

Welcome to the Era of Delayed Rewards for Language Agents: On Non-Verifiable Tasks

Anonymous ACL submission

Abstract

Agent tasks divide into *verifiable* (e.g., math, code) with immediate ground-truth rewards, and *non-verifiable* (e.g., marketing, policy, research communication) where rewards appear instant but are fundamentally *delayed* when considering users’ long-term goals. A user asking an LLM to write marketing copy may receive immediate output, but their true objective—readership, engagement, influence—unfolds over time through social propagation. Current approaches to non-verifiable tasks—self-refine, LLM-as-judge, multi-agent debate—rely on *instant feedback* that cannot capture this delayed, emergent value. We argue for a paradigm shift: **from instant feedback to delayed reward derivation via task-appropriate simulation environments**. Different tasks require different simulations: social media simulators for viral content, academic platforms for research communication, or policy debate forums for proposals. We validate this using OASIS, a scalable social media simulator with LLM-powered agents, comparing self-refine with our method MARFE (Multi-Agent Reward-Free Evolution). Across four tasks evaluated by three frontier LLM judges, MARFE achieves **58.3% win rate** versus baseline’s 41.7%, demonstrating that delayed social feedback provides superior signal for non-verifiable tasks.

1 Introduction

Agent tasks can be divided by *verifiability*. For verifiable tasks like mathematical reasoning (Wei et al., 2022; Shao et al., 2024) or code generation (OpenAI, 2023), outputs can be compared against ground truth to derive immediate rewards, enabling reinforcement learning (Shao et al., 2024; Schulman et al., 2017). But what about tasks where no ground truth exists?

Consider writing viral marketing content, crafting policy proposals, product announcements, or

sharing research findings on social media. These **non-verifiable tasks** may appear to have instant feedback—a judge can immediately rate “is this tweet good?” or “is this announcement compelling?” However, this misses a crucial insight: **users’ true objectives for non-verifiable tasks are inherently long-term**. A user asking an LLM to write marketing copy wants conversions and brand awareness over weeks. A user requesting a novel ultimately cares about readership, reviews, and influence—outcomes that unfold through social propagation over time. The immediate output quality is merely a proxy for these delayed, emergent rewards.

This reveals the **fundamental nature of non-verifiable tasks**: their true value is *socially constructed* and *temporally delayed*. A marketing tweet’s success depends not on intrinsic quality assessable at generation time, but on how it propagates through social networks over hours or days—who sees it, who engages, who converts (Tu and Neumann, 2022).

The Challenge: Instant Feedback Cannot Capture Delayed Value. Current approaches to non-verifiable tasks rely on *instant feedback*, fundamentally mismatched with the delayed nature of social value:

Self-Refine approaches (Madaan et al., 2023; Shinn et al., 2023) ask the model to critique and improve its own output. But a single perspective cannot anticipate how *diverse audiences* will react. When asked “is this tweet engaging?”, the model provides generic suggestions (“add more hooks”) without knowing which hooks resonate with developers versus executives, or what objections each segment might raise.

LLM-as-Judge (Zheng et al., 2023) provides external evaluation but remains *synchronous*—it scores content in isolation, unable to model how value *emerges* through propagation. A judge cannot

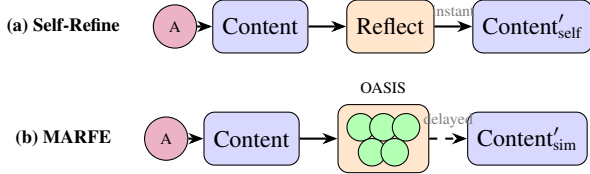


Figure 1: Comparison of feedback mechanisms. (a) Self-Refine: agent (A) generates content and reflects instantly. (b) MARFE: agent generates content, injects into OASIS where diverse users react over time, providing delayed feedback.

predict that a tweet will go viral because early tech adopters repost it to their networks, triggering a cascade.

Even **Multi-Agent Debate** (Du et al., 2023; Li et al., 2023), despite involving multiple perspectives, operates with *global information access*—all agents see all content simultaneously and respond instantly. This fundamentally differs from real social systems where: (1) users see *personalized feeds* filtered by recommendation algorithms; (2) content visibility *decays over time*; (3) engagement unfolds through *cascading reactions* across network connections. In debate, an “AI Skeptic” agent immediately sees and critiques content; in reality, skeptics might never see it if the algorithm doesn’t surface it to them, or might see it days later after it’s already gained momentum.

