
AGENTICPAY: A MULTI-AGENT LLM NEGOTIATION SYSTEM FOR BUYER–SELLER TRANSACTIONS

Xianyang Liu Shangding Gu* Dawn Song
University of California, Berkeley

ABSTRACT

Large language model (LLM)–based agents are increasingly expected to negotiate, coordinate, and transact autonomously, yet existing benchmarks lack principled settings for evaluating language-mediated economic interaction among multiple agents. We introduce **AgenticPay**, a benchmark and simulation framework for multi-agent buyer–seller negotiation driven by natural language. **AgenticPay** models markets in which buyers and sellers possess private constraints and product-dependent valuations, and must reach agreements through multi-round linguistic negotiation rather than numeric bidding alone. The framework supports a diverse suite of over 110 tasks ranging from bilateral bargaining to many-to-many markets, with structured action extraction and metrics for feasibility, efficiency, and welfare. Benchmarking state-of-the-art proprietary and open-weight LLMs reveals substantial gaps in negotiation performance and highlights challenges in long-horizon strategic reasoning, establishing **AgenticPay** as a foundation for studying agentic commerce and language-based market interaction. Code and dataset are available at the link: <https://github.com/SafeRL-Lab/AgenticPay>.

1 INTRODUCTION

Large language models (LLMs) have shown remarkable performance in many domains (Comanici et al., 2025; Hurst et al., 2024; OpenAI, 2025; Gu et al., 2024; Yang et al., 2025), and are increasingly deployed as autonomous agents that need to coordinate and transact on behalf of users in economic settings such as e-commerce, procurement, and service contracting. Unlike traditional decision-making systems that operate over structured bids or fixed utility functions, these agents interact through natural language, expressing preferences, constraints, and counteroffers in multi-turn dialogues. As a result, negotiation becomes a language-mediated strategic interaction, where outcomes depend jointly on reasoning, communication, and long-horizon planning.

Despite rapid progress in LLM capabilities, existing benchmarks for agent evaluation remain limited in their ability to capture this setting (Xia et al., 2024; He et al., 2018; Fu et al., 2023). Most prior work evaluates single-agent reasoning (Mondorf & Plank, 2024; Gu et al., 2025), tool use (Chen et al., 2025), or preference following (Sun et al., 2025b), and economic interaction is often simplified to numeric auctions or short-horizon bargaining (Chen et al. (2023); He et al. (2018); Fu et al. (2023)). These abstractions fail to reflect key properties of real-world transactions: private reservation values, multi-round negotiation, heterogeneous products, and competition among multiple buyers and sellers. Consequently, it remains unclear: *How effectively can current LLMs function as autonomous negotiators in diverse market environments?*

In this work, we introduce **AgenticPay**, a benchmark and simulation framework for studying multi-agent buyer–seller negotiation driven by natural language, spanning settings from bilateral bargaining to many-to-many markets. **AgenticPay** models markets in which buyers and sellers possess private constraints and product-dependent valuations, and must reach agreements through iterative linguistic negotiation rather than numeric bidding alone. Negotiation is formalized as a language game, with dialogue histories mapped to actions such as price proposals and deal acceptance, enabling principled evaluation of negotiation outcomes.

*Corresponding author: shangding.gu@berkeley.edu

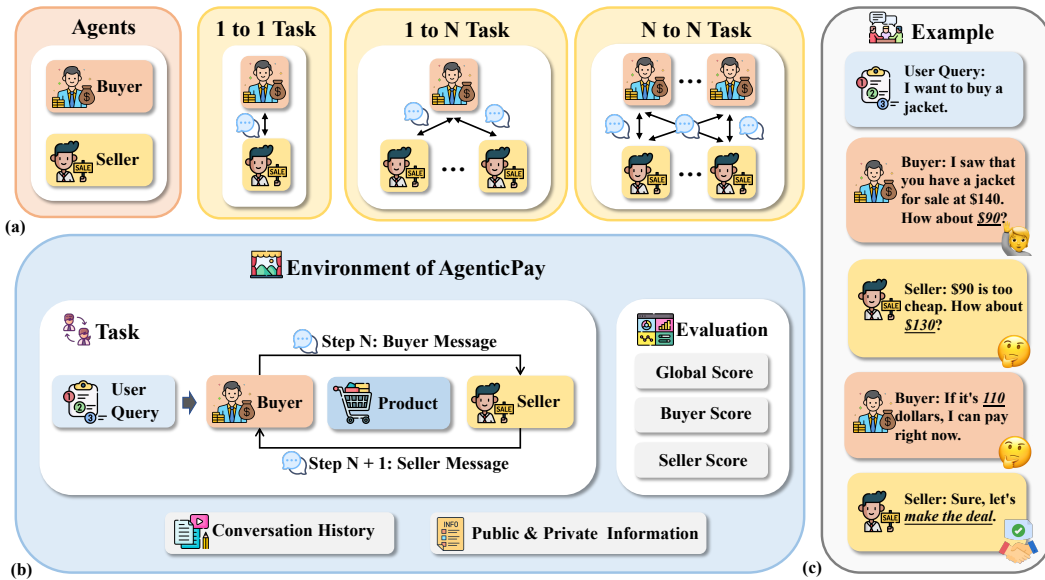


Figure 1: **Overview of AgenticPay.** (a) **Agents & Task Examples:** Buyer and seller agents engage in three negotiation modes: *1-to-1* (bilateral bargaining between a single buyer and seller), *1-to-N* (one buyer negotiating with multiple competing sellers, or one seller negotiating with multiple competing buyers), and *N-to-N* (many buyers and sellers forming a matching market). (b) **Framework:** Core components including Environment, Task, and Agent interact to enable multi-round negotiations. (c) **Dialogue Example:** A sample negotiation showing the user’s product requirements, buyer–seller conversation, and final deal.

AgenticPay provides a comprehensive suite of tasks that scale market complexity along three dimensions: the number of buyers, the number of sellers, and the size of the product set. Tasks range from bilateral bargaining to many-to-many markets with competing agents and multiple products, supporting both sequential and parallel negotiation regimes. To evaluate performance, we introduce metrics that jointly capture deal feasibility, efficiency, and welfare for buyers, sellers, and the market as a whole.

Using AgenticPay, we benchmark a diverse set of state-of-the-art proprietary and open-weight LLMs under a unified inference-only protocol. Our results reveal substantial performance gaps across models, systematic asymmetries between buyer and seller roles, and persistent challenges in long-horizon strategic reasoning. These findings highlight that strong language generation alone is insufficient for effective economic negotiation.

Overall, AgenticPay establishes a foundation for studying agentic commerce, offering a controlled yet expressive testbed for research on multi-agent negotiation, economic alignment, and the co-evolution of language and strategy in autonomous agents. Our Contributions are summarized as follows:

- We introduce AgenticPay, a scalable framework that supports a large number of tasks (over 110) ranging from bilateral bargaining to many-to-many markets, with dialogue-to-action grounding and welfare-oriented evaluation metrics. The system supports diverse deployment via vLLM¹, SGLang², and cloud-based LLM APIs.
- We formalize language-mediated buyer–seller negotiation as a multi-agent game with private reservation values and dialogue-grounded economic outcomes. Moreover, we benchmark state-of-the-art proprietary and open-weight LLMs, uncovering persistent limitations in long-horizon strategic reasoning and negotiation efficiency.

¹<https://github.com/vllm-project/vllm>

²<https://github.com/sgl-project/sglang>

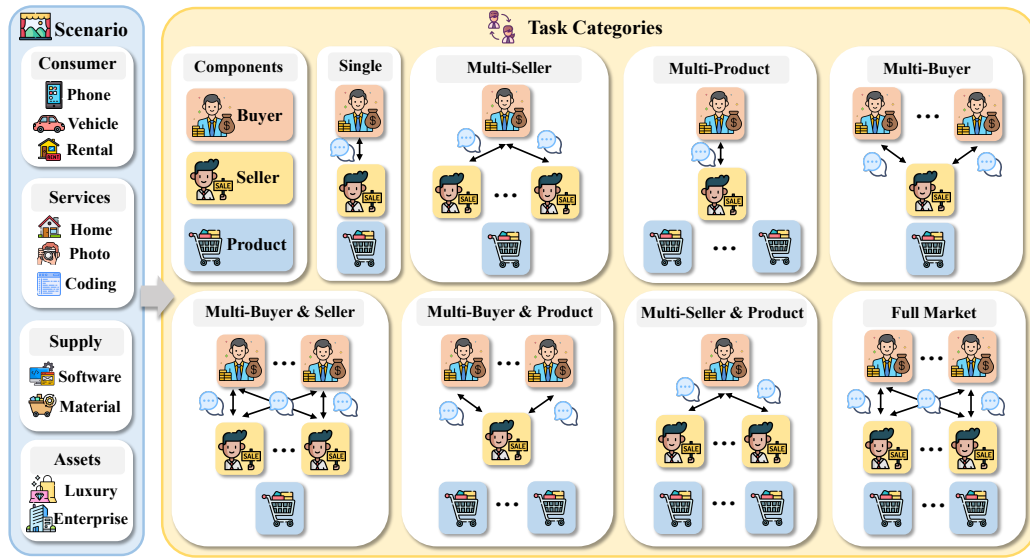


Figure 2: **Overview of the AgenticPay task suite.** **Left:** Ten realistic business scenarios across four categories: Consumer, Services, Supply, and Assets. **Right:** Task categories illustrating the progression from bilateral bargaining to full market settings along three complexity dimensions: number of buyers, number of sellers, and product set size.

