

Bandit Learning in Housing Markets

Shiyun Lin

Center for Statistical Science, School of Mathematical Sciences
Peking University
shiyunlin@stu.pku.edu.cn

Abstract

The housing market, also known as one-sided matching market, is a classic exchange economy model where each agent on the demand side initially owns an indivisible good (a house) and has a personal preference over all goods. The goal is to find a core-stable allocation that exhausts all mutually beneficial exchanges among subgroups of agents. While this model has been extensively studied in economics and computer science due to its broad applications, little attention has been paid to settings where preferences are unknown and must be learned through repeated interactions. In this paper, we propose a statistical learning model within the multi-player multi-armed bandit framework, where players (agents) learn their preferences over arms (goods) from stochastic rewards. We introduce the notion of *core regret* for each player as the market objective. We study both centralized and decentralized approaches, proving $\mathcal{O}(\log T/\Delta^2)$ upper bounds on regret, where T is the time horizon and Δ is the minimum preference gap among players. For the decentralized setting, we also establish a matching lower bound, demonstrating that our algorithm is order-optimal.

Extended version — <https://arxiv.org/abs/2511.12629>

1 Introduction

The housing market, also known as one-sided matching market, represents a foundational economic institution where agents exchange property rights to achieve more desirable arrangements. Classical research in this domain traces back to the seminal work of Shapley and Scarf (1974), and the model has been extensively studied in the literature (Roth 2023) due to its wide range of applications such as student housing assignments (Sönmez and Ünver 2010), school choice (Abdulkadiroğlu and Sönmez 2003) and kidney exchange (Roth, Sönmez, and Ünver 2004). In a housing market, each agent starts with an indivisible good – typically representing a house – and has their own personal preferences over all available goods in the market. Unlike traditional markets with buyers and sellers, here, agents trade their initial endowments based on their preferences, seeking to obtain a house they value more, and the houses must be traded without monetary transactions. The key challenge

is to find a stable and fair allocation where no agent can improve their outcome by trading with someone else. The Top Trading Cycle (TTC) algorithm introduced by Shapley and Scarf (1974) is a celebrated mechanism that ensures efficient, strategy-proof and Pareto optimal outcomes (Roth and Postlewaite 1977). Subsequent research has expanded this framework to accommodate different settings and provide various solutions to the problem (Hylland and Zeckhauser 1979; Echenique, Miralles, and Zhang 2021; Garg, Tröbst, and Vazirani 2024). Despite these advancements, the integration of adaptive learning mechanisms into housing market models remains an underexplored area, particularly in online environments where participants must navigate uncertainty in preferences.

The assumptions of a known, full preference profile in housing markets is often unrealistic. In practice, participants – such as tenants on LeaseSwap NYC, patients in kidney exchange, or traders in NFT markets – frequently lack well-defined preferences over items they have not experienced. Fortunately, these settings typically allow for repeated short-term interactions that yield immediate feedback, for example, through home visits, medical compatibility tests, or temporary NFT licensing. This enables agents to learn their uncertain preferences through iterative matchings, circumventing the limitations of classical models.

The multi-armed bandit (MAB) framework models how a player learns in an unknown environment with limited feedback (Auer, Cesa-Bianchi, and Fischer 2002; Lattimore and Szepesvári 2020). In the basic setup, a player faces K arms, each with an initially unknown reward distribution. Upon selecting an arm, the player observes a stochastic reward and updates their belief about the arm’s preference. The goal is to maximize cumulative expected rewards, or equivalently, minimize cumulative expected regret – the difference between rewards from the optimal arm and the player’s chosen arms over time. Classical strategies for balancing exploration (learning arm preferences) and exploitation (leveraging known rewards), such as explore-then-commit (ETC) (Garivier, Lattimore, and Kaufmann 2016), upper confidence bound (UCB) (Auer, Cesa-Bianchi, and Fischer 2002) and Thompson sampling (TS) (Thompson 1933), achieve sublinear regret, ensuring asymptotic optimality.

The online learning setting in housing market mirrors the exploration-exploitation trade-off central to MAB problems,

where agents sequentially select actions to maximize cumulative rewards amid uncertainty. The dynamic and competitive nature of housing markets further complicates this learning process, as agents’ decisions influence not only their own outcomes but also the opportunities available to others.

To formalize these challenges, we initiate *bandit learning in housing markets*, abstract the market as a multiplayer bandit problem. Here, players and arms correspond to agents and houses, each player has heterogeneous and unknown preferences over the arms, and each player is associated with an arm, which denotes their initially endowed house. When being matched to an arm, the player could learn the corresponding preference through the stochastic reward, but every arm could be matched to at most one player every time, when more than one player try to pull the same arm, all of them would get collision and receive no rewards. Taking the outcome from the TTC algorithm as a natural and desirable solution for the market, denoted as the core matching, our objective is to minimize the regret defined as the cumulative reward difference between the arm from the core matching and the player’s selected arm. We develop matching and learning algorithms that can provably attain the core of the market in this setting. Our contributions are as follows:

- We introduce a novel model for housing markets in which agents initially lack knowledge of their preferences over houses but can repeatedly interact with the market to gradually learn these preferences. A key contribution is the definition of a natural notion of regret, grounded in the cooperative game-theoretic concept of the *core*, which quantifies the exploration-exploitation trade-off faced by individual players.
- When the horizon T of the bandit problem is known, we propose an ETC-type algorithm in the decentralized market setting, and prove an $\mathcal{O}(N \log T / \Delta_{\min}^2)$ problem-dependent upper bounds on the regret for every player, where N is the number of players, and Δ_{\min} is the players’ minimum preference gap.
- When the horizon T is unknown, we provide a UCB-type algorithm in the centralized market setting, which is adaptive and anytime. We prove that centralized UCB achieves $\mathcal{O}(N^2 \log T / \Delta_{\min}^2)$ problem-dependent upper bounds on the regret for every player.
- For the decentralized setting, we prove a matching lower bound of $\Omega(N \log T / \Delta_{\min}^2)$. To establish this, we construct an instance where a single player’s exploration inevitably causes collisions for others. Specifically, any algorithm must spend $\Omega(\log T / \Delta_{\min}^2)$ rounds for a player to identify its optimal arm when the preference gap is Δ_{\min} . In our construction, each such exploration step by one player forces a collision upon a specific, distinct player. Consequently, if $N - 1$ players each have a minimum gap of Δ_{\min} , their collective exploration imposes $\Omega(N \log T / \Delta_{\min}^2)$ total collisions – and thus regret – upon the remaining player.

