

Preserving Condorcet Winners under Strategic Manipulation

Sirin Botan

Institute for Logic, Language and Computation
University of Amsterdam

Ulle Endriss

Institute for Logic, Language and Computation
University of Amsterdam

ABSTRACT

Condorcet extensions have long held a prominent place in social choice theory. A Condorcet extension will return the Condorcet winner as the unique winner whenever such an alternative exists. However, the definition of a Condorcet extension does not take into account possible manipulation by the voters. A profile where all agents vote truthfully may have a Condorcet winner, but this alternative may not end up in the set of winners if agents are acting strategically. Focusing on the class of tournament solutions, we show that many natural social choice functions in this class, such as the well-known Copeland and Slater rules, cannot guarantee the preservation of Condorcet winners when agents behave strategically. Our main result in this respect is an impossibility theorem that establishes that no tournament solution satisfying a very weak decisiveness requirement can provide such a guarantee. On the bright side, we identify several indecisive but otherwise attractive tournament solutions that do guarantee the preservation of Condorcet winners under strategic manipulation.

KEYWORDS

Voting; Condorcet Consistency; Strategic Manipulation

1 INTRODUCTION

By a seminal result in social choice theory, we know that every non-trivial resolute social choice function—or voting rule—is susceptible to strategic manipulation by voters [21, 33].

While there are instances of manipulation that arguably are not a cause for great concern—voting for your second choice rather than your top choice to avoid wasting your vote, for example—there are other instances that may seriously undermine the confidence we have in the end result of a vote. Our proposal in this paper is to define a more fine-grained strategyproofness axiom that dictates a specific kind of “undesirable” manipulation should not occur. We do not examine whether strategyproofness can be satisfied alongside a set of other desirable axioms, but whether the fact that strategyproofness fails can affect the “strength” of the other axiom(s). Our focus is on Condorcet extensions, and how failure of strategyproofness affects Condorcet consistency.

Suppose for a set of agents, that the preference profile resulting from everyone submitting their true preferences has a Condorcet winner, and suppose further that we use a Condorcet-consistent voting rule to aggregate their preferences. If everyone reports their true preferences, any Condorcet extension will return the Condorcet winner as the sole winner. However, not all such functions are strategyproof, and it may happen that the agents’ truthful preferences differ from their reported preferences. If an agent manipulates under such a rule—meaning she submits a preference other than her true one—it is possible that this will result in a reported profile without a Condorcet winner, despite the fact that the truthful

profile has one. In other words, some Condorcet extensions may fail to return the Condorcet winner of the truthful profile as the outcome. We examine exactly when the lack of strategyproofness affects whether we can trust that a Condorcet extension will give us all “true” Condorcet winners, even under the assumption that agents will vote strategically.

Related work. There are various methods found in the literature for circumventing the impossibility result due to Gibbard [21] and Satterthwaite [33]. One strategy is to consider a restricted domain for the social choice function, where strategyproofness is guaranteed. Among these domain restrictions, the best-known is the single-peaked domain of Black [6].¹ A more recent approach is to argue that the computational hardness of computing a successful manipulation strategy may be a barrier to manipulation [5, 14]. Similar results have been obtained by considering voters’ ignorance of others’ preference as an informational barrier [15, 28, 30].

Finally, while the Gibbard-Satterthwaite Theorem deals a blow to resolute social choice functions, irresolute functions come away relatively unscathed. While there are similar impossibility results for irresolute rules [3, 17, 20, 23, 31], these results differ in how they define manipulability as they by necessity must make assumptions about agents’ preferences over sets of alternatives. Duggan and Schwartz [17], for example, work with optimistic and pessimistic agents, while Gärdenfors [20] defines manipulability relative to his eponymous preference extension.

Although impossibilities abound, shifting focus away from resoluteness has also led to positive results regarding the strategyproofness of social choice functions. Gärdenfors [20] identifies two such strategyproof functions for the Gärdenfors extension—the one that returns all alternatives ranked first by at least one agent, and the one returning the Condorcet winner when one exists and the whole set of alternatives when one does not. Brandt [8] characterises the pairwise social choice functions that are strategyproof under the Kelly preference extension, among them the bipartisan set and the minimal covering set. Brandt and Brill [10] build on these results and find sufficient conditions for strategyproofness under the stronger Fishburn and Gärdenfors extensions as well, thereby identifying further social choice functions that are strategyproof for each of the three extensions.

More fine-grained approaches to strategyproofness have also been explored in the social choice literature. Sato [32] considers adjacency-strategyproofness, which requires agents to make large changes to their reported preference in order to successfully manipulate. For social welfare functions, Bossert and Sprumont [7] have obtained positive results for a form of strategyproofness that considers only manipulations resulting in outcomes located *between* an agent’s true preference order and the outcome when no

¹For a more thorough treatment of single-peakedness and many other domain restrictions, see Gaertner [19].

manipulation occurs. They find several social welfare functions, including the Kemeny and the Slater rules, that are strategyproof in this sense. More recently, Kruger and Terzopoulou [25] have examined the manipulability of scoring rules for agents with incomplete preferences. They distinguish between agents adding alternatives to their (incomplete) ranking, deleting alternatives, or swapping positions of alternatives in their preference order. They identify scoring rules that are strategyproof with respect to each of these forms of manipulation.

Requiring the preservation of Condorcet winners under strategic manipulation can also be seen as a stability axiom that tells us when—or under what social choice function—profiles with Condorcet winners are *stable*. Similar notions of stability exist in game theory, as well as game-theoretic examinations of voting. Trembling hand equilibria [34], for example, are Nash equilibria that are stable with regard to small “slips of the hand”. Trembling hand equilibria have also been studied in voting—Obraztsova et al. [27], for example, study trembling hand equilibria of Plurality voting—where such an equilibrium corresponds to a profile where no voter has an incentive to deviate, even under the assumption that others’ hands may tremble and slip.

Contribution. To distinguish between social choice functions that do not incentivise manipulation in profiles with a Condorcet winner and those that do, we introduce the notion of a *robust Condorcet extension*. We then show that no irresolute tournament solution that is *weakly resolute*—in the sense of returning a single winning alternative in at least one profile that does not have a Condorcet winner—can possibly be such a robust Condorcet extension. This class of weakly resolute rules includes the well-known Copeland [16] and Slater [35] social choice functions. Finally, we identify several attractive social choice functions that are robust Condorcet extensions (and thus fail weak resoluteness). This includes, in particular, the minimal extending set [9] and all of its coarsenings.

