

# Generalizing Obvious Dominance using the Sure-thing Principle

Chew Soo Hong\*      Wang Wenqian†

Sunday 19<sup>th</sup> June, 2022

## Abstract

We offer a class of solution concepts for dynamic games by allowing obvious dominance (OD, Li, 2017) to be facilitated by players' application of the sure-thing principle (STP, Savage, 1954) on relevant events in their contingent reasoning. In the resulting  $k$ -OD equilibrium, a lower  $k$  reflects a less demanding level of sophistication in applying STP. We present the class of gradual mechanisms, featuring dynamic information feedback, as a revelation principle for  $k$ -OD implementation. Applying  $k$ -OD to implement generalized median voter schemes, we show how an increasing sophistication level  $k$  extends the range of implementable social choice functions.

---

\*Southwestern University of Finance and Economics and National University of Singapore, 555 Liutai Street, Shuangliu District, Chengdu (e-mail: chew.soo hong@gmail.com)

†Hong Kong University of Science and Technology (Guangzhou) 2 South Huanshi Street, Nansha District, Guangzhou (e-mail: wqwang@ust.hk)

At the heart of the literature on mechanism design is the concept of incentive compatibility (Hurwicz, 1960, 1972), which stipulates that the mechanism provides incentives to its participants to reveal their private information truthfully. The concept of strategy-proofness has served as the gold standard of incentive compatibility (Satterthwaite, 1975) requiring that for each participant of the direct mechanism, revealing her private information is a dominant strategy. Despite the theoretical appeal of the concept of strategy-proofness, experimental studies have revealed difficulties with its dominant strategy implementation.<sup>1</sup>

There is substantial interest in understanding why agents participating in these mechanisms deviate from their dominant strategies (e.g., Cason and Plott, 2014; Dreyfuss et al., 2022), for which the literature on failure of contingent reasoning (e.g., Charness and Levin, 2009; Esponda and Vespa, 2014; Martínez-Marquina et al., 2019; Ngangoué and Weizsäcker, 2021) may offer a clue. One finding in this literature is that subjects are less likely to choose a theoretically inferior course of action when uncertainty is reduced or eliminated, e.g., postponing decision making until some relevant event is known to have happened. This motivates us to look for implementations of the same allocation rule with dynamic information feedback in order to reduce the difficulty in contingent reasoning and restore truthful behavior.

One recent influential contribution along this line is made by Li (2017) who offers a refinement of strategy-proofness by strengthening its underlying solution concept—dominance (dominant strategy equilibrium)—with what he calls *obvious dominance* (OD). One strategy is said to obviously dominate another when the worst outcome following the dominating strategy is no worse than the best outcome following the dominated strategy, both ranging across all possible contingencies after any information set where these two strategies initially diverge. In the context of mechanism design, if a social choice function can be implemented by a (possibly dynamic) game form with an equilibrium consisting of OD strategies for all agents, the social choice function is said to

---

<sup>1</sup>See, e.g., Kagel et al. (1987), Chen and Sönmez (2006), and Hakimov and Kübler (2019).

be *obviously dominant (OD) implementable* and the game form is called an *OD implementation* of that social choice function.

It has been recognized that OD is a particularly demanding concept. Due to the stringent standard for a dominant strategy to be obvious, many important strategy-proof social choice functions are OD implementable only in highly restrictive settings. For example, Li (2017), Ashlagi and Gonczarowski (2018), Thomas (2020), and Mandal and Roy (2022) study OD implementation of the two canonical matching rules—top trading cycles and deferred acceptance—and demonstrate that they may not be OD implementable in a matching market with 3 or more agents. Additionally, on the domain of single-peaked preference, according to Bade and Gonczarowski (2017) and Arribillaga et al. (2020), truth-telling in the majority rule for voting between two candidates is not an obviously dominant strategy.

There is a somewhat subtle issue in applying the idea of OD implementation to single-item ascending-price auctions. While Li (2017) motivates his solution concept with the experimentally observed difference in the proportion of type-revealing behavior between second-price and ascending-price private value auctions, the tie breaking rule adopted in his experiment—no winners in the event of a tie—is inefficient and different from the efficient tie-breaking rule commonly used, i.e., every remaining bidder wins with equal chance at the price of the tie (e.g., Kagel et al. 1987).<sup>2</sup> As it turns out, this seemingly minor difference in tie-breaking rule makes a significant theoretical difference—the ascending-price auction under the usual efficient tie-breaking rule is no longer an OD implementation.

To see why it is the case, suppose a bidder, Sophie, has private value  $v$  for the object being sold via an ascending-price auction. Her **optimal** strategy is to drop out when the next price level exceeds  $v$ , while she may entertain an **inferior** strategy to drop out earlier at price  $p$ .<sup>3</sup> Suppose the auction is

---

<sup>2</sup>In a second-price auction, a tie happens when there are multiple bids at the same highest price. In an ascending-price auction, a tie happens when multiple bidders leaves simultaneously at the last period.

<sup>3</sup>When there is no tie, a bidder wins the ascending-price auction at price  $p$  if she is the only bidder who chooses to stay in the auction. By dropping out at price  $p$ , a bidder avoids

continuing at the current price level  $p$  (i.e., there are multiple bidders, including Sophie, remaining in the auction). This is the point at which Sophie’s optimal strategy and inferior strategy first diverge. Sticking to the optimal strategy, when price eventually exceeds  $v$ , Sophie could lose the auction as long as there are still bidders staying in the auction. Following the inferior strategy, however, Sophie could possibly win the auction at price  $p$  in the event that all other remaining bidders also drop out and the tie is broken in Sophie’s favor (see Table 1 below).

Table 1: Ascending-price auction is not an OD implementation.

	Worst from optimal	Best from inferior
Among all possibility	(1) 0	(2) $v - p$

The fact that (1) < (2) demonstrates that the optimal strategy does not obviously dominate the inferior strategy in the ascending-price auction with an efficient tie-breaking rule.

An intermediate solution concept between dominance and OD extending the scope of implementable social choice functions would be valuable if it could retain some flavor of being more obvious and therefore easier for the agents to conform to the type-revealing strategy. In this paper, we set out to accomplish this.

In his *Foundations of Statistics*, Savage (1954) motivates the key postulate **P2** in his axiomatization of subjective expected utility by means of the following businessman example.

A businessman contemplates buying a certain piece of property. He considers the outcome of the next presidential election relevant. So, to clarify the matter to himself, he asks whether he would buy if he knew that the Democratic candidate were going to win, and decides that he would. Similarly, he considers whether he would buy if he knew that the Republican candidate were going to win,

---

the possibility of winning at the subsequent higher price levels.

and again finds that he would. Seeing that he would buy in either event, he decides that he should buy, even though he does not know which event obtains, or will obtain, as we would ordinarily say. It is all too seldom that a decision can be arrived at on the basis of this principle, but except possibly for the assumption of simple ordering, I know of no other extralogical principle governing decisions that finds such ready acceptance.

Savage describes the businessman’s decision making principle, what he calls **the sure-thing principle**, as being intuitive and appealing to cope with uncertainty and one which would find “ready acceptance”. To make the decision whether to buy a piece of property, the businessman seeks a relevant event (and together with its complement) such that buying would be a clearly better decision when both the event obtains and its complement obtains. While such a situation may be infrequent in real life, we show below how STP, when combined with Li’s (2017) idea of obviousness, addresses the robustness issue of OD implementation in the ascending-price auction under the efficient tie-breaking rule.

In the earlier discussion about Sophie’s choice between the optimal strategy and the inferior strategy, a naturally relevant event that may come to her mind is whether there would be a tie at price  $p$  if she drops out (i.e., all other remaining bidders choose to leave the auction at price  $p$ ). When there is such a tie at  $p$ , following the optimal strategy (staying in the auction) guarantees Sophie’s winning the auction at  $p$  which is no worse than the best outcome following the inferior strategy among all possibilities. Should this tie not occur, following the inferior strategy (dropping out of the auction early) simply loses the auction, which is no better than the worst outcome following the optimal strategy among all possibilities (see Table 2 below).

No matter whether there is a tie at  $p$ , the optimal strategy is obviously better than the inferior strategy. Thus, Sophie may adopt STP and recognize the dominance nature of the optimal strategy over the inferior strategy. In other words, the optimal strategy *obviously dominates* the inferior strategy *facilitated by STP* since  $(3) \geq (2)$  and  $(1) \geq (4)$ .

Table 2: Ascending-price auction is a STP-OD implementation.

	Worst from optimal	Best from inferior
Tie at $p$	(3) $v - p$	-
No Tie at $p$	-	(4) 0
Among all possibilities	(1) 0	(2) $v - p$

Our endeavor to generalize OD by incorporating a form of contingent reasoning captured by STP may seem surprising given that Li’s (2017) original idea is motivated by the failure of contingent reasoning. Notwithstanding the widely reported violations of STP in terms of Allais (1953) and Ellsberg (1961) style problems, Esponda and Vespa (2021) have found it to be effective as interventions in ameliorating failures in contingent reasoning.

In *The Republic*, Plato offers an observation about the penalty associated with not voting, “But the chief penalty is to be governed by someone worse if a man will not himself hold office and rule”. Here, he draws on the relevant scenario of an inferior candidate winning the election to offer an advice. In their laboratory study on voting, Esponda and Vespa (2021) employ a pair of complementary events to offer a contingent treatment, i.e., framing the decisions in a way facilitating the application of STP. On one event, the subject wins no matter what she votes, while on the other event, the subject’s vote is pivotal. Application of STP on them renders the dominant choice obvious.

After introducing some preliminaries of dynamic mechanism design to implement a (stochastic) social choice rule in Section 1, we provide in Section 2 a class of generalizations of the solution concepts of both OD and STP-OD, which we term as  $k$ -OD, by assuming that players are able to apply STP in more sophisticated manners. Instead of finding a single relevant event (together with its complement) to make dominance obvious, players can identify  $k$  relevant events forming a partition of the state space to make the comparison obvious using STP.

The following four sub-figures of Figure 1 give us an intuitive comparison among the solution concepts of dominance, OD, STP-OD, and  $k$ -OD, whose

formal definitions are given in Section 2.

