# Strategic Uncertainty and Equilibrium Selection in Stable Matching Mechanisms: Experimental Evidence

## Marco Castillo*  Ahrash Dianat**

**Abstract**

We present experimental evidence on the interplay between strategic uncertainty and equilibrium selection in stable matching mechanisms. In particular, we apply a version of risk-dominance to compare the riskiness of "truncation" against other strategies that secure against remaining unmatched. By keeping subjects' ordinal preferences fixed while changing their cardinal representation, our experimental treatments vary the risk-dominant prediction. We find that both truth-telling and truncation are played more often when they are risk-dominant. In both treatments, however, truncation strategies are played more often in later rounds of the experiment. Our results also shed light on several open questions in market design.

\* Department of Economics, Texas A&M University, marco.castillo@tamu.edu

\*\* Department of Economics, University of Essex, a.dianat@essex.ac.uk

# 1 Introduction

Strategic uncertainty plays a crucial role in games with multiple equilibria. The central insight is as follows: when there are multiple rationalizable actions, players often face a tension between profitability and safety. In $2 \times 2$ coordination games, Harsanyi and Selten (1988) developed the concept of *risk-dominance* to capture the intuition that one equilibrium may appear more or less strategically risky than another equilibrium. Lab experiments have documented the predictive power of risk-dominance in simple settings (Cooper et al., 1990; Van Huyck et al., 1990). More recently, variations of risk-dominance have been fruitfully applied to other domains. In the infinitely repeated prisoners' dilemma, for instance, there is experimental evidence that cooperation is more likely to be sustained when it is both an equilibrium action and a risk-dominant action (Dal Bó and Fréchette, 2011). In the context of continuous-time games, Calford and Oprea (2017) show that a particular measure of risk-dominance does an impressive job of organizing their experimental data.

In this paper, we apply a version of risk-dominance to stable matching mechanisms. In a stable matching mechanism, participants in a two-sided market report rank-order lists of their preferences over match partners to a central authority. The central authority then uses the reported preferences to calculate the final matching. Crucially, the final matching is stable with respect to the reported preferences.[1] This environment induces a *preference-revelation game* in which the strategy space is the set of all possible ordinal preference lists.

In particular, we investigate how varying the "riskiness" of preference misrepresentation affects selection among stable matchings. By automating the side of the market that has a dominant strategy, we are able to model the strategic environment as a coordination game. For a very general class of preferences, this coordination game has at least two symmetric and Pareto-ranked equilibria in pure strategies: an equilibrium in "truncation" strategies (i.e., removing less preferred match partners from the tail end of a preference list) and an equilibrium in "permutation" strategies (i.e., switching the order of match partners in a preference list). Although the truncation equilibrium yields a higher payoff, the presence of strategic uncertainty makes truncation strategies less appealing. Intuitively, truncation generates a trade-off between the likelihood of matching and the quality of match partner (conditional on matching). Thus, an agent who plays a truncation strategy opens herself up to the possibility of remaining unmatched for some profile of other agents' preference reports.[2] This insight allows us to apply a version of risk-dominance to compare the riskiness

---

[1] A matching is said to be stable if no agent prefers remaining unmatched to her current allocation and no pair of agents both prefer each other to their current allocations.

[2] Even after removing strategic uncertainty, there is also the possibility of remaining unmatched due to "over-truncation" (i.e., playing the wrong kind of truncation strategy). However, over-truncation is not

of truncation against other strategies that secure against remaining unmatched.

We then study this coordination game in the lab, using the simplest possible setting in which this strategic tension arises. Each experimental market consists of four agents: two firms (implemented by computers) and two workers (human subjects). Subjects play 20 rounds of the preference-revelation game induced by the firm-optimal stable mechanism, with random and anonymous re-pairing across rounds.[3] We use an ordinal constellation of preferences such that each subject has two stable match partners.[4] However, these ordinal preferences can be represented by different cardinal utilities. Indeed, we have full freedom in choosing the payoff difference between two ordered alternatives. In our experiment, we choose cardinal representations such that our treatments vary whether the criterion of risk-dominance selects truth-telling or truncation. When remaining unmatched is particularly costly (to be precisely defined later), then truth-telling is risk-dominant. When an agent has a strong intensity of preference for her first choice partner (to be precisely defined later), then truncation is risk-dominant.

We now preview our main results. Overall, we find that truth-telling is the modal strategy. However, we observe several notable patterns in the experimental data. First, both truth-telling and truncation are played more often when they are risk-dominant. This result is robust to whether the treatment effect is measured at the subject or the session level. Second, in both treatments, truncation strategies are played more often in later rounds of the experiment. Since truncation is a necessary component of all payoff-dominant equilibria, this suggests that the salience of payoff-dominance as a selection criterion increases with subject experience.

In our setting, the set of Nash equilibrium outcomes of the preference-revelation game is identical to the set of stable outcomes of the underlying two-sided matching market. When attention is confined to stable outcomes, the interests of the two sides of the market are opposed in a fundamental sense: the best stable matching for one side of the market is the worst stable matching for the other side of the market.[5] This makes equilibrium selection an important and relevant consideration for policymakers, who may have reasons to favor the welfare of one side of the market over another when designing matching institutions. An example of this is provided by the history of the National Resident Matching Program (NRMP), the entry-level labor market for American physicians. In May 1997, the NRMP

---

possible in our experimental set-up. In a related paper, Castillo and Dianat (2016) present evidence that laboratory subjects are less likely to truncate their preferences when the possibility of over-truncation exists.

[3]In this game, firms (computers) have a dominant strategy of truth-telling and workers (subjects) have incentives to misrepresent their preferences to influence the final outcome.

[4]That is, each worker can be matched to either firm at a stable matching.

[5]This is a consequence of the fact that the set of stable matchings has a lattice structure.

unanimously voted to alter the algorithm that was being used over concerns that the original design unduly favored hospitals at the expense of students.[6] Thus, our experimental results can be useful in predicting which stable matching is more likely to be implemented when there are multiple candidates for consideration. We discuss further implications for market design in more detail in Section 6 of the paper.

The argument for using laboratory experiments in this context is compelling. While data on participants' submitted rank-order lists may be available in field settings, participants' true preferences are unobserved. Under the assumption of truthful preference reporting, matching clearinghouses used in practice implement an extremal stable outcome (i.e., the most preferred stable outcome for one side of the market). But if agents strategically misrepresent their preferences in the field, then it is less clear which stable matching is implemented or even whether the final matching is stable. By allowing us to control for subjects' preferences, the laboratory setting is ideally suited for answering questions related to equilibrium selection.

Our paper naturally bridges several different strands of literature. First, there is a large body of work on the performance of matching mechanisms in the lab.[7] The current paper is most closely related to Castillo and Dianat (2016), which is an experimental investigation of truncation behavior in what approximates a decision-theoretic setting. To that end, Castillo and Dianat (2016) employ a restricted strategy space that yields a unique equilibrium in dominant strategies. The current design, on the other hand, introduces multiple Pareto-ranked equilibria and allows us to better understand the conditions that favor the implementation of one equilibrium over another. There are also experimental studies of matching mechanisms that investigate whether intensity of preference has implications for strategic behavior. For instance, Echenique et al. (2016) report that the cardinal representation of subjects' preferences has a significant effect on the stability of final outcomes, with instability more likely to arise in the presence of weak incentives. However, Klijn et al. (2013) find that subject behavior is fairly robust to changes in cardinal preferences in the Gale-Shapley mechanism but not the Boston mechanism. In addition, the prevalence of out-of-equilibrium truth-telling in our experiment mirrors earlier findings by Featherstone et al. (2019).

Second, there is a large experimental literature on behavior in coordination games.[8] Several of these studies focus on stag hunt games in an attempt to evaluate the merits of competing equilibrium selection principles. Our results are largely consistent with this

---

[6]Specifically, the NRMP switched from a version of the hospital-proposing deferred acceptance algorithm to a version of the student-proposing deferred acceptance algorithm.

[7]As far as we are aware, Harrison and McCabe (1989) is the first such study. Hakimov and Kübler (2020) provides a useful survey of the experimental literature on this topic.

[8]For a survey of this literature, see Camerer (2003).

literature, suggesting that equilibrium selection arguments may have significant generality and explanatory power across environments. In particular, our finding that the salience of payoff-dominance increases with subject experience mirrors that of Rankin et al. (2000), who show that subjects learn to gravitate toward payoff-dominance in a sequence of stag hunt games where cosmetic details (e.g., action labels, player labels, payoffs) are randomly perturbed.