Paper outline. The remainder of the paper is organised as follows. We introduce the framework and relevant literature in Section 2. In Section 3 we present an impossibility result for weakly resolute rules. In Section 4 we present a number of sufficient conditions for a Condorcet extension to be robust. We conclude in Section 5.

2 THE MODEL

In this section we introduce the framework and notation we will be using throughout the paper. Much of the material up to Section 2.6 is familiar from social choice theory [1, 12]. In Section 2.6, we introduce the novel notion of a robust Condorcet extension, the central concept of this paper.

2.1 Preference Profiles

Let A be a finite set of *alternatives*, and $N = \{1, \dots, n\}$ a finite set of *agents*. A *preference profile* $P = (>_1^P, \dots, >_n^P)$ is a vector of strict linear orders over A , where $>_i^P$ is the *preference relation* of agent i in the profile P . $\mathcal{L}(A)$ denotes the set of all linear orders over A , and $\mathcal{L}(A)^n$ denotes the set of all profiles for n agents.

For a profile P , \geq^P (with asymmetric part $>^P$) is the *majority relation* for P and is defined such that $a \geq^P a'$ if and only if $|\{i \in N \mid a >_i^P a'\}| \geq |\{i \in N \mid a' >_i^P a\}|$, for all $a, a' \in A$.

An alternative $a \in A$ is a *Condorcet winner* in profile P if it defeats every other alternative in a pairwise majority contest, meaning $a >^P a'$ for all $a' \in A \setminus \{a\}$. We define $\mathcal{D}_{\text{Condorcet}}$ as the set of profiles with a Condorcet winner.

For two profiles P and P' , and agent $i \in N$, we write $P =_{-i} P'$, and say they are i -variants, if $>_j^P = >_j^{P'}$ for all $j \in N \setminus \{i\}$.

We say a relation \geq over A is *complete* if for all $a, b \in A$ it is the case that $a \geq b$ or $b \geq a$. A relation $>$ is *connex* if $a > b$ or $b > a$ for all distinct $a, b \in A$. Note that the majority relation of any profile is complete, and any individual preference relation is connex.

2.2 Tournaments

A *tournament* T is a pair $(S, >^T)$, where S is a set of nodes and $>^T$ is an asymmetric and connex relation over S , which we call the *dominance relation* of the tournament. For a set S , we denote by $\mathcal{T}(S)$ all tournaments on S .

For a tournament $T = (S, >^T)$, we say an alternative $a \in S$ *dominates* $a' \in S$ in the tournament T if $a >^T a'$. The *dominion* of a in T is defined as $D_T(a) = \{a' \in S \mid a >^T a'\}$, the set of alternatives it dominates. The *dominators* of a in T is defined as $\bar{D}_T(a) = \{a' \in S \mid a' >^T a\}$, the set of alternatives that dominate it. For $S' \subseteq S$, we define the restriction $>_{S'}^T = \{(a, a') \in S' \times S' \mid a >^T a'\}$, which is $>^T$ restricted to the set S' . A *subtournament* of $T = (S, >^T)$ is a tournament $(S', >_{S'}^T)$ where $S' \subseteq S$. Thus, a subtournament of T is a subset of the nodes in T , and the edges between those nodes.

We say $\pi : S \rightarrow S'$ is an isomorphism between two tournaments $T = (S, >^T)$ and $T' = (S', >^{T'})$ if π is a bijection, and $a >^T a' \Leftrightarrow \pi(a) >^{T'} \pi(a')$ for all $(a, a') \in S \times S$.

We say a profile $P \in \mathcal{L}(A)^n$ *induces* tournament $T = (A, >^T)$ if $>^P = >^T$. So a profile induces a tournament if they range over the same alternatives, and the strict part of the majority relation is exactly the dominance relation of the tournament. Note that if a profile induces a tournament, the strict component of the majority relation of that profile must be connex. As tournaments do not speak about agents, we cannot directly talk about two tournaments being i -variants for some agent $i \in N$. Instead, we say two tournaments $T = (A, >^T)$ and $T' = (A, >^{T'})$ are *single-agent variants*, and write $T =_{-1} T'$, if there exist a set of agents N and profiles $P, P' \in \mathcal{L}(A)^n$ such that $P =_{-i} P'$ for some agent $i \in N$, and the profiles P and P' induce the tournaments T and T' , respectively.

We say an element $a \in S$ is the *Condorcet winner* of the tournament $T = (S, >^T)$ if $\bar{D}_T(a) = \emptyset$. This corresponds to the notion of a Condorcet winner of a profile; if a tournament has a Condorcet winner, that alternative will be the Condorcet winner in any profile that induces this tournament. We write $\mathcal{T}_{\text{Condorcet}}$ to mean the set of profiles that have a Condorcet winner.

2.3 Social Choice Functions

An irresolute *social choice function* (SCF) is a mapping from profiles to nonempty subsets of alternatives:

$$f : \mathcal{L}(A)^n \rightarrow 2^A \setminus \{\emptyset\}$$

To avoid having to break majority ties, we define social choice functions for odd n . Note, however, that while our functions are only defined for odd n , we do not require that the number of agents

is odd in all profiles. A SCF f is a *Condorcet extension*, or is *Condorcet-consistent*, if it returns (only) the Condorcet winner whenever one exists.

For irresolute social choice functions, the size of the set of winning alternatives is an important consideration. All things being equal, it is preferable that the SCF does not outsource the decision-making to a tie-breaking mechanism, but rather does most of the work of selecting a winner itself. More simply put, we would rather a SCF return small sets than large ones. As an example of a rule that returns quite large sets, take the rule that returns the Condorcet winner if one exists, and returns the whole set of alternatives otherwise. While this is clearly a Condorcet extension, it is a very *indecisive* rule, as it often results in many ties in the outcome.

A SCF is *resolute* if it always returns a singleton. In order to quantify how decisive an irresolute rule is, we define a weaker notion of resoluteness. We say f is *weakly resolute* if there exists a profile $P \in \mathcal{L}(A)^n \setminus \mathcal{D}_{\text{Condorcet}}$ for which $|f(P)| = 1$. So, a rule is weakly resolute if it *sometimes* returns a singleton for a profile without a Condorcet winner.