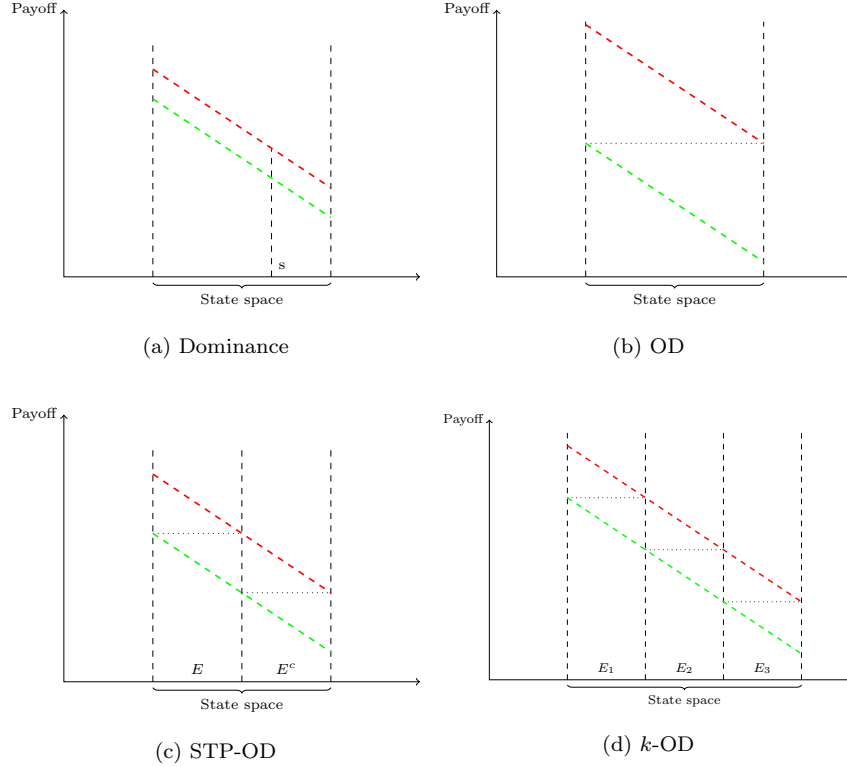


Figure 1: Comparison among Solution Concepts

In each sub-figure, two strategies, represented by the red (dark) and green (light) dashed lines, map possible states to payoffs with the red strategy dominating the green one.<sup>4</sup> To see that the red strategy is dominant in sub-figure (a), notice that its payoff is not worse than that of the green strategy at each possible state  $s$ . That the red strategy in sub-figure (b) is obviously dominant follows from observing that the worst payoff of the red strategy is not worse than the best payoff of the green strategy across all possible states. In sub-figure (c), STP-OD allows the agent to adopt the kind of contingent reasoning in the businessman example—the agent may find a relevant event  $E$  to make the comparison between two strategies obvious, i.e., on  $E$  the worst payoff of

<sup>4</sup>The co-movements between the red and green lines both downward sloping is only for expositional convenience.

the red strategy is no worse than the best possible payoff of the green strategy on the whole state space, and on  $E^c$  the best payoff of the green strategy is no better than the worst payoff of the red strategy on the whole state space. Finally, sub-figure (d) illustrates 3-OD in which STP is applied in a more sophisticated manner. The player identifies three relevant events  $E_1$ ,  $E_2$ , and  $E_3$  individually making the comparison between the red strategy and green strategy obvious in the sense that best payoff from the dominated strategy is no better than the worst outcome from the dominating strategy across all states in each event. It can be shown that OD is equivalent to 1-OD and STP-OD is equivalent to 2-OD (see Proposition 2 in Section 2).

Applying our solution concepts to mechanism design, we say that a dynamic game form is a  $k$ -OD implementation of a social choice function if the derived incomplete information game has a  $k$ -OD equilibrium delivering the same social outcome at every realization of the agents' type profile. When such a dynamic game form exists, we say that the social choice function is  $k$ -OD implementable. To study  $k$ -OD implementability, it is valuable to identify a canonical class of game forms such that as long as a strategy-proof social choice rule has a  $k$ -OD implementation, there exists one in this specific class. For strategy-proofness or dominance implementation, the classic revelation principle (Gibbard, 1973) comes in handy, providing a unique static mechanism—the direct mechanism—in which all agents simultaneously report their exact private types to the administrator. Note that direct mechanisms are without information feedback and truthful reporting may not be a  $k$ -OD equilibrium for a low sophistication level  $k$  that is achievable under some dynamic implementation.

A natural departure from the direct mechanism incorporating information feedback, we call the resulting dynamic game form a gradual mechanism, works as follows:<sup>5</sup>

Starting from the initial history, the administrator privately sends

---

<sup>5</sup>The idea of gradual mechanism is developed contemporaneously in this paper in parallel with our companion paper where it is called gradual revelation mechanism (Chew and Wang, 2022).



specific messages and forms to the active agents. Each form consists of a list of pairwise disjoint categories of the agent’s possible type and the agent can check the category to which she belongs. The message sent to each agent contains some information about how other agents have previously checked their forms.<sup>6</sup> The administrator keeps sending forms and messages and collecting returned forms until she collects enough information to arrive at a public outcome.

As it turns out, one extreme case of gradual mechanism, called round table mechanism (Mackenzie, 2020), in which the administrator *provides full information about how other agents checked their forms previously in each message*, serves as a revelation principle for OD implementation.

We show in Section 3 that the class of gradual mechanisms serves as a revelation principle for  $k$ -OD implementations (Theorem 1). Facilitated by the use of gradual mechanisms, we provide in Section 3 a necessary condition for  $k$ -OD implementability (Proposition 3) and demonstrate its usefulness in the application to anonymous generalized median voter schemes in Section 4. Section 5 offers a discussion of the relevant literature and some concluding remarks.

## 1 Preliminaries

There is a finite group of agents  $N$  who have interest on a set of public outcomes  $X$ . Each agent  $i \in N$  has a private type space  $\Theta_i$  whose element  $\theta_i$  corresponds to a complete and transitive preference ordering  $R(\theta_i)$  over outcomes in  $X$ . Conventionally, a type profile is written as  $\theta = (\theta_i, \theta_{-i}) \in \Theta = \prod_{i \in N} \Theta_i$  where  $\theta_{-i}$  is a type profile for agents other than  $i$ . The social planner wishes to implement a social choice function  $f : \Theta \rightarrow X$  that conditions a public outcome on agents’ type profile.

---

<sup>6</sup>Each history can be identified with an exact description of how each agent has checked a sequence of forms prior to this history. In the message to an agent, she is only informed that one history in her information set is reached.

To accommodate situations such as ascending price auction with a tie-breaking rule, we also consider stochastic social choice functions. Let  $\mathcal{A}_X$  be an algebra on  $X$ . A lottery is then a probability measure on  $(X, \mathcal{A}_X)$ , the space of which is denoted as  $\Delta(X)$ . A stochastic social choice function  $f$  maps a type profile in  $\Theta$  to a lottery in  $\Delta(X)$ .

We focus on implementation problems in which the main challenge for the social planner arises from private information, i.e., the type of any particular agent is known only to herself. To solve such a problem, the social planner designs a dynamic game form  $G$  with outcomes in  $X$  and invites agents, with any possible type profile  $\theta \in \Theta$ , to play the incomplete information game  $(G, \Theta)$ , in the hope that the game will deliver, in equilibrium, a public outcome stipulated by the social choice function.

Table 3: A list of notations in dynamic game forms

Name	Notation/Definition	Generic Element
Actions of $i \in N_0$	$A_i$	$a_i$
Action profiles	$A = \bigcup_{\emptyset \neq M \subset N_0} \prod_{i \in M} A_i$	$a = (a_i)_{i \in M}$
Histories	$\bar{H} \subseteq \bigcup_{t=0}^T A^t$ is a tree.	$h = \emptyset$ or $h = (h^{(1)}, \dots, h^{(t)})$
Precedence relation	$\preceq$	$h \preceq \bar{h}$
Terminal histories	$\preceq$ -maximal histories $Z \subseteq \bar{H}$	
Non-terminal histories	$H = \bar{H} \setminus Z$	
Outcome function	$\mathcal{X} : Z \rightarrow X$	
Active-player correspondence	$\mathbb{P} : H \rightarrow N_0$	
Active history of agent $i$	$H_i = \{h \in H : i \in \mathbb{P}(h)\}$	
Information sets	$\mathbf{H}_i$ is a partition of $H_i$	$\mathbf{h}_i$
Available actions at $\mathbf{h}_i$ or $h \in \mathbf{h}_i$	$A_i(h) = A_i(\mathbf{h}_i) \subseteq A_i$	
Space of interim strategies	$S_i$	$s_i : H_i \rightarrow A_i$

A dynamic game form  $G$  can be described by the following tuple

$$G = (\bar{H}, \{A_i, \mathbf{H}_i\}_{i \in N}, \mathcal{X})$$

and a randomized dynamic game form by

$$G = (\bar{H}, \{A_i, \mathbf{H}_i\}_{i \in N_0}, \mathcal{X}, m, \{\Omega, \mathcal{A}, \mu\})$$

where  $N_0 = N \cup \{0\}$  and agent 0 is the administrator who randomizes by adopting a mixed strategy  $m$  based on a randomization device (a probability space)  $(\Omega, \mathcal{A}, \mu)$ . We explicitly allow for simultaneous actions and model histories using sequences of action profiles (Osborne and Rubinstein, 1994; Battigalli et al., 2020). Table 3 lists some standard notations that will be used. A formal definition of dynamic game forms is offered in the appendix.

Starting from a non-terminal history  $h \in H$ , after specifying the profile of interim strategies  $s \in S = \prod_{i \in N} S_i$  and the realization  $\omega$  of administrator's random device, a unique terminal history  $Z(h, s, \omega)$  along with the associated outcome  $\mathcal{X}(h, s, \omega)$  will be determined. In these expressions, we will omit the history  $h$  if it is the initial history, i.e.,  $Z(\emptyset, s, \omega)$  will be denoted as  $Z(s, \omega)$ , and its associated outcome will simply be denoted as  $\mathcal{X}(s, \omega)$ .<sup>7</sup> Define  $H(\omega) = \{h \in \overline{H} : \exists s \in S \text{ s.t. } h \preceq Z(s, \omega)\}$  as the collection of histories that are consistent with the realization  $\omega$ . Similarly, define  $H(s_i) = \{h \in \overline{H} : \exists s_{-i} \in S_{-i} \text{ and } \omega \in \Omega \text{ s.t. } h \preceq Z(s_i, s_{-i}, \omega)\}$  to be the collection of histories that is consistent with the interim strategy  $s_i$  of agent  $i$ . Then  $H_i(s_i) = H(s_i) \cap H_i$  is its subset on which agent  $i$  is active. Define  $\mathbf{H}_i(s_i) = \{\mathbf{h}_i \in \mathbf{H}_i : \exists h \in H_i(s_i) \text{ s.t. } h \in \mathbf{h}_i\}$  as the collection of information sets of agent  $i$  consistent with her interim strategy  $s_i$ . Finally, define  $\mathbf{D}(s_i, s'_i)$  as the collection of  $\preceq$ -maximal information sets in  $\mathbf{H}_i(s_i) \cap \mathbf{H}_i(s'_i)$  which are the earliest information sets on which the two interim strategies  $s_i$  and  $s'_i$  diverge.

We use the combination of  $G$  and  $\Theta$ , i.e.,  $(G, \Theta)$ , to model an incomplete information game in which each agent  $i \in N$  only knows her own type.<sup>8</sup> Agent  $i$ 's (type) strategy in  $(G, \Theta)$  can be represented by  $\mathbb{S}_i : \Theta_i \rightarrow S_i$  (with  $\mathbb{S} = \prod_{i \in N} \mathbb{S}_i : \Theta \rightarrow S$  being a strategy profile) where  $\mathbb{S}_i(\theta_i) = s_i$  is the interim strategy agent  $i$  adopts when her private type is  $\theta_i$ .

We say that a social choice function  $f$  is implemented by  $(G, \Theta, \mathbb{S})$  if for each  $\theta \in \Theta$  it is the case that  $f(\theta) = \mathcal{X}(\mathbb{S}(\theta))$ . When  $f$  is a stochastic social choice function and  $G$  is randomized, we say  $(G, \Theta, \mathbb{S})$  implements  $f$  if for each

<sup>7</sup> $Z(h, s)$  and  $\mathcal{X}(h, s)$  can be defined for dynamic game form without randomization by dropping  $\omega$ . This also applies to subsequent notations.

<sup>8</sup>For solution concepts studied in this paper, players' belief about each other's types is irrelevant. Therefore, it is left unspecified.