Finally, there is a growing literature that applies risk-dominance or related concepts to other strategic environments. In particular, our work is methodologically related to a series of lab experiments that use strategic risk (as measured by the size of the basin of attraction) to predict behavior across a variety of settings: the centipede game (Healy, 2017), the infinitely repeated prisoner's dilemma (Dal Bó and Fréchette, 2011; Kartal and Müller, 2018), the finitely repeated prisoner's dilemma (Embrey et al., 2017), continuous-time games (Calford and Oprea, 2017), and dynamic games (Vespa and Wilson, 2016).

The rest of the paper is organized as follows. Section 2 introduces the necessary theoretical background, Section 3 presents our experimental design, Section 4 presents our experimental results, Section 5 evaluates the ability of related solution concepts to explain our results, Section 6 discusses broader implications for market design, and Section 7 concludes.

## 2   Theoretical Background

There are two finite and disjoint sets of agents of equal size: a set $F$ of firms and a set $W$ of workers. The preferences of worker $w \in W$ are represented by the von Neumann-Morgenstern utility function $u_w : F \cup \{w\} \to \mathbb{R}_+$, where $u_w(f) > 0$ is the utility she derives from matching with firm $f \in F$ and $u_w(w) = 0$ is the utility she derives from remaining single.[9] We assume that each function $u_w$ is one-to-one, such that it induces a strict preference ordering $P_w$ on the set $F$. We will refer to $P_w$ as the *preference list* of worker $w$. The preferences of the firms are defined similarly. We let $u = (u_i)_{i \in F \cup W}$ denote the profile of agents' utility functions and $P = (P_i)_{i \in F \cup W}$ denote the profile of agents' preference lists.

A *matching market* is a triple $(F, W, u)$. A matching is a function $\mu : F \cup W \to F \cup W$ such that:

1. for any $f \in F$, $\mu(f) \in W \cup \{f\}$

2. for any $w \in W$, $\mu(w) \in F \cup \{w\}$

---

[9]Although the classical results we present are usually framed in terms of ordinal preferences, the solution concept of risk-dominance is inherently cardinal. Thus, we assume cardinal preferences throughout the analysis.

3. for any $f \in F$, $w \in W$, $\mu(f) = w$ if and only if $\mu(w) = f$

A pair of agents $(f, w)$ is said to *block* a matching $\mu$ if they are not matched to one another at $\mu$ but they prefer each other to their assignments at $\mu$ (i.e., $u_w(f) > u_w(\mu(w))$ and $u_f(w) > u_f(\mu(f))$). A matching $\mu$ is *stable* if it is not blocked by any pair of agents. A firm $f$ and a worker $w$ are said to be *achievable* for each other in a matching market $(F, W, u)$ if they are matched to each other at some stable matching. A stable matching is called *firm-optimal* (*worker-pessimal*) if each firm is matched to her most preferred achievable worker (each worker is matched to her least preferred achievable firm). Similarly, a stable matching is called *worker-optimal* (*firm-pessimal*) if each worker is matched to her most preferred achievable firm (each firm is matched to her least preferred achievable worker). We denote the firm-optimal stable matching by $\mu_F$ and the worker-optimal stable matching by $\mu_W$. In other words, for each firm $f \in F$, each worker $w \in W$, and each stable matching $\mu$, we have that $u_f(\mu_F(f)) \geq u_f(\mu(f)) \geq u_f(\mu_W(f))$ and $u_w(\mu_W(w)) \geq u_w(\mu(w)) \geq u_w(\mu_F(w))$.

Let $\mathcal{M}$ denote the set of all possible matchings, $\mathcal{Q}$ denote the set of all possible preference profiles, and $\mathcal{Q}_i$ denote the set of all possible preference lists for agent $i \in F \cup W$. Let $\mu$, $Q$, and $Q_i$ denote arbitrary elements of the sets $\mathcal{M}$, $\mathcal{Q}$, and $\mathcal{Q}_i$, respectively. A mechanism is a function $\phi : \mathcal{Q} \to \mathcal{M}$ that assigns a matching to each preference profile. A mechanism $\phi$ that for each preference profile $Q$ produces a matching $\phi(Q)$ that is stable with respect to $Q$ is called a stable mechanism. If $\phi(Q)$ is the firm-optimal stable matching with respect to $Q$, then $\phi$ is called the firm-optimal stable mechanism. We denote the firm-optimal stable mechanism by $\phi_F$.

The firm-optimal stable mechanism can be modeled as a non-cooperative game in which the strategy space is the set of all possible ordinal preference lists. In this preference-revelation game, it is well-known that the firms have a dominant strategy of truth-telling (Dubins and Freedman, 1981). In markets with more than one stable matching, however, at least one worker will have an incentive to misrepresent her preferences to improve her match outcome (Gale and Sotomayor, 1985). Our goal is to characterize the different equilibria that can arise in this environment. To simplify our analysis, we define the *constrained preference-revelation game* induced by the firm-optimal stable mechanism:

**Definition 1.** *Consider a matching market $(F, W, u)$ in which $P = (P_i)_{i \in F \cup W}$ denotes the profile of agents' true preference lists and $Q = (Q_i)_{i \in F \cup W}$ denotes the profile of agents' reported preference lists. The **constrained preference-revelation game** induced by the firm-optimal stable mechanism $\phi_F$ is the preference-revelation game in which the firms are constrained to truth-telling. That is, for any profile of workers' reports $(Q_i)_{i \in W}$, the con-*

*strained preference-revelation game produces the matching* $\phi_F\big((P_i)_{i \in F}, (Q_i)_{i \in W}\big)$.[10]

We will find it useful to define two types of misrepresentation strategies for the workers:

**Definition 2.** *A **truncation** of a preference list $P_w$ containing $k$ firms is a list $Q_w$ containing $k' < k$ firms such that the $k'$ elements of $Q_w$ are the first $k'$ elements of $P_w$, in the same order.*[11]

**Definition 3.** *A **permutation** of a preference list $P_w$ is a list $Q_w \neq P_w$ that is not a truncation of $P_w$.*[12]

In other words, a truncation involves misrepresenting preferences by removing match partners from the tail end of a preference list, while a permutation involves misrepresenting preferences by switching the order of match partners in a preference list (regardless of the length of the list).

We now state and prove some basic results that are relevant to our experimental design. Throughout, we let $\mu_F$ and $\mu_W$ denote the firm-optimal and worker-optimal stable matchings with respect to the true preferences $P$.

**Proposition 1.** *Consider a matching market $(F, W, u)$ in which all agents have more than one achievable partner. In the constrained preference-revelation game induced by the firm-optimal stable mechanism $\phi_F$, there is a payoff-dominant equilibrium in which all workers play truncation strategies.*

*Proof.* See Appendix. □

We will refer to the equilibrium in which all workers play truncation strategies as the "symmetric" truncation equilibrium. However, one worker's truncation decision creates positive spillovers for other workers in the market (Ashlagi and Klijn, 2012; Coles and Shorrer, 2014). This implies that asymmetric equilibria also exist in which a subset of workers plays truncation strategies and the remaining workers report their true preferences, effectively free-riding on others' truncation behavior.

We now construct a "symmetric" permutation equilibrium in which all workers play permutation strategies.

---

[10]The assumption that firms play their dominant strategy in the firm-optimal stable mechanism is not entirely innocuous. While the firm-optimal stable mechanism is strategy-proof for the firms, Ashlagi and Gonczarowski (2018) show that it is not obviously strategy-proof in the sense of Li (2017). Furthermore, empirical studies by Rees-Jones (2018) and Hassidim et al. (2017) find that a small fraction of participants fail to play their dominant strategy in strategy-proof matching mechanisms.

[11]This definition is taken from Roth and Rothblum (1999). However, it has been slightly modified such that truthful preference revelation is no longer an "edge case" of a truncation strategy.

[12]The term "dropping strategy" is often used to refer to the act of removing a match partner from the middle of a preference list (rather than from the tail end of a preference list). According to our definitions, a dropping strategy would be classified as a permutation.

$$P_{f_1} = w_1, w_2 \quad P_{w_1} = f_2, f_1$$
$$P_{f_2} = w_2, w_1 \quad P_{w_2} = f_1, f_2$$

Table 1: The ordinal preferences used in the experiment. The firm-optimal stable matching is shown in red and the worker-optimal stable matching is shown in blue.

|  | $Truth$ | $Truncate$ | $Permute$ |
|---|---|---|---|
| $Truth$ | $v_2, v_2$ | $v_1, v_1$ | $v_2, v_2$ |
| $Truncate$ | $v_1, v_1$ | $v_1, v_1$ | $v_3, v_2$ |
| $Permute$ | $v_2, v_2$ | $v_2, v_3$ | $v_2, v_2$ |

Table 2: Normal-form representation of the constrained preference-revelation game ($v_1 > v_2 > v_3$). The Nash equilibrium that implements the firm-optimal stable matching is shown in red. The Nash equilibria that implement the worker-optimal stable matching are shown in blue.