For many social choice functions, we can directly compare how decisive they are relative to each other. A SCF f is a *refinement* of f' if for all profiles $P \in \mathcal{L}(A)^n$ it is the case that $f(P) \subseteq f'(P)$, meaning f always returns a subset of f' . If f is a refinement of f' , we say f' is a *coarsening* of f . If a rule is a refinement of another, it is clearly the more decisive of the two.

2.4 Tournament Solutions

A *tournament solution* is a mapping from tournaments to sets of alternatives, that does not distinguish between isomorphic tournaments:

$$F : \mathcal{T}(S) \rightarrow 2^S \setminus \{\emptyset\}$$

So $F(T') = \{\pi(a) \mid a \in F(T)\}$ if π is an isomorphism between T and T' . For ease of reading, we will sometimes write $F(>^T)$ to mean $F(S, >^T)$ when S is clear from context.

A social choice function f is *equivalent* to a tournament solution F if $f(P) = F(A, >^P)$ for all $P \in \mathcal{L}(A)^n$. Note that the majority relation of this profile P must be a strict order, as the SCF f is defined for odd n only. In a slight muddling of terminology, we will refer to social choice functions that are equivalent to tournament solutions as *tournament-solution SCFs*.

Tournament solutions roughly correspond to Fishburn's C1 functions [18], which require only the information in the majority graph to determine the winners. More precisely, tournament-solution SCFs correspond to *neutral* C1 functions,² as tournament solutions do not distinguish between isomorphic tournaments, and therefore do not favour any alternatives over others.

2.5 Extending Preferences

Because the rules we examine are irresolute—meaning they do not always return a single winner—we need to specify how agent preferences over alternatives are extended to preferences over *sets of alternatives*.

²A social choice function satisfies *Neutrality* if for any profile P and any permutation $\pi : A \rightarrow A$: it is the case that $f(\pi(P)) = \pi(f(P))$.

A preference extension e maps any given preference relation $>$ over alternatives in A , to a relation \geq^e (with strict part $>^e$) over sets of alternatives. We define two requirements for any preference extension e :

- $x > y$ implies $\{x\} >^e \{y\}$
- $X >^e Y$ implies that there exist some $x \in X$ and $y \in Y$ such that $x > y$ and $\{x, y\} \not\subseteq X \cap Y$

The first requirement simply dictates that e stays faithful to the agent's preferences when comparing singleton sets. The second requires that e does not extend the preferences in a way that completely disagrees with an agent's preferences over alternatives. Our first requirement corresponds to the Extension Rule of Barberà et al. [4], and the m -extension of Kruger and Terzopoulou [25], while our second requirement corresponds to their l -extension.

For an agent i with preference relation $>_i^P$ in profile P , we write $\geq_i^{P,e}$ to denote her preferences over sets of alternatives, extended according to e . We say an agent has *e -preferences*, if $\geq_i^{P,e}$ is her preference relation over sets of alternatives.

We will focus in particular on the *Gärdenfors preference extension*³ [20], which we refer to as g . For any two sets X and Y in $2^A \setminus \{\emptyset\}$, we have $X >^g Y$ if and only if one of the following three conditions is satisfied:

- (i) $X \subset Y$ and for all $x \in X$ and $y \in Y \setminus X$, we have $x > y$
- (ii) $Y \subset X$ and for all $x \in X \setminus Y$ and $y \in Y$, we have $x > y$
- (iii) Neither $X \subset Y$ nor $Y \subset X$, and for all $x \in X \setminus Y$ and $y \in Y \setminus X$, we have $x > y$

The Gärdenfors extension dictates that if one set is to be preferred over another, then any alternative added to the first set to reach the second should be preferred to the alternatives in the initial set. Similarly, those alternatives removed from the initial set to reach the new (and preferred) set, should be less preferred. The alternatives the two sets have in common are therefore not relevant, as the Gärdenfors extension only looks at how the two sets differ.

We say an agent has Gärdenfors-preferences if her preferences are extended to sets of alternatives according to the Gärdenfors extension.

2.6 Robust Condorcet Extensions

We say an irresolute social choice function f is *Condorcet-manipulable* by agent i in profile P , under preference extension e , if there exists another profile $P' =_{-i} P$ such that $f(P') >_i^{P,e} f(P)$ and $P \in \mathcal{D}_{\text{Condorcet}}$. We are now ready to present our central definition.

DEFINITION 1. A SCF f is a *robust Condorcet extension* under a preference extension e if f is Condorcet-consistent, and not Condorcet-manipulable in any profile $P \in \mathcal{D}_{\text{Condorcet}}$, by any $i \in N$ with e -preferences.

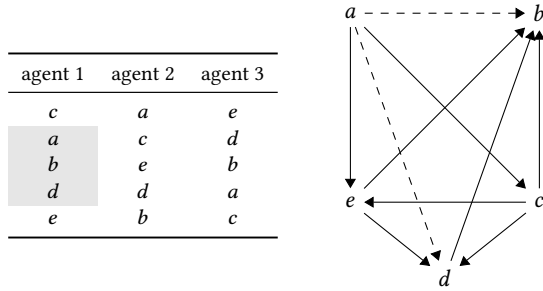
So a SCF is robust under a certain preference extension, if it is not Condorcet manipulable by any agent whose preferences over alternatives have been extended to sets of alternatives according to that extension.

While robustness is a weak strategyproofness requirement, it also speaks about how well a rule can preserve Condorcet winners.

³For a more thorough treatment of the Gärdenfors extension, as well as how it relates to other preference extensions in the literature, we refer to Brandt and Brill [10].

A robust Condorcet extension ensures that, if the truthful profile has a Condorcet winner, then it is a weakly dominant strategy for all agents to report their true preferences, thus ensuring that no Condorcet winner loses that designation because of strategic manipulation. A robust Condorcet extension therefore ensures that profiles with Condorcet winners are, in a sense, stable. We give an example of a Condorcet manipulation of the Copeland SCF,⁴ to demonstrate what failure of robustness looks like.