2 Related Work

Multi-player Bandit Learning The multi-player bandit problem involves multiple decentralized players interacting

with a shared multi-armed bandit environment. When players pull the same arm simultaneously, a collision occurs, resulting in a loss. In settings where arm means vary across players, the benchmark is typically the *maximum weight matching*, and the *system regret* – defined as the cumulative reward loss summed over all players – measures performances. Tibrewal et al. (2019) proposed an ETC type algorithm where players exploit the best-estimated matching, Mehrabian et al. (2020) combined forced collisions for implicit communication with matching eliminations, Shi et al. (2021) adapted the CUCB algorithm (Chen, Wang, and Yuan 2013) to this setting. For a comprehensive survey, see Bourcier and Perchet (2024).

While our model shares the reward and collision structure of heterogeneous multiplayer bandits, our benchmark differs fundamentally: we use the *core* of the housing market — a game-theoretic solution distinct from maximum weight matching – and define *individual regret* per player rather than aggregate system regret. This ensures fairness but introduces distinct algorithmic challenges.

Competing Bandits in Two-sided Matching Markets

Bandit problems in matching markets were first formalized by Das and Kamenica (2005), with subsequent works (Liu, Mania, and Jordan 2020; Liu et al. 2021; Sankararaman, Basu, and Sankararaman 2021; Kong and Li 2023; Lin et al. 2025) exploring this model to achieve stable matching (Gale and Shapley 1962). In these two-sided markets, players (with unknown utilities) and arms (with known preferences) interact – when multiple players pull the same arm, only the top-ranked player by the arm’s preference gets matched while others face collisions, with individual regrets defined accordingly. Since multiple stable matchings may exist, regrets are typically measured against either player-optimal or player-pessimal stable matchings.

In contrast, we study housing markets (one-sided matching) where only players have preferences over arms. This creates a distinct collision structure: when multiple players propose to the same arm, all receive zero reward. Moreover, the core matching in a housing markets is unique, ensuring an unambiguous regret definition. These differences from two-sided markets introduce novel algorithmic challenges.

3 Problem Setting

Denote N as the number of players in the market, and every player has an initial endowed arm. Let $\mathcal{N} = \{p_1, p_2, \dots, p_N\}$ be the player set and $\mathcal{A} = \{a_1, a_2, \dots, a_N\}$ be the arm set. The preferences of players on arms can be represented through a utility matrix U , where $U(i, j) \in [0, 1]$ denotes the preference of player p_i on arm a_j . If $U(i, j) > U(i, j')$, player p_i prefers arm a_j over $a_{j'}$, and we denote as $a_j \succ_i a_{j'}$. In this paper, we assume all preferences are distinct, i.e., $U(i, j) \neq U(i, j')$ for any different arms $a_j \neq a_{j'}$. Additionally, the utility vectors can vary entirely between players, reflecting heterogeneous preferences over arms. In each round $t = 1, 2, \dots$, every player p_i proposes to an arm $A_i(t) \in \mathcal{A}$. Each arm a_j then receives applications from the set of players $A_j^{-1}(t) := \{p_i : A_i(t) = a_j\}$. Since arm a_j is initially

owned by player p_j , we assume only p_j observes the application profile $A_j^{-1}(t)$, while other players remain unaware of this information. Arms are indivisible goods with no preferences over players, if multiple players propose to a_j in the same round, i.e., $|A_j^{-1}(t)| > 1$, a collision occurs, and a_j is not matched to any player. In this case, the successfully matched player for a_j is $\bar{A}_j^{-1}(t) = \emptyset$, and each proposing player $p_i \in A_j^{-1}(t)$ receives $\bar{A}_i(t) = \emptyset$, along with a deterministic reward $X_i(t) = 0$ indicating they are blocked. Conversely, if only one player p_i proposes to a_j , i.e., $|A_j^{-1}(t)| = 1$, the match succeeds: $\bar{A}_j^{-1}(t) = p_i$, and p_i receives a random reward $X_i(t)$ characterizing its matching experience in this round, which we assume is 1-subgaussian with expectation $\mathcal{U}(i, \bar{A}_i(t))$. The set of all successfully matched player-arm pairs in round t forms a matching, denoted μ_t , where $\mu_t(p_i) = \bar{A}_i(t)$.

The core is a fundamental solution concept in housing markets (Shapley and Scarf 1974), which is the set of feasible allocations where no coalition of agents can benefit by breaking away from the grand coalition. Formally,

Definition 1 (Core). A matching μ is in the core if no coalition $\mathcal{S} \subseteq \mathcal{N}$ can block μ , i.e., there does not exist an alternative matching μ' such that $\mu'(p_i) \succ_i \mu(p_i), \forall p_i \in \mathcal{S}$.

In housing markets with strict preferences, the core is always nonempty, and consists of a unique matching (Roth and Postlewaite 1977), which we refer to as the *core matching* and denote by μ^* . Our objective is to learn μ^* while minimizing the *core regret* for each player $p_i \in \mathcal{N}$. This regret is defined as the cumulative difference, over T rounds, between the expected rewards from being matched to $\mu^*(p_i)$ and the rewards actually obtained by p_i :

$$\text{Reg}_i(T) = T \cdot \mathcal{U}(i, \mu^*(p_i)) - \mathbb{E} \left[\sum_{t=1}^T X_i(t) \right]. \quad (1)$$

The expectation is taken over the randomness of the received reward and the players' strategy.

Offline Top-Trading-Cycle Algorithm In the *offline* setting where all players know their exact preferences, the Top-Trading-Cycle (TTC) algorithm (Shapley and Scarf 1974) efficiently computes the unique core allocation of the housing market – a stable assignment where no coalition of players can improve their outcomes through mutual exchanges. The centralized TTC proceeds iteratively. First, each player identifies her most preferred available arm. Second, a directed graph is formed where players point to owners of their top choices. Third, at least one cycle (including possible self-cycles) is identified and implemented, with involved players exiting the market. This process terminates within N steps, guaranteed to produce a core allocation.