**Proposition 2.** *Consider a matching market $(F, W, u)$ in which all agents have more than one achievable partner. In the constrained preference-revelation game induced by the firm-optimal stable mechanism $\phi_F$, there is a payoff-dominated equilibrium in which all workers play permutation strategies.*

*Proof.* See Appendix. □

Although all workers prefer the truncation equilibrium to the permutation equilibrium, truncation behavior introduces the possibility of remaining unmatched for some profiles of other agents' reported preferences. This exposure to the worst possible outcome is not present for the permutation strategy that we identify. To see this more clearly, it is instructive to consider the steps of the firm-proposing deferred acceptance algorithm.[13] A consequential truncation (i.e., a truncation that affects the final outcome) requires a worker to reject a proposal from an achievable firm. This rejection frees the firm to make other proposals, which may cause a chain of further rejections. If the truncating worker does not receive new proposals from this rejection chain, then she will remain single. The permutation strategy in Proposition 2 is inherently not consequential. It produces the same outcome as truth-telling: it merely prioritizes the least preferred achievable firm by elevating her to the top of a worker's preference ordering.

---

[13]The firm-proposing deferred acceptance algorithm is a procedure that generates the firm-optimal stable matching for any preference profile.

# 3    Experimental Design

Each experimental market consists of four agents: two firms (implemented by computers) and two workers (human subjects). The firms are automated to play their dominant strategy of truth-telling, while the workers are free to report any preference ordering. In each session, subjects play 20 rounds of the constrained preference-revelation game induced by the firm-optimal stable mechanism.[14] Subjects are randomly and anonymously re-paired at the start of each round.

Table 1 shows the ordinal preferences used across all 20 rounds of the experiment. With this constellation of preferences, there are two disjoint stable matchings (i.e., each agent has two achievable match partners). For convenience, we let $v_1$, $v_2$ and $v_3$ denote an agent's utilities from matching with her most preferred partner, her least preferred partner, and remaining single, respectively. Table 2 depicts the normal-form representation of the constrained preference-revelation game.[15] It should be noted that *Permute* combines two pure strategies that are strategically equivalent.[16,17]

The original Harsanyi and Selten (1988) criterion of risk-dominance concerns the pairwise comparison of Nash equilibria and is intended to capture the intuition of one equilibrium being more or less strategically risky than another equilibrium. In their formulation, constructed for $2 \times 2$ games, equilibrium $A$ risk-dominates equilibrium $B$ if $A$ has the larger Nash product (i.e., if the product of the two players' unilateral deviation losses are larger when moving from $A$ to $B$ than when moving from $B$ to $A$). However, generalizing the concept of risk-dominance presents complications.[18]

Furthermore, the matching environment that we consider (both the general setting and the specific case we take to the lab) has two important differences from the $2 \times 2$ games where risk-dominance is traditionally applied. First, the constrained preference-revelation game has asymmetric equilibria where agents choose different strategies (e.g., agent 1 reports her true preferences and agent 2 truncates her preferences). This makes it less meaningful to speak of "truncation" equilibria or "truth-telling" equilibria. Second, the constrained

---

[14]In the experiment, the firm-proposing deferred acceptance algorithm is used to illustrate to subjects how reported preferences map to final outcomes.

[15]When firms are unconstrained, there exist other equilibria in which firms play dominated strategies.

[16]For instance, consider the situation facing worker $w_1$ with preference list $P_{w_1} = f_2, f_1$. Both permutation strategies ($Q_{w_1} = f_1, f_2$ and $Q'_{w_1} = f_1$) yield the same outcome for all preference reports by the other player. More generally, the equivalence of different permutation strategies need not hold.

[17]Although the two permutation strategies are theoretically equivalent in our experimental set-up, they are not behaviorally equivalent. In the aggregate data, we find that the overwhelming majority of permutations (96%) involve submitting a full-length preference list.

[18]For instance, the binary relation imposed by risk-dominance can fail to be transitive. Morris et al. (1995) provide an example of a $3 \times 3$ game with three strict Nash equilibria in which the risk-dominance relationship is cyclical.
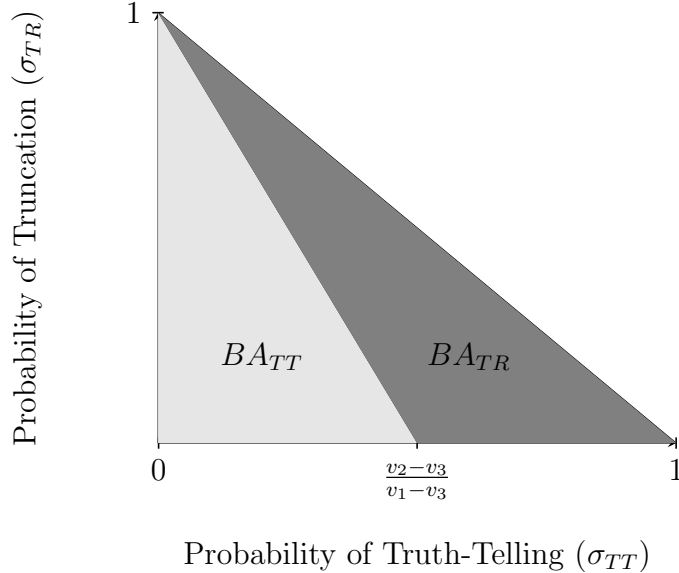
Figure 1: Basins of attraction of truth-telling ($BA_{TT}$) and truncation ($BA_{TR}$).

preference-revelation game does not have any strict Nash equilibria: at each equilibrium strategy profile, at least one agent has an alternative best-response. The implication is that all three classes of strategies we identify can be supported in equilibrium.

To resolve these tensions, we apply a version of risk-dominance to compare the risk-iness of different *strategies* rather than different *strategy profiles*. In particular, we use the size of the basin of attraction of a particular strategy to measure its level of strate-gic risk. For player $i$, the basin of attraction of strategy $a_i$ is the set of beliefs about other players' strategies such that playing $a_i$ is a best-response. More formally, $BA_{a_i} = \{\sigma_{-i} \in \Delta(A_{-i}) : u_i(a_i, \sigma_{-i}) \geq u_i(a_i', \sigma_{-i}) \text{ for all } a_i' \in A_i\}$. We will say that a strategy is risk-dominant if it has the largest basin of attraction. It should be noted that our notion of risk-dominance coincides with the Harsanyi and Selten (1988) criterion for $2 \times 2$ games and has the additional advantage of straightforwardly generalizing to any finite normal-form game.

We now apply our notion of risk-dominance to the constrained preference-revelation game. Suppose that all players play truth-telling with probability $\sigma_{TT}$, truncation with probability $\sigma_{TR}$, and permutation with probability $1 - \sigma_{TT} - \sigma_{TR}$. Then, player $i$'s expected payoffs from the three strategies are as follows:

$$u_{TT} = v_2 + \sigma_{TR}(v_1 - v_2),$$

$$u_{TR} = v_3 + (\sigma_{TT} + \sigma_{TR})(v_1 - v_3),$$

9

|                | Treatment | |
| --- | --- | --- |
|                | TT ($v_2 = 15$) | TR ($v_2 = 5$) |
| Truth-Telling  | risk dominant (RD) | neither PD nor RD |
| Truncation     | payoff dominant (PD) | PD and RD |

Table 3: Our experimental treatments vary the risk-dominant prediction. TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.

$$u_P = v_2.$$

It is easy to see that a player is indifferent between truth-telling and truncation when

$$\sigma_{TR} = 1 - \sigma_{TT}\left(\frac{v_1 - v_3}{v_2 - v_3}\right),$$

which yields the basins of attraction shown in Figure 1.[19] We will say that a strategy is risk-dominant if it has the larger basin of attraction. Truth-telling has the larger basin of attraction when $\frac{1}{2}\left(\frac{v_2 - v_3}{v_1 - v_3}\right) > \frac{1}{2}\left(1 - \frac{v_2 - v_3}{v_1 - v_3}\right)$, which yields $v_2 - v_3 > v_1 - v_2$. On the other hand, truncation has the larger basin of attraction when $\frac{1}{2}\left(1 - \frac{v_2 - v_3}{v_1 - v_3}\right) > \frac{1}{2}\left(\frac{v_2 - v_3}{v_1 - v_3}\right)$, which yields $v_1 - v_2 > v_2 - v_3$. Our characterization of risk-dominance neatly captures both the logic and the danger of playing a truncation strategy. When remaining unmatched is particularly costly (i.e., $v_2 - v_3 > v_1 - v_2$), then truth-telling is risk-dominant. When an agent has a strong intensity of preference for her first choice partner (i.e., $v_1 - v_2 > v_2 - v_3$), then truncation is risk-dominant.