$\theta \in \Theta$  and for each  $X' \in \mathcal{A}_X$ , let  $\Omega_{X'} = \{\omega \in \Omega : \mathcal{X}(\mathbb{S}(\theta), \omega) \in X'\}$ , then  $\Omega_{X'} \in \mathcal{A}$  and  $[f(\theta)](X') = \mu(\Omega_{X'})$ . That is, when agents act according to  $\mathbb{S}(\theta)$ , for each type profile  $\theta$ , the probability of obtaining any specific collection of outcomes  $X'$  is what is intended by the social choice function  $f(\theta)$ .

## 2 Obvious Dominance Facilitated by the Sure-thing Principle

To implement a (stochastic) social choice function  $f$  with a triple of  $(G, \Theta, \mathbb{S})$ , the social planner relies on the incomplete information game  $(G, \Theta)$  providing appropriate incentives for the agents to adopt the strategy profile  $\mathbb{S}$ . This is captured by  $\mathbb{S}$  being a proper solution of  $(G, \Theta)$ . For ease of comparison, we first provide the definition of (weak ) dominance as a comparison between two interim strategies under a specific type.<sup>9</sup>

**Definition 1.** *For any agent  $i \in N$  and any  $\theta_i \in \Theta_i$ , an interim strategy  $s_i$  **dominates**  $s'_i$  if it is the case that*

$$\mathcal{X}(s_i(\theta_i), s_{-i}, \omega) R(\theta_i) \mathcal{X}(s'_i, s_{-i}, \omega)$$

for all  $s_{-i} \in S_{-i}$  and any  $\omega \in \Omega$ .

A type strategy is **dominant** if the interim strategy it employs dominates any other interim strategies under each type. It is known that a social choice function  $f$  is dominance implementable if and only if it is strategy-proof, i.e.,  $f(\theta_i, \theta_{-i}) R(\theta_i) f(\theta'_i, \theta_{-i})$  for any  $\theta_i, \theta'_i \in \Theta_i$  and for each agent  $i \in N$ , and that the type strategy of truth telling is dominant in the direct mechanism where all agents simultaneously report their private types.

We are ready to introduce our solution concept of OD facilitated by STP (STP-OD). We first provide a definition of one interim strategy  $s_i$  of agent  $i$  obviously dominating another interim strategy  $s'_i$  facilitated by STP given a specific type  $\theta_i$ .

---

<sup>9</sup>The definitions in this section are given for the more complicated randomized dynamic games. They can be adapted to dynamic games without administrator's randomization by omitting  $\omega \in \Omega$ .

**Definition 2.** For any agent  $i$  and any  $\theta_i \in \Theta_i$ , an interim strategy  $s_i$  **obviously dominates  $s'_i$  facilitated by STP**, if for any  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ , there exists  $E \subseteq \mathbf{h}_i \times S_{-i} \times \Omega$  such that for all  $(h, s_{-i}, \omega), (h', s'_{-i}, \omega') \in \mathbf{h}_i \times S_{-i} \times \Omega$ , if  $(h, s_{-i}, \omega) \in E$  or  $(h', s'_{-i}, \omega') \in E^c$ , then

$$\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega').$$

The event  $E$  in Definition 2 is the **relevant event** on which sure-thing principle can be applied to render the dominance relation between  $s_i$  and  $s'_i$  obvious. Applying this definition to the incomplete information dynamic game  $(G, \Theta)$ , we provide the following definition of OD facilitated by STP for type strategies—a type strategy is STP-OD if it always employs an interim strategy that obviously dominates all other interim strategies facilitated by STP.

**Definition 3.**  $\mathbb{S}_i$  is **obviously dominant facilitated by STP**, if for any  $\theta_i \in \Theta_i$ , any  $\mathbf{h}_i \in \mathbf{H}_i(\mathbb{S}_i(\theta_i))$ , and any  $s'_i$  such that  $s'_i(\mathbf{h}_i) \neq [\mathbb{S}_i(\theta_i)](\mathbf{h}_i)$ , there exists  $E \subseteq \mathbf{h}_i \times S_{-i} \times \Omega$  such that for all  $(h, s_{-i}, \omega), (h', s'_{-i}, \omega') \in \mathbf{h}_i \times S_{-i} \times \Omega$ , if  $(h, s_{-i}, \omega) \in E$  or  $(h', s'_{-i}, \omega') \in E^c$ , then

$$\mathcal{X}(h, \mathbb{S}_i(\theta_i), s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega').$$

For any private type and at each information set of an agent, in comparing the dominating strategy and any deviating strategy, the uncertainty faced by the agent stems from three sources: (1) uncertainty in the history reached at this information set, (2) uncertainty in the profile of interim strategies adopted by other agents, and (3) uncertainty in the realization of administrator's randomization device.<sup>10</sup> STP-OD allows for obviousness in the comparison between the dominating strategy and the deviating one to be facilitated by STP in the sense that the agent could identify a relevant event such that the outcome from the dominating strategy when this event happens cannot be worse than any outcome following the deviating strategy across all possibilities and, should the complementary event happen, the outcome from the deviating

<sup>10</sup>If  $G$  is not randomized, the third source of uncertainty does not exist.

strategy cannot be better than any outcome following the dominating strategy across all possibilities. To link STP-OD with OD proposed in Li (2017), notice that by requiring  $E = \mathbf{h}_i \times S_{-i} \times \Omega$  in the definitions 2 and 3, STP-OD becomes OD.

The solution concept of STP-OD can be generalized further by allowing for more sophisticated applications of STP. In this sense, STP-OD in definitions 2 and 3 requires only the plainest application of STP and therefore the most stringent solution concept among such generalizations. Note that the following further generalization is given as a binary comparison between two interim strategies under a specific type of an agent. It is apparent that it enables a definition for type strategies as well.

**Definition 4.** For any agent  $i$  and any  $\theta_i \in \Theta_i$ , an interim strategy  $s_i$  **obviously dominates**  $s'_i$  **facilitated by STP with sophistication level  $k$** , if for any  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ , there exists a partition  $\{E_1, \dots, E_k\}$  of  $\mathbf{h}_i \times S_{-i} \times \Omega$  such that for all  $1 \leq l \leq k$ , if both  $(h, s_{-i}, \omega), (h', s'_{-i}, \omega') \in E_l$ , then

$$\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega').$$

From this definition, 1-OD is equivalent to Li's (2017) OD. We next develop an alternative definition to  $k$ -OD, based on the idea of eliminating obviously dominating states iteratively. This yields, in Proposition 1, an equivalence between these two definitions encompassing that between STP-OD and 2-OD. It also delivers a method in Proposition 2 to derive the lowest level of sophistication required by a dynamic mechanism.

In the definition of STP-OD, the relevant event  $E$  can be considered as a set of *obviously dominating states* in the sense that any outcome following the dominating strategy on  $E$  is weakly preferred to any outcome following the dominated strategy on  $S = \mathbf{h}_i \times S_{-i} \times \Omega$ . After *eliminating* these obviously dominating states, any outcome following the dominating strategy on the remaining  $E^c = S \setminus E$  is weakly preferred to any outcome following the dominated strategy. In the case of  $k$ -OD, eliminating obviously dominating states once may not suffice. Then, an iterated elimination of obviously dominating

states among the remaining state space may give rise eventually to a partition, on which a player can employ STP style contingent reasoning in place of state-by-state comparison between the dominating and dominated strategies.

**Definition 5.** For any agent  $i$  and any  $\theta_i \in \Theta_i$ , an interim strategy  $s_i$  **obviously dominates**  $s'_i$  **facilitated by STP with  $k$  iterations** if for any  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ , there exists a  $k$  partition  $\{E_1, \dots, E_k\}$  of  $\mathbf{h}_i \times S_{-i} \times \Omega$  such that  $(h, s_{-i}, \omega) \in E_l$  and  $(h', s'_{-i}, \omega') \in E_l \cup \dots \cup E_k$  imply

$$\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega')$$

for all  $1 \leq l \leq k$ .

In the above two definitions,  $E_1, \dots, E_k$  are the **relevant events** in players' application of STP to render the dominance relation obvious. An immediate implication of iterative- $k$ -OD is that  $(h', s'_{-i}, \omega') \in E_l$  and  $(h, s_{-i}, \omega) \in E_1 \cup \dots \cup E_l$  imply  $\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega')$  for all  $1 \leq l \leq k$ . This yields an alternative understanding of iterative- $k$ -OD in terms of iterated elimination of obviously dominated states. Moreover, if obviousness of dominance between two strategies is facilitated by STP with level  $k$ , it is also with level  $l$  when  $l > k$ .

**Proposition 1.** For any agent  $i$  and any  $\theta_i \in \Theta_i$ , an interim strategy  $s_i$   $k$ -obviously dominates  $s'_i$  if and only if  $s_i$  iterative- $k$ -obviously dominates  $s'_i$ .

*Proof.* We only show that when  $s_i$   $k$ -obviously dominates  $s'_i$  for  $\theta_i \in \Theta_i$ , it is also the case that  $s_i$  iterative- $k$ -obviously dominates  $s'_i$ . The opposite direction is straightforward. By definition, there exists a partition  $\{E_1, \dots, E_k\}$  of  $\mathbf{h}_i \times S_{-i} \times \Omega$  such that  $(h, s_{-i}, \omega) \in E_l$  and  $(h', s'_{-i}, \omega') \in E_l$  imply

$$\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega')$$

for all  $1 \leq l \leq k$ . Let  $E_0 = \emptyset$ , define a permutation function  $\sigma : \{0, 1, \dots, k\} \rightarrow \{0, 1, \dots, k\}$  with  $\sigma(0) = 0$  such that for each  $l \leq k$  there does not exist

$(h', s'_{-i}, \omega') \in \mathbf{h}_i \times S_{-i} \times \Omega \setminus \bigcup_{j=0}^{l-1} E_{\sigma(j)}$  such that

$$\mathcal{X}(h, s'_i, s_{-i}, \omega) P(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega')$$

for all  $(h, s_{-i}, \omega) \in E_{\sigma(l)}$  in which  $P(\theta_i)$  is the strict preference relation defined by  $R(\theta_i)$ . Notice that  $\sigma$  is always well defined and for any  $1 \leq l \leq k$ , any  $(h', s'_{-i}, \omega') \in E_{\sigma(l)} \cup \dots \cup E_{\sigma(k)}$ , there exists  $(h, s_{-i}, \omega) \in E_{\sigma(l)}$  such that

$$\mathcal{X}(h, s'_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega').$$

Observe that  $(h, s_{-i}, \omega) \in E_{\sigma(l)}$  and  $(h', s'_{-i}, \omega') \in E_{\sigma(l)} \cup \dots \cup E_{\sigma(k)}$  imply

$$\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega')$$

for all  $1 \leq l \leq k$ . Since  $\mathbf{h}_i$  is chosen arbitrarily in  $\mathbf{D}(s_i, s'_i)$ ,  $s_i$  iteratively- $k$ -obviously dominates  $s'_i$ .  $\square$

For any agent  $i$ , any  $\theta_i \in \Theta_i$ , and any information set  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ , we provide an iterated elimination procedure which is “greedy” in the sense that it eliminates, at each step, the biggest collection of obviously dominating states. Let  $S = \mathbf{h}_i \times S_{-i} \times \Omega$  and  $E_0^* = \emptyset$ . For each step  $l \geq 1$ , let

$$E_l^* = \left\{ (h, s_{-i}, \omega) \in S \setminus \bigcup_{t=0}^{l-1} E_t^* : \begin{array}{l} \mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega') \\ \text{for all } (h', s'_{-i}, \omega') \in S \setminus E_0^* \cup \dots \cup E_{l-1}^* \end{array} \right\}$$

When  $E_k^* \neq \emptyset$  and  $E_{k+1}^* = \emptyset$ , we say that the greedy iterative elimination of obviously dominating states delivers a partition with  $k$  cells for  $\theta_i$  and  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ . When comparing a dominating strategy with a dominated strategy, we can show that this greedy iteration procedure delivers the greatest lower bound for the required sophistication level in applying STP.