Our Insight: Rewards Emerge from Simulation Over Time. We argue that **rewards for non-verifiable tasks are emergent properties of output-environment interactions**—they cannot be computed instantly but must be *derived* by observing how content propagates through a realistic social system. This calls for a paradigm shift: **from instant feedback to delayed reward derivation.**

Our Proposal: Task-Appropriate Simulation Environments. We propose using simulation environments tailored to each task’s long-term objectives. **Different tasks require different simulations:** social media simulators for viral content, academic platforms for research sharing, debate forums for policy. We demonstrate using OASIS (Yang et al., 2024), which simulates realistic social dynamics: discrete time steps, personalized recommendation, and diverse user personas. By collecting feedback *after* propagation, we obtain rewards approximating real-world long-term value.

Our framework, MARFE (Multi-Agent Reward-Free Evolution), iteratively improves content by:

(1) generating output, (2) injecting it into OASIS simulation, (3) letting diverse users react over multiple time steps, and (4) summarizing delayed feedback to evolve the agent’s prompt. Evaluated across four tasks with three frontier LLM judges, MARFE achieves 58.3% win rate versus 41.7% for self-refine, demonstrating that delayed social feedback provides superior signal for non-verifiable content optimization.

2 Method

2.1 Problem Formulation

For non-verifiable tasks, the reward $R^*(o)$ for output o is not computable from o alone, but emerges from simulating o ’s propagation through a social system \mathcal{S} over time horizon T :

$$\hat{R}(o) = \text{Simulate}(o, \mathcal{S}, T) \quad (1)$$

where \mathcal{S} includes agents with diverse personas, a recommendation system, and temporal dynamics.

2.2 OASIS Simulation Environment

We use OASIS (Yang et al., 2024) (Open Agent Social Interaction Simulations), a scalable open-source social media simulator that combines LLM-powered agents with rule-based mechanisms to realistically model user behavior on platforms like Twitter and Reddit.

Discrete Time Steps. The simulation engine orchestrates time progression: each `env.step()` advances the clock, refreshes recommendation feeds, activates agents for observation and decision-making, and updates the environment based on executed actions.

Personalized Recommendation. OASIS implements interest-based recommendation using Twitter’s TWHINBert (Zhang et al., 2023). For each user u and post p , the score combines embedding similarity with logarithmic time decay:

$$\text{score}(p, u) = \cos(\mathbf{e}_p, \mathbf{e}_u) \cdot \log \frac{271.8 - \Delta t}{100} \quad (2)$$

where $\mathbf{e}_p, \mathbf{e}_u$ are embeddings of post content and user bio, and Δt is the time step difference. This models how older content receives diminishing visibility.

Diverse Action Space. Agents can perform 23 distinct operations including LIKE, REPOST, COMMENT, FOLLOW, and DO_NOTHING, enabling rich multi-faceted interactions that mirror real platform affordances.

LLM-Driven Agents. Each agent is a virtual user powered by LLMs, with unique profiles and decision processes. Agents observe their personalized feeds and use LLM reasoning to decide actions, creating authentic reaction patterns based on their personas.

2.3 MARFE Framework

Our **Multi-Agent Reward-Free Evolution (MARFE)** framework:

Algorithm 1 MARFE: Simulation-Grounded Evolution

```

1: Input: Task  $\tau$ , iterations  $N$ , sim steps  $S$ , agents  $K$ 
2: Initialize worker agent with base prompt  $P_0$ 
3: for  $i = 1$  to  $N$  do
4:    $o_i \leftarrow$  Worker.generate( $\tau, P_{i-1}$ )
5:   Inject  $o_i$  into OASIS as post
6:   for  $s = 1$  to  $S$  do
7:      $K$  persona agents observe and act (LLM-driven)
8:   end for
9:   Extract feedback: likes, comments, reposts
10:   $E_i \leftarrow$  Summarize(feedback)
11:   $P_i \leftarrow P_{i-1} + E_i$  (append experience)
12: end for
13: Output: Final content  $o_N$ 

```

The key innovation is that feedback comes from *delayed, emergent social dynamics* rather than instant self-critique.

3 Experiments

3.1 Setup

Tasks. We evaluate on four non-verifiable tasks: (1) **Marketing:** Twitter copy for ChatGPT Plus targeting developers and professionals; (2) **Research Sharing:** Twitter thread about the Transformer paper; (3) **Policy:** AI regulation proposal balancing safety and innovation; (4) **Product:** Vision Pro developer announcement.

Methods. Both methods use identical base prompts and 3 iterations:

- **Baseline:** Single-agent iterative refinement inspired by Self-Refine (Madaan et al., 2023). The agent generates content, then critiques and improves its own output each iteration using the same LLM, without external feedback.

Task	GPT-5.2	Gemini 3 pro	Claude 4.5 Sonnet
Marketing	MARFE	MARFE	MARFE
Research	MARFE	Base	MARFE
Policy	MARFE	Base	Base
Product	MARFE	Base	Base

Table 1: Pairwise winners by task and judge. MARFE wins 7/12.