For all agents, however, truth-telling constitutes the unique *protective* strategy in stable matching mechanisms (Barberà and Dutta, 1995).[20] This is because truth-telling is the only strategy that accomplishes the following two objectives: (1) it secures against the worst possible outcome (i.e., remaining single) and (2) it leads to the best possible outcome for some profile of other agents' preference reports.

Table 3 summarizes our experimental design. In our experiment, we fix the payoff from matching with the most preferred partner ($v_1 = 20$) and from the outside option of remaining single ($v_3 = 0$). Our treatments vary the risk-dominant prediction by manipulating the payoff from matching with the least preferred partner ($v_2$). For $v_2 \in (10, 20)$, truth-telling is risk-dominant (TT treatment). For $v_2 \in (0, 10)$, truncation is risk-dominant (TR treatment). In our experiment, we set $v_2 = 15$ in the TT treatment and $v_2 = 5$ in the TR treatment.

To secure comprehension, subjects are required to walk through a demonstration of the

---

[19]Since permutation is a weakly dominated strategy, it can never be a strict best-response.

[20]A protective strategy is a refinement of a maxmin strategy. Notice that while both *Truth* and *Permute* are maxmin strategies, *Truth* weakly dominates *Permute*.
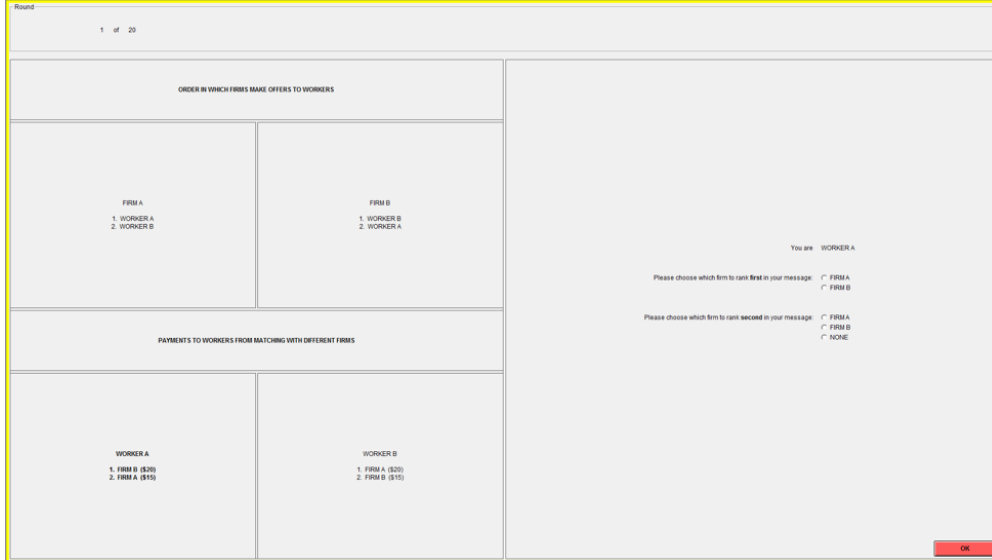
Figure 2: The experimental interface for the TT treatment.

deferred acceptance algorithm with a hypothetical set of reported preferences and to correctly answer a series of questions. In each round of the experiment, subjects observe the preferences of all market participants and they are then asked to report a preference ordering. At the end of a round, subjects receive feedback about the identity of their match partner (i.e., FIRM A, FIRM B, unmatched) and their payoff in that particular round. The experimental interface for the TT treatment is shown in Figure 2 and the experimental instructions for both treatments are provided in the Online Appendix.[21]

## 4 Experimental Results

The experimental sessions were run at the Experimental Social Science Laboratory (ESSL) at UC Irvine. A total of 120 subjects participated in the experiment (TT treatment: 64 subjects, TR treatment: 56 subjects). Each experimental session lasted approximately one hour. One of the 20 experimental rounds was randomly selected for subject payment. Average subject earnings were $21.75 (including a $7 show-up payment). The experiment was programmed and conducted with the software z-Tree (Fischbacher, 2007). We now present our main results.

**Result 1.** *Both truth-telling and truncation are played more often when they are risk-*

---

[21]To reduce experimenter demand effects, the terminology of preferences is never used in the experiment. A subject's true preference list is referred to as a "list of payments" and a subject's reported preference list is referred to as a "message."

|              | Round 1       |            | Round 20     |            |
| Strategy     | TT            | TR         | TT           | TR         |
| ------------ | ------------- | ---------- | ------------ | ---------- |
| Truth-Telling | 39 (61%)     | 39 (70%)   | 41 (64%)     | 26 (46%)   |
| Truncation   | 3 (5%)        | 4 (7%)     | 16 (25%)     | 30 (54%)   |
| Permutation  | 22 (34%)      | 13 (23%)   | 7 (11%)      | 0 (0%)     |
| Total        | 64 (100%)     | 56 (100%)  | 64 (100%)    | 56 (100%)  |

Table 4: Round 1 and Round 20 Choices by Treatment.
TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.

*dominant.*

## By round

We first demonstrate that this treatment effect is not present initially, but rather emerges with subject experience. Table 4 shows subject behavior in the first and last rounds of the experiment. In Round 1, truncation is rare in both treatments (TT: 5%, TR: 7%). In addition, a majority of subjects report their true preferences in both treatments (TT: 61%, TR: 70%). Neither of these treatment differences are statistically significant at conventional levels (truth-telling: two-sided t-test, $p = 0.323$; truncation: two-sided t-test, $p = 0.571$). In aggregate, we also find that initial behavior does not vary significantly across treatments (Fisher's exact test, $p = 0.394$).

However, there is a significant treatment effect at the end of the session (Fisher's exact test, $p < 0.001$). The observed treatment effect is consistent with subjects using risk-dominance as a selection criterion. In Round 20, subjects are more than twice as likely to play a truncation strategy when truncation is risk-dominant (TT: 25%, TR: 54%). At the same time, subjects are more likely to report their true preferences when truth-telling is risk-dominant (TT: 64%, TR: 46%). Both of these treatment differences are statistically significant at conventional levels (truth-telling: two-sided t-test, $p = 0.053$; truncation: two-sided t-test, $p = 0.001$).

## By session

An alternative way to measure a treatment effect is to consider experimental sessions (i.e., cohorts) as independent units of observation and rank the sessions by their average frequencies of different strategies. This procedure is shown in Table 5. The three TT sessions have higher average truth-telling rates than the three TR sessions, while the three TR sessions have higher average truncation rates than the three TT sessions. Using a non-parametric rank-sum test, we can reject the null hypothesis of no treatment difference for both cases

| Treatment | Frequency of Truth-Telling | Rank Score |
| --- | --- | --- |
| TT | 0.62 | 6 |
| TT | 0.60 | 5 |
| TT | 0.57 | 4 |
| TR | 0.52 | 3 |
| TR | 0.51 | 2 |
| TR | 0.46 | 1 |
| Treatment | Frequency of Truncation | Rank Score |
| TR | 0.47 | 6 |
| TR | 0.41 | 5 |
| TR | 0.36 | 4 |
| TT | 0.25 | 3 |
| TT | 0.19 | 2 |
| TT | 0.16 | 1 |
| Treatment | Frequency of Permutation | Rank Score |
| TT | 0.25 | 6 |
| TT | 0.22 | 5 |
| TT | 0.15 | 4 |
| TR | 0.12 | 3 |
| TR | 0.08 | 2 |
| TR | 0.08 | 1 |

Table 5: Ranking average frequencies of different strategies by session. TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.

($p = 0.050$). Truth-telling and truncation are both played significantly more often when they are the risk-dominant prediction.

## By subject

Each subject submits 20 rank-order lists during the course of the experiment (one in each experimental round). From this data, we can calculate a truth-telling, truncation, and permutation rate for each individual subject. The average subject-level truth-telling rates are 0.59 and 0.49 for the TT and TR treatments, respectively (two-sided t-test, $p = 0.114$). The average subject-level truncation rates are 0.20 and 0.42 for the TT and TR treatments, respectively (two-sided t-test, $p < 0.001$).

