# Condorcet's Principle and the Preference Reversal Paradox

DOMINIK PETERS, University of Oxford

We prove that every Condorcet-consistent voting rule can be manipulated by a voter who completely reverses their preference ranking, assuming that there are at least 4 alternatives. This corrects an error and improves a result of [Sanver, M. R., & Zwicker, W. S. (2009). One-way monotonicity as a form of strategy-proofness. Int J Game Theory 38(4), 553-574.] For the case of precisely 4 alternatives, we exactly characterise the number of voters for which this impossibility result can be proven. We also show analogues of our result for irresolute voting rules. We then leverage our result to state a strong form of the Gibbard–Satterthwaite Theorem.

## 1 INTRODUCTION

The Gibbard–Satterthwaite Theorem establishes that every non-trivial voting rule can be manipulated by voters through misrepresenting their preferences. In this paper, we will see that Condorcet extensions (voting rules that select the Condorcet winner if one exists) suffer from a particularly offensive failure of strategyproofness: all of them can be manipulated by a voter who completely reverses their preference ranking. For example, such a voting rule might designate $c$ to be the winning alternative if voter $i$ truthfully reports the ordering $a >_i b >_i c >_i d$, but choose $b$ as the winner if voter $i$ instead reports the ordering $d >_i c >_i b >_i a$. Since $i$ truthfully prefers $b$ to $c$, this is a successful manipulation, which one might consider surprising given that $i$ misreported every possible pairwise comparison. We will say that voting rules that are manipulable in this way suffer from the *preference reversal paradox*. While all Condorcet extensions exhibit this paradox, scoring rules (such a plurality and Borda's rule) are immune.

Preference reversal paradoxes were first introduced by Sanver and Zwicker [2009] in their study of monotonicity properties; they say that voting rules which avoid this paradox satisfy *half-way monotonicity*.[1] As Sanver and Zwicker [2009] show, half-way monotonicity is a weaker property than *participation*, an axiom stating that a voter cannot obtain a strictly better result by abstaining from an election; equivalently, participation says that voting truthfully guarantees a (weakly) better result than not voting at all. In a famous paper, Moulin [1988] showed that participation is incompatible with Condorcet-consistency, so that Condorcet extensions must suffer from the *no-show paradox* [Fishburn and Brams, 1983]. This result is often interpreted as showing that all Condorcet extensions are *manipulable* (through abstention). Notice, however, that this notion of manipulation (referring to electorates of different sizes) is quite different from the fixed-electorate manipulations that are the subject of the Gibbard–Satterthwaite Theorem, where a voter changes their preference ordering in some way [see also Núñez and Sanver, 2017]. We will see that half-way monotonicity, which is both weaker than participation and weaker than strategyproofness in the Gibbard–Satterthwaite sense, is already incompatible with Condorcet-consistency.

This result first appeared in Sanver and Zwicker [2009] who gave a proof that, for 4 or more alternatives and for sufficiently many voters, Condorcet extensions must fail half-way monotonicity. However, their proof contains an arithmetical mistake[2] that is non-trivial to fix. The proof technique also is only able to establish an impossibility for electorates containing a sufficiently large *even* number of voters. Further, their proof requires at least 702 voters to go through, this bound growing

---

[1] They chose this name because half-way monotonicity is a weaker version of their notion of *one-way* monotonicity.

[2] In the last paragraph of the proof or their Theorem 5.2, they calculate that $n^*(Q) = 30 + 8$, when in fact $n^*(Q) = 30 + 4 \cdot m! \gg 38$ which makes their "Condition M" inapplicable to profile $Q$. This problem was noticed by Wei Yu and Tokuei Higashino (Zwicker, private communication).

| $n =$ | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| participation | P | P | P | P | P | P | P | P | P | I | I | I | I | I | I | I | I | I | I | I | I | I | I | I |
| half-way monotonicity | P | P | P | P | P | P | P | P | P | P | P | P | I | P | I | P | I | P | I | P | I | I | I | I |

Possibility ▢   Impossibility ▨

Table 1. Numbers $n$ of voters for which Condorcet extensions can satisfy participation or half-way monotonicity, when there are exactly $m = 4$ alternatives.

exponentially as the number of alternatives increases,[3] which leaves open the question of whether the preference reversal paradox is a problem in practical voting situations with moderate numbers of voters.

We give a direct proof of the impossibility, treating the cases of electorates with odd and even numbers of voters separately. Our arguments require 15 voters for the odd case and 24 voters for the even case. These constant bounds hold for any number $m \geqslant 4$ of alternatives. Using computer-aided techniques we are able to show that these results are tight: for the case of precisely 4 alternatives, there exist Condorcet extensions satisfying half-way monotonicity for up to 13 voters and 22 voters, respectively. (For 3 alternatives, it is known that the maximin rule with some fixed tie-breaking is a Condorcet extension satisfying half-way monotonicity.)