**Proposition 2.** *For any agent  $i$  and any  $\theta_i \in \Theta_i$ , an interim strategy  $s_i$   $k$ -obviously dominates  $s'_i$  if and only if the greedy iterative elimination of obviously dominating states delivers a partition with no more than  $k$  cells for each  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ .*



*Proof.* It is trivial to show the “if” direction. Now, suppose  $s_i$   $k$ -obviously dominates  $s'_i$  while the greedy iterative elimination of obviously dominating states delivers a partition with  $k'$  (with  $k' > k$ ) cells for a specific  $\mathbf{h}_i \in \mathbf{D}(s_i, s'_i)$ . This implies that there exists a partition  $\{E_1, \dots, E_k\}$  of  $\mathbf{h}_i \times S_{-i} \times \Omega$  such that  $(h, s_{-i}, \omega) \in E_l$  and  $(h', s'_{-i}, \omega') \in E_l \cup \dots \cup E_k$  imply

$$\mathcal{X}(h, s_i, s_{-i}, \omega) R(\theta_i) \mathcal{X}(h', s'_i, s'_{-i}, \omega')$$

for all  $1 \leq l \leq k$ .

Notice that  $E_1 \subseteq E_1^*$  by the definition of  $E_1^*$ . Notice also that, inductively,  $E_1 \cup \dots \cup E_l \subseteq E_1^* \cup \dots \cup E_l^*$  implies that  $E_1 \cup \dots \cup E_{l+1} \subseteq E_1^* \cup \dots \cup E_{l+1}^*$  for all  $1 \leq l \leq k-1$ . Therefore,  $E_1 \cup \dots \cup E_k \subseteq E_1^* \cup \dots \cup E_k^* \subsetneq \mathbf{h}_i \times S_{-i} \times \Omega$ . This is a contradiction.  $\square$

Finally, we say (randomized)  $(G, \Theta, \mathbb{S})$  is a  **$k$ -OD implementation** of a (stochastic) social choice function  $f$  if  $\mathbb{S}$  consists of  $k$ -OD strategies for the incomplete information game  $(G, \Theta)$ . When it exists, we may omit mention of  $(G, \Theta, \mathbb{S})$  and say that  $f$  is  $k$ -OD implementable. If a social choice function  $f$  is  $k$ -OD implementable, it is also  $l$ -OD implementable for  $l > k$ . In other words, the class of  $k$ -OD implementable social choice functions enlarges as  $k$  increases. Should  $f$  be  $k$ -OD implementable but not  $(k-1)$ -OD implementable,  $k$  would be the minimal level of sophistication required for contingent reasoning facilitated by STP to simplify comparisons between dominating and dominated strategies in any dynamic implementation of  $f$ .

### 3 Gradual Mechanism and $k$ -OD implementation

As anticipated in the introduction, we develop a definition of (randomized) gradual mechanism to serve as a revelation principle for  $k$ -OD implementation of a (stochastic) social choice function. An example of gradual mechanism is provided in the next section in the context of the generalized median voter scheme.

For any history  $h \in \overline{H}$  of a dynamic game form  $G$ , we define  $h_i$  as the

sequence of actions player  $i$  has taken up to  $h$ . Formally, a gradual mechanism is a dynamic game form as defined below.

**Definition 6.** A *gradual mechanism*  $G$  for social choice function  $f$  is a dynamic game form such that:

1. For each agent  $i \in N$ , the collection of feasible actions are non-empty subsets of her possible types, i.e.,  $A_i = 2^{\Theta_i} \setminus \{\emptyset\}$ .
2. For each agent  $i \in N$  and any history  $h \in H_i$  where agent  $i$  is active, the available actions  $A_i(h)$  are pairwise disjoint subsets of  $\Theta_i$  whose union covers her previous reports, i.e.,  $a_i \cap a'_i = \emptyset$  for any  $a_i, a'_i \in A_i(h)$  and  $\bigcap h_i \subseteq \bigcup A_i(h)$ .<sup>11</sup>
3. The public outcome assigned to any terminal history  $z \in Z$  is aligned with the social choice function based on information accrued up to  $z$ , i.e.,  $\prod_{i \in N} \bigcap z_i \neq \emptyset$  for each  $z \in Z$  and for any  $\theta \in \prod_{i \in N} \bigcap z_i$ , it is the case that  $\mathcal{X}(z) = f(\theta)$ .

As we explicitly allow for simultaneous moves, gradual mechanisms encompass both direct mechanisms serving as the revelation principle for strategy-proof social choice functions (Gibbard, 1973) and round table mechanisms (i.e., gradual mechanisms with perfect information, see Mackenzie, 2020) which provide a revelation principle for OD implementable social choice functions. The following definition of a randomized gradual mechanism serves as a revelation principle for  $k$ -OD implementation of a stochastic social choice function in which the administrator randomly chooses and publicly announces a subsequent gradual mechanism based on the realization of her randomization device.

**Definition 7.** A *randomized gradual mechanism*  $G$  for stochastic social choice function  $f$  is a randomized dynamic game form such that:

---

<sup>11</sup>By the first condition,  $h_i$  is a sequence of subsets of agent  $i$ 's private type space whose intersection is denoted as  $\bigcap h_i$ . Note that  $\bigcup A_i(h) = \Theta_i$  if  $h_i = \emptyset$ , i.e., the whole type space is still possible when no action has been taken.

1. Administrator acts only at history  $\emptyset$  and the available actions for her are the realizations of her randomization device, i.e.,  $A_0 = \Omega$ , and  $m(\omega)(\emptyset) = \omega$ .
2. For any  $\omega \in \Omega$ , the subgame  $G^\omega$  starting at history  $(\omega)$  is a gradual mechanism for a social choice function  $f^\omega$  such that collectively for any type profile  $\theta \in \Theta$  and subset of public outcomes  $X' \in \mathcal{A}_X$ , it is the case that  $\mu(\{\omega \in \Omega : f^\omega(\theta) \in X'\}) = f(\theta)(X')$ .
3. Administrator's choice of gradual mechanism is commonly known by all agents, i.e., for any agent  $i$  and any of her information set  $\mathbf{h}_i$ , there exists  $\omega$  such that  $\mathbf{h}_i \subseteq H(\omega)$ .

In a (randomized) gradual mechanism, the type-revealing strategy  $\mathbb{T}_i$  for agent  $i \in N$  refers to a type strategy satisfying  $\theta_i \in \mathbb{T}_i(\theta_i)(h)$  for each type  $\theta_i$  and at each  $h \in H_i$  with  $\theta_i \in \bigcup A_i(h)$ . The following result guarantees that we can safely restrict our search for a (randomized) dynamic game form  $k$ -OD implementing a given (stochastic) social choice function within the class of (randomized) gradual mechanisms.

**Theorem 1.** *If a (stochastic) social choice function  $f$  is  $k$ -OD implementable, then it has a  $k$ -OD implementation with a (randomized) gradual mechanism and the corresponding type-revealing strategies.*

The theorem is proved using an adaptation of the pruning principle in (Li, 2017) to demonstrate that any  $k$ -OD implementation  $(G, \Theta, \mathbb{S})$  of a (stochastic) social choice function  $f$  can be transformed into  $(G^*, \Theta, \mathbb{T})$  where  $G^*$  is a (randomized) gradual mechanism and  $\mathbb{T}$  is the type-revealing strategy in  $G^*$  such that  $(G^*, \Theta, \mathbb{T})$  is also a  $k$ -OD implementation of  $f$ . The following pruning operation can transform, with respect to a strategy profile, a dynamic game form into a gradual mechanism (after relabeling its actions), in which the original strategy profile is also transformed into a profile of type-revealing strategies in the resulting gradual mechanism.

**Definition 8.** Let  $G = (\bar{H}, \{A_i, \mathbf{H}_i\}_{i \in N}, \mathcal{X})$ . The **pruning** of  $G$  with respect to a strategy profile  $\mathbb{S}$  of the incomplete information game  $(G, \Theta)$  is given by  $G^* = (\bar{H}^*, \{A_i^*, \mathbf{H}_i^*\}_{i \in N}, \mathcal{X}^*)$  such that:

1. In  $G^*$ , the set of feasible actions for each players is the same with that in  $G$ , i.e.,  $A_i^* = A_i$ .
2. The set of histories of  $G^*$  consists histories in  $G$  that can be reached for some  $\theta \in \Theta$ , i.e.,  $\bar{H}^* = \{h \in \bar{H} : h \preceq Z(\mathbb{S}(\theta)) \text{ for some } \theta \in \Theta\}$ .
3. In  $G^*$ ,  $\{\mathbf{H}_i^*\}_{i \in N}$  and  $\mathcal{X}^*$  are restrictions on the domain  $\bar{H}^*$  of their counterparts in  $G$ .

Let  $\mathbb{S}^*$  be a profile of type strategies in  $(G^*, \Theta)$  such that the same terminal history is reached for every type profile  $\theta$  with  $\mathbb{S}$  in  $(G, \Theta)$ , i.e.,  $Z(\mathbb{S}^*(\theta)) = Z(\mathbb{S}(\theta))$  for each  $\theta \in \Theta$ . This means that agent  $i$  adopting the interim strategy  $\mathbb{S}_i^*(\theta_i)$  makes the same decisions at each corresponding history when she adopts the interim strategy  $\mathbb{S}_i(\theta_i)$  in  $(G, \Theta)$ . Since the space of each agent's interim strategies shrinks after pruning, if  $\mathbb{S}$  is a  $k$ -OD strategy profile for  $(G, \Theta)$ , so is  $\mathbb{S}_i^*$  for  $(G^*, \Theta)$ . Each history  $h^*$  in  $\bar{H}^*$  is associated with a non-empty "rectangle"  $\prod_{i \in N} \Theta_i(h^*)$  in  $\Theta$  with  $\Theta_i(h^*) \subseteq \Theta_i$ , capturing the collection of type profiles for which  $h^*$  is reached when agents follow  $\mathbb{S}^*$ . Suppose agent  $i$  is active at  $h^*$ , each of her available action  $a_i \in A_i^*(h^*)$  can be identified by a subset  $a_i^* \subseteq \Theta_i(h^*)$  such that any immediate successor  $\bar{h}^*$  of  $h^*$  after agent  $i$  taking action  $a_i$  is associated with  $\Theta_i(\bar{h}^*) = a_i^*$ . We can relabel all actions in  $G^*$  as the corresponding subsets of the agents' type spaces. Doing so, we relabel  $G^*$  as a gradual mechanism and also relabel  $\mathbb{S}^*$  as a type-revealing strategy profile. Notice that the gradual mechanism derived in this process has an additional property: For each agent  $i$  and each  $h_i \in H_i$ , it is the case that  $\bigcup A_i(h) = \bigcap h_i$ .<sup>12</sup>

Two additional operations, de-randomization and pre-randomization,<sup>13</sup> are

---

<sup>12</sup>See our companion paper (Chew and Wang, 2022) for an interpretation of this property in terms of immediacy.

