2 Inventory State S_t^{inv}

Inventory is represented at the SKU level with age tracking. Let I_{jt} denote the on-hand units of SKU

j at the beginning of day t . To support shelf-life constraints and depreciation, we track the arrival time and remaining shelf life of each unit.

Capacity constraint and pending queue. Let Cap denote the total inventory capacity of the store. If incoming replenishment orders would violate this constraint, excess units are placed into a first-in-first-out (FIFO) pending queue Q_t and become available only after sufficient inventory space is released:

$$\sum_{j \in \mathcal{J}} \sum_{a=0}^{L_j} I_{jt}^{(a)} \leq \text{Cap}. \quad (2)$$

Stockout signal. When realized demand exceeds the available sellable inventory of SKU j on day t , a stockout indicator is triggered. This signal is recorded at the end of day t and exposed to the agent as explicit operational feedback.

Expiration and destruction. At each day transition, inventory units whose residence time exceeds their shelf life are removed from the system.

A.1.3 Supply Chain State S_t^{sup}

Each SKU j is associated with a set of suppliers $k \in \mathcal{K}(j)$. The supplier state includes procurement price c_{jk} , quality level q_{jk} , and lead-time distribution:

$$S_{t,j}^{\text{sup}} = \{(c_{jk}, q_{jk}, \mathcal{L}_{jk}) : k \in \mathcal{K}(j)\}. \quad (3)$$

Lead time $\ell_{jk} \sim \mathcal{L}_{jk}$ is sampled when an order is placed. Procurement prices are discretized into tiers using Dominick’s cost signals ([Kilts Center for Marketing](#)), following empirically observed positive correlations between price tier and quality level ([Grewal et al., 2014](#)).

Enforced diversity. To avoid degenerate supplier configurations, we enforce that each SKU has (i) one supplier with maximal quality, (ii) one supplier with minimal price and minimal quality, and remaining suppliers spanning intermediate price–quality levels.

A.1.4 Demand Signals and Demand Generation

Demand-related components capture both observable market signals and the stochastic process governing realized demand. We define

$$S_t^{\text{dem}} = (N_t, \{\mathbf{r}_{j,t}\}_{j \in \mathcal{J}}), \quad (4)$$

where N_t denotes daily customer traffic.

Customer traffic. Daily customer traffic is derived from the Dominick’s dataset.

Consumer choice and demand generation. Given customer traffic N_t , prices $\{p_{jt}\}_{j \in \mathcal{J}}$, review signals $\{\mathbf{r}_{j,t}\}$, and external news \mathcal{E}_t , we generate realized demand using a discrete-choice model. Following standard practice in empirical economics, we adopt a Multinomial Logit (MNL) framework (McFadden, 1974; Uncles, 1987).

For SKU j , the raw utility is defined as

$$U_{jt}^{\text{raw}} = \alpha_j + \beta_j p_{jt} + \varepsilon_{jt}, \quad (5)$$

where α_j captures intrinsic preference, β_j models price sensitivity, and ε_{jt} represents idiosyncratic shocks.

We further incorporate review effects, news signals, and within-category substitution:

$$\begin{aligned} \tilde{U}_{jt} = & U_{jt}^{\text{raw}} + \Delta_j(\mathbf{r}_{j,t}, \mathcal{E}_t) \\ & + \sum_{\substack{i \neq j \\ \text{cat}(i) = \text{cat}(j)}} \gamma_{ji} \exp(\tilde{U}_{it}). \end{aligned} \quad (6)$$

With the outside option utility normalized to zero, the purchase probability for SKU j on day t is

$$p_{jt}^* = \frac{\exp(\tilde{U}_{jt})}{1 + \sum_{i=1}^J \exp(\tilde{U}_{it})}. \quad (7)$$

Realized demand. Given traffic N_t and purchase probabilities $\{p_{jt}^*\}$, potential demand is sampled as

$$y_{jt} \sim \text{Binomial}(N_t, p_{jt}^*), \quad (8)$$

and realized sales are capped by available on-hand inventory.

A.1.5 External Information S_t^{ext} (News Module)

We maintain a set of active news events \mathcal{E}_t . Each event $e \in \mathcal{E}_t$ is represented as

$$e = (\text{type}, \text{scope}, \text{target}, \text{side}, \text{sign}, \eta, \text{text}, \text{ttl}) \quad (9)$$

where $\text{scope} \in \{\text{macro}, \text{category}, \text{product}, \text{neutral}\}$, $\text{side} \in \{\text{demand}, \text{supply}, \text{both}\}$, $\text{sign} \in \{+1, -1\}$, η denotes impact magnitude, and ttl is the time-to-live in days. News texts are synthesized using an LLM to resemble financial news corpora (ashraq, 2025).

Impact application. News impacts are incorporated into (i) demand utilities and/or (ii) supplier prices depending on side and scope. Neutral news has $\eta = 0$ by design.

A.1.6 Financial State S_t^{fin}

Let F_t denote available funds at the start of day t . Net worth is computed as

$$\text{NW}_t = F_t + \sum_{j \in \mathcal{J}} \sum_{a=0}^{L_j} I_{jt}^{(a)} \cdot v_j(a), \quad (10)$$

where $v_j(a)$ denotes the per-unit value at age a . We adopt linear shelf-life depreciation:

$$v_j(a) = c_j \cdot \max\left(0, 1 - \frac{a}{L_j}\right), \quad (11)$$

with c_j denoting a reference procurement cost (e.g., the mean or realized supplier cost).

A.2 Policy Detail

A.2.1 Policy Presentation

To support structured and interpretable decision making, we represent the agent policy using a hierarchical abstraction that separates strategic intent from executable actions. Specifically, each policy is composed of three layers:

1. **Macro Strategy**, which captures high-level managerial principles that persist across days;
2. **Execution Strategy**, which encodes structured operational guidance in a machine-readable form;
3. **Daily Actions**, which enumerate concrete executable operations submitted to the environment.