EXAMPLE 2.1. Suppose the profile below, along with the corresponding majority graph is the “truthful” profile, meaning all three agents have reported their true preferences. As alternative a is the Condorcet winner, Copeland will return a as the sole winner if all agents vote truthfully. Note however, that agent 1 has the ability to reverse the dashed edges (a, b) and (a, d) in the majority graph, by moving b and d above a in her own preference order, while keeping their relative ranking as it is. Agent 1 also has an incentive to do so, as this would result in a majority graph where c —her top choice—is the only alternative with a single incoming edge.



As the Copeland winner is the alternative with the smallest number of incoming edges in the majority graph, c would be the lone Copeland winner if agent 1 misreports her preferences, meaning, Copeland incentivises a Condorcet-manipulation in this profile. Δ

While the profile in Example 2.1 has a Condorcet winner, Copeland is not guaranteed to return this alternative as the winner (or even among them) unless we assume all agents vote truthfully. In the same profile, a robust Condorcet-extension would ensure no agent would have an incentive to misreport her preferences.

3 IMPOSSIBILITES

We present a first impossibility result, showing there is no perfect function that is robust under all preference extensions.

PROPOSITION 3.1. *No tournament-solution SCF is robust under all preference extensions.*

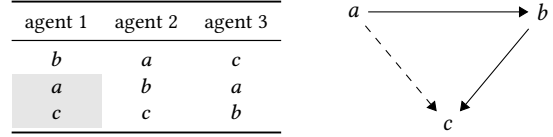
PROOF. Let $A = \{a, b, c\}$, $N = \{1, 2, 3\}$, and let f be a Condorcet-consistent SCF equivalent to the tournament solution F . Suppose agent 1’s preferences over sets of alternatives $\geq_1^{P,e}$ are such that $X \succ_1^{P,e} Y$ if and only if one of the following holds:

- $X = \{a, b, c\}$ and $Y = \{a\}$, or
- $X = \{x\}$ and $Y = \{y\}$ for some $x, y \in A$ such that $x > y$.

These preferences satisfy both our requirements for preference extensions, and are therefore a valid extension of \succ_1^P .

⁴The Copeland score of an alternative in a profile (for odd n) is the number of other alternatives it beats in a pairwise majority contest. The Copeland rule selects those alternatives with the highest Copeland scores [16].

Let P be the profile shown below, with the induced tournament T on the right. As f is a Condorcet extension, $f(P) = \{a\}$. Let $P' =_{-1} P$, where $b \succ_1^{P'} c \succ_1^{P'} a$, meaning agent 1 reverses the edge (a, c) in the induced tournament by reversing the order of these alternatives in her ranking. The tournament T' , induced by P' , consists of a 3-cycle.



As tournament solutions do not distinguish between isomorphic tournaments, $f(P') = F(T') = \{a, b, c\}$. As $\{a, b, c\} \succ_1^{P,e} \{a\}$, this would constitute a successful manipulation for agent 1, meaning f cannot be robust. \square

As there are no social choice functions that are robust for *all* preference extensions, we redirect our search to those that may be robust for *some* preference extension. We first recall a result by McGarvey [26], which we will use to prove the main result of this section. We include the proof for the sake of completeness.

THEOREM 3.2 (MCGARVEY, 1953). *Let A be a set of alternatives, and let \geq be a complete relation over A . Then there is a profile $P \in \mathcal{L}(A)^n$ for some even n such that $\geq^P = \geq$, and if $a > b$, there are $\frac{n}{2} + 2$ agents ranking a over b in P .*

PROOF. For a set of alternatives A and a relation \geq (with strict component $>$) over A , the profile P is constructed for an even number of agents $N = \{i_{ab}, j_{ab} \mid (a, b) \in >\}$ as follows. For every pair of alternatives such that $a > b$, there are two voters i_{ab} and j_{ab} , with the following preferences:

$$a \succ_{i_{ab}}^P b \succ_{i_{ab}}^P x_1 \succ_{i_{ab}}^P \cdots \succ_{i_{ab}}^P x_{|A|-2} \text{ and}$$

$$x_{|A|-2} \succ_{j_{ab}}^P \cdots \succ_{j_{ab}}^P x_1 \succ_{j_{ab}}^P a \succ_{j_{ab}}^P b,$$

Here $\{x_1, \dots, x_{|A|-2}\} = A \setminus \{a, b\}$. For each agent in $N \setminus \{i_{ab}, j_{ab}\}$ who prefers a over b , there will be exactly one corresponding agent who prefers b over a , meaning in the profile P exactly $\frac{n}{2} + 2$ agents prefer a to b . As this holds for any pair of alternatives, it is clear that $\geq^P = \geq$. \square

Note that while our statement of McGarvey’s Theorem is slightly stronger than in the original paper, the proof and the profile constructed in the proof remain the same.

We now show that weakly resolute rules fail robustness for all preference extensions, and further, that they are the only rules that do so.

THEOREM 3.3. *A tournament-solution SCF is weakly resolute if and only if it fails robustness under all preference extensions.*

PROOF. For the right-to-left direction we prove the contrapositive. That is, we suppose f is a tournament-solution SCF that fails weak resoluteness and show it must be robust under some preference extension. To see that this must be the case, note that

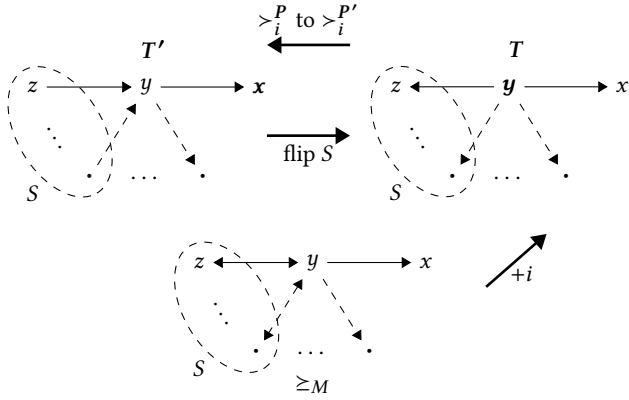


Figure 1: Tournaments T and T' —with winners marked in bold—and relation \geq_M involved in proof of Theorem 3.3. Ties are represented by bidirectional arrows.

any rule failing weak resoluteness never returns singletons outside $\mathcal{D}_{\text{Condorcet}}$. This means the preference extension ranging only over singletons would never (strictly) favour a larger set over the singleton set with the Condorcet winner. As f will always return a set larger than a singleton outside the Condorcet domain, Condorcet-manipulation with these preferences is not possible, thereby making f robust under this preference extension.