- **MARFE:** Injects content into OASIS with 5 user agents over 2 simulation steps, collecting delayed social feedback to guide evolution.

Evaluation. Following LMSys Arena methodology (Zheng et al., 2023), we use pairwise comparison with three frontier LLM judges (GPT-5.2, Gemini 3 Pro, Claude 4.5 Sonnet). Judges select winners based on persuasiveness, clarity, audience targeting, and engagement potential. This yields 12 total comparisons (4 tasks \times 3 judges).

3.2 Results

Table 1 shows pairwise evaluation results.

Method	W	L	Rate
Baseline	5	7	41.7%
MARFE (Ours)	7	5	58.3%

Table 2: Aggregate results (W=wins, L=losses).

Key Findings. MARFE achieves 58.3% win rate, outperforming Self-Refine baseline (41.7%). We report results from a single experimental run with 3 iterations per method. Notably, MARFE wins unanimously on Marketing (3/3 judges), suggesting social simulation feedback is particularly effective for audience-targeted content. Performance varies by task: MARFE excels on Marketing and Research but shows mixed results on Policy and Product, indicating task-dependent effectiveness.

3.3 Analysis

The simulation trajectories reveal concrete insights unavailable to Self-Refine:

Specific Objections Surfaced. In the Marketing task, simulated users raised “ad fatigue” (“Ugh, another ChatGPT ad”) and noted enterprise alternatives (“many companies provide team licenses”). By iteration 2, MARFE’s content directly addressed these: leading with pain points (“Tired of hitting ‘at capacity’ at 3 PM?”) rather than generic promotion, and offering concrete resources (“DM

230	me for 3 free custom GPTs”). The baseline, lacking this feedback, continued with generic improvements like “Join 500,000+ professionals.”	279
231		280
232		
233	Audience Mismatch Detection. Simulation revealed content reaching unintended audiences—	281
234	TikTok dancers commenting “I’m just here to	282
235	dance, not think about productivity tools.” This	283
236	prompted MARFE to add professional identity	284
237	signals (#DeveloperTools) and segment-specific	285
238	hooks. Self-Refine cannot detect such mismatches.	286
239		287
240	Actionable Structural Feedback. For Policy	288
241	content, users praised the “tiered, risk-based frame-	
242	work” as providing a “reusable mental model,” but	289
243	identified a gap: lack of “public/stakeholder in-	290
244	put mechanisms.” MARFE incorporated this in	291
245	subsequent iterations. For Product content, users	292
246	suggested “lead with the ‘so what’ for non-devs”	293
247	and use analogies like “Google Docs for the 3D	294
248	world.”	295
249	Engagement Pattern Learning. Across	
250	tasks, simulation revealed that ROI framing	$R(o) = \alpha \cdot \text{likes} + \beta \cdot \text{reposts} + \gamma \cdot \text{comments}$ (3)
251	(“\$20/month = 2 hours saved”), question hooks,	296
252	and debate-inviting content drove engagement.	
253	These patterns—discovered through delayed	Why Delayed Feedback Helps. Simulation re-
254	propagation—informed prompt evolution in ways	297
255	Self-Refine’s generic “add more hooks” advice can-	298
256	not.	299
257		300
258	4 Related Work	
259	Self-Improvement. Self-Refine (Madaan et al.,	Beyond Prompt Evolution. The paradigm ex-
260	2023) and Reflexion (Shinn et al., 2023) use self-	301
261	critique but lack external grounding. Our work	302
262	provides social simulation as external feedback.	303
263	Social Simulation. Generative Agents (Park	304
264	et al., 2023) and OASIS (Yang et al., 2024) simulate	305
265	realistic social dynamics. We repurpose these as	
266	reward derivation mechanisms for agent improve-	6 Conclusion
267	ment.	306
268	Reward Modeling. RLHF (Ouyang et al., 2022)	Non-verifiable tasks <i>appear</i> to have instant feed-
269	uses human preferences but is costly. RLAIIF (Lee	307
270	et al., 2023) provides instant rather than delayed	308
271	rewards. Our approach derives delayed rewards	309
272	from simulated social dynamics.	310
273		311
274	5 Discussion	312
275	The True Nature of Non-Verifiable Tasks. Non-	313
276	verifiable tasks <i>appear</i> to have instant feedback but	314
277	are <i>inherently delayed</i> when considering users’ ac-	315
278	tual goals. A user asking for a novel can immedi-	316
	ately judge quality, but their true objective is read-	317
	ership and influence over months. Instant feedback	
		Limitations
		318
		We evaluate on 4 tasks with 3 iterations each, which
		319
		320
		321
		322
		323
		324