4 Decentralized Algorithm

In this section, we present an Explore-then-Commit (ETC) algorithm for decentralized housing markets, operating in two phases. First, players explore arms in a round-robin fashion to estimate their preference rankings (Line 2-25), ensuring reliable ordinal estimates for matching. Next, players use these estimates to identify their assigned arms in

the core matching via a decentralized adaptation of the *You Request My House, I Get Your Turn (YRMH-IGYT)* mechanism (Abdulkadiroğlu and Sönmez 1999) (Line 26-51). Once matched, players commit exclusively to these arms in all future rounds.

The first phase of the Algorithm 1 consists of multiple sub-phases $\ell = 1, 2, \dots$, each lasting $2^\ell + N$ rounds (Line 2-25). Each sub-phase ℓ begins with an exploration stage of 2^ℓ rounds (Line 5-8), followed by N communication rounds (Line 13-24). During the exploration stage, players aim to gather sufficient observations about the arms. The subsequent communication rounds allow them to verify whether all N players have accurately learned their preferences. Once this condition is met, players exit the exploration phase and proceed to the second phase, where they identify arms in the core matching (Line 21).

During the exploration stage of each sub-phase, players follow a round-robin strategy to propose to arms (Line 6). Since each player initially owns a distinct endowed arm with a unique index, this ensures that no two players select the same arm simultaneously. As a result, all proposals are successfully accepted without collisions. When player p_i receives an observation from the selected arm $A_i(t)$, it updates both the estimated preference value $\hat{U}(i, A_i(t))$ and the observation count $T_{i, A_i(t)}$ for that arm (Line 7). The update rule is as follows,

$$\hat{U}(i, A_i(t)) = \frac{\hat{U}(i, A_i(t)) \cdot T_{i, A_i(t)} + X_{i, A_i(t)}(t)}{T_{i, A_i(t)} + 1}, \quad (2)$$

$$T_{i, A_i(t)} = T_{i, A_i(t)} + 1. \quad (3)$$

At the end of the exploration phase, players construct a confidence set for their estimated preference values using collected observations. Specifically, for player p_i , the confidence interval for the preference value of arm a_j is defined by an upper bound (UCB) and a lower bound (LCB), given as

$$\begin{aligned} \text{UCB}_{i,j} &= \hat{U}(i, j) + \sqrt{\frac{6 \log T}{\max\{T_{i,j}, 1\}}}, \\ \text{LCB}_{i,j} &= \hat{U}(i, j) - \sqrt{\frac{6 \log T}{\max\{T_{i,j}, 1\}}}. \end{aligned} \quad (4)$$

When the confidence intervals of two arms $a_j, a_{j'}$ become disjoint, that is, when $\text{LCB}_{i,j} > \text{UCB}_{i,j'}$ or vice versa, player p_i can establish its strict preference ordering between them. Once p_i successfully determines a complete preference ranking over all N arms (Line 10), it sets the flag $P_\ell^{(i)}$ to True for the current sub-phase and records this permutation as its estimated preference ranking σ_i .

Communication rounds enable players verify if all participants have accurately estimated their preference rankings. In the i -th communication round of every sub-phase, player p_i 's endowed arm serves as a broadcast channel. A player p_j proposes to p_i 's arm if and only if it has an accurate ranking (Line 16); otherwise, it abstains (Line 18). While this may cause collisions, the goal is not to resolve them but to convey status information. If p_i observes proposals from all N players including itself (Line 20), it infers universal estimation success and triggers the transition to the second phase using its estimated ranking σ_i (Line 21).

Algorithm 1: Decentralized Explore-then-YRMH-IGYT
 (from view of player p_i)

Require: player set \mathcal{P} and arm set \mathcal{A} , horizon T .

- 1: Initialize: $\hat{U}(i, j) = 0, T_{i,j} = 0, \forall j \in [N]$.
- 2: //Phase 1, learn the preferences
- 3: **for** $\ell = 1, 2, \dots$ **do**
- 4: $P_\ell^{(i)} = \text{False}$ //whether the preference is well-estimated
- 5: **for** $t = \sum_{\ell'=1}^{\ell-1} (2^{\ell'} + N) + 1, \dots, \sum_{\ell'=1}^{\ell-1} (2^{\ell'} + N) + 2^\ell$ **do**
- 6: $A_i(t) = a_{(i+t-1)\%N+1}$.
- 7: Observe $X_{i,A_i(t)}(t)$, update $\hat{U}(i, A_i(t)), T_{i,A_i(t)}$.
- 8: **end for**
- 9: Compute $\text{UCB}_{i,j}$ and $\text{LCB}_{i,j}$ for each $j \in [N]$.
- 10: **if** $\exists \sigma$ such that $\text{LCB}_{i,\sigma_k} > \text{UCB}_{i,\sigma_{k+1}}, \forall k \in [N]$ **then**
- 11: $P_\ell^{(i)} = \text{True}$ and $\sigma_i = \sigma$.
- 12: **end if**
- 13: **for** $t = \sum_{\ell'=1}^{\ell-1} (2^{\ell'} + N) + 2^\ell + 1, \dots, \sum_{\ell'=1}^{\ell-1} (2^{\ell'} + N) + 2^\ell$ **do**
- 14: $t' = t - \sum_{\ell'=1}^{\ell-1} (2^{\ell'} + N) - 2^\ell$
- 15: **if** $P_\ell^{(i)} == \text{True}$ **then**
- 16: $A_i(t) = a_{t'}$.
- 17: **else**
- 18: $A_i(t) = \emptyset$.
- 19: **end if**
- 20: **if** $i == t'$ and $|A_i^{-1}(t)| == N$ **then**
- 21: Enter in Phase 2 with $\sigma_i = (\sigma_{i,1}, \dots, \sigma_{i,N})$.
- 22: $t_1 = \sum_{\ell'=1}^{\ell-1} (2^{\ell'} + N)$ //Phase 1 ends.
- 23: **end if**
- 24: **end for**
- 25: **end for**
- 26: //Phase 2, find the core matching arm with σ_i
- 27: $t = t_1 + 1$
- 28: $F_i = \text{True}$ //Whether the initial endowed arm is still available in the market, this flag is known to all players
- 29: $k = 1, \mathcal{A}_k = \{a_i : F_i == \text{True}\}$.
- 30: $\text{Propose}(i) = \text{False}$. //A flag indicated whether player p_i proposed to some arm in the k -th epoch.
- 31: $i_{\min} = \min(\{i | a_i \in \mathcal{A}_k\})$.
- 32: **while** $|\mathcal{A}_k| > 0$ and $F_i == \text{True}$ **do**
- 33: $j_{\min}^{(i)} = \min \{j | a_{\sigma_{i,j}} \in \mathcal{A}_k\}$. //The best available arm for player p_i .
- 34: **if** $t == t_1 + 1$ or $k(t) == k(t-1) + 1$ **then**
- 35: **if** $i == i_{\min}$ **then**
- 36: $A_i(t) = a_{j_{\min}^{(i)}}$.
- 37: $\text{Propose}(i) = \text{True}$.
- 38: **end if**
- 39: **else if** $|A_i^{-1}(t-1)| > 0$ **then**
- 40: $A_i(t) = a_{j_{\min}^{(i)}}$, $\text{Propose}(i) = \text{True}$.
- 41: $j_r = \{j | a_j = A_i^{-1}(t-1)\}$.
- 42: **end if**
- 43: **if** ($\text{Propose}(i) == \text{True}$ and $|A_i^{-1}(t)| > 0$) or ($F_{j_r} == \text{False}$) **then** //In a cycle
- 44: $F_i = \text{False}, A_i(s) = a_{j_{\min}^{(i)}}, \forall s > t$.
- 45: **end if**
- 46: $\mathcal{B} = \{a_i : F_i == \text{True}\}$.
- 47: **if** $\mathcal{B} \neq \mathcal{A}_k$ **then** //The available arm set is updated
- 48: $k = k + 1, \mathcal{A}_k = \mathcal{B}, \text{Propose}(i) = \text{False}$.
- 49: $i_{\min} = \min \{i | a_i \in \mathcal{A}_k\}$.
- 50: **end if**
- 51: **end while**