Figure 3 shows the empirical cumulative distribution functions (CDFs) of subject-level behavior. When comparing the truncation CDFs, it is apparent that the first-order stochastic dominance relationship is consistent with the prediction of risk-dominance. In other words, the distribution of truncation rates from the TR treatment first-order stochastically dominates the distribution of truncation rates from the TT treatment (middle graph). We can also reject the null hypothesis that subject-level truncation rates have the same distribution function across treatments ($p = 0.006$, assessed with a Kolmogorov-Smirnov test). However, the first-order stochastic dominance relationship does not hold when comparing the truth-telling CDFs. Further, we fail to reject the null hypothesis that subject-level truth-telling rates from the two treatments come from the same theoretical distribution ($p = 0.136$, assessed with a Kolmogorov-Smirnov test).

Finally, we can use the CDFs to investigate the incidence of "purists" who always play a particular strategy. We find that very few subjects play the same strategy across all rounds of the experiment. Across both treatments, only 7% (8/120) of subjects consistently report their true preference list, 3% (3/120) of subjects consistently truncate their preference list, and 2% (2/120) of subjects consistently permute their preference list.

**Result 2.** *In both treatments, truncation strategies are played more often in later rounds of the experiment.*

We now investigate the effect of learning on subject behavior. Figure 4 shows the average frequencies of different strategies across rounds of the experiment. In the TT treatment, only 5% of subjects play a truncation strategy in Round 1 while 25% of subjects play a truncation strategy in Round 20. In the TR treatment, the corresponding numbers are 7% in Round 1 and 54% in Round 20. For both treatments, we reject the null hypothesis that truncation rates across the initial 10 rounds are the same as truncation rates across the final
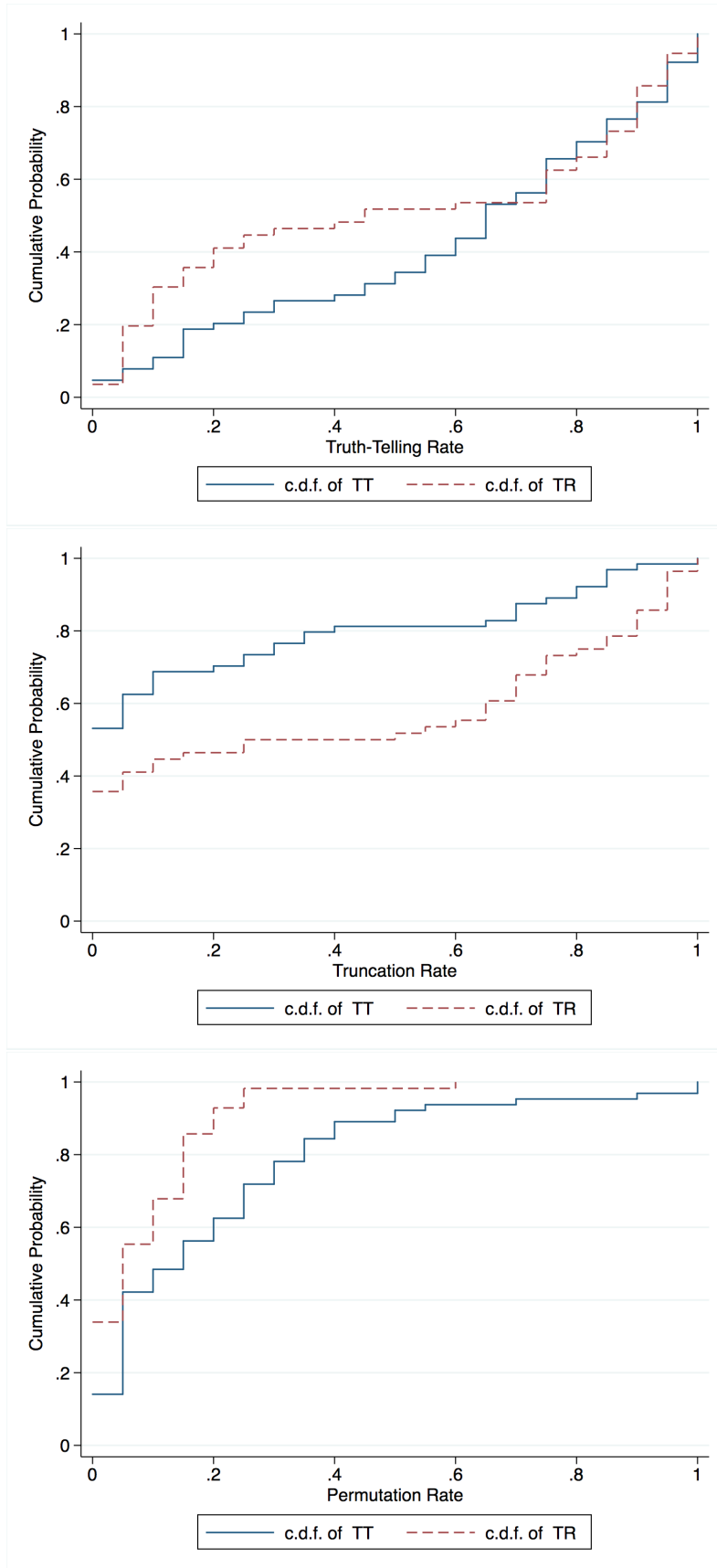
14

Figure 3: Empirical CDFs of subject-level behavior.
TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.
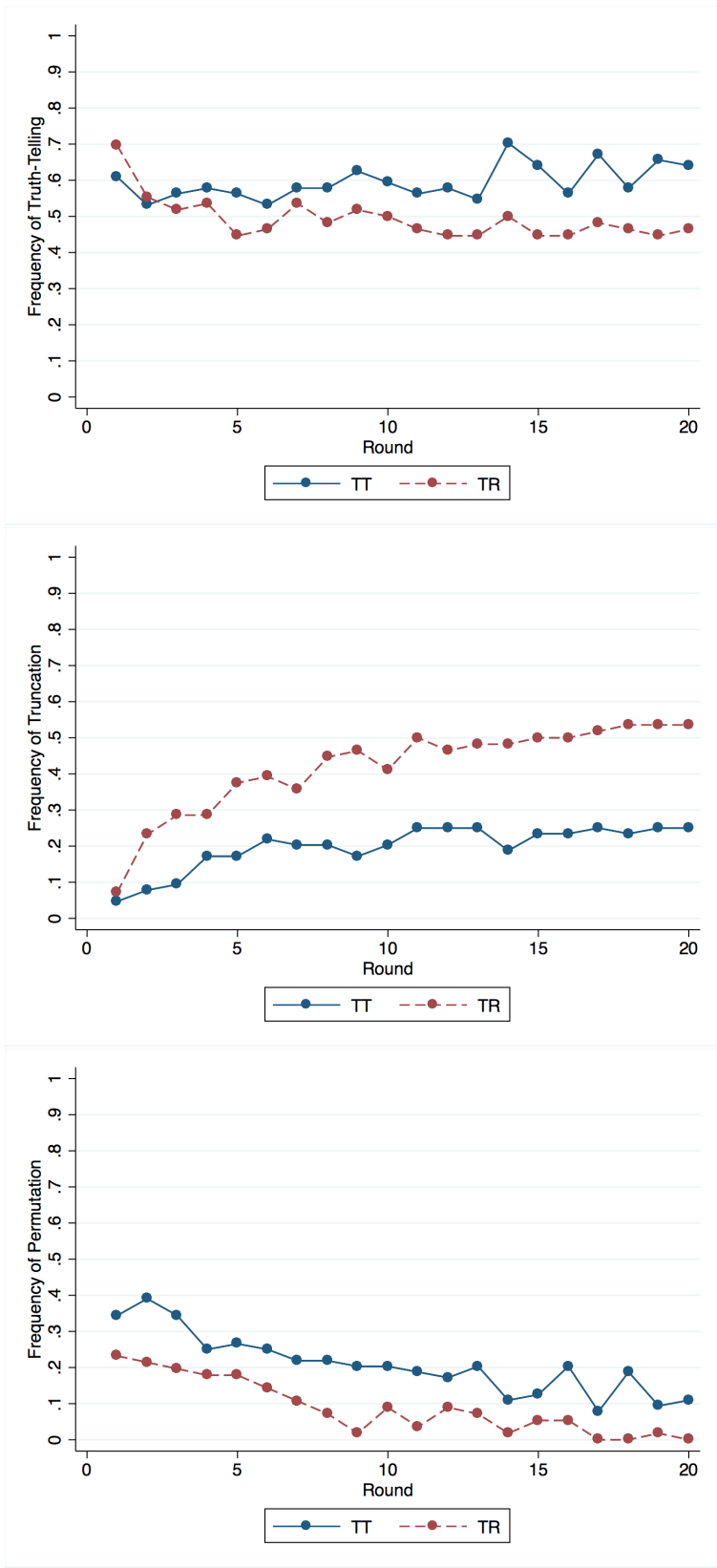
Figure 4: Average frequencies of different strategies across rounds of the experiment. TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.