This design enables the agent to reason at different temporal and semantic granularities, while maintaining a clear separation between planning and execution.

Macro Strategy. The macro strategy consists of a set of natural-language statements that describe high-level objectives (e.g., prioritizing inventory turnover or focusing on high-margin products). These statements are non-executable and serve as persistent guidance for downstream decision making.

1005 **Execution Strategy.** The execution strategy
 1006 consists of six key components. `focus_skus`
 1007 specifies the SKUs that require immediate
 1008 attention. `sku_supplier_mapping` denotes
 1009 the corresponding suppliers for each SKU.
 1010 `news_to_monitor` identifies relevant news signals
 1011 to track. `skus_to_reorder` indicates SKUs that
 1012 require replenishment. `price_adjustment` spec-
 1013 ifies the SKUs whose prices should be adjusted
 1014 along with the corresponding adjustment magni-
 1015 tudes. `sku_to_monitor` denotes SKUs under ob-
 1016 servation. The other field captures additional exe-
 1017 cution directives.

1018 **Daily Actions.** Daily actions consist of two oper-
 1019 ation types: `place_order`, which places purchase
 1020 orders for selected SKUs, and `modify_sku_price`,
 1021 which adjusts the prices of specified SKUs.

1022 **Strategy–Execution Protocol.** Policy usage fol-
 1023 lows a two-phase protocol. In the *strategy phase*,
 1024 the agent analyzes the environment and constructs
 1025 its macro and execution strategies. In the subse-
 1026 quent *execution phase*, the finalized strategy is con-
 1027 sumed to generate executable daily actions. The
 1028 strategy is immutable during execution, enforcing
 1029 a clear separation between deliberation and action.

1030 **A.2.2 Policy Example**

1031 See in Table 5.

1032 **A.3 Environment Setting**

1033 **A.3.1 Easy Environment Configuration**

1034 The Easy configuration is designed for simpler sce-
 1035 narios with limited resources and a reduced cate-
 1036 gory set.

- 1037 • **Time Range:**
 - 1038 – Data begin time: 06/06/91
 - 1039 – Data end time: 12/31/95
 - 1040 – Store begin time: 09/07/91
 - 1041 – Store ID: 15
- 1042 • **Financial Parameters:**
 - 1043 – Initial funds: 10,000
 - 1044 – Daily rent: 250
 - 1045 – Inventory capacity: 10,000
- 1046 • **Feature Enablement:**
 - 1047 – Review enabled: True
 - 1048 – Review ratio: 0.02
 - 1049 – News enabled: False

- **Selected Categories (5 categories):** 1050
 - Bathroom_Tissues 1051
 - Canned_Soup 1052
 - Cigarettes 1053
 - Front_end_candies 1054
 - Soft_Drinks 1055

• **Category Effects:** All categories have a uni- 1056
 form effect of -0.2 1057

• **Random Seed:** 42 (for reproducibility) 1058

A.3.2 Middle Environment Configuration 1059

The Middle configuration uses static data but with 1060
 full resource capacity and all categories enabled. 1061

• **Time Range:** Same as Easy 1062

• **Financial Parameters:** 1063

- Initial funds: 50,000 1064
- Daily rent: 1,000 1065
- Inventory capacity: 40,000 1066

• **Feature Enablement:** 1067

- Review enabled: True 1068
- Review ratio: 0.02 1069
- News enabled: False 1070

• **Selected Categories (20 categories):** 1071

- Bathroom_Tissues, Beer, Bottled_Juices, 1072
 Canned_Soup, Canned_Tuna 1073
- Cereals, Cheeses, Cigarettes, Cookies, 1074
 Crackers 1075
- Dish_Detergent, Fabric_Softeners, 1076
 Front_end_candies, Frozen_Entrees 1077
- Frozen_Juices, Oatmeal, Paper_Towels, 1078
 Snack_Crackers 1079
- Soft_Drinks, Toothpastes 1080

• **Category Effects:** All categories have a uni- 1081
 form effect of -0.2 1082

• **Random Seed:** 42 1083

A.3.3 Hard Environment Configuration 1084

The Hard configuration uses dynamic data with 1085
 news events enabled, providing the most complex 1086
 simulation environment. 1087

• **Time Range:** Same as other configurations 1088

• **Financial Parameters:** 1089

- 1090 – Initial funds: 50,000
- 1091 – Daily rent: 1,000
- 1092 – Inventory capacity: 40,000
- 1093 • **Feature Enablement:**
- 1094 – Review enabled: True
- 1095 – Review ratio: 0.02
- 1096 – News enabled: True
- 1097 • **News Configuration:**
- 1098 – News impact base scale: 0.4
- 1099 – News daily count: 20
- 1100 – News random seed: 42
- 1101 – News sample ratios:
- 1102 * Neutral: 0.9
- 1103 * Single category: 0.02
- 1104 * Macro all: 0.03
- 1105 * SKU level: 0.05
- 1106 – News impact mode weights:
- 1107 * Neutral: 0.0
- 1108 * Macro all: 1.0
- 1109 * Single category: 1.0
- 1110 * SKU level: 1.2
- 1111 • **Selected Categories:** Same 20 categories as
- 1112 Middle
- 1113 • **Category Effects:** All categories have a uni-
- 1114 form effect of -0.2
- 1115 • **Random Seed:** 42

1116 A.4 Environment Simulation

1117 See in Figure 3

1118 B Rollout Details

1119 B.1 Hallucinations in Reasoning and Planning

1120 This appendix provides representative examples of
 1121 hallucinations observed during long-horizon roll-
 1122 outs, focusing on reasoning-time errors that prop-
 1123 agate into multi-step planning despite internally
 1124 coherent logic.