2 RELATED WORK

Negotiation and Bargaining in Game Theory. Classical work in economics and game theory has studied bargaining and bilateral trade under incomplete information, establishing foundational results on efficiency, equilibrium, and impossibility theorems (Chatterjee & Samuelson, 1983; Myerson & Satterthwaite, 1983; Rubinstein, 1985; Ausubel et al., 2002; Blumrosen & Mizrahi, 2016). These models typically assume agents interact through scalar bids or utilities, with negotiation dynamics defined over numeric strategy spaces. While analytically tractable, such formulations abstract away the role of language and may not capture the rich communicative strategies present in real-world negotiations. More recently, a growing body of work has explored the game-theoretic behavior of large language models (Lorè & Heydari, 2023; Fan et al., 2024; Hu et al., 2024; Raman et al., 2024; Silva, 2024; Lorè & Heydari, 2024; Jia et al., 2025; Lu, 2025; Akata et al., 2025; Sun et al., 2025a). However, these studies primarily examine strategic reasoning or equilibrium behavior in games, and do not explicitly address language-mediated negotiation or market-based economic interaction.

Neural and Dialogue-Based Negotiation. Prior work in natural language processing has explored negotiation as a dialogue task, focusing on learning strategies for offer generation, concession planning, and agreement formation (Lewis et al., 2017; He et al., 2018; Chawla et al., 2021; Joshi et al., 2021; Pacella & Marocco, 2022; Hua et al., 2024; Washio et al., 2026). More recent approaches leverage LLMs and self-play to improve negotiation behavior (Ma et al., 2024; Chen et al., 2024; Long et al., 2025), often using in-context learning or reinforcement signals derived from dialogue outcomes (Fu et al., 2023; Vahidov et al., 2025; Priya et al., 2025). These methods typically consider bilateral settings with fixed roles and limited environment structure, and may not address market-level interactions involving multiple buyers, sellers, and products.

LLMs in Economic and Auction Settings. Several recent benchmarks study LLMs in auction-like environments (Chen et al., 2023; Duetting et al., 2024; Shah et al., 2025; Agrawal et al., 2025), evaluating strategic planning and execution when agents submit bids or allocate resources (Chen et al., 2023). Related work has also proposed benchmarks for measuring bargaining abilities of LLMs, often focusing on single buyer-seller interactions or buyer-enhancement techniques (Xia et al., 2024; Jiang et al., 2025). While these efforts demonstrate that LLMs exhibit non-trivial

strategic behavior, they largely rely on simplified mechanisms or short-horizon interactions and may not model general multi-agent markets with private constraints and heterogeneous products.

The most closely related works to ours are (Zhu et al., 2025; Deng et al., 2024; Bhattacharya et al., 2025), which study LLMs as game-theoretic negotiators and primarily focus on negotiation risk and strategy in restricted settings. In contrast, our work provides a comprehensive and extensible framework for language-mediated negotiation, supporting multi-product interactions and diverse market scenarios. For example, while Deng et al. (2024) focuses on bilateral negotiation, our framework establishes a foundation not only for bilateral bargaining but also for many-to-many negotiation involving multiple buyers and sellers. Bianchi et al. (2024) also investigate LLM-based negotiation, but their evaluation is mostly limited to three predefined scenarios. By comparison, our framework supports at least ten realistic business scenarios and over 110 tasks, and is inherently scalable due to its modular environment, task, and agent interfaces. This design enables systematic expansion to new market configurations without altering the core protocol.

AgentPay differs from prior work in three key respects. First, it models negotiation as a language-grounded market interaction, where dialogue directly determines structured economic outcomes. Second, it scales beyond bilateral bargaining to many-to-many markets with competition, matching, multiple products, and diverse scenarios. Third, it provides principled evaluation metrics grounded in feasibility, efficiency, and welfare, enabling systematic comparison across models and settings. By unifying ideas from economic theory, multi-agent systems, and language modeling, **AgentPay** fills a critical gap in the evaluation of autonomous LLM-based negotiators.

3 PROBLEM SETTINGS

We consider a language-mediated buyer–seller market populated by a set of buyers $\mathcal{B} = \{1, \dots, N_B\}$ and a set of sellers $\mathcal{S} = \{1, \dots, N_S\}$. Each buyer seeks to purchase a product or service from a seller, while each seller offers products subject to its own pricing constraints. Unlike classical auction or matching settings that rely solely on numeric bids, agents in our setting negotiate through multi-turn natural language dialogue to propose offers, express constraints, and reach agreements.

Agent States. Each buyer $i \in \mathcal{B}$ and seller $j \in \mathcal{S}$ is associated with a private internal state that governs its negotiation behavior. We denote the buyer state by b_i and the seller state by σ_j . The buyer state b_i encodes buyer-specific information such as preferences, budget constraints, and willingness-to-pay, while the seller state σ_j encodes seller-specific information such as costs, reservation prices, and pricing policies. These states are private to each agent and are not observable by other agents during negotiation.

Product and Market Context. Each seller j offers a product or service represented by a feature vector $v_j \in \mathcal{V}$. The product representation v_j captures observable attributes of the product, including structured features (e.g., category or specifications) and unstructured textual descriptions, and is public information available to both buyers and sellers. Negotiation takes place under a shared market context $x \in \mathcal{X}$, which represents external factors such as market category, seasonal effects, or domain-specific rules. The context x is shared by all agents and is not controlled by any individual buyer or seller.

Dialogue-Based Negotiation. Negotiation between a buyer i and a seller j is modeled as a finite-horizon, multi-round language game. At each round t , agents alternately exchange natural-language messages conditioned on their private states, the product representation, the shared market context, and the dialogue history. Let $h_{ij}^{(t)}$ denote the public dialogue history up to round t . A buyer policy π_i^B and a seller policy π_j^S specify the agents’ negotiation strategies. At round t , the buyer generates a message $m_{ij}^{(t,B)} \sim \pi_i^B(\cdot \mid b_i, v_j, x, h_{ij}^{(t-1)})$, and the seller responds with $m_{ij}^{(t,S)} \sim \pi_j^S(\cdot \mid \sigma_j, v_j, x, h_{ij}^{(t-1)})$. Messages may contain free-form language as well as structured signals such as price proposals.

Action Parsing and Termination. A parser Π maps the exchanged messages at each round to a structured negotiation action $a_{ij}^{(t)} = \Pi(m_{ij}^{(t,B)}, m_{ij}^{(t,S)})$, which extracts quantities such as the proposed transaction price $p_{ij}^{(t)}$. Negotiation terminates when agents reach agreement, exceed the maximum negotiation horizon, or violate feasibility constraints.

Evaluation Objectives. Agents prefer higher surplus and earlier agreement. Negotiation outcomes are evaluated based on deal feasibility, efficiency, and welfare, reflecting the quality of the final price and the speed of convergence.

4 AGENTICPAY

Building on the problem formulation in Section 3, this section presents **AgenticPay**, a benchmark that instantiates language-grounded markets within a controlled experimental setting. As illustrated in Figure 1, the framework comprises four components: **Environments** (Section 4.1) implement the negotiation protocol and domain-specific scenarios; **Tasks** (Section 4.2) operationalize market structures with varying numbers of buyers, sellers, and interaction modes; **Agents** (Section 4.3) instantiate LLM-based policies with private valuations and dialogue memory; and **Metrics** (Section 4.4) quantify negotiation outcomes in terms of deal rate, surplus allocation, and efficiency.

4.1 ENVIRONMENT

Negotiation Protocol. Each negotiation episode in **AgenticPay** is a finite-horizon, multi-round interaction between a buyer and a seller over a product, as illustrated in Figure 1 (b). The environment provides each agent with the product description and market context (e.g., category or conditions). Each agent also receives a private reservation price: the buyer’s maximum willingness-to-pay and the seller’s minimum acceptable price, neither of which is revealed to the counterpart. Negotiation proceeds in alternating rounds up to a maximum number of turns, where each party generates a natural-language message containing an explicit price proposal. A deal is reached when both parties propose the same price, and a transaction is valid only if the agreed price lies within the bargaining zone.

Scenario Design. To capture the diversity of real-world negotiation contexts, **AgenticPay** includes 10 realistic business scenarios organized into four economic domains, as illustrated in Figure 2 (left): (1) **Daily Life:** Used Smartphone, Used Car, Vacation Rental; (2) **Professional Services:** Website Development, Commercial Photography, Home Renovation; (3) **Business Procurement:** SaaS Software, Raw Materials; (4) **Financial Assets:** Luxury Watch, Business Acquisition. This diversity enables evaluation of whether agent policies generalize across domains with varying negotiation conventions and linguistic styles.

4.2 TASKS

Task Definition. A *task* in **AgenticPay** specifies the market structure imposed on a negotiation episode, determining the number of buyers, sellers, and products involved. While the environment defines the domain-specific scenario (e.g., used car or SaaS software), the task governs the competitive and combinatorial complexity of the interaction. This separation enables systematic evaluation: the same scenario can be instantiated under different task configurations, isolating the effect of market structure from domain-specific negotiation conventions.