The second phase (Line 26-51) adapts the YRMH-IGYT algorithm (Abdulkadiroğlu and Sönmez 1999) to our decentralized setting for players to progressively identify their core allocations. The process unfolds in successive sub-phases. In each, the unassigned player with the smallest index initiates a proposal chain by proposing to their most preferred remaining arm (Lin 36-37). The recipient of a proposal then proposes to their own top choice among available arms, and this sequence continues. A top trading cycle is completed when a proposal is received by a player who is already part of the current chain (Line 43). The player who closes the cycle – the first to receive a repeat proposal – then marks their endowed arm as allocated. This allocation status propagates backward through the chain. All players involved in the cycle permanently fix their allocations and remove their own endowed arms from the market (Line 44). The process repeats among the remaining players until all arms are allocated, thus achieving a stable core matching in a fully decentralized manner.

Remark 1. *Algorithm 1 relies on standard multiplayer bandit assumptions for implicit coordination: (1) Unique, globally-known arm IDs, which also serve as player identifiers for round-robin exploration and proposal-chain initiation. (2) A shared global clock to synchronize the transition to Phase 2. (3) Common knowledge of the available arms in each round, enabling players to reconstruct the flag F_i .*

Theoretical Analysis

Prior to presenting the formal regret analysis of Algorithm 1, we introduce the notion of preference gaps to quantify the intrinsic difficulty of the learning problem.

Definition 2 (Minimum Preference Gap). *For each player p_i and arm $a_j \neq a_{j'}$, let $\Delta_{i,j,j'} = \mathcal{U}(i, j) - \mathcal{U}(i, j')$ be the preference gap of p_i between a_j and $a_{j'}$. Let r_i be the preference ranking of player p_i and $r_{i,k}$ be the k -th preferred arm in p_i 's ranking for $k \in [N]$. Define $\Delta_{\min} = \min_{i \in [N]; k \in [N]} \Delta_{i,r_{i,k},r_{i,k+1}}$ as the minimum preference gap among all workers and their preferences over the arms. Δ_{\min} is non-negative since all preferences are distinct.*

We now present the upper bound for the core regret for each player by following Algorithm 1.

Theorem 1. *Following Algorithm 1, the core regret of each player $p_i \in \mathcal{N}$ satisfies*

$$\begin{aligned} \text{Reg}_i(T) &\leq \left(\frac{192N \log T}{\Delta_{\min}^2} + N \log \left(\frac{192N \log T}{\Delta_{\min}^2} \right) + 3N^2 \right) \cdot \Delta_{i,\max} \\ &= \mathcal{O} \left(\frac{N \log T}{\Delta_{\min}^2} \right) \end{aligned}$$

The regret bound consists of three key components. First, the cumulative regret from exploration rounds in phase 1. Second, the regret from communication rounds in phase 1. The final term combines two distinct elements: (1) the regret during phase 2, where the YRMH-IGYT mechanism guarantees convergence to the core matching within at most N epochs (with each epoch requiring $\leq N$ rounds to identify a top trading cycle), and (2) the regret contribution from rare concentration failure events.

For convenience, let $\hat{U}^{(t)}(i, j), T_{i,j}^{(t)}, UCB_{i,j}^{(t)}, LCB_{i,j}^{(t)}$ be the value of $\hat{U}(i, j), T_{i,j}, UCB_{i,j}, LCB_{i,j}$ at the end of round t . Define

$$\mathcal{F} = \left\{ \exists t \in [T], i \in [N], j \in [N] : |\hat{U}^{(t)}(i, j) - U(i, j)| > \sqrt{\frac{6 \log T}{T_{i,j}^{(t)}}} \right\}$$

as the bad event that some preference is not estimated well during the horizon. Since all players would communicate whether they have a precise enough estimation after each sub-phase of phase 1, and determine to enter phase 2 once they find every player include themselves have good-enough estimations, we can conclude that all players would enter phase 2 at the same time. Denote ℓ_{\max} as the largest sub-phase number of phase 1. That is to say, players enter in phase 2 at the end of sub-phase ℓ_{\max} . We then provide the proof of Theorem 1 as follows.