For the left-to-right direction, let A be our set of alternatives. Suppose f is a weakly resolute tournament-solution SCF, equivalent to a tournament solution F . We show it is possible for an agent to manipulate f from a profile with a Condorcet winner under an arbitrary preference extension e , meaning f cannot be robust under any preference extension.

We first define two tournaments T and T' , which we will show are single-agent variants. As F is equivalent to a weakly resolute SCF, there is some tournament $T' = (A, >^{T'}) \in \mathcal{T}(A) \setminus \mathcal{T}_{\text{Condorcet}}$ such that $F(T') = \{x\}$ for an alternative $x \in A$. As x is not a Condorcet winner in T' , there must be some $y \in A$ such that $y >^{T'} x$, and by the same reasoning, there must be at least one alternative $z \in A$ such that $z >^{T'} y$. We conclude that the nodes $\{x, y, z\}$ and the edges $(y, x), (z, y)$ must be present in T' . For a visual representation, see Figure 1.

Let $S = \overline{D}_{T'}(y)$ be the dominators of y in T' . We define a second tournament $T = (A, >^T)$, where $y >^T a$ for all $a \in S$, and $>^T$ agrees with $>^{T'}$ on all other pairs of alternatives. In other words, we simply reverse all incoming edges of y in T' to obtain T . Note that this makes y a Condorcet winner in T , meaning $F(T) = \{y\}$.

We now show that T and T' are single-agent variants. We start by constructing a profile P that induces T . To this end, consider a complete relation \geq_M (with strict component $>_M$, and symmetric component \sim_M) over A , such that $\geq_M = >^T \cup \{(a, y) \mid a \in S\}$. This means \geq_M and $>^T$ agree on all pairs of alternatives except those for which T and T' differ. In those cases, \geq_M gives a tie between the alternatives. By Theorem 3.2, we know there exists a profile $P^* = (>_1^{P^*}, \dots, >_n^{P^*})$ with majority relation \geq_M . Further, we know that we can construct P^* with an even number of agents n , such

that for any $a, a' \in A$, where $a >_M a'$, there are exactly $\frac{n}{2} + 2$ agents who prefer a to a' in P^* . We use P^* to construct the profile P . Let $P = (>_1^P, \dots, >_n^P, >_i^P)$, where $x >_i^P y >_i^P a$ for all $a \in A \setminus \{x, y\}$.

To see that P induces tournament T , note that for any pair of alternatives (a, a') , either

- (i) $a >_M a'$ —meaning $a >^T a'$, and $\frac{n}{2} + 2$ prefer a to a' in P^* , or
- (ii) $a \sim_M a'$ —meaning $a' = y$, and $a \in S$ (or vice versa).

If (i) is the case, a majority of agents in P will prefer a to a' regardless of agent i 's preferences; $\frac{n}{2} + 2$ agents still form a strict majority of $n + 1$ agents. If on the other hand (ii) is the case, we know from agent i 's preferences that $y >_i^P a$. As these alternatives were tied in P^* , adding agent i to the profile breaks these ties in favour of y , so a majority of agents in P will now prefer y to a .

This means the only differences between \geq_M and \geq^P relate to the same pairs on which \geq_M and \geq^T differ. As the changes agree with $>^T$, this makes $>^T = >^P$, meaning P induces T . As $F(T) = \{y\}$, we can conclude that $f(P) = \{y\}$.

It now remains to construct a profile P' such that $P =_{-i} P'$ and P' induces T' . Let $P' = (>_1^{P'}, \dots, >_n^{P'}, >_i^{P'})$, and $x >_i^{P'} a >_i^{P'} y$, for all $a \in A \setminus \{x, y\}$, meaning agent i moves y to the bottom of their ranking. Clearly, P' is an i -variant of P . In the tournament induced by P' , it must be the case that the edges (a, y) for all $a \in S$ are present as the majority on these alternatives is dictated by agent i (and all other edges remain as they were in T). As these edges correspond exactly to those where T and T' differ, P' must induce T' , and as $F(T') = \{x\}$ we can conclude $f(P') = \{x\}$.

Finally, let $\geq_i^{P', e}$ be agent i 's true preferences over sets of alternatives, extended according to e . It is immediately clear, as both outcomes are singletons and $x >_i^{P'} y$, that $f(P') >_i^{P', e} f(P)$. As P has a Condorcet winner, this constitutes a Condorcet-manipulation, meaning f cannot be robust under preference extension e . This concludes the proof. \square

We note that Theorem 3.3 applies to two of the most prominent Condorcet extensions—Copeland, and Slater.⁵

4 ROBUST TOURNAMENT SOLUTIONS

In this section, we present our robustness results for several tournament-solution SCF, and their coarsenings.

4.1 Relation to Kelly-Strategyproofness

Though some rules—the omninomination rule and the top cycle for example—have been shown to be strategyproof under Gärdenfors preferences [10, 20], Gärdenfors-strategyproofness—meaning no agent with Gärdenfors preferences can manipulate—is quite hard to attain. While strategyproofness proper for Gärdenfors preferences is not easily satisfied, there are several appealing tournament solution SCFs that have been shown to be strategyproof for the Kelly preference extension.

The *Kelly preference extension* [23]—which we will refer to as k —extends a (strict) preference relation $>$ as follows. For any two sets X and Y in $2^A \setminus \{\emptyset\}$, $X >^k Y$ if and only if $x > y$ for all $x \in X$ and all $y \in Y$. We say a SCF f is *Kelly-strategyproof* if no agent

⁵The result also extends to Slater's weighted counterpart, the Kemeny rule [24].

with Kelly preferences can manipulate successfully, i.e., if there are no agent $i \in N$ and profiles $P =_{-i} P'$ such that $f(P') \succ_i^{P,k} f(P)$. Gärdenfors-strategyproofness implies Kelly-strategyproofness, as the former must exclude more cases of manipulation to be satisfied. As robustness only requires taking into account comparisons where at least one singleton set is present, we can use strategyproofness results for Kelly preferences to show robustness for Gärdenfors preferences.