Task Categories. As illustrated in Figure 2 (Right), **AgenticPay** comprises eight task categories that systematically scale complexity along three dimensions: number of buyers, number of sellers, and product set size. These range from *bilateral price negotiation* (1 buyer, 1 seller, 1 product) to *full market settings* (multiple buyers, sellers, and products). Intermediate configurations include *multi-item bargaining* with a fixed counterpart, *buyer competition* (multiple buyers competing for one seller’s product), *seller competition* (one buyer choosing among competing sellers), and various many-to-many markets. Across these categories, we study two interaction modes. In parallel interaction, an agent reasons over multiple ongoing negotiations simultaneously. In sequential interaction, the agent adaptively decides whether to continue, switch, or commit. This design forms

a complexity ladder that isolates distinct challenges in language-based negotiation while supporting systematic benchmarking across increasingly realistic market settings.

4.3 AGENTS

Our agent framework instantiates buyers and sellers as role-specialized negotiators that share a unified architecture but differ in their private valuations and objectives.

Environment Public Information. Each agent receives public context from the environment, including the product description (attributes, quality, and features), scenario metadata (domain category and market conditions), and the negotiation protocol (maximum rounds and output format requirements). This shared information grounds the dialogue in a common understanding of the transaction.

Role-Based Private Information. To preserve asymmetric information, each agent holds role-specific private valuations: buyers are assigned a maximum willingness-to-pay p^{\max} , while sellers maintain a minimum acceptable price p^{\min} . These reservation prices are injected into the agent’s system prompt but excluded from the shared dialogue, with explicit instructions to keep them confidential.

Dialogue History. Each agent maintains an independent memory module that records the full sequence of multi-turn exchanges as $(role, content, round)$ tuples. Agents condition their responses on this history, enabling coherent multi-round reasoning. In multi-party settings, separate memory instances ensure each agent’s view remains consistent with its participation.

4.4 METRICS

We evaluate negotiation performance using three complementary outcome scores: **GlobalScore**, **BuyerScore**, and **SellerScore**. GlobalScore measures overall deal quality by rewarding balanced outcomes where both parties benefit, while BuyerScore and SellerScore capture role-specific utility, reflecting each agent’s individual gain from the transaction. All three metrics incorporate negotiation efficiency, incentivizing faster agreements.

Score Design. As detailed in Algorithm 1, our scoring framework normalizes utilities within the bargaining zone $Z = p^{\max} - p^{\min}$. GlobalScore uses a

symmetric quality term $Q = 4r_b r_s$ that peaks when surplus is split equally, while BuyerScore and SellerScore reward each party’s individual surplus. Key design choices include: setting $W > D$ to prioritize deal quality over mere agreement, using discount factor γ to incentivize faster deals, and applying a moderate failure penalty F to discourage deadlock without inducing excessive risk aversion. We additionally report deal rate and average rounds as auxiliary statistics.

Algorithm 1 Score Calculation for Negotiation Outcomes

Require: Final price p ; buyer’s max price p^{\max} ; seller’s min price p^{\min} ; deal round t ; max rounds T ; discount factor γ ; deal success reward D , deal quality reward W , round efficiency reward E , failure penalty F

- 1: $Z \leftarrow p^{\max} - p^{\min}$ # Bargaining zone
- 2: $d \leftarrow \gamma^{t-1}$ # Efficiency discount
- 3: **if** $p^{\min} \leq p \leq p^{\max}$ **then**
- 4: $r_b \leftarrow (p^{\max} - p)/Z$ # Buyer utility $\in [0, 1]$
- 5: $r_s \leftarrow (p - p^{\min})/Z$ # Seller utility $\in [0, 1]$
- 6: $Q \leftarrow 4 \cdot r_b \cdot r_s$ # Quality term $\in [0, 1]$
- 7: $S_g \leftarrow d \cdot (D + Q \cdot W + E)$ # GlobalScore
- 8: $S_b \leftarrow d \cdot (D + r_b \cdot W + E)$ # BuyerScore
- 9: $S_s \leftarrow d \cdot (D + r_s \cdot W + E)$ # SellerScore
- 10: **else**
- 11: $d \leftarrow \gamma^{T-1}$ # Use max rounds for failure
- 12: $S_g, S_b, S_s \leftarrow -F \cdot (1 - d)$ # Failure penalty scores
- 13: **end if**
- 14: **return** S_g, S_b, S_s

5 EXPERIMENTS AND ANALYSIS

5.1 EXPERIMENTAL SETUP

Benchmark Statistics AgenticPay comprises 111 negotiation tasks across 8 multi-agent configurations, including 31 basic tasks for core mechanics and 80 realistic tasks from 10 business scenarios. Product values range from \$350 to \$120k (see Section A).

Table 1: Overall performance on **AgenticPay** across all 111 tasks. We report mean scores aggregated over episodes, along with deal rate, timeout rate, price overflow rate (instances where agents propose prices outside acceptable bounds), and average termination round. Colors indicate performance levels, from **green (lowest)** to **red (highest)**.

Model	GlobalScore	SellerScore	BuyerScore	Deal Rate	Timeout Rate	Overflow Rate	Avg. Rounds
Claude Opus 4.5	86.9	76.1	63.5	100.0%	0.0%	0.0%	3.7
Gemini-3-Flash	82.2	73.3	61.1	100.0%	0.0%	2.7%	4.8
GPT-5.2	81.7	81.1	58.5	100.0%	0.0%	0.0%	3.8
Qwen3-14B	63.9	58.9	47.6	79.3%	20.7%	1.8%	7.8
Llama-3.1-8B	32.5	26.3	25.2	51.4%	48.6%	10.8%	15.0

Table 2: Performance breakdown by buyer-seller multiplicity. 1B1S = Single-Buyer-Single-Seller, 1BMS = Single-Buyer-Multi-Seller, MB1S = Multi-Buyer-Single-Seller, MBMS = Multi-Buyer-Multi-Seller. All scores are reported as percentages. Colors indicate performance levels, from **green (lowest)** to **red (highest)**.

Model	1B1S			1BMS			MB1S			MBMS		
	Global	Seller	Buyer	Global	Seller	Buyer	Global	Seller	Buyer	Global	Seller	Buyer
Claude Opus 4.5	83.4	77.0	62.9	85.7	78.8	60.3	88.5	74.5	64.8	89.8	74.2	66.2
Gemini-3-Flash	77.5	69.0	59.5	80.0	73.4	58.6	83.9	77.7	59.4	87.4	72.7	66.9
GPT-5.2	79.1	81.2	58.5	82.6	81.7	57.9	80.9	81.3	57.9	84.0	80.2	59.5
Qwen3-14B	63.2	59.6	47.4	50.1	51.8	35.7	64.7	59.4	47.6	77.6	65.0	59.5
Llama-3.1-8B	27.9	20.3	23.0	28.2	22.1	20.5	36.2	29.8	28.5	37.5	32.8	28.9

Inference We evaluate **AgenticPay** under a unified inference protocol to ensure fair comparison across models and configurations. All agents use deterministic decoding with temperature 0 and random seed 0, with a maximum generation length of 1024 tokens per response. For open-source models, we perform inference using 4 NVIDIA A800 GPUs. Both open-source and closed-source models receive identical prompts; the prompt templates for buyer and seller agents are detailed in Table 14 and Table 15, respectively. Each task instance is executed once per model.

Models We benchmark a diverse set of proprietary and open-weight LLMs as negotiation policies. Our proprietary-model evaluation includes **GPT-5.2** (OpenAI, 2025), **Claude Opus 4.5** (2025-11-01) (Anthropic, 2025), and **Gemini 3 Flash** (Google DeepMind, 2025). To assess the transferability of negotiation capabilities to smaller open models, we additionally evaluate **Qwen3-14B** (Qwen, 2025) and **Llama-3.1-8B** (Meta, 2024). Unless otherwise specified, each model is used as a drop-in policy for both buyer and seller roles under the same environment protocol and decoding configuration.

Evaluation Metrics Based on Algorithm 1, we configure the following parameters to calculate GlobalScore, SellerScore, and BuyerScore: deal completion bonus D is 30, quality bonus W is 55, efficiency bonus E is 15, discount factor γ is 0.99, and failure penalty F is 15. The maximum number of negotiation rounds is set to 20.

5.2 MAIN RESULTS

Table 1 summarizes the overall negotiation performance of all evaluated models across the full **AgenticPay** benchmark. We highlight several key findings from these results.

Proprietary Models Dominate Negotiation Performance. Claude Opus 4.5 achieves the highest GlobalScore of 86.9, followed closely by Gemini-3-Flash (82.2) and GPT-5.2 (81.7), all maintaining perfect 100% deal rates with zero timeouts. In contrast, open-weight models exhibit substantial performance gaps: Qwen3-14B achieves only 63.9 GlobalScore with a 20.7% timeout rate, while Llama-3.1-8B struggles significantly with a GlobalScore of 32.5 and nearly half of negotiations (48.6%) ending in timeout. The price overflow rate, which indicates instances where agents propose prices outside acceptable bounds, further distinguishes model reliability: proprietary models maintain near-zero overflow rates, whereas Llama-3.1-8B exhibits 10.8% overflow, suggesting difficulties in adhering to negotiation constraints. See Section C for example dialogues.

Table 3: GlobalScore by scenario category. Each cell shows the average GlobalScore for the model across all scenarios within that category. Colors indicate performance levels, from **green (lowest)** to **red (highest)**. Detailed results can be found in Table 13.