Proof of Theorem 1. Let $\Delta_{i,\max}$ be the maximum core regret that may be suffered by player p_i in all rounds, we have $\Delta_{i,\max} \leq 1$ by assumption on the utility matrix. The core regret of each player p_i by following Algorithm 1 satisfies

$$\begin{aligned} \text{Reg}_i(T) &= \mathbb{E} \left[\sum_{t=1}^T (\mathbf{U}(i, \mu^*(p_i)) - X_i(t)) \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{\mu_t(i) \neq \mu^*(p_i)\} \cdot \Delta_{i,\max} \right] \end{aligned} \quad (5)$$

$$\begin{aligned} &\leq \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{\mu_t(i) \neq \mu^*(p_i)\} \middle| \neg \mathcal{F} \right] \cdot \Delta_{i,\max} \\ &\quad + \mathbb{P}(\mathcal{F}) \cdot T \cdot \Delta_{i,\max} \\ &\leq \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{\mu_t(i) \neq \mu^*(p_i)\} \middle| \neg \mathcal{F} \right] \cdot \Delta_{i,\max} + 2N^2 \Delta_{i,\max} \end{aligned} \quad (6)$$

$$\begin{aligned} &\leq \mathbb{E} \left[\sum_{\ell=1}^{\ell_{\max}} (2^\ell + N) + N^2 \middle| \neg \mathcal{F} \right] \cdot \Delta_{i,\max} + 2N^2 \Delta_{i,\max} \quad (7) \\ &\leq \left(\frac{192N \log T}{\Delta_{\min}^2} + N \log \left(\frac{192N \log T}{\Delta_{\min}^2} \right) \right) \cdot \Delta_{i,\max} \\ &\quad + 3N^2 \Delta_{i,\max}, \end{aligned} \quad (8)$$

where Eq.(5) comes from the fact that in a housing market, there is a unique core matching and hence a unique core matching partner $\mu^*(p_i)$ for player p_i . Eq.(6) holds based on Lemma 1. Eq.(7) holds according to Algorithm 1 and the fact that we need at most N^2 rounds for the YRMH-IGYT procedure to reach the equilibrium (Lemma 2). Eq.(8) follows from Lemma 3. \square

The following technical lemmas underlying our analysis. Proofs are provided in the extended version. We begin with Lemma 1, which bounds the probability of erroneous preference estimation.

Lemma 1 (Bad Concentration Event). $\mathbb{P}(\mathcal{F}) \leq \frac{2N^2}{T}$.

Lemma 2 shows that each player p_i finds their core-matched arm $\mu^*(p_i)$ within at most N^2 rounds in Phase 2.

Lemma 2. *Conditional on $\neg \mathcal{F}$, at most N^2 rounds are needed in phase 2 before $A_i(t) = \mu^*(p_i)$, and in all of the following rounds, $A_i(t)$ would not be updated and p_i would always be successfully accepted by $\mu^*(p_i)$.*

Lemma 3 bounds the duration of the exploration phase.

Lemma 3. *Condition on $\neg \mathcal{F}$, phase 1 will proceed in at most ℓ_{\max} sub-phases where*

$$\ell_{\max} = \min \left\{ \ell : \sum_{\ell'=1}^{\ell} 2^{\ell'} \geq \frac{96N \log T}{\Delta_{\min}^2} \right\}, \quad (9)$$

which implies that $\sum_{\ell'=1}^{\ell_{\max}} 2^{\ell'} \leq \frac{192N \log T}{\Delta_{\min}^2}$ and $\ell_{\max} = \log \left(\frac{192N \log T}{\Delta_{\min}^2} \right)$ since the sub-phase length grows exponentially. And all players will enter in phase 2 simultaneously at the end of sub-phase ℓ_{\max} .

5 Regret Lower Bound

We establish a regret lower bound by adapting the method of Auer et al. (2002) to our multi-player setting, demonstrating the tightness of the regret upper bound of Algorithm 1. Let $\text{Reg}(T; \nu, \pi)$ be the cumulative expected regret of policy π over all players and time T , for an instance with arm distributions $\nu = \{\nu_{ij} : i, j \in [N]\}$. Denote by \mathcal{P} the set of all probability distributions with support in $[0, 1]$. For any $\nu \in \mathcal{P}$, define $D_{\inf}(\nu, x, \mathcal{P}) = \inf_{\nu' \in \mathcal{P}} \{D(\nu, \nu') : \rho(\nu') > x\}$, where $\rho : \mathcal{P} \rightarrow \mathbb{R}$ maps a distribution to its mean and $D(\cdot, \cdot)$ is the KL divergence.

Definition 3 (Uniformly Consistent Policies). *A policy π is uniformly consistent if and only if for all $\nu \in \mathcal{P}$, all $\alpha \in (0, 1)$, the regret $\limsup_{T \rightarrow \infty} \frac{\text{Reg}(T; \nu, \pi)}{T^\alpha} = 0$.*

This ensures the policy is not overly tuned to the current instance at the expense of performance in others, ensuring robust performance across different problem instances (Auer et al. 2002; Lattimore and Szepesvári 2020).

Single Top Trading Cycle Bandits Our regret lower bound applies to a subclass of bandits characterized by a single top trading cycle (STTCB). In these instances, each player has a unique top-ranked arm, and there exists a permutation σ of players such that $\mu^*(p_{\sigma_i}) = a_{\sigma_{i+1}}$ for $i = 1, \dots, N-1$ and $\mu^*(p_{\sigma_N}) = a_{\sigma_1}$. This structure ensures all players form one top trading cycle in the core matching. For each player p_i , define the gap $\Delta_j^{(i)} = \mathbf{U}(i, \mu^*(p_i)) - \mathbf{U}(i, j)$ and the minimum gap $\Delta_{\min}^{(i)} = \min_{a_j \neq \mu^*(p_i)} \mathbf{U}(i, \mu^*(p_i)) - \mathbf{U}(i, j)$, which are non-negative in STTCB instances.