10 rounds (TT: $\chi^2(1) = 13.8$, $p < 0.001$; TR: $\chi^2(1) = 34.5$, $p < 0.001$). Since this truncation time trend is present in both treatments, where risk-dominance yields different predictions but truncation remains a necessary component of all payoff-dominant equilibria, it suggests that payoff-dominance has increasing salience as a selection criterion in later rounds of the experiment.[22]

We also see that truth-telling rates are remarkably consistent across the experiment. In the TT treatment, the average frequency of truthful reporting ranges from a minimum of 53% (in Round 2) to a maximum of 70% (in Round 14). In the TR treatment, the minimum and maximum frequencies are 45% (in several rounds) and 70% (in Round 1). As before, we test for whether truth-telling rates across the initial 10 rounds are significantly different from truth-telling rates across the final 10 rounds. For the TT treatment, we fail to reject the null hypothesis of no time trend ($\chi^2(1) = 2.03$, $p = 0.155$). For the TR treatment, we can reject the null hypothesis of no time trend ($\chi^2(1) = 4.63$, $p = 0.031$). However, the latter result is attributable to the high level of truth-telling in Round 1 and disappears when Round 1 data is removed from the analysis ($\chi^2(1) = 2.17$, $p = 0.140$).

With respect to permutation strategies, we find a decrease in the TT treatment from 34% in Round 1 to 11% in Round 20. In the TR treatment, although 23% of subjects report a permuted preference list in Round 1, there are no permutations in Round 20. Again, for both treatments, we reject the null hypothesis that permutation rates across the initial 10 rounds are the same as permutation rates across the final 10 rounds (TT: $\chi^2(1) = 28.9$, $p < 0.001$; TR: $\chi^2(1) = 41.2$, $p < 0.001$). The decay in preference-list permutation, which constitutes a weakly dominated strategy, is consistent with an increase in subjects' strategic sophistication over the course of the experiment.

**Result 3.** *The worker-optimal (firm-optimal) stable matching is more likely to be implemented when truncation (truth-telling) is risk-dominant.*

Table 6 presents data on final outcomes. First, we find that stable matchings are the norm.[23] Across both treatments, only 6% (68/1,200) of markets fail to produce a stable matching. Second, it is clear that risk-dominance plays a crucial role in selecting among stable matchings. When truth-telling is risk-dominant, the firm-optimal stable matching is roughly twice as likely to be implemented (TT: 64%, TR: 34%). On the other hand, when truncation is risk-dominant, the worker-optimal stable matching is roughly twice as likely to be imple-

---

[22]In the context of stag hunt games, Rankin et al. (2000) also find that laboratory subjects focus on payoff-dominance rather than other solution concepts.

[23]There is a unique action profile that yields an unstable matching: a subject-pair where one subject plays a truncation strategy and the other subject plays a permutation strategy.

|  | Treatment | |
| Final Outcome | TT | TR |
| --- | --- | --- |
| Firm-Optimal Stable Matching | 409 (64%) | 193 (34%) |
| Worker-Optimal Stable Matching | 191 (30%) | 339 (61%) |
| Unstable Matching | 40 (6%) | 28 (5%) |
| Total | 640 (100%) | 560 (100%) |

Table 6: Empirical distribution of final outcomes.
TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.
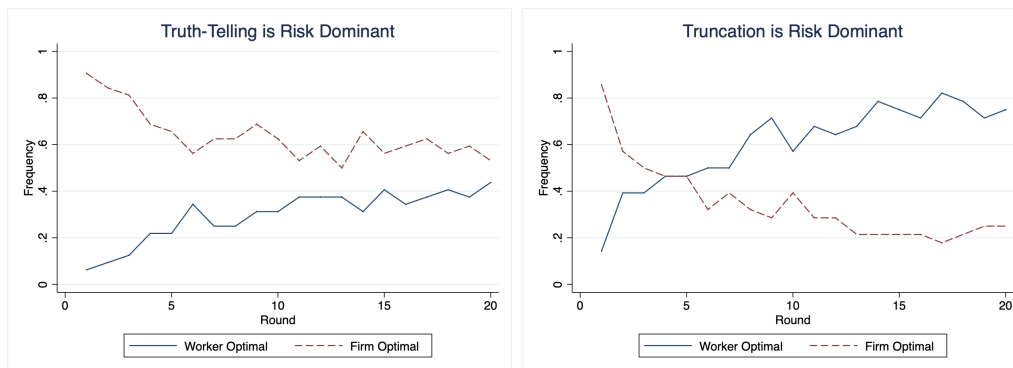


Figure 5: Frequency of stable matchings across rounds of the experiment.

mented (TR: 61%, TT: 30%). These treatment differences are statistically significant at conventional levels (two-sided t-test, $p < 0.001$ for both cases).

Figure 5 shows the frequency of stable matchings across rounds of the experiment. Unsurprisingly, there is also a tight connection between subject learning and selection among stable matchings. In particular, the worker-optimal (firm-optimal) stable matching is more (less) likely to be implemented in later rounds of the experiment. We find than an additional round of experience increases the probability of observing the worker-optimal stable matching by 1.6 percentage points in the TT treatment ($p = 0.037$) and by 2.6 percentage points in the TR treatment ($p = 0.004$).[24]

# 5    Related Solution Concepts

We have shown that a particular notion of strategic risk (based on the sizes of the basins of attraction of different strategies) is able to predict subjects' behavior in our experiment.

---

[24]For each treatment, we estimate an OLS regression of a dummy variable for whether the final outcome is the worker-optimal stable matching on the round of the experiment. Standard errors are clustered at the session level.

However, it is natural to wonder whether there are related solution concepts that are also capable of explaining our results. In this section, we discuss how well competing theories can organize our findings.

## 5.1 Protective Strategies

For all agents, truth-telling constitutes the unique *protective* strategy (Barberà and Dutta, 1995).[25] This is because truth-telling is the only strategy that accomplishes the following two objectives: (1) it secures against the worst possible outcome (i.e., remaining single) and (2) it leads to the best possible outcome for some profile of other agents' preference reports.

We find that truth-telling is the modal strategy in both treatments. A majority (59%) of submitted rank-order lists in the TT treatment and a plurality (49%) of submitted rank-order lists in the TR treatment correspond to subjects' true preferences. The prevalence of truthful behavior is consistent with the hypothesis that subjects use protective strategies. However, since truth-telling is the unique protective strategy in both treatments, theories in which subjects always use protective strategies would predict the absence of a treatment effect. As documented earlier, there is a significant treatment effect both at the session level and when comparing truth-telling rates in Round 20 of the experiment.

In principle, this finding leaves room for considerable heterogeneity in subject types: perhaps a fraction of subjects exclusively use protective strategies while the remaining subjects react to variations in strategic risk. However, as mentioned earlier, only 7% (8/120) of subjects consistently play truth-telling across all rounds. This fact, combined with the higher rate of truth-telling in the TT treatment, suggests that protective behavior fails to capture important aspects of the data.

## 5.2  $p$-Dominance

Morris et al. (1995) introduce $p$-dominance as a generalization of Harsanyi and Selten's risk-dominance from $2 \times 2$ games to many-action games. In their formulation, an action pair is $p$-dominant if each player's action is a best response to any belief that places at least probability $p$ on the other player taking her action in the pair. Since all Nash equilibria are 1-dominant in our setting, the concept of $p$-dominance does not offer a unique prediction and is also unable to account for any treatment effects that we observe in the experiment.

---

[25]A protective strategy is a refinement of a maxmin strategy. Notice that while both *Truth* and *Permute* are maxmin strategies, *Truth* weakly dominates *Permute*.

| **TT** | Truth | Truncation | Permutation |
|---|---|---|---|
| Simulation, s(1) = 10 | 62% | 4% | 34% |
| Simulation, s(1) = 100 | 62% | 4% | 34% |
| Experiment | 64% | 25% | 11% |
| **TR** | Truth | Truncation | Permutation |
| Simulation, s(1) = 10 | 65% | 16% | 19% |
| Simulation, s(1) = 100 | 66% | 15% | 19% |
| Experiment | 46% | 54% | 0% |

Table 7: Simulated and actual Round 20 behavior.
TT: truth-telling is risk-dominant. TR: truncation is risk-dominant.

## 5.3 $u$-Dominance

Kojima (2006) extends the concept of risk-dominance to more general finite games. They define an action to be $u$-dominant if it is a best response to any correlated action distribution where the number of opponents playing the action is uniformly distributed. In two-player games with more than two actions, an action is u-dominant if and only if the underlying action profile is 1/2-dominant in the sense of Morris et al. (1995). Thus, $u$-dominance offers no additional insight in our setting.