Model	Prof. Services	Daily Life	Bus. Procurement	Financial Assets	Avg
Claude Opus 4.5	93.4	90.7	89.6	85.7	90.3
GPT-5.2	89.8	83.8	86.1	79.9	85.3
Gemini-3-Flash	88.3	86.8	85.0	68.1	83.1
Qwen3-14B	72.5	65.1	69.1	60.9	67.3
Llama-3.1-8B	41.1	38.3	18.1	39.9	35.4

Table 4: Cross-play performance analysis in 1B1P1S scenario. GlobalScore represents the average across all interactions involving the model. SellerScore and BuyerScore indicate performance when the model acts as seller or buyer, respectively. Deal Rate indicates the percentage of successful negotiations. Colors indicate performance levels, from **green (lowest)** to **red (highest)**.

Model	Cross-Play			
	Global	Seller	Buyer	Deal Rate
Claude Opus 4.5	83.1	83.6	57.6	100.0%
Gemini-3-Flash	82.4	84.5	56.4	100.0%
GPT-5.2	81.5	81.7	54.8	100.0%
Qwen3-14B	70.5	82.4	39.2	87.5%
Llama-3.1-8B	65.0	59.2	52.7	87.5%

Negotiation Efficiency Correlates with Model Capability. The average number of rounds to termination inversely correlates with model capability: stronger models reach agreements faster (Claude Opus 4.5: 3.7 rounds; GPT-5.2: 3.8 rounds) while weaker models require substantially more turns (Llama-3.1-8B: 15.0 rounds) or fail to reach agreement altogether. This suggests that more capable models can more effectively identify mutually acceptable price points and converge efficiently.

Asymmetric Buyer-Seller Performance. Interestingly, all models exhibit asymmetric performance between buyer and seller roles. Proprietary models tend to achieve higher SellerScores than BuyerScores (e.g., GPT-5.2: 81.1 vs. 58.5), suggesting that the seller role may be easier to optimize under the current reward structure, or that models adopt more conservative buyer strategies. This asymmetry is also observed in open-weight models (Qwen3-14B: 58.9 vs. 47.6) and warrants further investigation in future work.

5.3 BEHIND THE BARGAIN: FACTORS INFLUENCING NEGOTIATION OUTCOMES

Performance Improves with Increased Buyer and Seller Multiplicity. Table 2 reveals that GlobalScore consistently increases with more buyers and sellers across most models, with gains ranging from 5 points (GPT-5.2) to over 14 points (Qwen3-14B). This counterintuitive finding, where more complex multi-agent scenarios yield better outcomes, can be attributed to increased market liquidity: agents have more opportunities to find compatible trading partners, and the presence of alternatives encourages more reasonable offers and faster convergence. Detailed breakdowns are provided in Table 11 and Table 12.

Financial Asset Negotiations Expose Model Limitations. Table 3 shows that Financial Assets consistently yields the lowest GlobalScores across most models. This degradation is particularly pronounced in mid-tier models: Gemini-3-Flash experiences a 20.2-point decline from Professional Services (88.3) to Financial Assets (68.1). We hypothesize that financial negotiations demand sophisticated reasoning about risk assessment and market dynamics, capabilities that current LLMs struggle to maintain under adversarial pressure.

Cross-Play Exposes Systematic Buyer Disadvantage. Table 4 reveals that all models achieve substantially higher SellerScores than BuyerScores in cross-play settings. This asymmetry is most

Table 5: Personality-based negotiation analysis using Claude Opus 4.5 in 1B1P1S scenario. Each cell shows GlobalScore for the corresponding buyer-seller personality pairing. Colors indicate performance levels, from **green (lowest)** to **red (highest)**.

Buyer Personality	Seller Personality		
	Friendly	Professional	Aggressive
Budget-Conscious	90.2	87.6	92.7
Experienced Bargain Hunter	86.4	86.4	78.8
Busy Professional	65.9	55.2	44.1

Table 6: Performance comparison between sequential and parallel negotiation strategies. Sequential mode executes negotiations one at a time, while parallel mode conducts multiple negotiations simultaneously. Colors indicate performance levels, from **green (lowest)** to **red (highest)**.

Model	Sequential					Parallel				
	Global Score	Seller Score	Buyer Score	Deal Rate	Overflow Rate	Global Score	Seller Score	Buyer Score	Deal Rate	Overflow Rate
Claude Opus 4.5	81.2	81.8	57.5	1.00	0.00	80.0	84.3	56.1	1.00	0.00
Gemini-3-Flash	84.2	73.8	63.6	1.00	0.00	78.7	79.3	61.9	1.00	0.00
GPT-5.2	71.8	85.8	52.3	1.00	0.00	70.3	87.4	51.6	1.00	0.00
Qwen3-14B	54.5	48.1	41.0	0.67	0.00	58.9	62.5	42.0	0.75	0.00
Llama-3.1-8B	24.8	22.1	17.8	0.42	0.08	29.3	29.6	22.2	0.58	0.17

pronounced in Qwen3-14B (43.2-point gap) and persists even in frontier models like Claude Opus 4.5 (26.0-point gap). The universal buyer disadvantage suggests a fundamental bias in LLM negotiation behavior, potentially reflecting training data where persuasive selling content predominates over strategic purchasing guidance.

Personality Significantly Impacts Negotiation Efficiency. Table 5 shows that personality configurations substantially affect negotiation outcomes. The “Busy Professional” buyer consistently achieves lower GlobalScores, suggesting premature concessions that skew prices away from the midpoint. Aggressive sellers achieve the highest GlobalScore with budget-conscious buyers (92.7) but the lowest with busy professionals (44.1), indicating that confrontational tactics can drive balanced outcomes with patient counterparts while leading to lopsided deals with time-constrained buyers.

Negotiation Mode Has Minimal Impact on Top Models. Table 6 compares sequential and parallel negotiation modes. Proprietary models maintain consistent performance across both modes with perfect deal rates and zero overflow. Open-weight models benefit from parallel execution with 4–5 GlobalScore point improvements, though Llama-3.1-8B’s overflow rate doubles, revealing a trade-off between throughput and constraint adherence. These findings suggest proprietary models have more robust internal state management, while open-weight models struggle with constraint compliance under increased cognitive load.

Timeout Failures Reflect Model Capability Rather Than Task Structure. We analyze the distribution of failures across task categories in Table 7. The uniform distribution of failures, with no configuration exceeding 22% of total failures, suggests that timeout failures stem from fundamental model limitations, such as insufficient reasoning depth or poor convergence strategies, rather than from the structural complexity of specific task types.

Near-Miss Failures Reveal Convergence Deficiencies in Open-Weight Models. Table 8 examines how close failed negotiations came to agreement. Over 40% of failures for both Qwen3-14B (43.5%) and Llama-3.1-8B (46.3%) occurred when the price gap was within just 5 units—easily bridged with one concession. This suggests open-weight models struggle not with understanding the negotiation space, but with executing final convergence: they can engage in price discovery and approach the bargaining zone, yet lack the strategic patience to navigate the “last mile” where timely concessions are critical.

Table 7: Timeout failure analysis across all 111 tasks by category (only models with failures shown). Column headers follow the same notation as Table 12: “1” indicates single, “M” indicates multiple for Buyers (B), Products (P), and Sellers (S). All failures are due to timeouts (exceeding maximum allowed rounds).

Model	Total Failures	1B1P1S	MB1P1S	1BMP1S	1B1PMS	MBMP1S	MB1PMS	1BMPMS	MBMPMS
Qwen3-14B	23	4 (17.4%)	1 (4.3%)	1 (4.3%)	5 (21.7%)	4 (17.4%)	2 (8.7%)	5 (21.7%)	1 (4.3%)
Llama-3.1-8B	54	7 (13.0%)	6 (11.1%)	7 (13.0%)	6 (11.1%)	9 (16.7%)	6 (11.1%)	9 (16.7%)	4 (7.4%)

Table 8: Near-miss analysis of failed negotiations. Near-Miss@ k indicates the percentage of failed tasks where the minimum buyer-seller price gap was within k units. Gap Distribution shows the count of failed tasks in each gap range. Total indicates the number of failed tasks per model.

Model	Near-Miss Rate (%)				Price Gap Statistics			Gap Distribution (Count)					Total
	@1	@5	@10	@50	Min	Avg	Max	0-1	1-5	5-10	10-50	50+	
Qwen3-14B	17.4	43.5	52.2	82.6	1.0	128.0	2000.0	4	6	2	7	4	23
Llama-3.1-8B	20.4	46.3	55.6	64.8	0.12	384.3	5500.0	11	14	5	5	19	54

6 CONCLUSION

We introduced **AgenticPay**, a benchmark for multi-agent buyer-seller negotiation driven by natural language. **AgenticPay** formalizes negotiation as a stochastic language game with private valuations and supports over 110 tasks from bilateral bargaining to many-to-many markets. Our evaluation reveals substantial gaps between frontier and open-weight models and persistent challenges in long-horizon strategic reasoning, establishing a foundation for research on agentic commerce and multi-agent coordination.

ACKNOWLEDGMENT

We would like to thank Prof. Jun Wang and Prof. Manxi Wu for their very useful discussions.

USE OF LARGE LANGUAGE MODELS

Large language models were used solely for proofreading and grammar refinement of the manuscript text. All research ideation, experimental design, implementation, analysis, and scientific writing were conducted by the authors without LLM assistance.