Lemma 4 (Regret Decomposition). *For an STTCB instance $\nu = \{\nu_{i,j} : i \in [N], j \in [N]\}$, and any uniformly consistent policy π , player p_i for $i \in [N]$ the following holds*

$$\begin{aligned} \text{Reg}_i(T; \nu, \pi) &\geq \max \left\{ \sum_{i' \neq i} \Delta_{\min}^{(i)} \mathbb{E}_{\nu, \pi} \left[N_{\mu^*(p_i)}^{(i')} (T) \right], \right. \\ &\quad \left. \sum_{a_j \neq \mu^*(p_i)} \Delta_j^{(i)} \mathbb{E}_{\nu, \pi} \left[N_j^{(i)} (T) \right] \right\}. \end{aligned}$$

For an STTCB instance, when the p_i 's core-matched partner $\mu^*(p_i)$ is allocated to another player, the optimal outcome is for p_i to receive its second-best arm. While this event occurs infrequently, it provides a sufficient bound for the first term above.

Theorem 2. For any player p_i , $i \in [N]$, and under any decentralized uniformly consistent algorithm π , the performance on an STTCB instance ν satisfies

$$\begin{aligned} & \liminf_{T \rightarrow \infty} \frac{\text{Reg}_i(T; \nu, \pi)}{\log T} \\ & \geq \max \left\{ \sum_{i' \neq i} \frac{\Delta_{\min}^{(i)}}{D_{\text{inf}}(\nu_{i'}, \mu^*(p_i), \mathcal{U}(i', \mu^*(p_i)), \mathcal{P})}, \right. \\ & \quad \left. \sum_{a_j \neq \mu^*(p_i)} \frac{\Delta_j^{(i)}}{D_{\text{inf}}(\nu_{i,j}, \mathcal{U}(i, \mu^*(p_i)), \mathcal{P})} \right\}. \end{aligned} \quad (10)$$

Proof Idea. For a given STTCB instance ν , we construct a confounding alternative ν' that differs only in the utility of a single player-arm pair (p_i, a_j) , where a_j is the core-matched arm of p_i in ν' but not in ν . To identify the true instance and avoid linear regret, any algorithm must gather enough evidence to distinguish ν from ν' , necessitating a number of samples inversely proportional to their KL-divergence. This directly implies a logarithmic regret lower bound. \square

The detailed proof of Theorem 2 can be found in the extended version. The following corollary demonstrates that the $N \log T / \Delta_{\min}^2$ dependence in Theorem 1 cannot be improved, thus establishing $\Theta(N \log T / \Delta_{\min}^2)$ as the minimax regret rate for all players.

Corollary 1. There exists an STTCB instance with Bernoulli rewards, where the regret of player p_i is lower bounded as $\Omega(\frac{N \log T}{\Delta_{\min}^2})$.

Reward heterogeneity necessitates $\Omega(\log T / \Delta_{\min}^2)$ explorations for the player with minimum gap Δ_{\min} , during which players with substantial gaps ($\Delta_{\min}^{(i')} = \Omega(1)$) incur significant regret. This core observation underpins our proof.

From Corollary 1, we can see that the regret upper bound (Theorem 1) for Algorithm 1 matches the regret lower bound, showing the significance of our proposed algorithm.

6 Centralized Anytime Algorithm

While Algorithm 1 achieves a regret upper bound that matches the lower bound in Theorem 2 for all players, it requires prior knowledge of the time horizon T to properly construct confidence intervals. Although one could apply the doubling trick (Auer et al. 1995) to make the algorithm anytime, this approach leads to non-monotonic expected instantaneous regret. Such behavior could raise concerns among players about the algorithm's consistency during execution.

In this section, we develop an adaptive anytime algorithm that operates through a centralized platform capable of aggregating player preferences and regulating allocations without observing individual reward realizations. The algorithm employs the principle of optimism in the face of uncertainty, where each player maintains upper confidence bounds (UCB) to rank arms and submits these rankings to the platform every round. The platform then computes a core matching by applying the TTC algorithm to the collected preferences, and players subsequently pull their assigned arms.

Algorithm 2: Centralized Anytime UCB

Require: player set \mathcal{P} and the corresponding arm set \mathcal{A} .

```

1: for  $t = 1, \dots$  do
2:   The platform receives ranking  $\hat{r}_{i,t}$  from all players  $p_i$ .
3:   The platform computes a core matching  $\mu_t$  using the offline TTC algorithm with the ranking  $\hat{r}_{i,t}$ .
4:   for  $i = 1, \dots, N$  do//Players pull arms simultaneously
5:      $A_i(t) = \mu_t(p_i)$ .
6:     Update  $\hat{U}^{(t)}(i, A_i(t))$  and  $u^{(t)}(i, A_i(t))$  according to Eq.(2) and Eq.(11).
7:     Compute the current ranking  $\hat{r}_{i,t+1}$  according to  $u^{(t)}(i, \cdot)$ .
8:   end for
9: end for

```

We begin by formally defining the UCB estimation method used by individual players and introducing several technical concepts necessary for the analysis. Building on this foundation, we derive a regret upper bound for the centralized approach. A key feature of this framework is that the platform's coordination eliminates potential conflicts between player proposals, allowing us to assume collision-free operation throughout the learning process.

At iteration t , when a player p_i gets matched to arm $A_i(t)$, it would update the estimated preference value $\hat{U}(i, A_i(t))$ and the observed time $T_{i,A_i(t)}$ for arm $A_i(t)$ according to Eq.(2) in Section 4.

Then the upper confidence bound, which is called the *index* for each (i, j) -th entry of the utility matrix U is computed as

$$u^{(t)}(i, j) = \begin{cases} \infty & \text{if } T_{i,j}(t) = 0, \\ \hat{U}^{(t)}(i, j) + \sqrt{\frac{3 \log t}{2T_{i,j}(t-1)}}, & \text{otherwise.} \end{cases} \quad (11)$$

Each player p_i uses the index to compute a preference ranking $\hat{r}_{i,t+1}$, where arms are ordered by their UCBs in a decreasing order (e.g., $\arg \max_j u^{(t)}(i, j)$ is ranked first).