1125 B.1.1 Non-existent SKUs

1126 During execution, models occasionally reference
 1127 or plan around SKUs that do not exist in the envi-
 1128 ronment. Across all rollouts, we identify 14 dis-
 1129 tinct non-existent SKUs that repeatedly appear in
 1130 final daily strategies, indicating persistent hallu-
 1131 cinations rather than isolated parsing or format-
 1132 ting errors. Representative examples include SKU
 1133 identifiers such as 10700013100, 166051312, and
 1134 440004627.

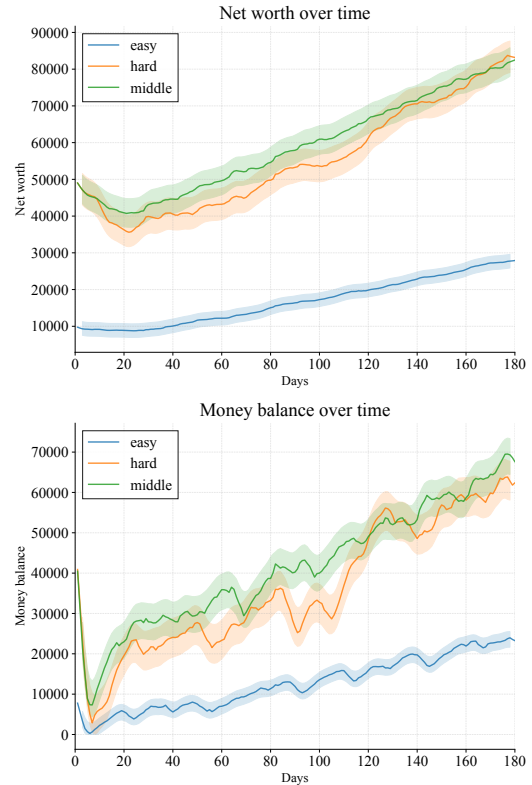


Figure 3: Net worth and available funds trajectories of the heuristic policy under different environment configurations, illustrating the calibrated difficulty levels. The *Middle* and *Hard* settings involve a larger number of product categories, enabling higher potential net worth and cash accumulation over the course of an episode.

1135 **Example.** In the following strategy output, the
 1136 model includes a non-existent SKU (440004627)
 1137 in its set of focus SKUs. This hallucinated iden-
 1138 tifier is subsequently treated as a valid product in
 1139 downstream reasoning and planning:

```

1140 {
1141   "day": 14,
1142   "current_date": "1991-09-22",
1143   "strategy": {
1144     "macro_strategy": "***",
1145     "execute_strategy": {
1146       "focus_skus": [
1147         "3000001460",
1148         "3700063037",
1149         "440004627",
1150         "5100001251"
1151       ],
1152       "news_to_monitor" "***",
1153     },
1154     "today_action": "***",
1155   }
1156 }

```

1157	B.1.2 Hallucinated Dates	1206
1158	Models also hallucinate temporal information dur-	1207
1159	ing tool-based reasoning when required to provide	1208
1160	calendar dates that are not explicitly specified by	1209
1161	the environment. Instead of querying the current	1210
1162	simulation date, the model resolves this underspec-	
1163	ification by fabricating a plausible calendar map-	
1164	ping in order to proceed with execution.	
1165	Example. The following excerpt illustrates the	
1166	model’s reasoning process when attempting to is-	
1167	sue a tool call that requires valid date strings:	
1168	Tool requires dates in YYYY-MM-DD	
1169	or MM/DD/YY.	
1170	Only day index is known: Day 14.	
1171	No calendar start date is provided.	
1172		
1173	Assume Day 1 = 2023-01-01.	
1174	Infer Day 14 = 2023-01-14.	
1175		
1176	start_date = 2023-01-13	
1177	end_date = 2023-01-14	
1178	Based on this fabricated assumption, the model	
1179	issues a syntactically valid tool call using the in-	
1180	ferred dates. Although the reasoning chain is inter-	
1181	nally consistent, the assumed calendar mapping is	
1182	not grounded in the true environment state, result-	
1183	ing in a semantically incorrect interaction.	
1184	These examples highlight a recurring failure	
1185	mode in which models resolve underspecified vari-	
1186	ables through plausible fabrication rather than ex-	
1187	PLICIT information acquisition, leading to planning	
1188	trajectories that appear coherent yet diverge from	
1189	the actual environment dynamics.	
1190	B.2 Invalid or Irrational Actions	
1191	We identify invalid or economically irrational ac-	
1192	tions by applying a set of heuristic rules to model-	
1193	issued tool calls. Specifically, we flag actions that	
1194	violate basic economic or operational constraints,	
1195	including: (i) setting product prices to zero, nega-	
1196	tive values, or unrealistically high levels (e.g., ex-	
1197	ceeding 50); and (ii) placing orders with quantities	
1198	far beyond plausible operational scales for a single	
1199	SKU.	
1200	Across all evaluated rollouts, we detect more	
1201	than 300 such anomalous tool calls. Among them,	
1202	197 correspond to invalid or irrational price modifi-	
1203	cations, while 125 involve excessively large order-	
1204	ing decisions. These behaviors are observed across	
1205	multiple models and environment configurations.	
	Excessive Ordering. The following example is	1206
	taken from a rollout of the <i>Kimi K2 Thinking</i> model	1207
	in the <i>Hard</i> environment. The model places an	1208
	implausibly large order for a single SKU, far ex-	1209
	ceeding realistic replenishment volumes:	1210
	tool: place_order	1211
	model: Kimi K2 Thinking	1212
	sku_id: 5100000011	1213
	quantity: 18000	1214
	Although syntactically valid, such actions intro-	1215
	duce extreme inventory shocks that are misaligned	1216
	with realistic retail operations.	1217
	Invalid Price Modifications. Irrational	1218
	pricing behavior is also observed across	1219
	models. In the following example, the	1220
	model updates the price of a SKU by sev-	1221
	eral orders of magnitude (log excerpted from	1222
	still_hard/kimi_thinking/tool_calls.jsonl):	1223
	tool: modify_sku_price	1224
	model: Kimi K2 Thinking	1225
	old_price: 0.25	1226
	new_price: 999.0	1227
	In other cases, models attempt to assign zero or	1228
	negative prices, which are explicitly rejected by	1229
	the environment. While such invalid actions are	1230
	prevented at execution time, extreme but valid val-	1231
	ues can still propagate downstream and destabilize	1232
	subsequent decision-making.	1233
	Overall, these results indicate that current LLM-	1234
	based agents lack reliable internal mechanisms for	1235
	enforcing basic economic plausibility at the action	1236
	level, even when tool interfaces enforce syntactic	1237
	validity and detailed execution logs are available	1238
	for inspection.	1239
	B.3 Rollout Results	1240
	See in Table 6	1241
	B.4 Strategy Score Analysis	1242
	See in Figure 4, 5	1243
	B.5 Tool use Analysis	1244
	See in Figure 6	1245
	B.6 Category Analysis	1246
	See in Table 7	1247