PROPOSITION 4.1. *If a Condorcet-consistent SCF f is Kelly-strategyproof, then it is a robust Condorcet extension under Gärdenfors preferences.*

PROOF. Suppose we have a rule f that is Kelly-strategyproof. That is, for any two profiles P and P' , and any agent $i \in N$, if $P' =_{-i} P$, then $f(P') \not\succeq_i^{P,k} f(P)$. Suppose P has a Condorcet winner, meaning $f(P) = \{a\}$ for some $a \in A$. If this is the case, either $f(P') = f(P)$, or there is some $a' \in f(P')$ such that $a \succ_i^P a'$. If $f(P') = f(P)$ then agent i cannot strictly prefer one outcome over the other under any preference extension, so $f(P') \not\succeq_i^{P,g} f(P)$.

If there is some $a' \in f(P')$ s.t. $a \succ_i^P a'$, then, as $a' \in f(P') \setminus f(P)$, it is immediate from the definition of the Gärdenfors extension that $f(P') \not\succeq_i^{P,g} f(P)$. \square

We therefore get robustness “for free” for Condorcet extensions known to be Kelly-strategyproof. Among these are social choice functions that are not fully strategyproof for Gärdenfors preferences, such as the bipartisan set and the minimal covering set [10].

4.2 Minimal Extending Set & Beyond

Before we present our positive robustness results, we need to define the Banks set and the minimal extending set.

A tournament $T' = (S', \succ^{T'})$ is a *maximal transitive subtournament* of $T = (S, \succ^T)$ if

- (i) T' is a subtournament of T ,
- (ii) $\succ^{T'}$ is a transitive relation, and
- (iii) there is no other transitive subtournament $(S'', \succ^{T''})$ of T such that $S' \subset S''$.

We write \hat{T} to denote the set of all maximal transitive subtournaments of tournament T , and $\text{top}(\succ)$ to denote the maximal element of the strict linear order \succ . Note that if a tournament T has a Condorcet winner, it will be the maximal element of *all* maximal transitive subtournaments.⁶

The *Banks set* [2] is the set of maximal elements of all maximal transitive subtournaments of a tournament:

$$\text{BA}(T) = \{\text{top}(\succ_S^T) \mid (S, \succ_S^T) \in \hat{T}\}.$$

Because the Condorcet winner will top all maximal transitive subtournaments, Banks is a Condorcet extension.

A set $S \subseteq A$ is a *BA-stable set* of a tournament T if $a \notin \text{BA}(S \cup \{a\}, \succ_{S \cup \{a\}}^T)$ for all $a \in A \setminus S$. A BA-stable set of a tournament T is *minimal* if there is no BA-stable set S' of T such that $S' \subset S$. The

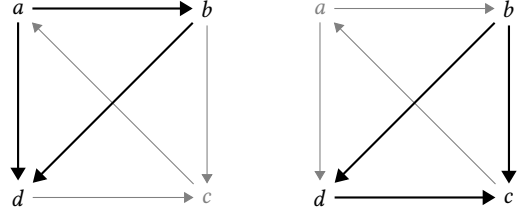
⁶As the existence of a Condorcet winner does not imply no cycles are present, there may indeed be several maximal transitive subtournaments.

minimal extending set $\text{ME}(T)$ [9] of a tournament T is the union of all minimal BA-stable sets of T :

$$\text{ME}(T) = \bigcup \{S \subseteq A \mid S \text{ is a minimal BA-stable set of } T\}.$$

We give an example to shed some light on these definitions.

EXAMPLE 4.1. In the tournament T below, the two maximal transitive subtournaments are indicated using black edges. It is clear that the subtournaments are transitive, and they are both maximal; adding the last alternative will break transitivity. From examining these subtournaments, we can see that $\text{BA}(T) = \{a, b\}$.



The tournament has two minimal BA-stable sets: $\{a, b, d\}$ —because $c \notin \text{BA}(T)$, and $\{a, b, c\}$ —because $d \notin \text{BA}(T)$. ME will output the union of these sets: $\text{ME}(T) = \{a, b, c, d\}$. Note that the set $\{b, c, d\}$ is not BA-stable, as $a \in \text{BA}(T)$. \triangle

ME is one of several tournament solutions defined based on this notion of stability. The top cycle for example, is the union of all minimal CNL-stable sets [9], where CNL is the tournament solution returning the set of all Condorcet nonlosers—meaning all alternatives with at least one outgoing edge. We will use BA and ME to refer both to the tournaments solutions above, and the equivalent social choice functions.

The minimal extending set is not Kelly-strategyproof as Kelly-strategyproofness of tournament solutions implies set-monotonicity, a strong monotonicity axiom [8]. As ME is not a monotonic rule, it also fails the stronger set-monotonicity axiom. However, we show it is still robust under Gärdenfors preferences, and extend this result to all coarsenings of ME.

THEOREM 4.2. *ME is a robust Condorcet extension under Gärdenfors preferences.*

PROOF. For a set of alternatives A , and a set of agents N , let $P =_{-i} P'$ be i -variant profiles for an agent $i \in N$. Let P be such that $x \in A$ is the Condorcet winner in P . Let $T = (A, \succ^T)$ and $T' = (A, \succ^{T'})$ be the (single-agent variant) tournaments induced by P and P' , respectively.

We assume $\text{ME}(T') \neq \text{ME}(T)$.⁷ Because of this, we know $\overline{D}_{T'}(x)$ is nonempty, as the two outcomes cannot differ if x remains a Condorcet winner in T' . As $P =_{-i} P'$, any changes going from T to T' must be counter to agent i 's preferences. This implies $x \succ_i^P a$ for all $a \in \overline{D}_{T'}(x)$. So, all alternatives in $\overline{D}_{T'}(x)$ are worse than x according to agent i .

We want to show that there is some minimal BA-stable set S of T' , such that $S \cap \overline{D}_{T'}(x) \neq \emptyset$. This would guarantee the existence of an alternative $a \in \overline{D}_{T'}(x)$ that is also in $\text{ME}(T')$, precluding agent i from preferring this outcome to $\text{ME}(T)$.

⁷If no such i -variants exist, robustness of the rule would immediately follow.

So suppose for contradiction that S is a minimal BA-stable set of T' such that $S \cap \overline{D_{T'}(x)} = \emptyset$. The only way this can be the case is if $S \subseteq D_{T'}(x) \cup \{x\}$. We consider two cases.