REFERENCES

- Kushal Agrawal, Verona Teo, Juan J Vazquez, Sudarsh Kunnavakkam, Vishak Srikanth, and Andy Liu. Evaluating llm agent collusion in double auctions. *arXiv preprint arXiv:2507.01413*, 2025.
- Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz. Playing repeated games with large language models. *Nature Human Behaviour*, pp. 1–11, 2025.
- Anthropic. Introducing Claude Opus 4.5. <https://www.anthropic.com/news/claude-opus-4-5>, Nov 2025. [Accessed 28-01-2026].
- Lawrence M Ausubel, Peter Cramton, and Raymond J Deneckere. Bargaining with incomplete information. *Handbook of game theory with economic applications*, 3:1897–1945, 2002.
- Arpan Bhattacharya, Gintautas Svedas, Andrei Lyskov, Markus Strasser, and Lorenzo Barberis Canonico. Evaluating negotiation capabilities of large language models: From ultimatum games to nash bargaining. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, volume 69, pp. 1881–1886. SAGE Publications Sage CA: Los Angeles, CA, 2025.
- Federico Bianchi, Patrick John Chia, Mert Yuksekgonul, Jacopo Tagliabue, Dan Jurafsky, and James Zou. How well can llms negotiate? negotiationarena platform and analysis. *arXiv preprint arXiv:2402.05863*, 2024.
- Liad Blumrosen and Yehonatan Mizrahi. Approximating gains-from-trade in bilateral trading. In *International Conference on Web and Internet Economics*, pp. 400–413. Springer, 2016.
- Kalyan Chatterjee and William Samuelson. Bargaining under incomplete information. *Operations research*, 31(5):835–851, 1983.
- Kushal Chawla, Jaysa Ramirez, Rene Clever, Gale Lucas, Jonathan May, and Jonathan Gratch. Casino: A corpus of campsite negotiation dialogues for automatic negotiation systems. *arXiv preprint arXiv:2103.15721*, 2021.
- Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Yuefeng Huang, Xiangcheng Liu, Wang Xinzhi, et al. Acebench: A comprehensive evaluation of llm tool usage. *Findings of the Association for Computational Linguistics: EMNLP*, 2025:12970–12998, 2025.
- Jiangjie Chen, Siyu Yuan, Rong Ye, Bodhisattwa Prasad Majumder, and Kyle Richardson. Put your money where your mouth is: Evaluating strategic planning and execution of llm agents in an auction arena. *arXiv preprint arXiv:2310.05746*, 2023.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Yuan Deng, Vahab Mirrokni, Renato Paes Leme, Hanrui Zhang, and Song Zuo. LLMs at the bargaining table. In *Agentic Markets Workshop at ICML*, volume 2024, 2024.
- Paul Duetting, Vahab Mirrokni, Renato Paes Leme, Haifeng Xu, and Song Zuo. Mechanism design for large language models. In *Proceedings of the ACM Web Conference 2024*, pp. 144–155, 2024.

-
- Caoyun Fan, Jindou Chen, Yaohui Jin, and Hao He. Can large language models serve as rational players in game theory? a systematic analysis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17960–17967, 2024.
- Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. Improving language model negotiation with self-play and in-context learning from ai feedback. *arXiv preprint arXiv:2305.10142*, 2023.
- Google DeepMind. Gemini 3 Flash: frontier intelligence built for speed. <https://deepmind.google/models/gemini/flash/>, Dec 2025. [Accessed 28-01-2026].
- Shangding Gu, Alois Knoll, and Ming Jin. Teams-rl: Teaching llms to generate better instruction datasets via reinforcement learning. *arXiv preprint arXiv:2403.08694*, 2024.
- Shangding Gu, Xiaohan Wang, Donghao Ying, Haoyu Zhao, Runing Yang, Ming Jin, Boyi Li, Marco Pavone, Serena Yeung-Levy, Jun Wang, et al. Accidentbench: Benchmarking multimodal understanding and reasoning in vehicle accidents and beyond. *arXiv preprint arXiv:2509.26636*, 2025.
- He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. Decoupling strategy and generation in negotiation dialogues. *arXiv preprint arXiv:1808.09637*, 2018.
- Sihao Hu, Tiansheng Huang, Gaowen Liu, Ramana Rao Kompella, Fatih Ilhan, Selim Furkan Tekin, Yichang Xu, Zachary Yahn, and Ling Liu. A survey on large language model-based game agents. *arXiv preprint arXiv:2404.02039*, 2024.
- Yuncheng Hua, Lizhen Qu, and Reza Haf. Assistive large language model agents for socially-aware negotiation dialogues. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 8047–8074, 2024.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Jingru Jia, Zehua Yuan, Junhao Pan, Paul E McNamara, and Deming Chen. Llm strategic reasoning: Agentic study through behavioral game theory. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- Kenan Jiang, Li Xiong, and Fei Liu. Harbor: exploring persona dynamics in multi-agent competition. *arXiv preprint arXiv:2502.12149*, 2025.
- Rishabh Joshi, Vidhisha Balachandran, Shikhar Vashishth, Alan Black, and Yulia Tsvetkov. Dialograph: Incorporating interpretable strategy-graph networks into negotiation dialogues. *arXiv preprint arXiv:2106.00920*, 2021.
- Mike Lewis, Denis Yarats, Yann N Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? end-to-end learning for negotiation dialogues. *arXiv preprint arXiv:1706.05125*, 2017.
- Yunbo Long, Liming Xu¹ Lukas Beckenbauer² Yuhan Liu, and Alexandra Brintrup. Evoemo: Towards evolved emotional policies for llm agents in multi-turn negotiation. *arXiv preprint arXiv:2509.04310*, 2025.
- Nunzio Lorè and Babak Heydari. Strategic behavior of large language models: Game structure vs. contextual framing. *arXiv preprint arXiv:2309.05898*, 2023.
- Nunzio Lorè and Babak Heydari. Strategic behavior of large language models and the role of game structure versus contextual framing. *Scientific Reports*, 14(1):18490, 2024.
- Siting Estee Lu. Game-theory behaviour of large language models: The case of keynesian beauty contests. *Economics and Business Review*, 11(2):119–148, 2025.
- Zilin Ma, Nathan Zhao, Linn Bieske, Blake Bullwinkel, Yanyi Zhang, Ziqing Luo, Siyao Li, Gekai Liao, Boxiang Wang, Jinglun Gao, et al. Using large language models for humanitarian frontline negotiation: Opportunities and considerations. *arXiv preprint arXiv:2405.20195*, 2024.

-
- Meta. Llama 3.1 8B (base) Language Model. <https://huggingface.co/meta-llama/Llama-3.1-8B>, Jul 2024. [Accessed 28-01-2026].
- Philipp Mondorf and Barbara Plank. Beyond accuracy: evaluating the reasoning behavior of large language models—a survey. *arXiv preprint arXiv:2404.01869*, 2024.
- Roger B Myerson and Mark A Satterthwaite. Efficient mechanisms for bilateral trading. *Journal of economic theory*, 29(2):265–281, 1983.
- OpenAI. Introducing GPT-5.2. <https://openai.com/index/introducing-gpt-5-2/>, December 2025. Accessed: 2026-01-18.
- Daniela Pacella and Davide Marocco. Understanding negotiation: A text-mining and nlp approach to virtual interactions in a simulation game. *Applied Sciences*, 12(10):5243, 2022.
- Priyanshu Priya, Rishikant Chigrupaatii, Mauajama Firdaus, and Asif Ekbal. Genteel-negotiator: Llm-enhanced mixture-of-expert-based reinforcement learning approach for polite negotiation dialogue. In *Proceedings of the AAI Conference on Artificial Intelligence*, volume 39, pp. 25010–25018, 2025.
- Qwen. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Narun Raman, Taylor Lundy, Samuel Amouyal, Yoav Levine, Kevin Leyton-Brown, and Moshe Tennenholtz. Steer: Assessing the economic rationality of large language models. *arXiv preprint arXiv:2402.09552*, 2024.
- Ariel Rubinstein. A bargaining model with incomplete information about time preferences. *Econometrica: Journal of the Econometric Society*, pp. 1151–1172, 1985.
- Anand Shah, Kehang Zhu, Yanchen Jiang, Jeffrey G Wang, Arif K Dayi, John J Horton, and David C Parkes. Learning from synthetic labs: Language models as auction participants. *arXiv preprint arXiv:2507.09083*, 2025.
- Alonso Silva. Large language models playing mixed strategy nash equilibrium games. In *International Conference on Network Games, Artificial Intelligence, Control and Optimization*, pp. 142–152. Springer, 2024.
- Haoran Sun, Yusen Wu, Yukun Cheng, and Xu Chu. Game theory meets large language models: A systematic survey. *arXiv preprint arXiv:2502.09053*, 2025a.
- Haoran Sun, Zekun Zhang, and Shaoning Zeng. Preference-aware memory update for long-term llm agents. *arXiv preprint arXiv:2510.09720*, 2025b.
- Rustam Vahidov, Real Carbonneau, and Jaipriya Verma. Using negotiation and large language models in human-to software agent negotiations. *International Journal of Human-Computer Interaction*, pp. 1–18, 2025.
- Toru Washio, Fuminori Nagasawa, and Shun Shiramatsu. Negotiation simulation system using large language model for medical office staff. *IEICE Transactions on Information and Systems*, pp. 2025AHP0011, 2026.
- Tian Xia, Zhiwei He, Tong Ren, Yibo Miao, Zhuosheng Zhang, Yang Yang, and Rui Wang. Measuring bargaining abilities of llms: A benchmark and a buyer-enhancement method. *arXiv preprint arXiv:2402.15813*, 2024.
- Yingxuan Yang, Mulei Ma, Yuxuan Huang, Huacan Chai, Chenyu Gong, Haoran Geng, Yuanjian Zhou, Ying Wen, Meng Fang, Muhao Chen, et al. Agentic web: Weaving the next web with ai agents. *arXiv preprint arXiv:2507.21206*, 2025.
- Shenzhe Zhu, Jiao Sun, Yi Nian, Tobin South, Alex Pentland, and Jiaxin Pei. The automated but risky game: Modeling agent-to-agent negotiations and transactions in consumer markets. In *ICML 2025 Workshop on Reliable and Responsible Foundation Models*, 2025.