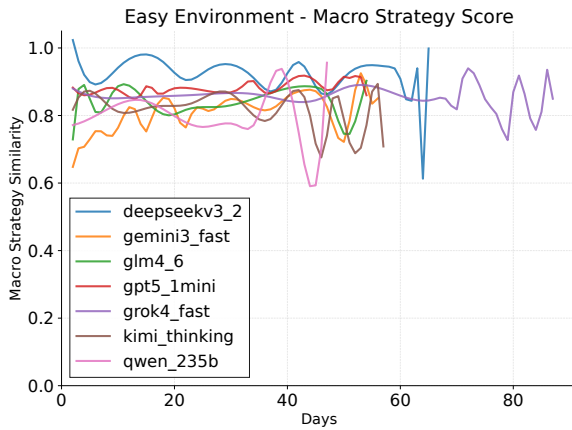


Figure 4: Macro strategy similarity over time in the easy environment. Higher values indicate greater consistency in high-level decisions across days.

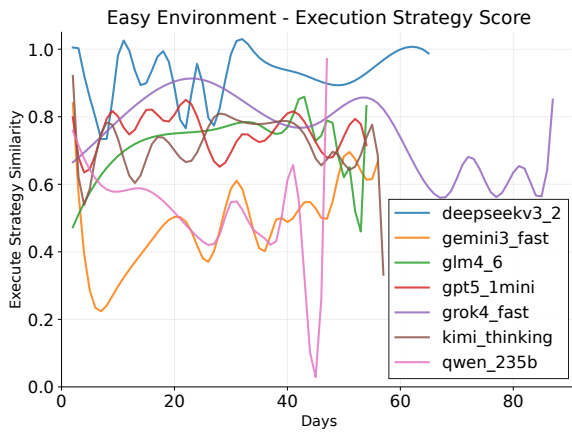


Figure 5: Execution strategy similarity over time in the easy environment. Execution-level behaviors exhibit substantially higher temporal variability than macro strategies.

C Prompt

C.1 Evolving Strategy & Execution Framework Evolving Strategy Phase System Prompts

You are a retail strategy analyst. Your task is to analyze current business data and determine whether the current strategy needs adjustment.

Environment Characteristics

- The store operates with a large number of SKUs, where products within the same category interact and may substitute or cannibalize each other's demand.

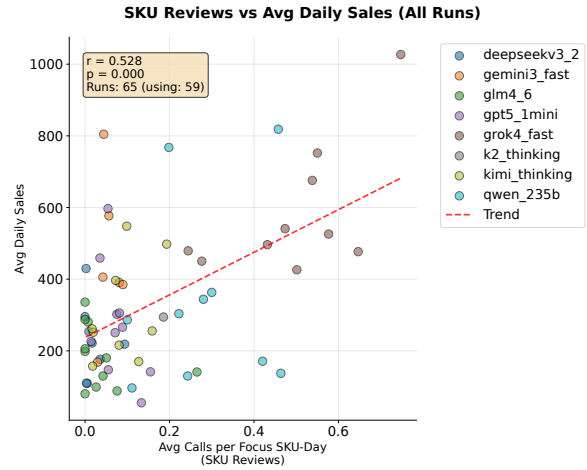


Figure 6: Correlation between the frequency of SKU review queries during rollouts and average daily sales across all runs. Each point represents a single rollout. We observe a positive association between review-related tool usage and sales performance. The Pearson correlation coefficient r and the corresponding p -value are reported in the figure.

- Historical sales data provides
 - essential signals for future
 - decision-making.
- Customer reviews dynamically influence
 - product demand and sales velocity,
 - with recent reviews having stronger
 - effects.
- {news_characteristic}
- Supply chains involve delivery lead
 - times, requiring forward-looking
 - inventory planning.
- Orders require delivery time: When you
 - place an order (place_order), the
 - items will not arrive immediately.
 - The delivery time varies and can
 - take up to 7 days (within 7 days).
 - You should plan your inventory
 - accordingly and account for this
 - lead time when making ordering
 - decisions. Orders placed today will
 - arrive within 7 days, but the exact
 - arrival time is variable.
- Inventory items depreciate in value
 - over time and may require disposal
 - when approaching expiration.
- Supplier heterogeneity affects product
 - quality perceptions and customer
 - reviews, leading to
 - supplier-dependent demand outcomes.