Both our positive and our negative results were proved with the help of SAT solvers, using a technique introduced by Geist and Endriss [2011] and Tang and Lin [2009]. For a recent survey, see the book chapter by Geist and Peters [2017]. The general approach is to produce a formula of propositional logic whose models correspond to voting rules that are Condorcet-consistent and half-way monotonic. We can then pass this formula to a SAT solver. If the formula is satisfiable, we obtain an example of a good voting rule; if it is unsatisfiable, we have an impossibility result. In the unsatisfiable case, using an idea of Brandt and Geist [2016], we can then extract a *minimal unsatisfiable set* (MUS) which can often be translated into a human-readable impossibility proof. This technique was used to prove impossibility results for Fishburn-strategyproofness in majoritarian social choice functions [Brandt and Geist, 2016], for Fishburn-participation in the same setting [Brandl et al., 2015], for the no-show paradox [Brandt, Geist, and Peters, 2016], and for probabilistic social choice rules [Brandl et al., 2016]. Since our techniques are variations of the technique of Brandt, Geist, and Peters [2016], we will only give a brief overview of the method in this paper.

As mentioned above, our theorem implies Moulin's result for participation. Brandt, Geist, and Peters [2016] recently showed that, for 4 alternatives, Moulin's impossibility requires 12 voters to go through, while there exists a Condorcet extension satisfying participation for up to 11 voters. This gives us a rough but intriguing way to compare the relative strengths of participation and half-way monotonicity (see Table 1); we can see that half-way monotonicity is weaker than participation, but not by much.

In Section 6, we will discuss some extensions of this result. First, we consider *irresolute* voting rules which return a *set* of alternatives; we check whether in this more general model we can guarantee half-way monotonicity for larger numbers of voters (the answer turns out to be *no*). Then we consider the *strong* preference reversal paradox, which occurs when a voter can cause

---

[3]The large number arises because the proof uses several copies of the full profile containing a copy of each of the $m!$ preference orders. Fixing the arithmetical error described above tends to necessitate using many more voters than this (Zwicker, private communication).

their *most*-preferred alternative to win by reversing their preferences. We show that most, but not all, Condorcet extensions exhibit this strong paradox.

Finally, in Section 7, we combine our results with a theorem of Campbell and Kelly [2003] to give a strengthened version of the Gibbard–Satterthwaite Theorem. This version claims that every non-trivial voting rule is either *needlessly* or *egregiously* manipulable. This gives a more explicit description of the types of manipulations that are sufficient to obtain an impossibility.

## 2  DEFINITIONS

A *linear order* $\geq$ is a complete, antisymmetric, transitive binary relation over $A$. We write $>$ for the strict (irreflexive) part of $\geq$. The set of all linear orders over $A$ is denoted by $A!$. The *reverse* $>^{\text{rev}}$ of a linear order $>$ is defined by $a >^{\text{rev}} b \iff b > a$ for all $a, b \in A$.

Let $N = \{1, \ldots, n\}$ be a finite set of $n$ *voters*, and let $A$ be a finite set of $m$ *alternatives*. Often, we will consider the case of precisely 4 alternatives, when $A = \{a, b, c, d\}$. A *profile* $P$ is a function assigning to every $i \in N$ a linear ordering $\geq_i$ of the alternatives. Thus, the set of profiles is $A!^N$. A (resolute) *voting rule* is a function $f : A!^N \to A$ that assigns a winning alternative $f(P) \in A$ to every profile $P \in A!^N$.

Given a profile $P$, we say that $a \in A$ is the (unique) *Condorcet winner* if $|\{i \in N : a >_i b\}| > |\{i \in N : b >_i a\}|$ for all $b \in A \setminus \{a\}$. Thus, a Condorcet winner wins against every other alternative in a pairwise majority comparison. We say that a voting rule $f$ is a *Condorcet extension* if $f$ selects the Condorcet winner whenever one exists.

Given a profile $P \in A!^N$ where $i \in N$, we write $P_{-i} := P|_{N \setminus \{i\}}$ for the profile obtained from $P$ by removing voter $i$. We also write $(P_{-i}, >_i') := P_{-i} \cup \{(i, >_i')\}$ for the profile obtained from $P$ by replacing $i$'s vote by $>_i'$.

*Definition 2.1.* A voting rule $f$ satisfies *half-way monotonicity* if

$$f(P_{-i}, >_i) \geq_i f(P_{-i}, >_i^{\text{rev}}) \quad \text{for all profiles } P \in A!^N \text{ and all voters } i \in N.$$

Thus, voters weakly prefer voting truthfully to voting the reverse of their preferences. If a rule violates half-way monotonicity, we say that it suffers from the *preference reversal paradox*.