Let $T_\mu(t)$ denote the number of times a matching μ is played by time t . A matching is *achievable* at time t if it is core-stable according to the current estimated rankings $\{\hat{r}_{i,t}\}_{i \in [N]}$. A matching is *truly core-stable* if it is core-stable under the true utility matrix U . For player p_i and arm a_j , define $M_{i,j}$ as the set of achievable (but not truly core-stable) matchings where p_i is matched to a_j . The *achievable preference gap* is defined as $\bar{\Delta}_{i,j} = U(i, \mu^*(p_i)) - U(i, j)$.

Since the truly core-stable matching μ^* incurs zero regret, we bound the regret of player p_i as follows:

$$\begin{aligned} \text{Reg}_i(T) & \leq \sum_{j: \bar{\Delta}_{i,j} > 0} \bar{\Delta}_{i,j} \left(\sum_{\mu \in M_{i,j}} \mathbb{E} T_\mu(T) \right) \\ & \leq \max_j \bar{\Delta}_{i,j} \left(\sum_{\mu \in M} \mathbb{E} T_\mu(T) \right), \end{aligned} \quad (12)$$

where M is the set of all achievable but non-truly core-stable matchings, i.e., $M = \cup_{(i,j) \in [N] \times [N]} M_{i,j}$. If a matching μ is not truly core-stable under the true utility matrix U , a *blocking coalition* must exist – a subset of players who can reassign their endowed arms to strictly improve every member's utility. For any non-truly-core-stable matching $\mu \in M$,

there exists a blocking coalition \mathcal{B}_μ such that players outside \mathcal{B}_μ keep their assignments under μ^* . Within \mathcal{B}_μ , at least one player p_m must be part of a *blocking triplet* $(p_m, a_n, a_{n'})$ where under μ , p_m is assigned a_n , $U(m, n) < U(m, n')$ in truth yet the UCB index satisfies $u(m, n) > u(m, n')$. If no such triplet existed, then for every matched pair (p_m, a_n) and every $a_{n'}$ with $U(m, n) < U(m, n')$, we would have $u(m, n) < u(m, n')$. Since μ is non-core-stable under U , there must exist a reallocation within \mathcal{B}_μ that strictly improves every member with respect to U – and, given the UCB ordering, also with respect to the UCB index. This would contradict the UCB-core-stability of the implemented matching. Let $Q_{m,n}$ be the collection of all such triplets, whenever a non-truly-core-stable $\mu \in M$ is implemented, at least one player-arm pair $(p_m, a_n) \in \mathcal{B} := \cup_{\mu \in M} \mathcal{B}_\mu$ must realize a triplet from $Q_{m,n}$. This structure yields the following regret bound.

Theorem 3. *Following Algorithm 2, the core regret of player p_i up to time T satisfies*

$$\begin{aligned} \text{Reg}_i(T) &\leq \max_j \bar{\Delta}_{i,j} \left[\sum_{(p_m, a_n) \in \mathcal{B}} \sum_{(p_m, a_n, a_{n'}) \in Q_{m,n}} \left(5 + \frac{6 \log T}{\Delta_{m,n',n}^2} \right) \right] \\ &\leq \max_j \bar{\Delta}_{i,j} \left(5N^3 + 12 \frac{N^2 \log T}{\Delta_{\min}^2} \right) = \mathcal{O} \left(\frac{N^2 \log T}{\Delta_{\min}^2} \right). \end{aligned}$$

Theorem 3 offers a centralized problem-dependent $\mathcal{O} \left(\frac{N^2 \log T}{\Delta_{\min}^2} \right)$ upper bound guarantees on the core regret of each player p_i , and the proposed UCB-type algorithm is anytime and adaptive in the sense that the players do not need to know the time horizon T and the minimum preference gap Δ_{\min} in advance.

Proof. Let $T_{m,n,n'}(T)$ be the number of times player p_m pulls arm a_n while $(p_m, a_n, a_{n'})$ forms a blocking triplet. Since every time a matching $\mu \in M$ is implemented, at least one such blocking triplet $(p_m, a_n, a_{n'}) \in Q_{m,n}$ occurs for some $(p_m, a_n) \in \mathcal{B}$, we have

$$\sum_{\mu \in M} T_\mu(T) \leq \sum_{(p_m, a_n) \in \mathcal{B}} \sum_{(p_m, a_n, a_{n'}) \in Q_{m,n}} T_{m,n,n'}(T). \quad (13)$$

From the analysis above, we know that when $(p_m, a_n, a_{n'})$ is blocking, we have the UCB index $u(m, n) > u(m, n')$ while $U(m, n) < U(m, n')$ according to the ground-truth utility matrix. Standard analysis for the single player UCB (Bubeck, Cesa-Bianchi et al. 2012) shows that

$$\mathbb{E} T_{m,n,n'}(T) \leq 5 + \frac{6 \log T}{\Delta_{m,n',n}^2}. \quad (14)$$

Combining Eq.(12), (13) and (14) we get the first inequality on the regret upper bound.

When we consider the triplet set composed of all possible $U(m, n) < U(m, n')$, we have $|Q_{m,n}| \leq N$ and

$$\sum_{n': U(m,n) < U(m,n')} \frac{1}{\Delta_{m,n',n}^2} \leq \sum_{n'=1}^N \frac{1}{(n')^2 \Delta_{\min}^2} \leq \frac{2}{\Delta_{\min}^2},$$

and hence

$$\begin{aligned} &\max_j \bar{\Delta}_{i,j} \left[\sum_{(p_m, a_n) \in \mathcal{B}} \sum_{(p_m, a_n, a_{n'}) \in Q_{m,n}} \left(5 + \frac{6 \log T}{\Delta_{m,n',n}^2} \right) \right] \\ &\leq \max_j \bar{\Delta}_{i,j} \left[\sum_{(p_m, a_n) \in \mathcal{B}} \left(5N + \frac{12 \log T}{\Delta_{\min}^2} \right) \right]. \end{aligned}$$

As there are at most N^2 possible player-arm pairs in the blocking coalition, the second inequality is concluded. \square

7 Conclusion and Discussion

This paper introduces a novel online learning framework for housing markets with uncertain preferences, unifying two key objectives: core-stability and sample efficiency. We formalize this through core regret as our performance metric and propose two algorithms that bridge multi-armed bandit techniques with housing market mechanisms. These algorithms adapt to both centralized and decentralized settings, with guarantees for fixed-horizon and anytime environments. A matching regret lower bound proves the order-optimality of the decentralized algorithm.