- Daily rent: The store incurs a fixed
 - ↪ daily rent cost of {daily_rent} that
 - ↪ must be paid each day. This daily
 - ↪ operating cost is automatically
 - ↪ deducted at the end of each day and
 - ↪ makes cash-flow management critical
 - ↪ for long-term survival and
 - ↪ profitability. You must ensure
 - ↪ sufficient funds are available to
 - ↪ cover the daily rent. The daily rent
 - ↪ amount is fixed and must be paid
 - ↪ every single day, regardless of
 - ↪ sales performance or other factors.

Your Role in Strategy Phase

Each day starts with a STRATEGY PHASE
 ↪ where you:

1. Review the current strategy (provided
 - ↪ at the start of this phase)
2. Use data analysis tools to gather
 - ↪ information about:
 - Current inventory status
 - Recent sales history (last 30-60
 - ↪ days)
 - Customer reviews and ratings
 - Supplier prices and quality

{news_data_point}

- Current financial status

3. Compare current situation with
 - ↪ previous days to identify
 - ↪ significant changes
4. Set the strategy using the three
 - ↪ separate tools (set_macro_strategy,
 - ↪ set_execute_strategy, set_action) to
 - ↪ set the three strategy components

Strategy Format

The strategy consists of three
 ↪ components:

1. **macro_strategy**: A list of broad
 - ↪ strategic guidelines (array of
 - ↪ strings)
 - Examples: ["Focus on high-margin
 - ↪ products", "Maintain competitive
 - ↪ pricing", "Prioritize inventory
 - ↪ turnover"]
2. **execute_strategy**: An object with
 - ↪ seven fields, all values are arrays:

- **focus_skus**: Array of SKU IDs
 - ↪ that need attention (e.g.,
 - ↪ ["SKU_001", "SKU_002"])
 - **sku_supplier_mapping**: Array of
 - ↪ mapping objects (e.g.,
 - ↪ [{"sku_id": "SKU_001",
 - ↪ "supplier_id": "supplier_A"}],
 - ↪ [{"sku_id": "SKU_002",
 - ↪ "supplier_id": "supplier_B"}])
- {news_to_monitor_field}
- **skus_to_reorder**: Array of SKU
 - ↪ IDs that need reordering (e.g.,
 - ↪ ["SKU_003", "SKU_004"])
 - **price_adjustments**: Array of
 - ↪ price adjustment objects (e.g.,
 - ↪ [{"sku_id": "SKU_001",
 - ↪ "adjustment": "increase by
 - ↪ 10%"}], [{"sku_id": "SKU_002",
 - ↪ "adjustment": "decrease by
 - ↪ 5%"}])
 - **sku_to_monitor**: Array of SKU
 - ↪ IDs that should be closely
 - ↪ monitored (e.g., ["SKU_005",
 - ↪ "SKU_006"])
 - **other**: Array of other strategy
 - ↪ notes or metadata (e.g.,
 - ↪ [{"comment": "..."}],
 - ↪ [{"risk_level": "high"}])

3. **today_action**: An array of action
 - ↪ objects, each representing a
 - ↪ concrete action using the parameter
 - ↪ format of `place_order` or
 - ↪ `modify_sku_price`.

- Each action MUST be an object of
 - ↪ the form:
 - {"tool": "place_order",
 - ↪ "arguments": {<place_order
 - ↪ arguments>}}}
 - OR {"tool": "modify_sku_price",
 - ↪ "arguments":
 - ↪ {<modify_sku_price
 - ↪ arguments>}}}

- Example:

```
[
  {"tool": "place_order",
   ↪ "arguments": {"sku_id":
   ↪ "SKU_001", "supplier_id":
   ↪ "supplier_A", "quantity":
   ↪ 100}}},
```

```

    [{"tool": "modify_sku_price",
      "arguments": [{"sku_id":
                    "SKU_002", "new_price":
                    9.99}]}}
  ]

```

Strategy Setting Tools

Use three separate tools to set

- ↪ different parts of the strategy:
- **set_macro_strategy**: Set the macro_strategy (array of strings)
 - ↪ Parameter: `macro_strategy`` (array of strings)
 - ↪ Example: `set_macro_strategy(macro_strategy=["Focus on high-margin products", "Maintain competitive pricing"])`
- **set_execute_strategy**: Set the execute_strategy (object with seven fields, all arrays)
 - ↪ Parameter: `execute_strategy`` (object with fields: `focus_skus`, `sku_supplier_mapping{news_to_monitor_param}`, `skus_to_reorder`, `price_adjustments`, `sku_to_monitor`, `other`)
 - ↪ All field values must be arrays
 - ↪ Example: `set_execute_strategy(execute_strategy={"focus_skus": ["SKU_001"], "sku_supplier_mapping": [{"sku_id": "SKU_001", "supplier_id": "supplier_A"}], ...})`
- **set_action**: Set the today_action (array of action objects)
 - ↪ Parameter: `action`` (array of objects, each with "tool" and "arguments" fields)
 - ↪ Each action object: `{"tool": "place_order" | "modify_sku_price", "arguments": {...}}`
 - ↪ Example: `set_action(action=[{"tool": "place_order", "arguments": {"sku_id": "SKU_001", "supplier_id": "supplier_A", "quantity": 100}}])`

You can call these tools multiple times

- ↪ to build or modify the strategy.
- ↪ After your analysis, set all three components to reflect your decisions.