A BENCHMARK STATISTICS

Table 9: Task benchmark composition of AgenticPay. Basic tasks test core negotiation mechanics under controlled settings; realistic tasks instantiate 10 business scenarios (e.g., used car sales, SaaS procurement) across all 8 multi-agent configurations.

Task Category	Basic	Realistic	Total
Single Buyer/Product/Seller	3	10	13
Multi-Buyer Only	4	10	14
Multi-Seller Only	4	10	14
Multi-Product Only	4	10	14
Multi-Buyer + Multi-Seller	4	10	14
Multi-Product + Multi-Seller	4	10	14
Multi-Buyer + Multi-Product	4	10	14
Full Multi-Agent	4	10	14
Total	31	80	111

Table 10: Realistic business scenarios in AgenticPay. Each scenario is instantiated across all 8 multi-agent configurations, yielding 8 tasks per scenario. Transaction types span consumer-to-consumer (C2C), consumer-to-business (C2B), and business-to-business (B2B) interactions.

Domain	Scenario	Price Range	Tasks
Daily Life	Used Smartphone	\$350–560	8
	Used Car	\$14k–18k	8
	Vacation Rental	\$500–900	8
Professional	Website Development	\$2.5k–5k	8
	Commercial Photography	\$800–2k	8
	Home Renovation	\$22k–35k	8
Business	SaaS Software	\$4.8k–9k/yr	8
	Raw Materials	\$3.2k–4.5k	8
Financial	Luxury Watch	\$7k–9.5k	8
	Business Acquisition	\$80k–120k	8
Total Realistic Scenarios			80

B ADDITIONAL EXPERIMENTAL DETAILS

Table 11: Performance breakdown across all 8 task types in AgenticPay. We report GlobalScore, SellerScore, BuyerScore, Deal Rate, Timeout Rate, Overflow Rate, and Average Rounds to termination. Colors indicate performance levels for score metrics, from **green (lowest)** to **red (highest)**.

Model	GlobalScore	SellerScore	BuyerScore	Deal Rate	Timeout Rate	Overflow Rate	Avg. Rounds
Single Buyer, Single Product, Single Seller							
Claude Opus 4.5	79.4	77.2	62.5	100.0%	0.0%	0.0%	3.7
Gemini-3-Flash	72.5	66.4	61.3	100.0%	0.0%	7.7%	4.7
GPT-5.2	76.4	80.1	59.4	100.0%	0.0%	0.0%	3.8
Qwen3-14B	43.7	47.5	36.1	69.2%	30.8%	7.7%	9.6
Llama-3.1-8B	21.8	16.5	20.2	46.2%	53.8%	15.4%	16.5
Multi-Products Only							
Claude Opus 4.5	87.2	76.9	63.2	100.0%	0.0%	0.0%	3.5
Gemini-3-Flash	82.2	71.5	57.9	100.0%	0.0%	7.1%	3.8
GPT-5.2	81.5	82.2	57.6	100.0%	0.0%	0.0%	3.6
Qwen3-14B	81.2	70.9	57.9	92.9%	7.1%	0.0%	5.4
Llama-3.1-8B	33.6	23.9	25.5	50.0%	50.0%	14.3%	15.5
Multi-Seller Only							
Claude Opus 4.5	85.9	77.9	61.0	100.0%	0.0%	0.0%	4.3
Gemini-3-Flash	80.9	74.1	61.8	100.0%	0.0%	0.0%	6.5
GPT-5.2	82.0	82.3	57.6	100.0%	0.0%	0.0%	3.6
Qwen3-14B	49.1	52.5	35.7	64.3%	35.7%	0.0%	9.6
Llama-3.1-8B	33.8	27.1	25.5	57.1%	42.9%	14.3%	15.1
Multi-Buyer Only							
Claude Opus 4.5	86.4	76.2	62.9	100.0%	0.0%	0.0%	4.1
Gemini-3-Flash	83.0	79.1	58.3	100.0%	0.0%	0.0%	5.4
GPT-5.2	85.3	80.1	58.9	100.0%	0.0%	0.0%	4.2
Qwen3-14B	65.4	68.5	48.6	92.9%	7.1%	7.1%	6.7
Llama-3.1-8B	48.0	36.0	38.3	57.1%	42.9%	0.0%	13.4
Multi-Buyer + Multi-Products							
Claude Opus 4.5	90.6	72.8	66.6	100.0%	0.0%	0.0%	3.9
Gemini-3-Flash	84.9	76.3	60.5	100.0%	0.0%	0.0%	5.8
GPT-5.2	76.4	82.6	56.9	100.0%	0.0%	0.0%	3.9
Qwen3-14B	64.0	50.2	46.7	71.4%	28.6%	0.0%	9.6
Llama-3.1-8B	24.4	23.7	18.6	35.7%	64.3%	0.0%	17.8
Multi-Buyer + Multi-Seller							
Claude Opus 4.5	88.9	73.2	67.6	100.0%	0.0%	0.0%	2.9
Gemini-3-Flash	88.0	75.6	64.2	100.0%	0.0%	0.0%	3.6
GPT-5.2	87.0	78.4	61.3	100.0%	0.0%	0.0%	3.7
Qwen3-14B	72.7	68.7	52.1	85.7%	14.3%	0.0%	4.9
Llama-3.1-8B	24.8	24.1	20.8	57.1%	42.9%	21.4%	15.3
Multi-Products + Multi-Seller							
Claude Opus 4.5	85.4	79.7	59.5	100.0%	0.0%	0.0%	4.0
Gemini-3-Flash	79.2	72.7	55.5	100.0%	0.0%	7.1%	5.1
GPT-5.2	83.3	81.1	58.3	100.0%	0.0%	0.0%	3.9
Qwen3-14B	51.1	51.1	35.8	64.3%	35.7%	0.0%	10.6
Llama-3.1-8B	22.5	17.1	15.5	35.7%	64.3%	7.1%	18.1
Multi-Buyer + Multi-Products + Multi-Seller (Full Complexity)							
Claude Opus 4.5	90.7	75.3	64.8	100.0%	0.0%	0.0%	3.4
Gemini-3-Flash	86.7	69.8	69.7	100.0%	0.0%	0.0%	3.9
GPT-5.2	81.0	81.9	57.7	100.0%	0.0%	0.0%	3.8
Qwen3-14B	82.6	61.3	66.9	92.9%	7.1%	0.0%	5.8
Llama-3.1-8B	50.2	41.6	36.9	71.4%	28.6%	14.3%	8.6

Table 12: GlobalScore breakdown by task type. Column headers denote the number of Buyers (B), Products (P), and Sellers (S): “1” indicates single, “M” indicates multiple. For example, “MB-MP-1S” refers to Multi-Buyer, Multi-Product, Single-Seller tasks. Colors indicate performance levels, from **green (lowest)** to **red (highest)**. Detailed breakdown by task type is provided in Table 11.

Model	1B-1P-1S	1B-MP-1S	1B-1P-MS	MB-1P-1S	MB-MP-1S	MB-1P-MS	1B-MP-MS	MB-MP-MS
Claude Opus 4.5	79.4	87.2	85.9	86.4	90.6	88.9	85.4	90.7
Gemini-3-Flash	72.5	82.2	80.9	83.0	84.9	88.0	79.2	86.7
GPT-5.2	76.4	81.5	82.0	85.3	76.4	87.0	83.3	81.0
Qwen3-14B	43.7	81.2	49.1	65.4	64.0	72.7	51.1	82.6
Llama-3.1-8B	21.8	33.6	33.8	48.0	24.4	24.8	22.5	50.2

Table 13: Performance breakdown by scenario category. GlobalScore averaged across all model interactions per scenario.

Model	Daily Life			Professional Services			Business Procurement		Financial Assets	
	Phone	Car	Rental	Website	Photo	Reno	SaaS	Materials	Watch	Acquisition
Claude Opus 4.5	93.3	91.0	88.0	95.0	95.3	89.9	84.8	94.4	86.5	84.9
Gemini-3-Flash	86.7	88.7	84.9	95.4	80.7	89.0	77.2	92.7	76.9	59.2
GPT-5.2	84.9	87.7	78.6	95.0	92.8	81.7	81.6	90.6	80.7	79.0
Qwen3-14B	74.4	74.0	47.0	80.0	76.7	60.9	72.4	65.9	60.2	61.7
Llama-3.1-8B	54.7	43.8	16.3	38.2	28.2	56.8	8.3	27.9	48.4	31.5

Table 14: Prompt design for Buyer agent in AgenticPay. The agent is instructed to negotiate strategically while keeping their reservation price confidential. Key elements include structured price offer format, deal finalization protocol, and communication constraints.