## 3  RELATIONSHIP TO PARTICIPATION

*Participation* is a property of voting rules that assign outcomes to profiles with varying numbers of voters. Let us define a *variable-electorate* voting rule as a function that assigns a winning alternative to every profile defined on some finite electorate $N \subseteq \mathbb{N}$. If $N$ is an electorate with $i \notin N$, $>_i$ is some linear order, and $P$ is a profile on $N$, then we define $P + (i, >_i)$ to be the profile obtained by letting voter $i$ join $P$. Then we say that a variable-electorate voting rule $f$ satisfies *participation* if for all electorates $N$, all voters $i \notin N$, and all preference orders $>_i$, we have $f(P + (i, >_i)) \geq_i f(P)$. In other words, voters always weakly prefer joining an election.

It turns out that participation is a stronger requirement than half-way monotonicity. This was shown by Sanver and Zwicker [2009, Theorem 4.1] using a proof that established several interrelated implications among their monotonicity axioms. Here, we give a direct proof of this implication.

LEMMA 3.1 (SANVER AND ZWICKER, 2009). *If a variable-electorate voting rule $f$ satisfies participation, then $f$ satisfies half-way monotonicity.*

PROOF. The key idea is that the reversal of a vote $>_i$ is equivalent to $>_i$ leaving the election and $>_i^{\text{rev}}$ joining it. Let $P \in A!^N$ be a profile and let $i \in N$ be a voter with preferences $\geq_i$ in $P$. Consider the profile $P_{-i}$ with $i$ removed. By participation, we have $f(P) \geq_i f(P_{-i})$. Also by participation, we have $f(P_{-i}, >_i^{\text{rev}}) \geq_i^{\text{rev}} f(P_{-i})$. Putting these together, and using the definition of the reverse of an

order, we have

$$f(P) \succcurlyeq_i f(P_{-i}) \succcurlyeq_i f(P_{-i}, >_i^{\text{rev}}).$$

Thus, using transitivity, we have verified half-way monotonicity.                    □

Interestingly, to deduce half-way monotonicity for electorates of $n$ voters, we only require participation to hold between electorates of size $n-1$ and $n$. Núñez and Sanver [2017] also prove the implication of Lemma 3.1 by proposing an intermediate "Condition $\lambda$" that is implied by participation and that implies half-way monotonicity.

## 4 METHOD

To obtain the possibility and impossibility results of the next section, we used the computer-aided technique introduced by Geist and Endriss [2011] and Tang and Lin [2009]. In this section, we will give a brief overview of the basic ideas. For a more detailed discussion of the method, see the survey by Geist and Peters [2017].

We begin by translating our question (of whether a Condorcet extension satisfying half-way monotonicity exists) into *propositional logic*. To do so, we fix a set $A$ of $m$ alternatives and a set $N$ of $n$ voters. We then explicitly enumerate the set $A!^N$ of profiles, and introduce propositional variables $x_{P,a}$ for each profile $P \in A!^N$ and each alternative $a \in A$. The intended meaning of the variables is

$$x_{P,a} \text{ is set true} \iff f(P) = a,$$

where $f$ is a voting rule. To pin down this meaning, we produce a propositional formula $\varphi$ in conjunctive normal form (CNF) by encoding three classes of constraints:

- *functionality* of $f$, i.e., that $f(P) = a$ for exactly one alternative, which means that there is *at least* one and *at most* one such alternative:

$$\varphi_{\text{functionality}} \equiv \bigwedge_{P \in A!^N} \left( \left( \bigvee_{a \in A} x_{P,a} \right) \wedge \bigwedge_{a \neq b \in A} (\neg x_{P,a} \vee \neg x_{P,b}) \right)$$

- *Condorcet-consistency*: for $a \in A$, let $C_a \subseteq A!^N$ be the set of profiles in which $a$ is the Condorcet winner.

$$\varphi_{\text{Condorcet}} \equiv \bigwedge_{a \in A} \bigwedge_{P \in C_a} x_{P,a}$$

- *half-way monotonicity*: if $a, b \in A$ are such that $a >_i b$ for voter $i$ in profile $P \in A!^N$, and $f(P) = b$, then $f(P_{-i}, >_i^{\text{rev}}) \neq a$.

$$\varphi_{\text{half-way monotonicity}} \equiv \bigwedge_{i \in N} \bigwedge_{P \in A!^N} \bigwedge_{\substack{a,b \in A \\ a >_i b}} (\neg v_{P,b} \vee \neg v_{(P_{-i}, >_i^{\text{rev}}),a}).$$

Putting these formulas together, we obtain $\varphi \equiv \varphi_{\text{functionality}} \wedge \varphi_{\text{Condorcet}} \wedge \varphi_{\text{half-way monotonicity}}$. Then it is clear that each true/false assignment to the propositional variables that satisfies $\varphi$ induces a voting rule $f$ which is Condorcet-consistent and satisfies half-way monotonicity.