Case 1: Suppose $x \notin S$. As x dominates all alternatives in $D_{T'}(x)$, it will dominate all alternatives in S , as $S \subseteq D_{T'}(x)$. This means x is a Condorcet winner in the tournament $(S \cup \{x\}, >_{S \cup \{x\}}^{T'})$, and thus, $x \in \text{BA}(>_{S \cup \{x\}}^{T'})$. So S cannot be a BA-stable set, contradicting our assumption that it is a minimal one.

Case 2: Suppose instead $x \in S$. To reach our contradiction, we want to show there exists an alternative $a \in \overline{D_{T'}(x)}$ such that $a \in \text{BA}(>_{S \cup \{a\}}^{T'})$ —which would imply S is not BA-stable.

We use an algorithm proposed by Hudry [22] to find such an alternative $a \in \text{BA}(>_{S \cup \{a\}}^{T'})$. We start at step 1 with a transitive subtournament of $(S \cup \{a\}, >_{S \cup \{a\}}^{T'})$. Let $S_1 = (\{x, a\}, >_{\{x, a\}}^{T'})$, for some $a \in \overline{D_{T'}(x)}$. We label all remaining elements of S —which are all elements of $D_{T'}(x)$ —in any order from 2 to $|S|$. At step k , we look at the alternative labelled k , and add it to the tournament S_{k-1} to create S_k , if it does not break transitivity to do so. As a dominates x , and x dominates all $a' \in D_{T'}(x)$, adding any alternative outside the dominion of a will break transitivity, as it will create a 3-cycle. Thus, at any step, an alternative $a' \in D_{T'}(x)$ will only be added to the tournament if $a >^{T'} a'$. When the algorithm terminates after iterating through all alternatives, we will be left with a subtournament $S_{|S|}$ of $(S \cup \{a\}, >_{S \cup \{a\}}^{T'})$. It is easy to see the resulting tournament will be transitive, and it will indeed be a maximal transitive subtournament of $(S \cup \{a\}, >_{S \cup \{a\}}^{T'})$, as no further alternatives can be added to the tournament without breaking transitivity. Importantly, the maximal element of the resulting subtournament will be a , meaning $a \in \text{BA}(>_{S \cup \{a\}}^{T'})$. Thus, S cannot be an BA-stable set, which contradicts our assumption that it is a minimal one.

As we have shown that no subset of $D_{T'}(x) \cup \{x\}$ can be a BA-stable set of T' , any minimal BA-stable set must contain at least one element of $\overline{D_{T'}(x)}$, meaning it cannot be the case that $\text{ME}(T') >_i^{P, g} \text{ME}(T)$. \square

In terms of decisiveness, ME is among the more decisive tournament solutions that fail weak resoluteness, as it is a refinement of several prominent tournament solutions, including the top cycle, the Banks set [13] and the uncovered set (see Brandt et al. [11] for definitions and further rules that fall into this category).

We now show that these coarsenings of ME inherit the robustness property.

LEMMA 4.3. *If a Condorcet-consistent SCF f is robust under Gärdenfors preferences, then all Condorcet-consistent coarsenings of f are robust under Gärdenfors preferences.*

PROOF. Let f be a SCF that is robust under Gärdenfors preferences, and let f' be a Condorcet-consistent coarsening of f . Let P be a profile with a Condorcet winner a . Note that $f(P) = f'(P) = \{a\}$ as they are both Condorcet extensions.

Suppose P' is an i -variant of P for some agent $i \in N$. Because f is robust under Gärdenfors preferences, either (i) there must be some $a' \in f(P')$ such that $a >_i^P a'$, or (ii) $f(P) = f(P')$.

If (i) is the case, then as $f(P') \subseteq f'(P')$, a' is also an element of $f'(P')$. As $f'(P) = \{a\}$, we know $a' \in f'(P') \setminus f'(P)$, meaning by definition of the Gärdenfors extension that it cannot be the case that $f'(P') >_i^{P, g} f'(P)$.

If (ii) is the case, we know a must also be the Condorcet winner in P' as f cannot satisfy weak resoluteness if it is robust under any preference extension, and therefore does not return singletons outside the Condorcet domain. Since f' is also a Condorcet extension, we know $f'(P') = \{a\}$, meaning $f'(P') \not>_i^{P, g} f'(P)$. \square

COROLLARY 4.4. *Condorcet-consistent coarsenings of ME are robust under Gärdenfors preferences.*

Corollary 4.4 follows from Lemma 4.3 and Theorem 4.2, and it establishes the robustness of the Banks set and the uncovered set, neither of which is Kelly-strategyproof. Note that Corollary 4.4 is not restricted to tournament-solution SCFs, but holds for all Condorcet-consistent social choice functions.

5 CONCLUSION

We have introduced the strategyproofness-related notion of a robust Condorcet extension. We have argued that Condorcet extensions that are robust are preferable to those that are not, as we can trust that they will return true Condorcet winners when they exist. We have introduced an axiom—weak resoluteness—and shown that no weakly resolute tournament solution can be a robust Condorcet extension. Finally, we have shown that the minimal extending set is a robust Condorcet extension under Gärdenfors preferences, and have extended this result to all coarsenings of ME.

We have argued that in lieu of searching for fully strategyproof rules, a fruitful endeavor is to explore immunity against more specific manipulations that may interact with, and compromise, other desirable properties satisfied by manipulable social choice functions. We have scratched the surface in this paper, but have contained our exploration to robustness of irresolute rules in general, and tournament solutions in particular. These are, of course, only a small class of all Condorcet extensions, and it remains to be seen if similar results can be obtained for other classes. Finally, there are many other arenas of social choice theory open to similar notions of strategyproofness. While proportionality and strategyproofness cannot be satisfied by the same multiwinner approval voting rule for example [29], one might want to ask to what degree proportionality is affected by the failure of strategyproofness.