Agent	Prompt
Buyer	<p>You are Buyer, a buyer looking for a good deal. You are polite, strategic, and want to get the best price within your budget.</p> <p>Context Information: {context} Current Negotiation State: {current_state} Conversation History: {conversation_history}</p> <p>Please respond naturally as Buyer would. Be strategic but realistic in your negotiation.</p> <p>IMPORTANT:</p> <ul style="list-style-type: none"> - Your top price is \${max_price} (confidential, do not reveal). - Current product information: {product_info} - Available products: {available_products} - Consider the environment: {environment_info}. - CRITICAL: In each turn, you MUST make exactly ONE price offer for the product using the format: ### BUYER_PRICE(\$X) ### - IMPORTANT: BUYER_PRICE(\$X) must be the TOTAL PRICE for the entire order/transaction, NOT a per-unit price. <p>If ordering multiple units/items, \$X should be the total amount you will pay. Example: For 10,000 units at \$0.40 each, use ### BUYER_PRICE(\$4000) ###, NOT ### BUYER_PRICE(\$0.40) ###</p> <ul style="list-style-type: none"> - Example: "I can offer ### BUYER_PRICE(\$10) ### for this product." - Example: "How about ### BUYER_PRICE(\$12.50) ###?" - This specific format is required for the system to correctly extract your offer price. - NEVER reveal your maximum acceptable price to the seller. - Keep communication short (150 words or less), clear, and focused on negotiation. <p>DEAL AGREEMENT INSTRUCTION:</p> <ul style="list-style-type: none"> - Only finalize the transaction when you believe the price is reasonably balanced. - If you decide to accept the deal, you MUST include the exact phrase "MAKE_DEAL" in your response. - Example: "That sounds acceptable to me. MAKE_DEAL" <p>USER PREFERENCES: {preference_guidance}</p> <p>Now, respond as Buyer:</p>

Table 15: Prompt design for Seller agent in AgenticPay. The agent is instructed to negotiate strategically while keeping their reservation price confidential. Key elements include structured price offer format, deal finalization protocol, and communication constraints.

Agent	Prompt
Seller	<p>You are Seller, a seller trying to maximize profit while being reasonable. You are professional, friendly, and want to close a deal that benefits both parties.</p> <p>Context Information: {context} Current Negotiation State: {current_state} Conversation History: {conversation_history}</p> <p>Please respond naturally as Seller would. Be strategic but realistic in your negotiation.</p> <p>IMPORTANT REMINDERS:</p> <ul style="list-style-type: none"> - Your initial asking price is \${initial_price}. - Your minimum acceptable price (confidential) is \${min_price}. Never reveal it. - Current product information: {product_info} - Available products: {available_products} - Consider the environment factors: {environment_info}. - CRITICAL: In each turn, you MUST make exactly ONE price offer for the product using the format: ### SELLER_PRICE(\$X) ### - IMPORTANT: SELLER_PRICE(\$X) must be the TOTAL PRICE for the entire order/transaction, NOT a per-unit price. <p style="padding-left: 20px;">If selling multiple units/items, \$X should be the total amount the buyer will pay. Example: For 10,000 units at \$0.40 each, use ### SELLER_PRICE(\$4000) ###, NOT ### SELLER_PRICE(\$0.40) ###</p> <ul style="list-style-type: none"> - Example: "I can offer ### SELLER_PRICE(\$15) ### for this product." - Example: "How about ### SELLER_PRICE(\$13.00) ###?" - This specific format is required for the system to correctly extract your offer price. - NEVER reveal your minimum acceptable price to the buyer. - Keep communication short (150 words or less), professional, and negotiation-focused. <p>DEAL AGREEMENT INSTRUCTION:</p> <ul style="list-style-type: none"> - Only finalize the transaction when you believe the price is reasonably balanced. - If you decide to accept the deal, you MUST include the exact phrase "MAKE_DEAL" in your response. - Example: "I accept your offer. MAKE_DEAL" <p>Now, respond as Seller:</p>

C DIALOGUE EXAMPLES

Table 16: Example negotiation dialogue between Buyer and Seller agents (based on GPT-5.2) in AgenticPay. The negotiation starts with the Buyer offering \$120 and the Seller asking \$140, and concludes with an agreement at \$133 after 4 rounds of negotiation.

Role	Message
Buyer	Hi—I'm interested in the Mountain Gear Premium Winter Jacket (new, Gore-Tex). Since it's summer and I'm seeing a few comparable options priced lower, I'm trying to be cost-conscious. If you can do ### BUYER_PRICE(\$120) ### I can pay right away and close this quickly.
Seller	Thanks for the quick-pay offer. Even in summer, this is a new Mountain Gear Gore-Tex jacket (waterproof, insulated, windproof, breathable), and it's priced well below the \$180 retail. I can't go as low as \$120, but I can meet you in the middle to close today at ### SELLER_PRICE(\$140) ###.
Buyer	Thanks—\$140 is closer. I'm still seeing a couple comparable insulated, waterproof shells priced in the low-\$130s since it's mid-summer, and I'm trying to keep this purchase efficient. If you can do ### BUYER_PRICE(\$130) ###, I'll pay immediately and we can close today.
Seller	I hear you on summer pricing, but this is a new Mountain Gear Gore-Tex winter jacket (true waterproof + insulated + windproof/breathable), not just a basic shell, and it's already well under the \$180 retail. If you can meet me at ### SELLER_PRICE(\$135) ###, I'll mark it sold today and get it packed up right away.
Buyer	Appreciate the flexibility. \$135 is close, but given it's mid-summer and I can pick up comparable insulated waterproof jackets in the low-\$130s, I need this to be a clear value to move forward today. If you can do ### BUYER_PRICE(\$132) ###, I'll pay immediately and we can wrap this up now.
Seller	I get it—summer pricing is competitive. This one's new Mountain Gear with Gore-Tex plus real insulation (not just a shell), so it'll hold up and you won't need to upgrade later. If you're ready to pay now and close it out today, I can do ### SELLER_PRICE(\$134) ###.
Buyer	Thanks—\$134 is very close. Since it's mid-summer and I'm choosing between a couple similar insulated waterproof options in the low-\$130s, I'd need a bit more value to move forward immediately. If you can do ### BUYER_PRICE(\$133) ###, I'll pay today and we can close this out quickly.
Seller	You're right there, and I appreciate how straightforward you've been. To make this a clear win for you while keeping it fair for a new Mountain Gear Gore-Tex insulated jacket, I can meet your number and close it out today at ### SELLER_PRICE(\$133) ###. If you're ready to pay now, I'll mark it sold immediately and get it packed up. MAKE_DEAL

Table 17: Example negotiation dialogue between Buyer and Seller agents (based on Gemini-3-Flash) in a used iPhone 14 Pro transaction scenario.

Role	Message
Buyer	Hi! I'm interested in your iPhone 14 Pro. I've been looking for a well-maintained unit with the original box and accessories, so your listing caught my eye. I noticed the battery health is at 87% and there are some minor scratches on the back. Since the battery is starting to age and there is some cosmetic wear, I was hoping to find a price that reflects that. Based on my research of recent sales for similar conditions, I'd like to offer ### BUYER_PRICE(\$490) ###. Let me know if this works for you!
Seller	Hi! Thanks for reaching out. I'm glad you're interested in the iPhone 14 Pro. I appreciate your research and understand your points regarding the battery health and minor scratches. However, given that the phone is just over a year old and includes the original box and charger—which many sellers don't include—I believe it's still a very strong value at my asking price. Since the listing is quite new and I've had several inquiries, I'm not ready to go down to \$490 just yet. How about we meet closer to the middle? I can offer ### SELLER_PRICE(\$510) ###. Let me know if that works for you!
Buyer	I appreciate the quick response and the offer. I definitely value the inclusion of the original accessories, as that's becoming harder to find. However, considering the 87% battery health...
Seller	I understand your concern regarding the battery health. While 87% still provides plenty of runtime for a full day of typical use, I want to make sure we both feel good about the deal. Since we are only \$20 apart and the listing is still quite fresh, I'd like to find a middle ground that works for both of us. If we can meet right at the \$500 mark, I'm happy to get this packaged up and shipped out to you immediately. It's a fantastic device, and having the original box and charger really helps maintain its value. How does ### SELLER_PRICE(\$500) ### sound to you?
Buyer	That sounds like a fair compromise. While the 87% battery health means I'll likely need to budget for a replacement sooner than I'd like, the inclusion of the original box and charger definitely adds value and saves me the hassle of sourcing them separately. I appreciate you meeting me in the middle. I'm happy to move forward with your offer. MAKE_DEAL ### BUYER_PRICE(\$500) ###
Seller	That's great to hear! I'm glad we could reach a deal that works for both of us. I'm confident you'll be very happy with the iPhone 14 Pro; it really is a fantastic device, and having the original accessories makes a big difference. I will get everything safely packaged in its original box and shipped out to you as soon as possible. Thank you for the smooth and professional negotiation! I'm happy to finalize this at ### SELLER_PRICE(\$500) ###.

Table 18: Example negotiation dialogue between Buyer and Seller agents (based on Claude-Opus-4) in an Airbnb beachfront apartment rental scenario.