## 5.4 Reinforcement Learning

We now investigate whether subject behavior can be described by a basic model of reinforcement learning, a low-cognitive theory of decision-making in which agents are more likely to repeat actions that yield more favorable outcomes. In our experiment, at the end of a round, each subject only receives feedback about her own payoff (but not the payoff of the other player nor the action taken by the other player). Since the theory of reinforcement learning only requires that an agent update her play based on information about her own payoff, it can be directly tested with our experimental data.

We use the one-parameter reinforcement learning model from Erev and Roth (1998). Each player $n$ has a propensity to play strategy $k$ in round $t$ of the experiment, denoted by $q_{nk}(t)$. The reinforcement from receiving a payoff $x \in \mathbb{R}_+$ is represented by the function $R(x) = x$. If player $n$ plays strategy $k$ in round $t$ and receives a payoff of $x$, propensities are then updated as follows:

$$q_{nj}(t+1) = \begin{cases} q_{nj}(t) + R(x) & j = k \\ q_{nj}(t) & j \neq k. \end{cases} \tag{1}$$

The probability that player $n$ chooses strategy $k$ in round $t$ is given by

$$p_{nk}(t) = \frac{q_{nk}(t)}{\sum q_{nj}(t)}, \tag{2}$$

where the sum is taken over all the pure strategies. The single parameter of the model is the sum of the initial propensities, which we assume is the same for all players and we denote by

$$s(1) = s_n(1) = \sum q_{nj}(1). \tag{3}$$

The parameter $s(1)$ is often called the *strength* of the initial propensities, since it can substantially affect the speed of learning.

To determine whether the one-parameter reinforcement learning model is consistent with our experimental data, we conduct simulations for both treatments and for two different values of the free parameter $s(1)$.[26] For each of these four cases, we conduct 1000 simulations where two players interact for 20 rounds. We set the initial propensities such that the Round 1 choice probabilities (in the simulation) mirror the Round 1 distribution of strategies (in the experiment). To simulate play in later rounds, propensities are updated according to Equation (1) and strategies are randomly determined according to Equation (2).

Table 7 reports the simulated Round 20 behavior (averaged across both players and all 1000 simulations) alongside the actual Round 20 behavior for both treatments. We find that reinforcement learning fails to capture important patterns in the data. First, the simulations predict very little difference in Round 20 truth-telling rates across the two treatments. Second, in both treatments, the simulations under-predict the Round 20 truncation rates and over-predict the Round 20 permutation rates. In other words, reinforcement learning is unable to fully account for both the treatment effects and the learning effects that we find in the experiment.

# 6 Implications for Market Design

Our experiment can help shed light on several open questions in market design. In a seminal paper, Roth and Peranson (1999) conduct computational experiments using NRMP submitted rank-order lists from 1987 and 1993-1996. They find that only 0.1% of applicants would have received a different match from the applicant-proposing and hospital-proposing versions of the deferred acceptance algorithm. Assuming that the NRMP rank-order lists accurately reflect participants' true preferences, this exercise suggests that the applicant-optimal and

---

[26]We use $s(1) = 10$ and $s(1) = 100$.

| Average Number of Applicants | 2010 | 2012 | 2014 | 2016 | 2018 |
|---|---|---|---|---|---|
| Interviewed | 85 | 89 | 96 | 94 | 95 |
| Ranked | 66 | 66 | 77 | 80 | 82 |

Table 8: NRMP Program Director Survey Reports (across all medical specialties).

hospital-optimal stable matchings coincide.[27] This result has profound implications for market design: if there is a unique stable matching in these environments, then there is no incentive for market participants to behave strategically.

However, the Roth and Peranson (1999) result depends crucially on the assumption of truthful preference reporting. By credibly documenting the use of non-truth-telling strategies, our experiment suggests a degree of caution in this regard. In our experimental markets, there are two disjoint stable matchings with respect to the true (i.e., induced) preferences. In our experimental data, however, 72% of markets have a unique stable matching with respect to the reported preferences (TT treatment: 67%, TR treatment: 77%). In other words, the assumption of truthful preference reporting can lead to the false conclusion that a matching market has a unique stable matching when in fact there is a multiplicity of stable matchings. It is also noteworthy how robust this perception can be when confronted with reported preference data. In our experiment, any strategy profile in which at least one subject in a pair deviates from truth-telling is sufficient to produce a unique stable matching with respect to the reported preferences. Thus, a large range of reporting behavior is consistent with the empirical regularity of a singleton core.

In addition, previous theoretical work on strategic behavior in matching markets has largely focused on truncation strategies (Coles and Shorrer, 2014; Roth and Rothblum, 1999). There are two common arguments for this emphasis. First, truncation can be profitably implemented even with incomplete information about other agents' preferences (Coles and Shorrer, 2014; Roth and Rothblum, 1999). Second, it is sufficient to restrict attention to truncation when considering the space of profitable misrepresentation strategies (Roth and Vate, 1991). In other words, any agent who can improve her match partner by deviating from truth-telling can do so by submitting a truncation of her true preferences.

However, the empirical content of truncation strategies remains an open question. Field data are insufficient to address this issue: while centralized matching clearinghouses may provide access to participants' submitted rank-order lists, participants' true preferences are

---

[27]In fact, the computational exercise of Roth and Peranson (1999) has generated a literature on "core convergence" in matching models (e.g. Immorlica and Mahdian, 2005; Kojima and Pathak, 2009; Lee, 2016). Under certain conditions, these papers show that the set of stable matchings shrinks as the size of the market increases.

unobserved. Our experiment provides preliminary support for the empirical relevance of truncation strategies and helps highlight the conditions that foster truncation behavior. Our results suggest that truncation should be more likely to arise in settings where participants either have a strong intensity of preference for top-ranked alternatives or have previous experience with the particular mechanism that is being used. Although real-world matching protocols often more closely resemble one-shot games for many participants, hospitals in the NRMP usually do have considerable experience based on their participation in the match process in previous years. The NRMP conducts regular surveys of the directors of all hospital programs participating in the residency match.[28] Data from recent survey reports are shown in Table 8. In recent years, between 15-26% of interviewed applicants have not been included in hospitals' submitted rankings. While not conclusive, the NRMP survey data suggests that truncation strategies may play an important role in field settings.

# 7    Conclusion

In this paper, we investigate the interplay between strategic uncertainty and equilibrium selection in the context of stable matching mechanisms. We report three main findings from a laboratory experiment: (1) truth-telling is the most common strategy, (2) both truth-telling and truncation are played more often when they are risk-dominant, and (3) in both treatments, truncation strategies are played more often in later rounds of the experiment. The final point suggests that the salience of payoff-dominance as a selection criterion increases with subject experience.

It is worthwhile to briefly discuss the interpretation of our findings. We have shown that a simple measure of strategic uncertainty - the size of the basin of attraction of a strategy - is able to explain important patterns in our experimental data. However, we do not necessarily take the view that subjects are explicitly using this measure in their calculations. Rather, it is more plausible that subjects are reacting to variations in strategic risk using a convenient behavioral shortcut (e.g., the size of cardinal payoff differences). We do not believe this is a shortcoming of our implementation of risk-dominance: "as if" interpretations of behavior have a long tradition in economic theory dating back to Friedman (1953).[29]

A related question is to what extent our findings can inform the design of real-world matching markets. Given the large number of participants in centralized matching clearinghouses, it can be computationally infeasible for market designers to apply our notion of

---

[28]The NRMP Program Director Survey Reports can be found at the following website: http://www.nrmp.org/report-archives.

[29]Calford and Oprea (2017) make a similar point in the context of continuous-time games in the lab.

risk-dominance in a literal sense. Yet, we believe our experimental results can still provide guidance on the conditions under which we can expect different strategies to prevail. As discussed earlier, our results suggest that truncation is more likely to occur in settings characterized by strong preference intensities for top-ranked alternatives and truth-telling is more likely to occur in settings where remaining unmatched is particularly costly. More broadly, our experiment shows that we cannot necessarily expect uniform behavior across different settings. Thus, having an accurate model of the underlying preference structure can be crucial for interpreting reported preference data. Identifying the most common patterns of naturally-occurring preferences in the field is an important gap for empirical research to fill.

There are other promising avenues for further research. As it stands, there is no consensus on how to appropriately generalize concepts such as risk-dominance from $2 \times 2$ games to games with either additional players or larger action spaces. Theoretical work can hopefully bridge this gap. Further, the theory of two-sided matching is largely silent on the question of which stable matching will arise in markets with a multiplicity of stable matchings. Our experiment indicates that cardinal payoff differences can play a role in this selection process, but more empirical work is needed before a unified theory of selection among stable matchings can be developed.