Next, we write down $\varphi$ in a text file in the standardised DIMACS format, and pass this formula to a *SAT solver*, that is, a computer program which checks whether a given propositional formula is satisfiable or unsatisfiable. Despite this decision problem being NP-complete, modern SAT solvers such as lingeling [Biere, 2013] or glucose [Audemard and Simon, 2009] can often solve even large formulas in a relatively short time.

For our choice of $\varphi$, a satisfiability result gives us an example of a Condorcet extension which avoids the preference reversal paradox. An unsatisfiability result implies an impossibility. In the

former case, the SAT solver will return a satisfying assignment, which induces a specific voting rule $f$ which is available as a look-up table. In the latter case, the SAT solver will merely report "UNSAT". It would be desirable to obtain a proof of this claim. While many solvers are able to produce an unsatisfiability proof in a computer-readable format, these proofs can be very big (a recent result in Ramsey theory required 200 TB [Heule et al., 2016]) and uninsightful. Following Brandt and Geist [2016], we use a technique that is sometimes able to produce short and human-readable proofs. To do so, we obtain a *minimal unsatisfiable set* (MUS), which is a subselection of the clauses of $\varphi$ which is already unsatisfiable. If we can find a small enough MUS using tools such as MUSer2 [Belov and Marques-Silva, 2012] or MARCO [Liffiton et al., 2015], we can then translate this MUS into a human-readable proof.

For more details on this approach, we refer to Geist and Peters [2017] and to Brandt, Geist, and Peters [2016].

## 5 IMPOSSIBILITY RESULTS

In this section, we will present our main results. Our technique is inductive: we give positive and negative results for a specific number $n$ of voters, and then use the following lemma to conclude that positive results also hold for smaller $n$ and negative results hold for larger $n$, as long as parity is preserved. While the parity of $n$ is immaterial as to whether participation can be satisfied by a variable-electorate voting rule defined on electorates up to size $n$, we will see that half-way monotonicity is less restrictive on Condorcet extensions defined for even electorates.

LEMMA 5.1 (INDUCTION STEP). *Fix a number $m$ of alternatives, and let $n \geqslant 1$. If there exists a Condorcet extension defined on electorates with $n + 2$ voters which satisfies half-way monotonicity, then there also exists a Condorcet extension for $n$ voters satisfying half-way monotonicity.*

PROOF. Fix some linear order $>_*$ over $A$. Suppose $|N| = n$, and suppose $f_{n+2}$ is a Condorcet extension satisfying half-way monotonicity, defined for the electorate $N \cup \{i, j\}$. Then define the voting rule $f_n$ on the electorate $N$ by

$$f_n(P) := f_{n+2}(P + (i, >_*) + (j, >_*^{\mathrm{rev}})) \text{ for all profiles } P \in A!^N.$$

Then the voting rule $f_n$ is Condorcet-consistent: if a profile $P \in A!^N$ admits a Condorcet winner, then this alternative remains the Condorcet winner after adding two completely opposed orders to $P$, since this operation does not change the majority margins. Further, $f_n$ satisfies half-way monotonicity, since any successful manipulation through preference reversal for $f_n$ can also be pulled off for $f_{n+2}$. □

Contrapositively, this lemma implies that an incompatibility result between Condorcet-consistency and half-way monotonicity for $n$ voters also applies to $n + 2k$ voters, for each $k \geqslant 0$. Thus, in our impossibility proofs below, we only need to handle the base case for $n = 15$ and $n = 24$, respectively.

Before we present the proofs, let us have a look at our main positive result.

PROPOSITION 5.2 (POSSIBILITIES). *For $m = 4$ alternatives, and for either $n = 13$ or $n = 22$ voters, there exists a Condorcet extension satisfying half-way monotonicity.*

This result was obtained by running a SAT solver to decide the satisfiability of a suitable (large) formula of propositional logic as described in Section 4. The major downside of this technique is that the voting rules of Proposition 5.2 are only available as *look-up tables*. Both of the voting rules mentioned are C2 functions in Fishburn's classification, i.e., they only depend on the majority

margins of the input profile.[4] The only available description of these voting rules are text files indicating, for every weighted majority tournament, which alternative is to be selected. For the case of $n = 22$ voters, the uncompressed file has a size of 1.7GB. As with the voting rules found by Brandt, Geist, and Peters [2016], it would be desirable to find such rules that have a more concise description.

Now let us move on to our negative results. The proof diagrams in Figures 1 and 2 give a graphical representation of the proof steps involved. An arc from $P$ to $P'$ labelled "*reverse* $2\,cabd$" is interpreted as "if the voting rule chooses $a$ or $b$ at $P$, then the rule must also choose $a$ or $b$ at $P'$ by half-way monotonicity". The profiles at the leafs all admit a Condorcet winner, which leads to a contradiction. The general proof strategy of our impossibility proofs is a follows: we identify an initial profile $P_0$, and iterate through each possible value of $f(P_0) \in A$. Assuming that $f(P_0) = x$, say, will then, by half-way monotonicity, imply restrictions on the possible values that $f$ can take at profiles obtained from $P_0$ by reversing some of the votes. In particular, it will imply that $f$ must not pick the Condorcet winner at some of these profiles, contradicting $f$ being a Condorcet extension.