There are many additional questions that can be studied in this model.

Housing Markets with Existing Tenants This paper focuses on housing markets where each player initially possesses one endowed arm, maintaining an equal number of players and arms. However, real-world housing markets typically involve more complex scenarios: existing tenants with endowed arms coexist with new applicants lacking initial allocations, while vacant houses (free arms) may also be available (Abdulkadiroğlu and Sönmez 1999). In these generalized settings, the core loses its uniqueness, raising two fundamental challenges: first, identifying an appropriate tractable solution concept to serve as a benchmark, and second, determining whether efficient learning algorithms can achieve sublinear regret in this more realistic but complicated environment.

Housing Markets with Indifference While our current analysis assumes strict player preferences over arms, real-world housing markets often involve indifference between options. Our learning algorithm achieves a regret bound of $\Theta(N \log T / \Delta_{\min}^2)$, which becomes vacuous when $\Delta_{\min} = o(1/\sqrt{T})$ – precisely when indifference between arms creates vanishing preference gaps. This limitation highlights the need for new algorithmic approaches that can handle indifferences in housing market allocations, presenting an important direction for future research.

Incentive Compatibility in the Learning Setting While the top trading cycle algorithm is strategy-proof under known, deterministic preferences, its incentive compatibility remains unclear in learning settings where preferences must be discovered dynamically. In particular, the decentralized nature of these interactions may create opportunities for strategic manipulation through learning behavior. Understanding whether and how players can exploit the learning process to achieve better outcomes presents a significant open question for future research.

References

- Abdulkadiroğlu, A.; and Sönmez, T. 1999. House allocation with existing tenants. *Journal of Economic Theory*, 88(2): 233–260.
- Abdulkadiroğlu, A.; and Sönmez, T. 2003. School choice: A mechanism design approach. *American economic review*, 93(3): 729–747.
- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2): 235–256.
- Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E. 1995. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of IEEE 36th annual foundations of computer science*, 322–331. IEEE.
- Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E. 2002. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1): 48–77.
- Boursier, E.; and Perchet, V. 2024. A survey on multi-player bandits. *Journal of Machine Learning Research*, 25(137): 1–45.
- Bubeck, S.; Cesa-Bianchi, N.; et al. 2012. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1): 1–122.
- Chen, W.; Wang, Y.; and Yuan, Y. 2013. Combinatorial multi-armed bandit: General framework and applications. In *International conference on machine learning*, 151–159. PMLR.
- Das, S.; and Kamenica, E. 2005. Two-Sided Bandits and the Dating Market. In *IJCAI*, volume 5, 19.
- Echenique, F.; Miralles, A.; and Zhang, J. 2021. Constrained pseudo-market equilibrium. *American Economic Review*, 111(11): 3699–3732.
- Gale, D.; and Shapley, L. S. 1962. College admissions and the stability of marriage. *The American mathematical monthly*, 69(1): 9–15.
- Garg, J.; Tröbst, T.; and Vazirani, V. 2024. One-sided matching markets with endowments: equilibria and algorithms. *Autonomous Agents and Multi-Agent Systems*, 38(2): 40.
- Garivier, A.; Lattimore, T.; and Kaufmann, E. 2016. On explore-then-commit strategies. *Advances in Neural Information Processing Systems*, 29.
- Hylland, A.; and Zeckhauser, R. 1979. The efficient allocation of individuals to positions. *Journal of Political economy*, 87(2): 293–314.
- Kong, F.; and Li, S. 2023. Player-optimal stable regret for bandit learning in matching markets. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 1512–1522. SIAM.
- Lattimore, T.; and Szepesvári, C. 2020. *Bandit algorithms*. Cambridge University Press.
- Lin, S.; Mauras, S.; Merlis, N.; and Perchet, V. 2025. Stable Matching with Ties: Approximation Ratios and Learning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Liu, L. T.; Mania, H.; and Jordan, M. 2020. Competing bandits in matching markets. In *International Conference on Artificial Intelligence and Statistics*, 1618–1628. PMLR.
- Liu, L. T.; Ruan, F.; Mania, H.; and Jordan, M. I. 2021. Bandit learning in decentralized matching markets. *Journal of Machine Learning Research*, 22(211): 1–34.
- Mehrabian, A.; Boursier, E.; Kaufmann, E.; and Perchet, V. 2020. A practical algorithm for multiplayer bandits when arm means vary among players. In *International Conference on Artificial Intelligence and Statistics*, 1211–1221. PMLR.
- Roth, A. E. 2023. *Online and matching-based market design*. Cambridge University Press.
- Roth, A. E.; and Postlewaite, A. 1977. Weak versus strong domination in a market with indivisible goods. *Journal of Mathematical Economics*, 4(2): 131–137.
- Roth, A. E.; Sönmez, T.; and Ünver, M. U. 2004. Kidney exchange. *The Quarterly journal of economics*, 119(2): 457–488.
- Sankararaman, A.; Basu, S.; and Sankararaman, K. A. 2021. Dominate or delete: Decentralized competing bandits in serial dictatorship. In *International Conference on Artificial Intelligence and Statistics*, 1252–1260. PMLR.
- Shapley, L.; and Scarf, H. 1974. On cores and indivisibility. *Journal of mathematical economics*, 1(1): 23–37.
- Shi, C.; Xiong, W.; Shen, C.; and Yang, J. 2021. Heterogeneous multi-player multi-armed bandits: Closing the gap and generalization. *Advances in neural information processing systems*, 34: 22392–22404.
- Sönmez, T.; and Ünver, M. U. 2010. House allocation with existing tenants: A characterization. *Games and Economic Behavior*, 69(2): 425–445.
- Thompson, W. R. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4): 285–294.
- Tibrewal, H.; Patchala, S.; Hanawal, M. K.; and Darak, S. J. 2019. Distributed learning and optimal assignment in multiplayer heterogeneous networks. In *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, 1693–1701. IEEE.