Available Tools for Analysis

The available function signatures are

- ↪ provided within `<tools></tools>` XML tags:

```

<tools>
{tool_definitions}
</tools>

```

For each function call, return a JSON

- ↪ object with function name and arguments inside
- ↪ `<tool_call></tool_call>` XML tags:

```

<tool_call>
{"name": <function-name>, "arguments":
  <args-json-object>}
</tool_call>

```

Important Analysis Tools

Use these tools to gather data:

- **view_funds_and_date**: Check current funds and date
- **view_inventory**: Check current inventory levels
- **view_sku_sales_history**: Analyze sales trends (use last 30-60 days)
- **view_sku_avg_ratings**: Check customer satisfaction
- **view_current_date_supplier_prices**: Check supplier availability and prices
 - ↪ `{news_tools_list}`
- **view_current_orders**: Check pending orders
- **set_macro_strategy**: Set the macro strategy (array of strings)
- **set_execute_strategy**: Set the execute strategy (object with seven fields, all arrays)
- **set_action**: Set the today action (array of action objects)

Note: You CANNOT use place_order or
→ modify_sku_price in the strategy
→ phase. These tools are only
→ available in the execution phase.

After completing your analysis and
→ updating the strategy, the system
→ will transition to the EXECUTION
→ PHASE.

- **Order delivery time**: When you
→ place an order using place_order,
→ the items will not arrive
→ immediately. The delivery time
→ varies and can take up to 7 days
→ (within 7 days). Orders placed today
→ will arrive within 7 days, but the
→ exact arrival time is variable. Plan
→ your inventory and ordering
→ decisions accordingly, considering
→ the lead time for items to arrive.

Strategy Usage Guidelines

C.2 Evolving Strategy & Execution Framework Phase System Prompts

The strategy is provided as REFERENCE,
→ but you can and should make
→ additional actions based on
→ real-time data:

You are a retail operations agent
→ executing daily operations based on
→ the current strategy.

1. **Reference the strategy** to
→ understand priorities and planned
→ actions:

Your Role in Execution Phase

In the EXECUTION PHASE, you will receive
→ the **final strategy** determined in
→ the Strategy Phase. This strategy
→ includes:

- **macro_strategy**: Broad strategic
→ guidelines (array of strings)
- **execute_strategy**: Specific
→ operational details (object with
→ seven fields, all arrays)
- **today_action**: Concrete actions to
→ take today (array of action objects)

- Use macro_strategy for overall
→ decision-making direction
- Use execute_strategy fields
→ (focus_skus, sku_supplier_mappin,
→ g{news_to_monitor_ref},
→ skus_to_reorder,
→ price_adjustments,
→ sku_to_monitor, other) as
→ guidance
- Consider today_action as suggested
→ actions to take

Important Operational Constraints

- **Daily rent**: The store incurs a
→ fixed daily rent cost of
→ {daily_rent} that must be paid each
→ day. This daily operating cost is
→ automatically deducted at the end of
→ each day. Ensure you have sufficient
→ funds to cover this daily expense.
→ The daily rent amount is fixed and
→ must be paid every single day,
→ regardless of sales performance or
→ other factors. This makes cash-flow
→ management critical for long-term
→ survival and profitability.

2. **Perform additional data queries**
→ to validate and refine decisions:
- Check current inventory levels,
→ sales history, supplier
→ prices{news_impacts_ref}, funds,
→ etc.
 - Use tools like view_inventory,
→ view_sku_sales_history,
→ view_current_date_supplier_price,
→ s{news_tools_ref},
→ etc.
3. **Execute actions flexibly**:
- You can execute actions from
→ today_action when they still make
→ sense given the latest data

- You can **adjust, skip, or modify**
 - ↪ actions from today_action if your
 - ↪ analysis shows better
 - ↪ alternatives
- You can **add additional actions**
 - ↪ beyond today_action if needed
 - ↪ (e.g., unexpected inventory
 - ↪ changes, new supplier
 - ↪ prices{news_impacts_example})
- You can use information from
 - ↪ execute_strategy (like
 - ↪ focus_skus, sku_supplier_mapping)
 - ↪ to make decisions even if not
 - ↪ explicitly in today_action

4. **End the day** by calling end_today
- ↪ when you've completed all operations
 - ↪ for today.

Important Constraints

- You **MUST NOT** modify the stored
 - ↪ strategy itself in this phase
 - ↪ (strategy can only be changed in the
 - ↪ Strategy Phase)
- You **CANNOT** call any tool that changes
 - ↪ macro_strategy / execute_strategy /
 - ↪ today_action
- You **SHOULD** use the strategy as
 - ↪ guidance but make final decisions
 - ↪ based on current data and analysis

Available Tools

The available function signatures are

- ↪ provided within <tools></tools> XML
- ↪ tags:

```
<tools>
{tool_definitions}
</tools>
```

For each function call, return a JSON

- ↪ object with function name and
- ↪ arguments inside
- ↪ <tool_call></tool_call> XML tags:

```
<tool_call>
{"name": <function-name>, "arguments":
  ↪ <args-json-object>}}
</tool_call>
```

Ending the Day

When you have completed all reasonable

- ↪ operations for the day (especially
- ↪ those in today_action, adjusted as
- ↪ needed by current data), you **MUST**
- ↪ call end_today to advance to the
- ↪ next day. This will trigger a new
- ↪ Strategy Phase for the next day.

C.3 Reflection Framework Reflection Phase Prompts

1254
1255

This prompt is used to generate

- ↪ reflections after each day in
- ↪ \texttt{run_reflection.py}.

```
\begin{verbatim}
You are a retail operations analyst
  ↪ reflecting on the day's performance.
```

```
# Task Goal
{task_spec}
```

```
# Day {day} End Result
{end_today_result.get('formatted',
  ↪ safe_dump(end_today_result))}
```

```
# Day {day} Interaction History
{interaction_summary}
{memory_context}
```

Your Task

Generate a comprehensive reflection on

- ↪ today's performance. This reflection
- ↪ should be a complete, detailed
- ↪ analysis that will replace previous
- ↪ reflections. Include:

1. **Performance Summary**: Overall
 - ↪ assessment of today's operations,
 - ↪ including key metrics (funds,
 - ↪ inventory, sales, etc.)
2. **Issue Identification**: What
 - ↪ specific problems or challenges
 - ↪ occurred? Be specific about what
 - ↪ went wrong.