<sup>13</sup>See Ashlagi and Gonczarowski (2018) for de-randomization and Pycia and Troyan (2021) for pre-randomization.

needed to transform a randomized implementation of a stochastic social choice function to a randomized gradual mechanism. These operations shift all of administrator's randomization to the initial history. To arrive at a randomized gradual mechanism, let each subgame starting from a realization of administrator's random signal be replaced by its pruned gradual mechanism. Formal definitions of these two operations in our framework and the corresponding part of the proof of Theorem 1 are provided in the appendix.

Theorem 1 helps develop the following idea of a starting point to find the collection of  $k$ -OD implementations of a social choice function and a way to check if they exist. We begin with investigating the direct mechanism in which all agents simultaneously report their types and check for each agent  $i \in N$  whether there exists some information about her private type (modeled by a partition  $\Pi_i$  of her type space  $\Theta_i$ ) for which type revelation is  $k$ -OD, the meaning of which we elaborate below. We say that **type revelation is  $k$ -OD between  $\Pi_i$ 's two distinct cells  $[\theta_i]$  and  $[\tilde{\theta}_i]$**  if for any  $\theta'_i \in [\theta_i]$  (or any  $\tilde{\theta}'_i \in [\tilde{\theta}_i]$ ), whenever it is the true type, reporting it truthfully in the direct mechanism  $k$ -obviously dominates reporting any type  $\tilde{\theta}'_i \in [\tilde{\theta}_i]$  (or any type  $\theta'_i \in [\theta_i]$ ). We say that **type revelation is  $k$ -OD among  $\Pi_i$**  if type revelation between any two cells of  $\Pi_i$  is  $k$ -OD. Potentially, information  $\Pi_i$  of agent  $i$  can be transmitted in a gradual mechanism and applied by some other agents, making their type-revealing decisions easier.

From another perspective, for two types  $\theta_i$  and  $\theta'_i$  belonging to the same cell, some information about other agents' type profile may be needed to make type-revelation between them  $k$ -OD. These understandings of what information may be available or required could facilitate the search process by suggesting what to attempt at the initial history where some information must be transmitted among the agents to render the overall type-revealing strategy profile  $k$ -OD. Note that if type revelation is  $k$ -OD among  $\Pi_i$  in the direct mechanism, this is also the case among any  $\Pi'_i$  that is coarser than  $\Pi_i$ . Note also that if type revelation is  $k$ -OD among  $\Pi_i^1$  and among  $\Pi_i^2$  in the direct mechanism, this is also the case among  $\Pi_i$  which is the coarsest common refinement of  $\Pi_i^1$  and  $\Pi_i^2$ .

The discussion above yields the following necessary condition for  $k$ -OD implementability.

**Proposition 3.** *If  $f$  is  $k$ -OD implementable, then there exist some agent  $i \in N$  and some  $E \subseteq \Theta_i$  such that type revelation between  $E$  and  $E^c$  is  $k$ -OD in the direct mechanism.*

*Proof.* For any given gradual mechanism  $G$  implementing social choice function  $f$ , we first identify  $\Theta_i$  as a subset of  $S_i$  for each agent  $i$ . Let  $\theta_i \in \Theta_i$  correspond to a  $s_i \in S_i$  such that  $\theta_i \in s_i(\mathbf{h}_i)$  for all  $\mathbf{h}_i \in \mathbf{H}_i$  satisfying  $\theta_i \in \bigcup A_i(\mathbf{h}_i)$ . Notice that  $\mathcal{X}(\theta) = f(\theta)$  under this identification.

Suppose gradual mechanism  $G$  is a  $k$ -OD implementation of  $f$  and let agent  $i$  be active at the initial history. By definition of  $k$ -OD implementation, for any  $a, a' \in A_i(\emptyset)$ , any  $\theta_i \in a$ , and any  $\theta'_i \in a'$ , there exists a partition  $\{E_1, \dots, E_k\}$  of  $\prod_{j \neq i} S_j$  such that  $\mathcal{X}(\theta_i, s_i) R(\theta_i) \mathcal{X}(\theta'_i, s'_i)$  if  $s_{-i}, s'_{-i} \in E_l$  for all  $1 \leq l \leq k$ , in which  $\theta_i$  and  $\theta'_i$  are two interim strategies identified previously. Notice that  $\{E_1, \dots, E_k\}$  defines a partition  $\{F_1, \dots, F_k\}$  on  $\Theta_{-i}$  by  $F_l = E_l \cap \Theta_{-i}$  for all  $1 \leq l \leq k$ . Let  $\{E, E^c\}$  be any binary partition of  $\Theta_i$  that is a coarsening of  $A_i(\emptyset)$ . We have shown, in the direct mechanism, that type revelation between them is  $k$ -OD.  $\square$

Intuitively, suppose agent  $i$  is active at the initial history of a gradual mechanism which  $k$ -OD implements social choice function  $f$  with available actions  $A_i(\emptyset) = \Pi_i$ . By definition, the type-revealing strategy  $k$ -obviously dominates any deviating type strategy at the initial history. In the direct mechanism, agent  $i$  faces less uncertainty given that other agents cannot condition their actions on her action. Therefore type revelation must be  $k$ -OD for agent  $i$  among  $\Pi_i$  in the direct mechanism.

Moreover,  $k$ -OD has a hereditary property: A  $k$ -OD implementable social choice function will remain  $k$ -OD implementable on a smaller domain of type profiles. This is true since the gradual mechanism that  $k$ -OD implements  $f$  on the domain of  $\Theta$  also  $k$ -OD implements  $f$  on the domain of  $\tilde{\Theta}$ .

## 4 Generalized Median Voter Scheme

Consider the problem of voting on the domain of single-peaked preferences. There is a finite and linearly ordered outcome space  $X$  with a top (best) alternative for each agent such that alternatives that are further away from the best are progressively less preferred by the agent. This framework captures situations such as possible levels or locations of a public good, temperatures in a meeting room, and platforms of political parties. The OD implementation of these social choice functions is studied in Arribillaga et al. (2020).<sup>14</sup>

Formally, we write  $X = \{x_1, \dots, x_M\}$  with subscripts representing a linear order on outcomes. Every agent  $i$  has the same type space  $\Theta_i = \mathcal{R}$  in which each preference relation  $R$  satisfies single-peakedness, i.e., there is a *peak*  $t(R) \in \{1, \dots, M\}$  such that (i)  $x_{t(R)} R x_a$  and (ii)  $a < b \leq t(R)$  or  $t(R) \leq b < a$  implies  $x_b P x_a$  (where  $P$  refers to the strict preference) for all  $a, b \in \{1, \dots, M\}$ . In addition to strategy-proofness, we consider the following two properties of a social choice function  $f$ . We say that  $f$  is *onto* if for each outcome  $x_a \in X$ , there is a type profile  $R \in \mathcal{R}^N$  such that  $f(R) = x_a$  and that  $f$  is *anonymous* if  $f(R) = f(\tilde{R})$  for any type profile  $R \in \mathcal{R}^N$  and any  $\tilde{R}$  that is a permutation of  $R$ .

On the domain of single-peaked preferences, it is known that a social choice function  $f$  is strategy-proof and onto if and only if it is a generalized median voter scheme (Moulin, 1980; Barberà et al., 1993). We present here results pertaining to STP-OD implementation of generalized median voter schemes with two levels and  $k$ -OD implementation of anonymous generalized median voter schemes.

When voting between two candidates, i.e.,  $X = \{x_1, x_2\}$ , generalized median voter schemes can be described simply by a single family of minimal winning coalitions  $\mathcal{C}$ , i.e., a collection of subsets of agents such that  $S \not\subseteq T$  for any two distinct  $S, T \in \mathcal{C}$ , such that  $x_1$  is the public outcome unless the subset of subjects who prefer  $x_2$  contains a committee in  $\mathcal{C}$ , in which case  $x_2$  will be

---

<sup>14</sup>Also see Bade and Gonczarowski (2017) for a similar setting with infinite number of levels.

the outcome.

Let us first verify that type revelation is not OD in the direct mechanism. Suppose agent  $i$  does not form a singleton winning committee, i.e.,  $\{i\} \notin \mathcal{C}$ , and that there exists a winning committee  $S$  such that  $i \notin S$ . Truthful reporting, for each of the two types of agent  $i$ , could yield the less preferred public outcome while misrepresentation may deliver the preferred outcome.

To see that the type revelation is STP-OD in the direct mechanism, consider first the case when agent  $i$  prefers  $x_2$  over  $x_1$ . Agent  $i$  may find the event  $E$  defined by the existence of committee in  $\mathcal{C}$  voting for  $x_2$  to be relevant. When there is such a committee, agent  $i$ 's voting will not change the outcome  $x_2$  which is the preferred one. It is therefore obvious that truthful reporting is weakly better. Should this not be the case, voting for  $x_1$  will result in the less preferred outcome  $x_1$ . We can do a similar analysis for the case of agent  $i$  preferring  $x_1$  over  $x_2$  by applying STP to the same event.

**Observation 1.** *When there are two levels, the direct mechanism STP-OD implements the generalized median voter scheme.*

For the more general case of  $M$  levels of a public good  $X = \{x_1, \dots, x_M\}$ , it is known that a generalized median voter scheme is anonymous if and only if it can be described by an increasing sequence of voting quotas  $1 = T_1 \leq T_2 \leq \dots \leq T_M \leq T_{M+1} = N$  such that  $x_k$  is the public outcome if there are at least  $T_k$  agents whose peaks are not less than  $k$ , i.e.,  $\#\{i \in N : t(R_i) \geq k\} \geq T_k$ , while the number of agents whose peaks are greater than  $k$  is less than  $T_{k+1}$ , i.e.,  $\#\{i \in N : t(R_i) > k\} < T_{k+1}$ . From Arribillaga et al. (2020), we know that the anonymous a generalized median voter scheme is OD implementable if and only if the voting quotas are either 1 or  $N$ .

Consider  $k$ -OD implementation of generalized median voter schemes. Suppose  $l$  is the highest level with voting quota being 1 and there are  $(k-1)$  levels with quotas strictly between 1 and  $N$ , then it must be that  $1 = T_l < T_{l+1} \leq \dots \leq T_{l+k-1} < T_{l+k} = N$ . We can show that such a generalized median voter scheme is  $k$ -OD implementable via the following gradual mechanism.

1. At the initial history, administrator sends a form to each agent with



the following three categories concerning her type  $R$ : (a)  $t(R) \leq l$ ; (b)  $l < t(R) < l + k$ ; (c)  $t(R) \geq l + k$ .

2. In the following two cases, administrator have collected enough information to determine the final level after collecting the returned forms.

(a) When there are less than  $T_{l+1}$  agents but at least one checks category  $l < t(R) < l + k$  or category  $t(R) \geq l + k$ , administrator has collected enough information to set the level at  $l$ .

(b) When, for some  $1 < k' < k$ , there are no less than  $T_{l+k'}$  agents checking category  $t(R) \geq l + k$  and less than  $T_{l+k'+1}$  agents checking category  $t(R) \geq l + k$  or category  $l < t(R) < l + k$ , administrator has collected enough information to set the level at  $l + k'$ .

3. Should administrator not be able to set the level in Step 2, it must fall into the following three cases.

(a) When all agents check category  $t(R) \leq l$ , administrator does not have enough information to determine the level among  $\{1, \dots, l\}$ .