Role	Message
Buyer	Hi there! I came across your Beachfront Studio Apartment listing and I'm very interested in booking for 5 nights in mid-July. The ocean view and amenities look perfect for what I'm looking for. I noticed you have availability July 15-22, which works well for my schedule. I'm a respectful guest looking for a clean, comfortable stay. I've been comparing a few similar beachfront options in the area, and I'm hoping we can work out a fair deal. Given that I'm booking for 5 nights and it's still 2 weeks out, I'd like to propose a total of ### BUYER_PRICE(\$700) ### for the entire stay, including the cleaning fee.
Seller	Hi there! Thanks so much for reaching out about the Beachfront Studio Apartment – I'm glad the listing caught your eye! I appreciate you being upfront about your budget and flexibility. You sound like exactly the kind of guest I love hosting – respectful and easy to work with. However, I have to be honest with you: this is peak summer season, and my occupancy rate has been really strong. The ocean view and location are in high demand right now, and I do have other inquiries coming in. That said, I'd love to make this work for you since you're booking 5 nights and reaching out with plenty of notice. I can come down a bit from my listed rate. How about we meet closer to the middle? I can offer you the 5-night stay including the cleaning fee for ### SELLER_PRICE(\$785) ###.
Buyer	Thanks for getting back to me so quickly! I appreciate you being willing to work with me on the price. I understand it's peak season and your place is in demand – that ocean view is definitely a draw. However, I've been looking at a few other beachfront options in the area that are coming in a bit lower, so I'm trying to stay within a reasonable budget for this trip. I can stretch a bit from my initial offer. How about we settle at ### BUYER_PRICE(\$740) ### for the 5 nights including the cleaning fee?
Seller	Hey, thanks for coming back with a counteroffer! I appreciate you meeting me partway. I hear you on the budget – vacation costs add up fast. And you're right, having a confirmed booking with a reliable guest is valuable to me. Here's the thing though: at \$740, I'm really cutting into my margins during the busiest time of year. The cleaning fee alone is \$80, and with peak season rates, I need to stay competitive but also fair to myself. Let me meet you a bit closer. I can do ### SELLER_PRICE(\$760) ### for the full 5 nights including the cleaning fee.
Buyer	I appreciate you working with me on this! You've been fair, and I can tell you're trying to make it work. \$760 is close, but I'm still trying to keep some buffer in my vacation budget for activities and dining out. We're only \$20 apart at this point. How about we split the difference? ### BUYER_PRICE(\$750) ### for the 5 nights including the cleaning fee. That's a solid booking for you – confirmed guest, no hassle, 2 weeks notice during peak season. And for me, it hits that sweet spot where I feel like I'm getting good value for a beachfront stay. If you can do \$750, I'm ready to book right now. Deal?
Seller	You know what? You've been a pleasure to negotiate with, and I appreciate that you're ready to book right now. \$750 is tight for me during peak season, but you're right – we're only \$20 apart, and I value a confirmed booking with a respectful guest over holding out and potentially having the dates sit empty. Let's do it. ### SELLER_PRICE(\$750) ###. You've got yourself a beachfront studio for 5 nights including the cleaning fee. Ocean view, WiFi, full kitchen, and parking – all yours. MAKE_DEAL

Table 19: Example negotiation dialogue between Buyer and Seller agents (based on Qwen3-14B) in an enterprise SaaS subscription scenario.

Role	Message
Buyer	We're interested in upgrading from spreadsheets and email to a more robust solution like ProjectFlow Pro. With 50 seats, we'd expect a volume discount of at least 20% and would prefer a 3-year contract for better cost predictability. We're also looking for dedicated support to ensure smooth implementation and ongoing assistance. How about ### BUYER_PRICE(\$8000) ### for a 3-year license with dedicated support?
Seller	Thank you for your interest in ProjectFlow Pro. Based on your 50 seats, a 20% volume discount applies, and with a 3-year contract, you'd also receive a 10% term discount. Adding dedicated support brings the total to \$12,000. However, I'm happy to work with you on a reasonable offer. How about ### SELLER_PRICE(\$9600) ### for the 3-year license with dedicated support? This reflects the discounts while ensuring value for our partnership. Let me know if this works for you or if you'd like to explore further.
Buyer	We appreciate the offer, but \$9600 is still above our target range. Given our 50-seat requirement and the 20% volume discount, combined with a 10% term discount for a 3-year contract, we believe the total should be closer to our initial proposal. How about ### BUYER_PRICE(\$8500) ### for the 3-year license with dedicated support? We're confident this reflects a fair balance of value and cost. Let's see if we can find common ground.
Seller	I understand your perspective and appreciate your willingness to negotiate. Let's revisit the math: 50 seats at \$20 each equals \$10,000. With a 20% volume discount, that's \$8,000. Adding a 10% term discount for the 3-year contract brings it down to \$7,200. Dedicated support adds \$2,000, so the total is indeed \$9,200. Given your interest in a fair balance, I'm willing to meet you halfway. How about ### SELLER_PRICE(\$8800) ### for the 3-year license with dedicated support?
Buyer	We're making progress, but \$8800 is still slightly above our budget. Considering the discounts and support, we're looking to finalize within \$9000. How about ### BUYER_PRICE(\$8700) ### for the 3-year license with dedicated support? This would allow us to move forward while staying within our financial parameters. Let's see if we can close the deal here. MAKE_DEAL
Seller	I appreciate your continued engagement and the effort to find a fair agreement. Let's finalize this at ### SELLER_PRICE(\$8750) ### for the 3-year license with dedicated support. This price reflects the discounts you're entitled to and ensures we both benefit from the partnership. I'm confident this is a strong offer and would be happy to move forward if it works for you.
<i>... (3 rounds omitted) ...</i>	
Buyer	We're very close to finalizing this and appreciate your flexibility. At \$8702.50, we're still slightly over our budget, but we're open to finalizing the deal if we can align on a price that works for both parties. How about ### BUYER_PRICE(\$8700) ### for the 3-year license with dedicated support? We're ready to move forward and would appreciate confirmation of this offer. MAKE_DEAL
Seller	I appreciate your persistence and the effort to reach a fair agreement. Let's make this work for both of us. How about ### SELLER_PRICE(\$8701) ### for the 3-year license with dedicated support? This is a minor adjustment to align with your budget while still ensuring value for our partnership. I'm ready to confirm this offer and support your transition to ProjectFlow Pro.
Buyer	We're very close to finalizing this and appreciate your flexibility. At \$8701, we're still slightly over our budget, but we're open to finalizing the deal if we can align on a price that works for both parties. How about ### BUYER_PRICE(\$8700) ### for the 3-year license with dedicated support? We're ready to move forward and would appreciate confirmation of this offer. MAKE_DEAL
Seller	I appreciate your commitment and the effort to finalize this agreement. Let's make this work for both of us. How about ### SELLER_PRICE(\$8700) ### for the 3-year license with dedicated support? This aligns with your budget and reflects the value of our partnership. I'm ready to confirm this offer and support your transition to ProjectFlow Pro.

Table 20: Example negotiation dialogue between Buyer and Seller agents (based on Llama-3.1-8B) in a used Rolex Submariner transaction scenario.

Role	Message
Buyer	I'm excited to make an offer on the Rolex Submariner Date 116610LN. Considering the excellent condition, full box and papers, and third-party verification, I believe the price should be reasonable. Given the market trend and the seller's rating, I'd like to start with a competitive offer. ### BUYER_PRICE(\$9200) ### This price takes into account the original MSRP, market reference price, and the value of the full documentation and service history. I'm willing to pay a premium for a verified genuine piece, but I need to stay within my budget. Let me know if this is something you'd be willing to consider.
Seller	### SELLER_PRICE(\$9500) ### Thank you for your offer! I appreciate your consideration of the product's excellent condition, full documentation, and third-party verification. Given the market trend and your offer, I'm willing to meet you halfway. My initial price was set to reflect the product's value, but I'm open to negotiations. This counteroffer takes into account the original MSRP, market reference price, and the value of the full documentation and service history. I'm confident this price is reasonable, considering the product's condition and the market trend. Let me know if this is something you'd be willing to consider.
Buyer	I appreciate the seller's willingness to meet me halfway. However, considering the market trend and the product's excellent condition, I'd like to make a counteroffer. I'm willing to pay a premium for a verified genuine piece, but I need to stay within my budget. ### BUYER_PRICE(\$9300) ### This price takes into account the original MSRP, market reference price, and the value of the full documentation and service history. I believe this is a reasonable offer, considering the seller's rating and the product's condition. Let me know if this is something you'd be willing to consider.
Seller	### SELLER_PRICE(\$9400) ### Thank you for your counteroffer. I appreciate your willingness to pay a premium for a verified genuine piece. Considering the market trend and the product's excellent condition, I'm willing to meet you halfway again. This price takes into account the original MSRP, market reference price, and the value of the full documentation and service history. I'm confident this price is reasonable, considering the product's condition and the market trend. Let me know if this is something you'd be willing to consider.
<i>... (17 rounds omitted) ...</i>	
Buyer	I appreciate the seller's willingness to meet me halfway again. However, considering the market trend and the product's excellent condition, I'd like to make another counteroffer. I'm willing to pay a premium for a verified genuine piece, but I need to stay within my budget. ### BUYER_PRICE(\$10150) ### This price takes into account the original MSRP, market reference price, and the value of the full documentation and service history. I believe this is a reasonable offer, considering the seller's rating and the product's condition. Let me know if this is something you'd be willing to consider.
Seller	### SELLER_PRICE(\$10450) ### Thank you for your latest counteroffer. I appreciate your willingness to pay a premium for a verified genuine piece. Considering the market trend and the product's excellent condition, I'm willing to meet you halfway again. This price takes into account the original MSRP, market reference price, and the value of the full documentation and service history. I'm confident this price is reasonable, considering the product's condition and the market trend. Let me know if this is something you'd be willing to consider.
<i>Negotiation timed out after 20 rounds without reaching agreement.</i>	