REFERENCES

- [1] Kenneth J. Arrow, Amartya Sen, and Kotaro Suzumura. 2010. *Handbook of Social Choice and Welfare*. Vol. 2. Elsevier.
- [2] Jeffrey S. Banks. 1985. Sophisticated Voting Outcomes and Agenda Control. *Social Choice and Welfare* 1, 4 (1985), 295–306.
- [3] Salvador Barberà. 1977. Manipulation of Social Decision Functions. *Journal of Economic Theory* 15, 2 (1977), 266–278.
- [4] Salvador Barberà, Walter Bossert, and Prasanta K. Pattanaik. 2004. Ranking Sets of Objects. In *Handbook of Utility Theory*. Springer, 893–977.
- [5] John J. Bartholdi, Craig A. Tovey, and Michael A. Trick. 1989. The Computational Difficulty of Manipulating an Election. *Social Choice and Welfare* 6, 3 (1989), 227–241.
- [6] Duncan Black. 1948. On the Rationale of Group Decision-making. *Journal of Political Economy* 56, 1 (1948), 23–34.
- [7] Walter Bossert and Yves Sprumont. 2014. Strategy-proof Preference Aggregation: Possibilities and Characterizations. *Games and Economic Behavior* 85 (2014), 109–126.

- [8] Felix Brandt. 2011. Group-Strategyproof Irresolute Social Choice Functions. In *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI-2011)*, 79–84.
- [9] Felix Brandt. 2011. Minimal Stable Sets in Tournaments. *Journal of Economic Theory* 146, 4 (2011), 1481–1499.
- [10] Felix Brandt and Markus Brill. 2011. Necessary and Sufficient Conditions for the Strategyproofness of Irresolute Social Choice Functions. In *Proceedings of the 13th Conference on Theoretical Aspects of Rationality and Knowledge (TARK-2011)*.
- [11] Felix Brandt, Markus Brill, and Paul Harrenstein. 2016. Tournament Solutions. In *Handbook of Computational Social Choice*, F. Brandt, V. Conitzer, U. Endriss, J. Lang, and A. D. Procaccia (Eds.). Cambridge University Press, Chapter 3.
- [12] Felix Brandt, Vincent Conitzer, Ulle Endriss, Jérôme Lang, and Ariel D Procaccia. 2016. *Handbook of Computational Social Choice*. Cambridge University Press.
- [13] Felix Brandt, Paul Harrenstein, and Hans Georg Seedig. 2017. Minimal Extending Sets in Tournaments. *Mathematical Social Sciences* 87 (2017), 55–63.
- [14] Vincent Conitzer and Toby Walsh. 2016. Barriers to Manipulation in Voting. In *Handbook of Computational Social Choice*, F. Brandt, V. Conitzer, U. Endriss, J. Lang, and A. D. Procaccia (Eds.). Cambridge University Press, Chapter 6.
- [15] Vincent Conitzer, Toby Walsh, and Lirong Xia. 2011. Dominating Manipulations in Voting with Partial Information. In *Proceedings of the 25th AAAI Conference on Artificial Intelligence (AAAI-2011)*.
- [16] Arthur H. Copeland. 1951. A “Reasonable” Social Welfare Function. (1951). University of Michigan Seminar on Applications of Mathematics to the Social Sciences.
- [17] John Duggan and Thomas Schwartz. 2000. Strategic Manipulability Without Resoluteness or Shared Beliefs: Gibbard-Satterthwaite Generalized. *Social Choice and Welfare* 17, 1 (2000), 85–93.
- [18] Peter C. Fishburn. 1977. Condorcet Social Choice Functions. *SIAM J. Appl. Math.* 33, 3 (1977), 469–489.
- [19] Wulf Gaertner. 2001. *Domain Conditions in Social Choice Theory*. Cambridge University Press.
- [20] Peter Gärdenfors. 1976. Manipulation of Social Choice Functions. *Journal of Economic Theory* 13, 2 (1976), 217–228.
- [21] Allan Gibbard. 1973. Manipulation of Voting Schemes: A General Result. *Econometrica* 41, 4 (1973), 587–601.
- [22] Olivier Hudry. 2004. A Note on “Banks Winners in Tournaments are Difficult to Recognize” by G.J. Woeginger. *Social Choice and Welfare* 23, 1 (2004), 113–114.
- [23] Jerry S. Kelly. 1977. Strategy-proofness and Social Choice Functions without Singlevaluedness. *Econometrica* 45, 2 (1977), 439–446.
- [24] John G. Kemeny. 1959. Mathematics without Numbers. *Daedalus* 88, 4 (1959), 577–591.
- [25] Justin Kruger and Zoi Terzopoulou. 2020. Strategic Manipulation with Incomplete Preferences: Possibilities and Impossibilities for Positional Scoring Rules. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS-2020)*.
- [26] David C. McGarvey. 1953. A Theorem on the Construction of Voting Paradoxes. *Econometrica* 21, 4 (1953), 608–610.
- [27] Svetlana Obratsova, Zinovi Rabinovich, Edith Elkind, Maria Polukarov, and Nicholas R. Jennings. 2016. Trembling Hand Equilibria of Plurality Voting. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI-2016)*.
- [28] Martin J. Osborne and Ariel Rubinstein. 2003. Sampling Equilibrium, with an Application to Strategic Voting. *Games and Economic Behavior* 45, 2 (2003), 434–441.
- [29] Dominik Peters. 2018. Proportionality and Strategyproofness in Multiwinner Elections. In *Proceedings of the 17th International Conference on Autonomous Agents and Multiagent Systems, (AAMAS-2018)*, 1549–1557.
- [30] Annemieke Reijngoud and Ulle Endriss. 2012. Voter Response to Iterated Poll Information. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems (AAMAS-2012)*, 635–644.
- [31] Shin Sato. 2008. On Strategy-proof Social Choice Correspondences. *Social Choice and Welfare* 31, 2 (2008), 331–343.
- [32] Shin Sato. 2013. A Sufficient Condition for the Equivalence of Strategy-proofness and Nonmanipulability by Preferences Adjacent to the Sincere one. *Journal of Economic Theory* 148, 1 (2013), 259–278.
- [33] Mark Allen Satterthwaite. 1975. Strategy-proofness and Arrow’s Conditions: Existence and Correspondence Theorems for Voting Procedures and Social Welfare Functions. *Journal of Economic Theory* 10, 2 (1975), 187–217.
- [34] Reinhard Selten. 1975. Reexamination of the Perfectness Concept for Equilibrium Points in Extensive Games. *Journal of Game Theory* 103 (1975), 25–55.
- [35] Patrick Slater. 1961. Inconsistencies in a Schedule of Paired Comparisons. *Biometrika* 48, 3–4 (1961), 303–312.