(b) When all agents check category  $t(R) \geq l + k$ , administrator does not have enough information to determine the the level among  $\{l + k, \dots, M\}$ .

(c) When, for some  $1 < k' < k'' < k$ , there are no less than  $T_{l+k'}$  and more than  $T_{l+k'-1}$  agents checking category  $t(R) \geq l + k$  and no less than  $T_{l+k''}$  while less than  $T_{l+k''+1}$  agents checking category  $t(R) \geq l + k$  or category  $l < t(R) < l + k$ , administrator does not have enough information to determine the level among  $\{l + k', \dots, l + k''\}$ .

4. For case 3(a) (or 3(b) respectively), administrator proceeds as follows.

(a) Sends each agent a form containing two categories (i)  $t(R) = l$  (or  $t(R) = l + k$  in case 3(b)) and (ii)  $t(R) < l$  (or  $t(R) > l + k$  in case 3(b)). Informs each agent that all agents have previously checked  $t(R) \leq l$  (or  $t(R) \geq l + k$  in case 3(b)).

- (b) When at least one agent checks  $t(R) = l$  (or  $t(R) = l + k$  in case 3(b)), administrator has collected enough information to set the level at  $l$  (or  $l + k$  in case 3(b)).
  - (c) Otherwise, iteratively for  $h \geq 0$ , suppose the administrator cannot set the level to be  $l - h$  (or  $l + k + h$  in case 3(b)), then each agent is sent a form containing two categories (i)  $t(R) = l - h - 1$  (or  $t(R) = l + k + h + 1$  in case 3(b)) and (ii)  $t(R) < l - h - 1$  (or  $t(R) > l + k + h + 1$  in case 3(b)) and informed that all agents have previously checked  $t(R) < l - h$  (or  $t(R) > l + k + h$  in case 3(b)).
  - (d) When at least one agent checks  $t(R) = l - h - 1$  (or  $t(R) = l + k + h + 1$  in case 3(b)), administrator has collected enough information to set the level at  $l - h - 1$  (or  $l + k + h + 1$  in case 3(b)).
5. In case 3(c), each agent who has checked category  $l < t(R) < l + k$  receives a form with  $k - 1$  categories consisting of  $t(R) = l + h$  for one  $1 \leq h \leq k - 1$  and is informed that more information is needed to determine the level among  $\{l + 1, \dots, l + k - 1\}$ . After collecting the returned forms, administrator has enough information to set the level.

In addition,  $k$  is lowest level of sophistication that is achievable in any dynamic implementation of such an anonymous generalized median voter scheme. Summarizing, we have the following proposition.

**Proposition 4.** *An anonymous generalized median voter scheme is  $k$ -OD implementable if and only if there exist at most  $k - 1$  levels with voting quotas strictly between 1 and  $N$ .*

*Proof.* Observe that for each agent  $i$  and for each information set  $\mathbf{h}_i$ , the number of overlapping levels that could follow any two actions  $a, a'$  available at  $\mathbf{h}_i$  is at most  $k$ . This is obvious for information sets after the initial history. For the initial history, notice that (i) the possible levels after choosing category  $t(R) \leq l$  are among  $\{1, \dots, l + k - 1\}$ ; (ii) the possible levels after choosing category  $l < t(R) < l + k$  are among  $\{l, \dots, l + k - 1\}$ ; and (iii) the possible levels after choosing category  $t(R) \geq l + k$  are among  $\{l, \dots, M\}$ .

Since the overlapping levels are less than  $k$  after any two actions at each information sets, given any true type  $\theta_i$ , the greedy iterative elimination of obviously dominating states delivers a partition with no more than  $k$  cells. By Proposition 2, type-revealing strategy is  $k$ -OD and therefore anonymous generalized median voter schemes with  $k - 1$  levels of quotas being strictly 1 and  $N$  is  $k$ -OD implementable.

We next show that such an anonymous generalized median voter scheme is not a  $(k - 1)$ -OD implementation. For each agent  $i$ , let  $\tilde{\Theta}_i$  be the subset of  $\Theta_i$  in which  $l \leq t(R) \leq l + k$ . Consider the following table summarizing some outcome relevant events of the direct mechanism in which the action  $l + h$  corresponds to reporting any  $R$  with  $t(R) = l + h$ .

Events Actions	$E_l^1$	$E_l^2$	$\dots$	$E_{l+h-1}^2$	$E_{l+h}^1$	$E_{l+h}^2$	$\dots$	$E_{l+k-1}^1$	$E_{l+k-1}^2$
$l$	$l$	$l$	$\dots$	$l+h-1$	$l+h$	$l+h$	$\dots$	$l+k-1$	$l+k-1$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$l+h$	$l$	$l+1$	$\dots$	$l+h$	$l+h$	$l+h$	$\dots$	$l+k-1$	$l+k-1$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$l+k$	$l$	$l+1$	$\dots$	$l+h$	$l+h$	$l+h+1$	$\dots$	$l+k-1$	$l+k$

For each  $1 \leq h \leq k - 1$ ,  $E_{l+h}^1$  is the event in which there are no less than  $T_{l+h}$  agents reporting  $t(R) \geq l + h$  while there are less than  $T_{l+h+1} - 2$  agents reporting  $t(R) > l + h$ ;  $E_{l+h}^2$  is the event in which there are no less than  $T_{l+h}$  agents reporting  $t(R) \geq l + h$  while there are exactly  $T_{l+h+1} - 1$  agents reporting  $t(R) > l + h$  with at least one agent reporting  $t(R) = l + h + 1$ . Notice that each of these events is non-empty and any pair of them is disjoint.

Observe that for any two actions  $l + h$  and  $l + h'$  ( $0 \leq h, h' \leq k$ ), since the outcomes are identical on events  $E_{l+k'}^1$  for all  $0 \leq k' \leq k - 1$ , greedy iterated elimination will deliver  $k$  cells. By propositions 2, 3, and the hereditary property, the anonymous generalized median voter scheme is not  $(k - 1)$ -OD implementable.  $\square$

The above proposition demonstrates how  $k$ -OD extends the range of im-

plementable anonymous generalized median voter schemes as the level of sophistication in the application of STP increases.

## 5 Discussion and Conclusion

In the implementation literature, there has been a longstanding concern on the need for robustness of mechanisms in terms of actual performance. In environments with complete information, this need may arise from flaws in reasoning and mistakes in the specification of the agents' preferences, knowledge, or situation. In this regard, Moore (1992) highlights the value of simplicity for the mechanism itself with a parallel message for implementation theory that mechanisms which do not allow for mis-specifications in the agent's preference, knowledge, and situation may perform poorly in actual applications, especially if they place undue demands in terms of attention, cognition and strategic calculations.<sup>15</sup>

In an incomplete information setting, Bergemann and Morris (2005) investigate the question of robustness by relaxing the common knowledge assumption among players and designer by studying mechanism design on richer type spaces. This inspires Börgers and Li (2019) to propose a new simplicity criterion such that participants can deduce their optimal strategy using only their first-order beliefs about other players' preferences.<sup>16</sup>

Li's (2017) definition of OD, which has ushered in a novel direction of research on simplicity of mechanisms, has inspired the present paper along with Pycia and Troyan (2021) and Zhang and Levin (2021).

Relying on an exogenously given partition of the state space, Zhang and Levin (2017, 2021) provide a generalization of OD, named *partition obvious*

---

<sup>15</sup>In terms of iterative elimination of dominated strategies in a normal game, Glazer and Rubinstein (1996) offer a more formal sense of what may constitute simplicity by showing how a dynamic game can be viewed as a guide, and thus conclude that calculating the subgame perfect equilibrium outcome in a dynamic game is simpler. The sense of this need to address the issue of complexity is echoed in Maskin and Sjöström's (2002) recognition of the incidence of bounded rationality in developing mechanisms that are more forgiving of departures from full-blown "homo game theoreticus".

<sup>16</sup>De Clippel et al. (2019) offers a nonequilibrium model of robustness involving level-k reasoning.

*dominance*, in which the decision maker compares the worst outcome of one act (strategy) with the best outcome of another act (strategy) on each cell. When applied to dynamic mechanism design, the exogenously given partition is updated by removing possibilities in each cell of the partition that are ruled out by the emergence of new information. The major difference between partition obvious dominance and  $k$ -OD is that the relevant partition in our definition emerges endogenously from the agents' application of STP and that this partition may vary at different information sets in comparing different strategies. This being said, partition obvious dominance can be viewed as a special case of  $k$ -OD requiring an exogenously given "fixed" partition.

Pycia and Troyan (2021) offer a model of simplicity which delivers a sequence of solution concepts, also indexed by a natural number, capturing the number of forward-looking steps rather than planning for the entire future of a game.<sup>17</sup> We apply Savage's (1954) sure-thing principle to generalize OD with  $k$ -OD serving as a bridge between dominance and OD. In going beyond OD, Pycia and Troyan's limited foresight model is motivated by the proverb "you can cross the bridge when you come to it" also appearing in Savage (1954).

We develop in Section 3 a definition of a (randomized) gradual mechanism to serve as a revelation principle for  $k$ -OD implementation. This proof is straightforward once we have developed a framework capable of representing the design domain of dynamic mechanisms. A similar observation is stated in Myerson (1989) in his demonstration of the direct revelation principle for Bayesian implementation. Another shared feature with the proof of direct revelation principle is a pruning operation which removes actions never used by the agents. Li (2017) provides this pruning operation for dynamic mechanisms which has appeared in several follow up studies, such as Ashlagi and Gonczarowski (2018), Bade and Gonczarowski (2017), Mackenzie (2020), and Pycia and Troyan (2021).

Facilitated by gradual mechanism, we provide in Section 3 a necessary condition for  $k$ -OD implementability and demonstrate its usefulness, together

---

<sup>17</sup>Relatedly, Catonini and Xue (2021) study simplicity by applying a one-step foresight model.

with the hereditary property of  $k$ -OD implementability, in the application to anonymous generalized median voter schemes in Section 4. There has been several papers studying the OD implementability of specific social choice rules offering new insights arising from obviousness of the truth-telling strategy in dynamic implementation.<sup>18</sup> Given the sense in which  $k$ -OD may capture a bound on the required level of sophistication in applying STP for contingent reasoning, it may hold promise for further research to explore its applicability in a range of mechanisms such as matching, auctions, and voting.

The ability of dynamic mechanisms to make dominant strategy obvious stems from the dynamic information flow which brings other benefits, such as simplicity (Bó and Hakimov, 2021; Pycia and Troyan, 2021), credibility (Akbarpour and Li, 2020), and privacy preservation (Mackenzie and Zhou, 2020; Haupt and Hitzig, 2022). Many of the dynamic mechanisms proposed in the literature, including the round table mechanisms (Mackenzie, 2020), millipedes games (Pycia and Troyan, 2021), pick-an-object mechanisms (Bó and Hakimov, 2021), and menu mechanisms (Mackenzie and Zhou, 2020), are in effect special cases of gradual mechanisms.<sup>19</sup> In subsequent research, besides serving as revelation principle for  $k$ -OD implementations, the gradual mechanism defined in this paper could serve as a canonical class of dynamic mechanisms to accommodate a rich range of properties associated with the underlying information flow (Chew and Wang, 2022). Moreover, it would be valuable to exploit, following Golowich and Li (2021), the class of gradual mechanism towards a computationally efficient way to check whether a given social choice rule is  $k$ -OD implementable.