As we noted, we will treat the cases of odd and even electorates separately, since the induction step of Lemma 5.1 only works in steps of two. Let us start with the odd case.

THEOREM 5.3 (ODD ELECTORATES). *For $m \geqslant 4$ alternatives and odd $n \geqslant 15$, there does not exist a Condorcet extension satisfying half-way monotonicity.*

*Proof* By Lemma 5.1, we only need to handle the case with $n = 15$. Write $A = \{a, b, c, d\} \cup X$, where $X = \{x_1, \ldots, x_{m-4}\}$. Suppose there exists a half-way monotonic Condorcet extension $f$ for 15 voters. Consider the 15-voter profile $P_0$ depicted on the right. The column numbers indicate how many voters submit a given ordering; for example, there are exactly 3 voters in $P_0$ with the ordering $a \succ b \succ d \succ c \succ X$; let us abbreviate this ordering as "$abdc$", and similarly for other voters. The $X$ at the bottom of each vote should be replaced by an arbitrary ordering of the

| 1 | 3 | 3 | 4 | 2 | 2 |
|---|---|---|---|---|---|
| $a$ | $a$ | $b$ | $c$ | $d$ | $d$ |
| $b$ | $b$ | $d$ | $a$ | $c$ | $c$ |
| $c$ | $d$ | $c$ | $b$ | $a$ | $b$ |
| $d$ | $c$ | $a$ | $d$ | $b$ | $a$ |
| $X$ | $X$ | $X$ | $X$ | $X$ | $X$ |

alternatives in $X$. Our proof is by case analysis on the value of $f(P_0)$, arriving at a contradiction in each case.

Suppose first that $f(P_0) \in \{a, b\} \cup X$. Let $P_1$ be the profile after one $dcba$ voter reverses their preferences in $P_0$. By half-way monotonicity, we have $f(P_1) \in \{a, b\} \cup X$. Suppose that $f(P_1) \in \{a\} \cup X$. Let $P_2$ be the profile after two $bdca$ voters reverse their preferences in $P_1$. By half-way monotonicity, we have $f(P_2) \in \{a\} \cup X$; however $c$ is the Condorcet winner in $P_2$, contradicting Condorcet-consistency of $f$. Thus $f(P_1) = b$. Let $P_3$ be the profile obtained from $P_1$ after one $dcab$ voter and two $cabd$ voters reverse their preferences. By half-way monotonicity, we have $f(P_3) \in \{b, d\}$. However, $a$ is the Condorcet winner in $P_3$, a contradiction.

Thus $f(P_0) \in \{c, d\}$. Let $P_4$ be the profile obtained from $P_0$ by reversing an $abcd$ voter. By half-way monotonicity, $f(P_4) \in \{c, d\}$. Suppose $f(P_4) = d$. Let $P_5$ be the profile obtained from $P_4$ by reversing two $cabd$ voters; then $f(P_5) = d$. But $b$ is the Condorcet winner at $P_5$, a contradiction. Hence $f(P_4) = c$. Let $P_6$ be the profile obtained from $P_4$ by reversing three $abdc$ voters; then $f(P_6) = c$. But $d$ is the Condorcet winner at $P_6$, a contradiction. □

The proof above was obtained with the help of computers, and in particular by using SAT solvers, in a way similar to the technique described by Brandt, Geist, and Peters [2016]. In particular, our

---

[4]There will also exist other example functions that satisfy our axioms but are not C2; restricting attention to C2 functions allows our computer search approach to be tractable. We do not have an explanation for why this restriction still allows for tight bounds.
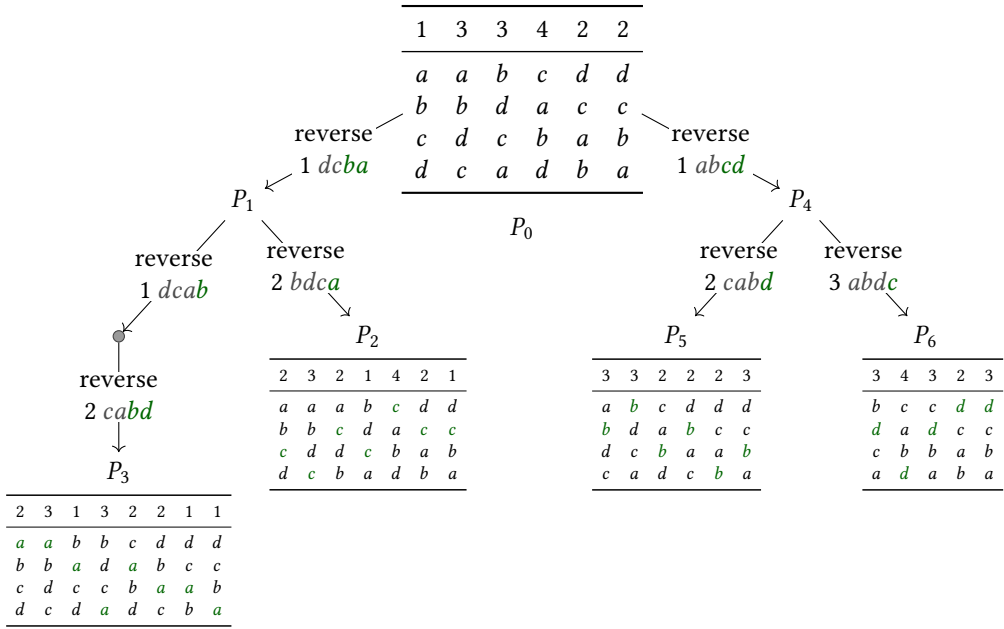
$P_0$:

| 1 | 3 | 3 | 4 | 2 | 2 |
|---|---|---|---|---|---|
| a | a | b | c | d | d |
| b | b | d | a | c | c |
| c | d | c | b | a | b |
| d | c | a | d | b | a |

reverse 1 dcba → $P_1$  reverse 1 abcd → $P_4$

$P_1$: reverse 1 dcab ; reverse 2 bdca → $P_2$

reverse 2 cabd → $P_3$

$P_2$:

| 2 | 3 | 2 | 1 | 4 | 2 | 1 |
|---|---|---|---|---|---|---|
| a | a | a | b | c | d | d |
| b | b | c | d | a | c | c |
| c | d | d | c | b | a | b |
| d | c | b | a | d | b | a |

$P_3$:

| 2 | 3 | 1 | 3 | 2 | 2 | 1 | 1 |
|---|---|---|---|---|---|---|---|
| a | a | b | b | c | d | d | d |
| b | b | a | d | a | b | c | c |
| c | d | c | c | b | a | a | b |
| d | c | d | a | d | c | b | a |

$P_4$: reverse 2 cabd → $P_5$ ; reverse 3 abdc → $P_6$

$P_5$:

| 3 | 3 | 2 | 2 | 2 | 3 |
|---|---|---|---|---|---|
| a | b | c | d | d | d |
| b | d | a | b | c | c |
| d | c | b | a | a | b |
| c | a | d | c | b | a |

$P_6$:

| 3 | 4 | 3 | 2 | 3 |
|---|---|---|---|---|
| b | c | c | d | d |
| d | a | d | c | c |
| c | b | b | a | b |
| a | d | a | b | a |

Fig. 1. Proof diagram of the proof of Theorem 5.3.

$P_0$:

| 2 | 4 | 6 | 6 | 4 | 2 |
|---|---|---|---|---|---|
| a | a | b | c | d | d |
| b | b | d | a | c | c |
| c | d | c | b | a | b |
| d | c | a | d | b | a |

reverse 2 dcba → $P_1$  reverse 2 abcd → $P_4$

$P_1$: reverse 2 dcab ; reverse 3 bdca → $P_2$

reverse 3 cabd → $P_3$

$P_2$:

| 4 | 4 | 3 | 3 | 6 | 4 |
|---|---|---|---|---|---|
| a | a | a | b | c | d |
| b | b | c | d | a | c |
| c | d | d | c | b | a |
| d | c | b | a | d | b |

$P_3$:

| 4 | 4 | 2 | 6 | 3 | 3 | 2 |
|---|---|---|---|---|---|---|
| a | a | b | b | c | d | d |
| b | b | a | d | a | b | c |
| c | d | c | c | b | a | a |
| d | c | d | a | d | c | b |

$P_4$: reverse 3 cabd → $P_5$ ; reverse 2 abdc → $P_6$

reverse 3 bdca → $P_6$

$P_5$:

| 4 | 6 | 3 | 3 | 4 | 4 |
|---|---|---|---|---|---|
| a | b | c | d | d | d |
| b | d | a | b | c | c |
| d | c | b | a | a | b |
| c | a | d | c | b | a |

$P_6$:

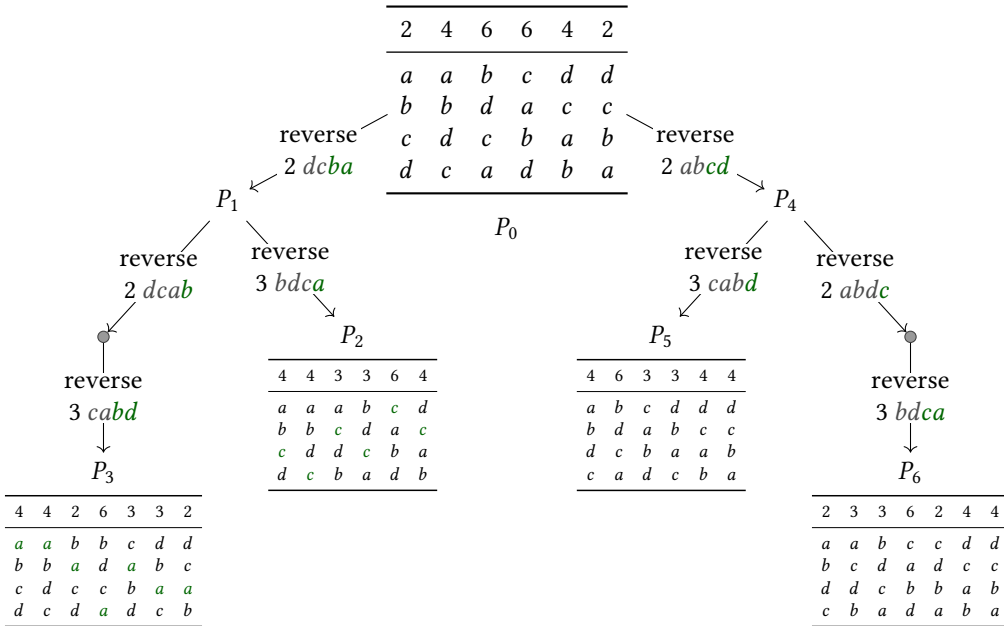| 2 | 3 | 3 | 6 | 2 | 4 | 4 |
|---|---|---|---|---|---|---|
| a | a | b | c | c | d | d |
| b | c | d | a | d | c | c |
| d | d | c | b | b | a | b |
| c | b | a | d | a | b | a |

Fig. 2. Proof diagram of the proof of Theorem 5.4.

search was aided by only considering profiles made up of the about 6–10 preference orders that appear in their proofs of the no-show paradox.

The bound on $n$ for even electorates is significantly higher than for odd ones. Intuitively, the reason is that Condorcet-consistency is less demanding in even electorates, since there are 'fewer' Condorcet winners because they need to beat every other alternative by a majority margin of at least 2.

THEOREM 5.4 (EVEN ELECTORATES). *For $m \geqslant 4$ alternatives and even $n \geqslant 24$, there does not exist a Condorcet extension satisfying half-way monotonicity.*

*Proof* By Lemma 5.1, we only need to handle the case with $n = 24$. Write $A = \{a, b, c, d\} \cup X$, where $X = \{x_1, \ldots, x_{m-4}\}$. Suppose there exists a half-way monotonic Condorcet extension $f$ for 24 voters. Consider the 24-voter profile $P_0$ depicted on the right. The column numbers indicate how many voters submit a given ordering; for example, there are exactly 4 voters in $P_0$ with the ordering $a > b > d > c > X$. The $X$ at the bottom should be replaced by an arbitrary ordering of the alternatives in $X$. Our proof is by case analysis on the value of $f(P_0)$, arriving at a contradiction in each case.

| 2 | 4 | 6 | 6 | 4 | 2 |
|---|---|---|---|---|---|
| $a$ | $a$ | $b$ | $c$ | $d$ | $d$ |
| $b$ | $b$ | $d$ | $a$ | $c$ | $c$ |
| $c$ | $d$ | $c$ | $b$ | $a$ | $b$ |
| $d$ | $c$ | $a$ | $d$ | $b$ | $a$ |
| $X$ | $X$ | $X$ | $X$ | $X$ | $X$ |

Suppose first that $f(P_0) \in \{a, b\} \cup X$. Let $P_1$ be the profile after two $dcba$ voters reverses their preferences in $P_0$. By half-way monotonicity, we have $f(P_1) \in \{a, b\} \cup X$. Suppose that $f(P_1) \in \{a\} \cup X$. Let $P_2$ be the profile after three $bdca$ voters reverse their preferences in $P_1$. By half-way monotonicity, we have $f(P_2) \in \{a\} \cup X$; however $c$ is the Condorcet winner in $P_2$, contradicting Condorcet-consistency of $f$. Thus $f(P_1) = b$. Let $P_3$ be the profile obtained from $P_1$ after two $dcab$ voter and three $cabd$ voters reverse their preferences. By half-way monotonicity, we have $f(P_3) \in \{b, d\}$. However, $a$ is the Condorcet winner in $P_3$, a contradiction.