325	networks) remain future work. LLM-simulated	limits of mathematical reasoning in open language	378
326	personas may not perfectly mimic human behav-	models. <i>Preprint</i> , arXiv:2402.03300.	379
327	ior. Simulation is more computationally expensive	Noah Shinn, Federico Cassano, Edward Berman, Ash-	380
328	than self-refinement. Future work should explore	win Gopinath, Karthik Narasimhan, and Shunyu Yao.	381
329	diverse simulation environments, human studies to	2023. Reflexion: Language agents with verbal rein-	382
330	validate fidelity, and RL-based training with simu-	forcement learning . <i>Preprint</i> , arXiv:2303.11366.	383
331	lation rewards.	Sijing Tu and Stefan Neumann. 2022. A viral marketing-	384
332	References	based model for opinion dynamics in online social	385
333	Yilun Du, Shuang Li, Antonio Torralba, Joshua B.	networks . <i>Preprint</i> , arXiv:2202.03573.	386
334	Tenenbaum, and Igor Mordatch. 2023. Improving	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	387
335	factuality and reasoning in language models through	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	388
336	multiagent debate . <i>Preprint</i> , arXiv:2305.14325.	Denny Zhou. 2022. Chain-of-thought prompting elic-	389
337	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas	its reasoning in large language models . <i>Preprint</i> ,	390
338	Mesnard, Johan Ferret, Kellie Lu, Colton Bishop,	arXiv:2201.11903.	391
339	Ethan Hall, Victor Carbune, Abhinav Rastogi, and	Ziyi Yang, Zaibin Zhang, Zirui Zheng, Yuxian Jiang,	392
340	Sushant Prakash. 2023. Rlaif vs. rlhf: Scaling re-	Ziyue Gan, Zhiyu Wang, Zijian Ling, Jinsong Chen,	393
341	inforcement learning from human feedback with ai	Martz Ma, Bowen Dong, Prateek Gupta, Shuyue	394
342	feedback . <i>Preprint</i> , arXiv:2309.00267.	Hu, Zhenfei Yin, Guohao Li, Xu Jia, Lijun Wang,	395
343	Guohao Li, Hasan Abed Al Kader Hammoud, Hani	Bernard Ghanem, Huchuan Lu, Chaochao Lu, and	396
344	Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023.	4 others. 2024. Oasis: Open agent social interac-	397
345	Camel: Communicative agents for "mind" explo-	tion simulations with one million agents . <i>Preprint</i> ,	398
346	ration of large language model society . <i>Preprint</i> ,	arXiv:2411.11581.	399
347	arXiv:2303.17760.	Xinyang Zhang, Yury Malkov, Omar Florez, Serim	400
348	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Park, Brian McWilliams, Jiawei Han, and Ahmed	401
349	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	El-Kishky. 2023. Twin-bert: A socially-enriched	402
350	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	pre-trained language model for multilingual tweet	403
351	Shashank Gupta, Bodhisattwa Prasad Majumder,	representations . <i>Preprint</i> , arXiv:2209.07562.	404
352	Katherine Hermann, Sean Welleck, Amir Yazdan-	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	405
353	bakhsh, and Peter Clark. 2023. Self-refine: It-	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	406
354	erative refinement with self-feedback . <i>Preprint</i> ,	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	407
355	arXiv:2303.17651.	Joseph E. Gonzalez, and Ion Stoica. 2023. Judg-	408
356	OpenAI. 2023. Gpt-4 technical report . <i>Preprint</i> ,	ing llm-as-a-judge with mt-bench and chatbot arena .	409
357	arXiv:2303.08774.	<i>Preprint</i> , arXiv:2306.05685.	410
358	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-		
359	roll L. Wainwright, Pamela Mishkin, Chong Zhang,		
360	Sandhini Agarwal, Katarina Slama, Alex Ray, John		
361	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,		
362	Maddie Simens, Amanda Askell, Peter Welinder,		
363	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.		
364	Training language models to follow instructions with		
365	human feedback . <i>Preprint</i> , arXiv:2203.02155.		
366	Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai,		
367	Meredith Ringel Morris, Percy Liang, and Michael S.		
368	Bernstein. 2023. Generative agents: Interac-		
369	tive simulacra of human behavior . <i>Preprint</i> ,		
370	arXiv:2304.03442.		
371	John Schulman, Filip Wolski, Prafulla Dhariwal,		
372	Alec Radford, and Oleg Klimov. 2017. Prox-		
373	imal policy optimization algorithms . <i>Preprint</i> ,		
374	arXiv:1707.06347.		
375	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,		
376	Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu,		
377	and Daya Guo. 2024. Deepseekmath: Pushing the		