As an intermediate solution concept between dominance and OD, there is value in testing the performance of  $k$ -OD implementation. A useful test may involve contingent versus non-contingent framing (Esponda and Vespa, 2021). Building on the reported effectiveness in the provision of advice about the

---

<sup>18</sup>In addition to the ones mentioned in the introduction, there are others such as Bade and Gonczarowski (2017), Bade (2019), Ferraioli and Ventre (2021), Tsakas (2019), and Troyan and Morrill (2020)

<sup>19</sup>In pick-an-object mechanisms and menu mechanisms, the administrator always collects information about agents' top preference in a menu.

value of truthful reporting (Masuda et al., 2022), we can enrich the dynamic mechanism under consideration towards prompting subjects into engaging in the type of contingent thinking implicit in STP and observe their effectiveness in bringing about greater incidence of truthful revelation.<sup>20</sup>

## References

- Akbarpour, M. and S. Li (2020). Credible Auctions: A Trilemma. *Econometrica* 88, 425–467.
- Allais, M. (1953). Le Comportement de l’Homme Rationnel devant le Risque: Critique des Postulats et Axiomes de l’Ecole Americaine. *Econometrica* 21, 503–546.
- Arribillaga, R. P., J. Massó, and A. Neme (2020). On Obvious Strategy-proofness and Single-peakedness. *Journal of Economic Theory* 186, 104992.
- Ashlagi, I. and Y. A. Gonczarowski (2018). Stable Matching Mechanisms Are Not Obviously Strategy-proof. *Journal of Economic Theory* 177, 405–425.
- Bade, S. (2019). Matching with Single-peaked Preferences. *Journal of Economic Theory* 180, 81–99.
- Bade, S. and Y. A. Gonczarowski (2017). Gibbard-Satterthwaite Success Stories and Obvious Strategyproofness. Working Paper.
- Barberà, S., F. Gul, and E. Stacchetti (1993). Generalized Median Voter Schemes and Committees. *Journal of Economic Theory* 61, 262–289.
- Battigalli, P., P. Leonetti, and F. Maccheroni (2020). Behavioral Equivalence of Extensive Game Structures. *Games and Economic Behavior* 121, 533–547.
- Bergemann, D. and S. Morris (2005). Robust Mechanism Design. *Econometrica* 73, 1771–1813.
- Bó, I. and R. Hakimov (2021). Pick-an-object Mechanisms. Working Paper.

---

<sup>20</sup>In Chen et al.’s (2021) study of the robustness of their simultaneous report mechanisms, each subject completes a screening quiz containing contingencies which may arise in the experiment. Like Masuda et al. (2022), they also provide advice to subjects about reports consistent with their true values being in their material interests.

- Börger, T. and J. Li (2019). Strategically Simple Mechanisms. *Econometrica* 87, 2003–2035.
- Cason, T. N. and C. R. Plott (2014). Misconceptions and Game Form Recognition: Challenges to Theories of Revealed Preference and Framing. *Journal of Political Economy* 122, 1235–1270.
- Catonini, E. and J. Xue (2021). Local Dominance. Working Paper.
- Charness, G. and D. Levin (2009). The Origin of the Winner’s Curse: A Laboratory Study. *American Economic Journal: Microeconomics* 1, 207–236.
- Chen, Y. and T. Sönmez (2006). School Choice: An Experimental Study. *Journal of Economic Theory* 127, 202–231.
- Chen, Y.-C., R. Holden, T. Kunimoto, Y. Sun, and T. Wilkening (2021). Getting Dynamic Implementation to Work. Working Paper.
- Chew, S. H. and W. Wang (2022). Information Design of Dynamic Mechanisms. Working Paper.
- De Clippel, G., R. Saran, and R. Serrano (2019). Level- $k$  Mechanism Design. *Review of Economic Studies* 86, 1207–1227.
- Dreyfuss, B., O. Heffetz, and M. Rabin (2022). Expectations-Based Loss Aversion May Help Explain Seemingly Dominated Choices in Strategy-Proof Mechanisms. *American Economic Journal: Microeconomics*. Forthcoming.
- Ellsberg, D. (1961). Risk, Ambiguity, and the Savage Axioms. *Quarterly Journal of Economics* 75, 643–669.
- Esponda, I. and E. Vespa (2014). Hypothetical Thinking and Information Extraction in the Laboratory. *American Economic Journal: Microeconomics* 6, 180–202.
- Esponda, I. and E. Vespa (2021). Contingent Preferences and the Sure-Thing Principle: Revisiting Classic Anomalies in the Laboratory. Working Paper.
- Ferraioli, D. and C. Ventre (2021). Approximation Guarantee of OSP Mechanisms: The Case of Machine Scheduling and Facility Location. *Algorithmica* 83, 695–725.



- Gibbard, A. (1973). Manipulation of Voting Schemes: A General Result. *Econometrica* 41, 587–601.
- Glazer, J. and A. Rubinstein (1996). An Extensive Game as a Guide for Solving a Normal Game. *Journal of Economic Theory* 70, 32–42.
- Golowich, L. and S. Li (2021). On the Computational Properties of Obviously Strategy-Proof Mechanisms. Working Paper.
- Hakimov, R. and D. Kübler (2019). Experiments on Matching Markets: A Survey. Working Paper.
- Haupt, A. and Z. Hitzig (2022). Contextually Private Implementation. Working Paper.
- Hurwicz, L. (1960). Optimality and Informational Efficiency in Resource Allocation Processes. In K. J. Arrow, S. Karlin, and P. Suppes (Eds.), *Mathematical methods in the social sciences*, pp. 27–46. Stanford, CA: Stanford University Press.
- Hurwicz, L. (1972). On Informationally Decentralized Systems. In C. B. McGuire and R. Rad (Eds.), *Decision and organization: A volume in Honor of J. Marschak*, pp. 297–336. Amsterdam: North-Holland.
- Kagel, J. H., R. M. Harstad, and D. Levin (1987). Information Impact and Allocation Rules in Auctions with Affiliated Private Values: A Laboratory Study. *Econometrica* 55, 1275–1304.
- Li, S. (2017). Obviously Strategy-Proof Mechanisms. *American Economic Review* 107, 3257–3287.
- Mackenzie, A. (2020). A Revelation Principle for Obviously Strategy-proof Implementation. *Games and Economic Behavior* 124, 512–533.
- Mackenzie, A. and Y. Zhou (2020). Menu Mechanisms. Working Paper.
- Mandal, P. and S. Roy (2022). Obviously Strategy-Proof Implementation of Assignment Rules: a New Characterization. *International Economic Review* 63, 261–290.
- Martínez-Marquina, A., M. Niederle, and E. Vespa (2019). Failures in Contingent Reasoning: The Role of Uncertainty. *American Economic Review* 109, 3437–3474.

- Maskin, E. S. and T. Sjöström (2002). Implementation Theory. In K. J. Arrow, A. K. Sen, and K. Suzumura (Eds.), *Handbook of Social Choice and Welfare*, Volume 1, pp. 237–288. Amsterdam: Elsevier.
- Masuda, T., R. Mikami, T. Sakai, S. Serizawa, and T. Wakayama (2022). The Net Effect of Advice on Strategy-proof Mechanisms: An Experiment for the Vickrey Auction. *Experimental Economics* 25, 902–941.
- Moore, J. (1992). Implementation, Contracts, and Renegotiation in Environments with Complete Information. In J.-J. Laffont (Ed.), *Advances in Economic Theory*, Volume 1, pp. 182–282. Cambridge, MA: Cambridge University Press.
- Moulin, H. (1980). On Strategy-proofness and Single Peakedness. *Public Choice* 35, 437–455.
- Myerson, R. B. (1989). Mechanism Design. In J. Eatwell, M. Milgate, and P. Newman (Eds.), *Allocation, Information and Markets*, pp. 191–206. New York, W.W. Norton: Springer.
- Ngangoué, M. K. and G. Weizsäcker (2021). Learning from Unrealized versus Realized Prices. *American Economic Journal: Microeconomics* 13, 174–201.
- Osborne, M. J. and A. Rubinstein (1994). *A Course in Game Theory*. Cambridge: MIT Press.
- Pycia, M. and P. Troyan (2021). A Theory of Simplicity in Games and Mechanism Design. Working Paper.
- Satterthwaite, M. A. (1975). Strategy-proofness and Arrow’s Conditions: Existence and Correspondence Theorems for Voting Procedures and Social Welfare Functions. *Journal of Economic Theory* 10, 187–217.
- Savage, L. J. (1954). *The Foundations of Statistics*. New York: John Wiley & Sons.
- Thomas, C. (2020). Classification of Priorities such that Deferred Acceptance is Obviously Strategyproof. Working Paper.
- Troyan, P. and T. Morrill (2020). Obvious Manipulations. *Journal of Economic Theory* 185, 104970.
- Tsakas, E. (2019). Obvious Belief Elicitation. *Games and Economic Behavior* 118, 374–381.

Zhang, L. and D. Levin (2017). Bounded Rationality and Robust Mechanism Design: An Axiomatic Approach. *American Economic Review* 107, 235–239.

Zhang, L. and D. Levin (2021). Partition Obvious Preference and Mechanism Design: Theory and Experiment. Working Paper.

## Appendix

In this appendix, we provide the formal definitions of (randomized) dynamic game forms and the two randomization operations—de-randomization and pre-randomization.

- *Players.* In addition to agents in  $N$ , there is one additional player 0 representing the mechanism administrator who brings randomization to the mechanism.
- *Actions.* For each agent  $i \in N$ ,  $A_i$  is a nonempty action set. Denote the set of action profiles by

$$A = \{\emptyset\} \cup \bigcup_{\emptyset \subsetneq M \subseteq N_0} \prod_{i \in M} A_i$$

in which the empty action profile  $\emptyset$  is introduced only to simplify exposition.

- Pick an action profile  $a \in \prod_{i \in M} A_i \subseteq A$ . Suppose  $i \in M$ . Let  $a_i$  denote the action of agent  $i$  in  $a$  and  $a_{-i}$  the action profile of other agents in  $a$ . Otherwise, suppose  $i \notin M$ . Then  $a_i = \emptyset$  and  $a_{-i} = a$ .
- For each  $T > 0$ , let  $A^T$  denote the collection of histories of length  $T$  with a generic history being denoted by  $h = (h^{(1)}, \dots, h^{(T)})$  in which  $h^{(T)}$  is also referred to as  $h^{(-1)}$ . Let  $A^0 = \{\emptyset\}$  be the singleton set of the empty history.
- Let  $A^{<\mathbb{N}} = \bigcup_{T \in \mathbb{N}} A^T$  denote the collection of all histories of finite length.