# References

Itai Ashlagi and Yannai A Gonczarowski. Stable matching mechanisms are not obviously strategy-proof. *Journal of Economic Theory*, 177:405–425, 2018.

Itai Ashlagi and Flip Klijn. Manipulability in matching markets: conflict and coincidence of interests. *Social Choice and Welfare*, 39(1):23–33, 2012.

Salvador Barberà and Bhaskar Dutta. Protective behavior in matching models. *Games and Economic Behavior*, 8(2):281–296, 1995.

Evan Calford and Ryan Oprea. Continuity, inertia, and strategic uncertainty: A test of the theory of continuous time games. *Econometrica*, 85(3):915–935, 2017.

Colin Camerer. *Behavioral game theory: Experiments in strategic interaction*. Princeton University Press, 2003.

Marco Castillo and Ahrash Dianat. Truncation strategies in two-sided matching markets: Theory and experiment. *Games and Economic Behavior*, 2016.

Peter A Coles and Ran I Shorrer. Optimal truncation in matching markets. *Games and Economic Behavior*, 87:591–615, 2014.

Russell W Cooper, Douglas V DeJong, Robert Forsythe, and Thomas W Ross. Selection criteria in coordination games: Some experimental results. *American Economic Review*, 80(1):218–233, 1990.

Pedro Dal Bó and Guillaume R Fréchette. The evolution of cooperation in infinitely repeated games: Experimental evidence. *American Economic Review*, 101(1):411–429, 2011.

Lester E Dubins and David A Freedman. Machiavelli and the gale-shapley algorithm. *American Mathematical Monthly*, 88(7):485–494, 1981.

Federico Echenique, Alistair J Wilson, and Leeat Yariv. Clearinghouses for two-sided matching: An experimental study. *Quantitative Economics*, 2016.

Matthew Embrey, Guillaume R Fréchette, and Sevgi Yuksel. Cooperation in the finitely repeated prisoners dilemma. *The Quarterly Journal of Economics*, 133(1):509–551, 2017.

Ido Erev and Alvin E Roth. Predicting how people play games: Reinforcement learning in experimental games with unique, mixed strategy equilibria. *American Economic Review*, pages 848–881, 1998.

Clayton R Featherstone, Eric Mayefsky, and Colin D Sullivan. Learning to manipulate: Out-of-equilibrium truth-telling in matching markets. *Working Paper*, 2019.

Urs Fischbacher. z-tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, 10(2):171–178, 2007.

Milton Friedman. *Essays in Positive Economics*. University of Chicago Press, 1953.

David Gale and Marilda Sotomayor. Ms. machiavelli and the stable matching problem. *American Mathematical Monthly*, 92(4):261–268, 1985.

Rustamdjan Hakimov and Dorothea Kübler. Experiments on centralized school choice and college admissions: A survey. *Experimental Economics*, 2020.

Glenn W Harrison and Kevin A McCabe. Stability and preference distortion in resource matching: An experimental study of the marriage market. *Research in Experimental Economics*, 8, 1989.

John C Harsanyi and Reinhard Selten. *A general theory of equilibrium selection in games*. MIT Press, 1988.

Avinatan Hassidim, Déborah Marciano, Assaf Romm, and Ran I Shorrer. The mechanism is truthful, why aren't you? *American Economic Review*, 107(5):220–24, 2017.

Paul J Healy. Epistemic experiments: Utilities, beliefs, and irrational play. *Working Paper*, 2017.

Nicole Immorlica and Mohammad Mahdian. Marriage, honesty, and stability. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 53–62, 2005.

Melis Kartal and Wieland Müller. A new approach to the analysis of cooperation under the shadow of the future: Theory and experimental evidence. *Working Paper*, 2018.

Flip Klijn, Joana Pais, and Marc Vorsatz. Preference intensities and risk aversion in school choice: A laboratory experiment. *Experimental Economics*, 16(1):1–22, 2013.

Fuhito Kojima. Risk-dominance and perfect foresight dynamics in n-player games. *Journal of Economic Theory*, 128(1):255–273, 2006.

Fuhito Kojima and Parag A Pathak. Incentives and stability in large two-sided matching markets. *American Economic Review*, 99(3):608–627, 2009.

SangMok Lee. Incentive compatibility of large centralized matching markets. *Review of Economic Studies*, 84(1):444–463, 2016.

Shengwu Li. Obviously strategy-proof mechanisms. *American Economic Review*, 107(11): 3257–87, 2017.

Stephen Morris, Rafael Rob, and Hyun Song Shin. p-dominance and belief potential. *Econometrica*, pages 145–157, 1995.

Frederick W Rankin, John B Van Huyck, and Raymond C Battalio. Strategic similarity and emergent conventions: Evidence from similar stag hunt games. *Games and Economic Behavior*, 32(2):315–337, 2000.

Alex Rees-Jones. Suboptimal behavior in strategy-proof mechanisms: Evidence from the residency match. *Games and Economic Behavior*, 108:317–330, 2018.

Alvin E Roth and Elliott Peranson. The redesign of the matching market for american physicians: Some engineering aspects of economic design. *American Economic Review*, 89 (4):748–780, 1999.

Alvin E Roth and Uriel G Rothblum. Truncation strategies in matching markets - in search of advice for participants. *Econometrica*, 67(1):21–43, 1999.

Alvin E Roth and Marilda A Oliveira Sotomayor. *Two-sided matching: A study in game-theoretic modeling and analysis*. Cambridge University Press, 1992.

Alvin E Roth and John H Vande Vate. Incentives in two-sided matching with random stable mechanisms. *Economic Theory*, 1(1):31–44, 1991.

John B Van Huyck, Raymond C Battalio, and Richard O Beil. Tacit coordination games, strategic uncertainty, and coordination failure. *American Economic Review*, 80(1):234–248, 1990.

Emanuel Vespa and Alistair J Wilson. Experimenting with equilibrium selection in dynamic games. *Working Paper*, 2016.

# Appendix

## Proof of Proposition 1

*Proof.* Let $T$ denote the profile of reported preferences in which each worker $w$ truncates her preference list by removing all firms ranked below $\mu_W(w)$. By Theorem 4.17 of Roth and Sotomayor (1992), $T$ is a Nash equilibrium and it produces the matching $\mu_W$. Suppose another Nash equilibrium $Q$ payoff-dominates $T$. Let $\mu$ denote the matching that is produced by $Q$. Since $Q$ is a Nash equilibrium, we know by Theorem 4.16 of Roth and Sotomayor (1992) that the matching $\mu$ is also stable with respect to the true preferences $P$. Since $Q$ payoff-dominates $T$, we know that $u_w(\mu(w)) > u_w(\mu_W(w))$ for all $w \in W$. We have arrived at a contradiction, since $\mu_W$ is the W-optimal stable matching with respect to $P$. Thus, there is no other Nash equilibrium that payoff-dominates $T$. We conclude that $T$ is payoff-dominant. $\qquad\square$

## Proof of Proposition 2

*Proof.* Let $Q$ denote the profile of reported preferences in which each worker $w$ reports a preference list $Q_w$ that ranks $\mu_F(w)$ in the first position (regardless of the length of the list). Each preference list $Q_w$ is clearly a permutation since $\mu_F(w)$ is not at the head of any worker's true preference list.[30] It is straightforward to see that $Q$ produces the matching $\mu_F$.

We argue that the profile of reported preferences $Q$ constitutes a Nash equilibrium.[31] To see this, suppose that $Q$ is not a Nash equilibrium. Then, there exists some worker $w$ who can deviate and report a preference list $Q'_w$, which leads to a new profile of reported preferences $Q' = (Q_{-w}, Q'_w)$ and a new matching $\mu'$ such that $u_w(\mu'(w)) > u_w(\mu_F(w))$. Let $f = \mu'(w)$. Then firm $f$ must have been matched to a worker she prefers to $w$ at $\mu_F$, otherwise $(f, w)$ would have blocked the matching $\mu_F$ under the true preferences $P$. But now firm $f$ and worker $\mu_F(f)$ block the matching $\mu'$ under the reported preferences $Q'$, which is a contradiction. Therefore, $Q$ is a Nash equilibrium. Furthermore, $Q$ is payoff-dominated by the truncation equilibrium constructed in Proposition 1. $\qquad\square$

---

[30]If $\mu_F(w)$ were at the head of any worker's true preference list, then this contradicts the assumption that all workers have more than one achievable partner.

[31]The proof of this claim closely follows the proof of Theorem 4.15 in Roth and Sotomayor (1992). The only difference is that we allow the preference lists in $Q$ to be of any length.