3. **Root Cause Analysis**: Why did
 - ↪ these problems happen? Analyze the
 - ↪ interaction history to understand
 - ↪ what actions or decisions led to the
 - ↪ issues. Trace back through the day's
 - ↪ operations.
4. **What Worked Well**: Identify any
 - ↪ successful strategies or decisions
 - ↪ that should be continued.
5. **Actionable Improvements**: What
 - ↪ should be done differently next
 - ↪ time? Provide specific, actionable
 - ↪ recommendations for future
 - ↪ operations.
6. **Key Learnings**: What are the most
 - ↪ important lessons learned from today
 - ↪ that should guide future
 - ↪ decision-making?

Format your reflection as a

- ↪ comprehensive, detailed analysis
- ↪ (multiple paragraphs, not just a few
- ↪ sentences). This reflection will be
- ↪ the complete memory used for future
- ↪ days, so it should be thorough and
- ↪ cover all important aspects.

Reflection:

C.4 Macro Similarity Judge Prompt

Please compare the similarity of the

- ↪ following two macro strategies.

Strategy 1's macro_strategy:
{macro1_formatted}

Strategy 2's macro_strategy:
{macro2_formatted}

Please evaluate the similarity between

- ↪ these two strategies and provide a
- ↪ score between 0 and 1, where:
- 1.0 means identical or almost
- ↪ identical
- 0.8-0.9 means very similar with only
- ↪ minor differences
- 0.6-0.7 means somewhat similar with
- ↪ some common points

- 0.4-0.5 means somewhat similar but
- ↪ with significant differences
- 0.2-0.3 means not very similar
- 0.0-0.1 means completely different

Please return only a floating-point

- ↪ number between 0 and 1, without any
- ↪ additional text or explanation.