- An action profile  $a$  and the corresponding sequence  $(a)$  containing only  $a$  are used interchangeably. The empty history  $\emptyset$  and any sequence of empty action profiles are used interchangeably. For history  $h = (h^{(1)}, \dots, h^{(T)})$  consisting of  $S$  non-empty action profiles and some empty action profiles, let  $\underline{h} = (\underline{h}^{(1)}, \dots, \underline{h}^{(S)})$  be the corresponding history without empty action profiles, i.e.,  $s$ -th non-empty action profiles in  $h$  equals  $\underline{h}^{(s)}$  for each  $1 \leq s \leq S$ . We will use  $h$  and  $\underline{h}$  interchangeably.
  - There is a precedence relation  $\preceq$  on  $A^{<\mathbb{N}}$ , i.e.,  $\underline{h} \preceq h$  (reads  $\underline{h}$  is a predecessor of  $h$  or  $h$  is a successor of  $\underline{h}$ ) if  $\underline{h} \in A^S$  and  $h \in A^T$  such that  $S = 0$  or that  $0 < S \leq T$  and  $\underline{h}^{(s)} = h^{(s)}$  for any  $1 \leq s \leq S$ .
  - Let  $h_1, \dots, h_m \in A^{<\mathbb{N}}$  be  $m$  sequences of action profiles, define  $(h_1, \dots, h_m) \in A^{<\mathbb{N}}$  by concatenation.
  - If  $\underline{h} \preceq h$  in  $A^{<\mathbb{N}}$  such that  $\underline{h} \in A^T$ ,  $h \in A^{T+1}$ , and  $h^{(-1)} \neq \emptyset$ , we say  $\underline{h}$  is an immediate predecessor of  $h$  or  $h$  is an immediate successor of  $\underline{h}$ . Note that a nonempty history  $h$  has a unique immediate predecessor.
  - Let  $h = (h^{(1)}, \dots, h^{(T)})$ , define  $h_i = (h_i^{(1)}, \dots, h_i^{(T)})$  and  $h_{-i} = (h_{-i}^{(1)}, \dots, h_{-i}^{(T)})$ .
- *Histories.* The set of histories  $\overline{H}$  is modeled by a tree of finite length in  $A^{<\mathbb{N}}$ , i.e., a subset of  $A^{<\mathbb{N}}$  such that (i) there exists  $\mathcal{T}$  such that  $\overline{H} \subseteq \bigcup_{T=0}^{\mathcal{T}} A^T$ , (ii)  $\emptyset \in \overline{H}$ , and (iii) for any  $h \in \overline{H}$  such that  $h \neq \emptyset$ , the immediate predecessor of  $h$  is in  $\overline{H}$ . We make use of the following assumptions and notations.
    - Denote the set of terminal histories by  $Z$  and non-terminal histories by  $H$ .
    - For any non-terminal history  $h \in H$ , denote by  $\sigma(h)$  the collection of its immediate successors.
    - $\overline{H}$  satisfies the following property: there exists an active-player correspondence  $\mathbb{P} : H \rightarrow N_0$  such that for any non-terminal history

$h \in H$  and any  $a \in A$  satisfying  $(h, a) \in \sigma(h)$ , it is the case that  $a \in \prod_{i \in \mathbb{P}(h)} A_i$ . Therefore,  $\mathbb{P}(\cdot)$  depicts the players that are simultaneously active at a particular non-terminal history.

- Let  $H_i = \{h \in H : i \in \mathbb{P}(h)\}$  represent the collection of histories on which player  $i$  is active. For each  $h \in H$  and each  $i \in \mathbb{P}(h)$ , define  $A_i(h) = \{a_i \in A_i : (h, a) \in \overline{H} \text{ for some } a \in A\}$  as the collection of available actions for player  $i$  at history  $h$ .
- $\overline{H}$  satisfies the following property: for any  $h \in H$  and any  $a \in \prod_{i \in \mathbb{P}(h)} A_i(h)$ , we have  $(h, a) \in \overline{H}$ . The two assumptions here depict decentralized decision making.

- *Information Structure.* For each  $i \in N_0$ ,  $\mathbf{H}_i$  is a partition of  $H_i$  whose elements are information sets of player  $i$ , with a generic information set being denoted by  $\mathbf{h}_i$ . We introduce further assumptions, notations, and observations below.

- Assume that for any  $\mathbf{h}_i \in \mathbf{H}_i$  and any  $h, \tilde{h} \in \mathbf{h}_i$ , we have  $A_i(h) = A_i(\tilde{h})$ . Then  $A(\mathbf{h}_i) = A_i(h)$  for which  $h \in \mathbf{h}_i$  is well defined.
- Assume that the game form  $G$  has perfect recall. Formally, for any  $h, \tilde{h} \in \mathbf{h}_i$ , we have  $h_i = \tilde{h}_i$  and for any  $\underline{h} \preceq h$  with  $\underline{h} \in H_i$ , there exists  $\tilde{\underline{h}} \preceq \tilde{h}$  such that  $\underline{h}$  and  $\tilde{\underline{h}}$  are in the same information set of agent  $i$ .
- Given perfect recall, an ordering on  $\mathbf{H}_i$  can be defined,  $\underline{\mathbf{h}}_i \preceq \mathbf{h}_i$ , if there exist  $\underline{h} \in \underline{\mathbf{h}}_i$  and  $h \in \mathbf{h}_i$  such that  $\underline{h} \preceq h$ .
- Administrator has perfect information, i.e.,  $\mathbf{H}_0$  consists of singleton information sets.

- *Outcomes.*  $\mathcal{X} : Z \rightarrow X$  assigns each terminal history a public outcome.

Given a game form  $G$  and any of its non-terminal history  $h \in H$ , we can define the subgame form of  $G$  starting from  $h$  by restricting various components of  $G$  to its truncated successors in  $\overline{H}$ , i.e., to  $\{h' \in A^{<\mathbb{N}} : (h, h') \in \overline{H}\}$ .

An interim strategy  $s_i : H_i \rightarrow A_i$  of agent  $i \in N$  specifies an available action  $s_i(h) \in A_i(h)$  for each history  $h \in H_i$  such that  $s_i(h) = s_i(\tilde{h})$  if  $h, \tilde{h}$  belong to the same information set. Therefore,  $s_i : \mathbf{H}_i \rightarrow A_i$  is well defined. We use  $S_i$  to denote the set of interim strategies for agent  $i$  and use  $s \in S$  and  $s_{-i} \in S_{-i}$  to denote the profile of interim strategies for all agents in  $N$  and that for those other than agent  $i$  respectively.

*Administrator's randomization* is modeled by a probability space  $(\Omega, \mathcal{A}, \mu)$  providing a randomization device for the administrator to adopt a mixed strategy  $m$ . In particular, upon receiving  $\omega \in \Omega$ , the administrator adopts the interim strategy  $m(\omega) : H_0 \rightarrow A_0$  which specifies an action at every history the administrator is active. The two randomization operations are formally defined as follows.

**Definition 9.** Let  $G = (\overline{H}, \{A_i, \mathbf{H}_i\}_{i \in N_0}, \mathcal{X}, m, \{\Omega, \mathcal{A}, \mu\})$ , then for each  $\omega \in \Omega$ , the  $\omega$ -derandomization of  $G$  is  $G^\omega = (\overline{H}^\omega, \{A_i^\omega, \mathbf{H}_i^\omega\}_{i \in N}, \mathcal{X}^\omega)$  such that:

1. The set of feasible actions for each agent in  $G^\omega$  is the same with that in  $G$ , i.e.,  $A_i^\omega = A_i$ .
2.  $\overline{H}^\omega$  is the collection of histories that are consistent with administrator acting according to  $m(\omega)$  but with administrator's actions removed, i.e.,  $\overline{H}^\omega = \{h_{-0} : h \in H(\omega)\}$ .
3.  $\{\mathbf{H}_i^\omega\}_{i \in N}$  and  $\mathcal{X}^\omega$  are restrictions of their counterparts in  $G$  on the domain  $\overline{H}^\omega$ .

Note that both  $\{\mathbf{H}_i^\omega\}_{i \in N}$  and  $\mathcal{X}^\omega$  are well defined since for each agent  $i$  and each history  $h^\omega \in H_i^\omega$ , there exists a unique  $h \in H_i \cap H(\omega)$  such that  $h_{-0} = h^\omega$ .

**Definition 10.** Let  $G = (\overline{H}, \{A_i, \mathbf{H}_i\}_{i \in N_0}, \mathcal{X}, m, \{\Omega, \mathcal{A}, \mu\})$ , then the pre-randomization of  $G$  is  $G^* = (\overline{H}^*, \{A_i^*, \mathbf{H}_i^*\}_{i \in N_0}, \mathcal{X}^*, m^*, \{\Omega, \mathcal{A}, \mu\})$  such that:

1. In  $G^*$ , the administrator is active only at the initial history whose available actions are the realizations of the same randomization device as in  $G$ , i.e.,  $H_0^* = \{\emptyset\}$ ,  $A_0^* = \Omega$ , and  $m^*(\omega)(\emptyset) = \omega$ .

2. For each  $\omega \in \Omega$ , the subgame in  $G^*$  following history  $\omega \in \overline{H}^*$  is the  $\omega$ -derandomization  $G^\omega$  of the original  $G$ .
3. The action of the administrator is known by all agent, i.e., for all agent  $i$  and all information set  $\mathbf{h}_i$ , there exists  $\omega$  such that  $\mathbf{h}_i \subseteq H^*(\omega)$ .

*Proof of Theorem 1.* Consider a randomized dynamic game form  $G$ . Let  $G^*$  be its pre-randomization. Observe that each history  $h^* \in \overline{H}^*$  corresponds to a unique  $m(h^*) = h \in \overline{H}$  such that  $h_{-0}^* = h_{-0}$ . Suppose agent  $i$  is active at  $h^*$ , then she will be active at  $h = m(h^*)$  as well and that the sets of available actions for agent  $i$  at  $h$  and at  $h^*$  are the same. Based on this observation, there is a surjective mapping from  $S_i$  to  $S_i^*$  (the set of interim strategies in  $G$  and  $G^*$  respectively) for each agent  $i$ , in which  $s_i \in S_i$  of the original game form  $G$  is mapped to an interim strategy  $s_i^*$  such that  $s_i^*(h^*) = s_i(h)$  for any  $h^* \in H_i^*$  and  $h \in H_i$  such that  $h_{-0}^* = h_{-0}$ . Notice that  $\mathcal{X}^*(s^*, \omega) = \mathcal{X}(s, \omega)$  when  $s$  is mapped to  $s^*$  agent by agent.

Let  $\mathbb{S}$  be a strategy profile in the incomplete information game  $(G, \Theta)$ , then we can define a strategy profile  $\mathbb{S}^*$  for  $(G^*, \Theta)$  by letting  $\mathbb{S}_i^*(\theta_i)(h^*) = \mathbb{S}_i(\theta_i)(h)$  for each agent  $i$ , each private type  $\theta_i$ , and each non-terminal history  $h^*$  such that  $h_{-0}^* = h_{-0}$ . In short,  $\mathbb{S}^*$  defines a strategy profile in  $G^*$  that acts in the same way as  $\mathbb{S}$  in  $G$  at corresponding histories.

To demonstrate that  $\mathbb{S}$  being a  $k$ -OD strategy profile in  $G$  implies that  $\mathbb{S}^*$  being a  $k$ -OD strategy profile in  $G^*$ , note that an information set  $\mathbf{h}_i^*$  of agent  $i$  in  $G^*$  corresponds to a unique information set  $\mathbf{h}_i$  in  $G$  defined by the correspondence relation between histories mentioned in the previous paragraph. Suppose two strategies  $\mathbb{S}^*(\theta_i)(\mathbf{h}_i^*)$  and  $s_i^*$  deviate on  $\mathbf{h}_i^*$ , then  $\mathbb{S}(\theta_i)(\mathbf{h}_i)$  deviate from any  $s_i \in S_i$  that maps to  $s_i^*$  on  $\mathbf{h}_i$ . Now, notice that the space of uncertainties  $\mathbf{h}_i^* \times S_{-i}^* \times \Omega$  in  $G^*$  shrinks in each of the three sources compared with that in  $G$ .

Next, do pruning on each subgame form following  $(\omega)$  in  $G^*$  and we will have randomized gradual mechanism  $k$ -OD implementing the same stochastic social choice function.  $\square$