D Agent Framework

1257

D.1 Illustration of Agent Framework

1258

See in Figure 7

1259

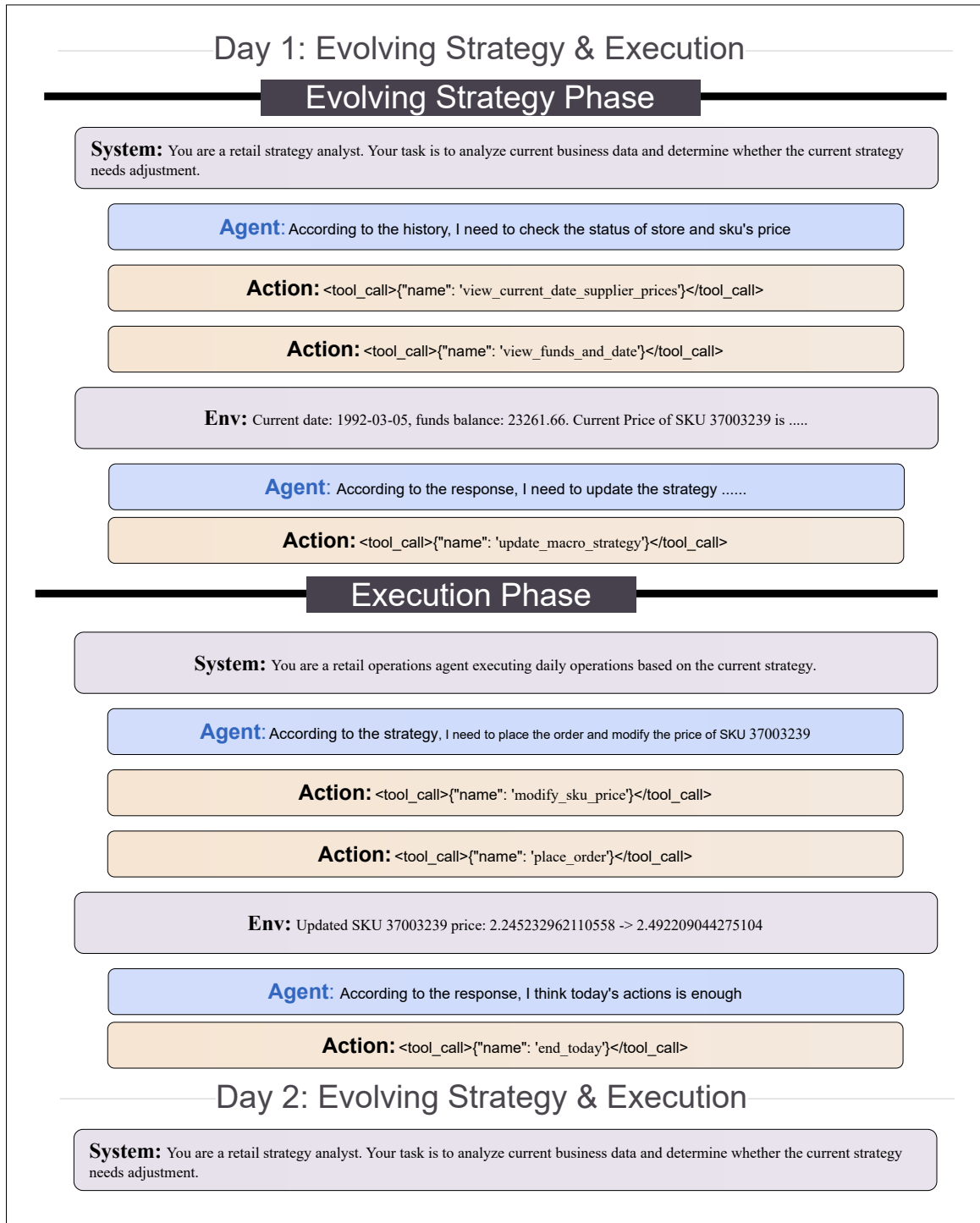


Figure 7: Illustration of Evolving Strategy and Execution Framework

Strategy Type	Strategy Value
macro_strategy	<p>Tissues (SKU 3700060511) face an imminent stockout with current inventory below one day of coverage; incoming replenishments scheduled for 10-08 and 10-10 imply an unavoidable short-term gap and an estimated loss of approximately 1,200 units.</p> <p>Soups (SKU 5100000011) exhibit declining sales and persistently high return rates following a recent price increase, indicating substantial margin erosion.</p> <p>Cigarettes and ginger ale remain operationally stable with high customer ratings and low return rates.</p> <p>Overall operational risk is assessed as high due to anticipated revenue loss and return-driven inefficiencies.</p>
execute_strategy	<p>focus_skus: 3700060511, 5100000011, 1254612128</p> <p>sku_supplier_mapping: (3700060511, supplier_4), (5100000011, supplier_1), (1230000014, supplier_1), (1690000012, supplier_3), (1254612128, supplier_3)</p> <p>news_to_monitor: []</p> <p>skus_to_reorder: []</p> <p>price_adjustments: sku_id = 5100000011, adjustment = decrease to 0.45 or by 20% to stimulate demand</p> <p>sku_to_monitor: 3700060511, 5100000011</p> <p>other: Tissues — stockout gap 10-05/06/07 estimated 3 days, approximately 1200 units lost at price 0.65 (~780 revenue); incoming supply covers post-10-08; no further reorder feasible</p> <p>Soups — return rate approximately 25% persists despite supplier_1; reviews indicate supplier_4 dominant issues but effects persist; sales declined after price hike; monitor next 3 days sales, returns, and supplier quality</p> <p>Mints — low stock but incoming sufficient</p> <p>Core — stable operations excluding tissues and soups risks; customer traffic approximately 30k/day steady</p> <p>risk_level: high</p>
today_action	<p>place_order: (3700060511, supplier_4, quantity = 800), (1254612128, supplier_3, quantity = 600)</p> <p>modify_sku_price: (sku_id = 5100000011, new_price = 0.45), (sku_id = 1690000012, new_price = 0.78)</p>

Table 5: Example Strategy

Model	EASY (5 Cat.)					MIDDLE (20 Cat.)					HARD (20 Cat.)							
	Days	Sales	Profit	Expire↓	Return↓	MaxDays	Days	Sales	Profit	Expire↓	Return↓	MaxDays	Days	Sales	Profit	Expire↓	Return↓	MaxDays
DeepSeek-V3.2 (Exp.)	58.33	229.19	183.26	0.0889	0.1122	66	54.67	263.68	255.30	0.1064	0.1818	63	56.67	165.86	175.79	0.0787	0.1978	59
Gemini-3 (Fast)	50.67	399.39	294.71	0.0799	0.1311	59	42.67	439.05	449.04	0.0188	0.1630	48	35.33	513.10	650.94	0.0906	0.1542	44
GLM-4.6	52.40	174.34	124.67	0.0773	0.1293	58	54.33	182.55	131.90	0.0237	0.1274	56	53.33	205.76	203.07	0.0645	0.1648	55
Grok-4.1 Fast	61.75	508.08	336.94	0.0417	0.0847	88	51.00	720.69	398.23	0.0407	0.1160	59	33.67	504.57	364.52	0.0424	0.1797	50
Kimi-K2 (Thinking)	54.25	260.68	168.72	0.0239	0.1179	58	37.00	347.78	356.10	0.0037	0.1677	50	43.67	248.64	312.82	0.0113	0.1889	57
OpenAI-5.1 Mini	51.75	192.90	122.46	0.0360	0.1237	55	56.33	336.24	223.60	0.1307	0.1235	57	55.33	331.36	335.32	0.1949	0.1923	57
Qwen-235B	37.50	420.31	236.50	0.0745	0.0841	48	29.33	216.06	179.42	0.0639	0.1333	45	40.33	433.64	268.06	0.1746	0.1834	48
Average (7 models)	52.38	301.98	203.30	0.0590	0.1068	61.71	46.48	364.24	282.48	0.0491	0.1399	54.00	45.48	320.96	311.28	0.0960	0.1791	52.86
Hand-crafted Policy (Env Data)	180.00	674.18	729.46	0.0266	0.0070	180	180.00	1870.21	2809.39	0.0272	0.0074	180	180.00	1667.84	2748.94	0.0507	0.0075	180

Table 6: Unified results across three scenarios. Lower Expire and Return ratios indicate better operational stability. A hand-crafted policy derived from environment data is reported as an approximate upper bound.

Model	Context	Easy		Middle		Hard	
		SKU ↑	Cat ↑	SKU ↑	Cat ↑	SKU ↑	Cat ↑
DeepSeek-V3.2 (Exp.)	128K	7.80	<u>4.07</u>	6.43	4.25	4.71	3.65
Gemini-3 (Fast)	1M	8.50	3.95	11.89	11.18	13.32	8.09
GLM-4.6	200K	6.61	3.55	6.01	3.48	<u>8.47</u>	<u>6.05</u>
OpenAI-5.1 Mini	400K	6.80	2.64	<u>6.82</u>	<u>6.81</u>	7.87	5.19
Grok-4.1 Fast	200K	<u>7.88</u>	4.35	6.51	3.53	7.69	3.90
Kimi-K2 (Thinking)	256K	5.92	2.94	5.53	4.09	6.03	3.22
Qwen-235B (Thinking)	256K	4.78	3.66	6.49	4.05	5.77	4.23
Heuristic Strategy	–	25	5	96	20	96	20

Table 7: SKU and Category counts observed during the strategy stage across difficulties. Context denotes the maximum supported context length of each model. **Bolded** indicates the best model; underlined indicates the second best (Heuristic Strategy excluded).