Thus $f(P_0) \in \{c, d\}$. Let $P_4$ be the profile obtained from $P_0$ by reversing two $abcd$ voters. By half-way monotonicity, $f(P_4) \in \{c, d\}$. Suppose $f(P_4) = d$. Let $P_5$ be the profile obtained from $P_4$ by reversing three $cabd$ voters; then $f(P_5) = d$. But $b$ is the Condorcet winner at $P_5$, a contradiction. Hence $f(P_4) = c$. Let $P_6$ be the profile obtained from $P_4$ by reversing two $abdc$ and three $bdca$ voters; then $f(P_6) = c$. But $d$ is the Condorcet winner at $P_6$, a contradiction.                                                                                          □

One may wonder whether it is a coincidence that our cut-off for half-way monotonicity in even electorates ($n = 24$) is double the cut-off for participation ($n = 12$). The answer is no, as suggested by the proof of Theorem 4.1(3) of Sanver and Zwicker [2009], which (roughly) shows that half-way monotonicity for $2n$ voters implies participation for $n$ voters, at least in the presence of homogeneity and reversal cancellation. In fact, we have obtained the proof of Theorem 5.4 by taking an impossibility proof for the no-show paradox for $n = 12$, and doubling all the profiles involved in the proof.

## 6   EXTENSIONS: IRRESOLUTE RULES AND STRONG PARADOXES

### 6.1   Irresolute Voting Rules

A *set-valued* (or *irresolute*) voting rule is a function $F : A!^N \to 2^A \setminus \{\emptyset\}$ that assigns to every profile $P$ a non-empty subset $F(P) \subseteq A$ of winning alternatives. The usual interpretation is that the ties will later be broken by some other mechanism. A set-valued voting rule $F$ is a *Condorcet extension* if it uniquely selects the Condorcet winner if one exists; thus, if $x$ is the Condorcet winner of the profile $P$, then $F(P) = \{x\}$. One can define analogues of half-way monotonicity for set-valued voting rules in several ways (see Sanver and Zwicker, 2012). Here, we will focus on the approach using *set extensions*, where voters' preferences over alternatives are lifted to preferences

over *sets* of alternatives. In particular, following Jimeno et al. [2009], we focus on the *optimistic* and the *pessimistic* set extensions that are the subject of Duggan and Schwartz's [2000] impossibility theorem. An optimist prefers sets with better most-preferred alternative, while a pessimist prefers sets with better least-preferred alternative. If $X = \{a, d\}$ and $Y = \{b, c\}$, then an optimist with preferences $a > b > c > d$ would prefer $X$ to $Y$, while a pessimist with the same underlying preferences would prefer $Y$ to $X$.

Given a set $X \subseteq A$, let us write $\max_{\succcurlyeq_i} X$ (resp. $\min_{\succcurlyeq_i} X$) for the most-preferred (resp. least-preferred) alternative in $X$ according to $\succcurlyeq_i$, so that for all $x \in X$ we have $\max_{\succcurlyeq_i} X \succcurlyeq_i x \succcurlyeq_i \min_{\succcurlyeq_i} X$. This allows us to define variants of half-way monotonicity for these set extensions:

*Definition 6.1.* A set-valued voting rule $F$ satisfies *optimistic half-way monotonicity* if

$$\max_{\succcurlyeq_i} F(P_{-i}, \succ_i) \succcurlyeq_i \max_{\succcurlyeq_i} F(P_{-i}, \succ_i^{\mathrm{rev}}) \quad \text{for all profiles } P \text{ and all voters } i.$$

A set-valued voting rule $F$ satisfies *pessimistic half-way monotonicity* if

$$\min_{\succcurlyeq_i} F(P_{-i}, \succ_i) \succcurlyeq_i \min_{\succcurlyeq_i} F(P_{-i}, \succ_i^{\mathrm{rev}}) \quad \text{for all profiles } P \text{ and all voters } i.$$

One might hope that dropping the requirement of resoluteness makes it easier for Condorcet extensions to satisfy half-way monotonicity. Indeed, this happens for participation: As Brandt, Geist, and Peters [2016] show, there are Condorcet extensions satisfying optimistic participation for 16 voters and pessimistic participation for 13 voters, while the limit is 11 voters for resolute rules. For half-way monotonicity, surprisingly, it turns out that giving up resoluteness buys us nothing:

THEOREM 6.2. *For $m \geqslant 4$ alternatives and odd $n \geqslant 15$ or even $n \geqslant 24$, there does not exist a set-valued Condorcet extension satisfying either pessimistic or optimistic half-way monotonicity.*

Why does the move to the irresolute setting not allow us larger bounds on $n$, when it does for participation? An intuitive reason is suggested by the proof of Lemma 3.1, where we showed that participation implies half-way monotonicity by decomposing a preference reversal into a voter *leaving* the electorate and the reverse voter *joining* the electorate. Now, a *pessimist* reversing their preferences can be decomposed into the pessimist leaving, and an optimist with reverse preferences joining. Thus, pessimistic half-way monotonicity is related to the *conjunction* of optimistic and pessimistic participation, and neither of the latter properties alone implies pessimistic half-way monotonicity. As Brandt, Geist, and Peters [2016] find, imposing this conjunction of properties does not allow for larger bounds in the participation setting as well.

PROOF OF THEOREM 6.2. For pessimistic half-way monotonicity, we can follow the proofs of Theorems 5.3 and 5.4 almost verbatim. By way of example, let us translate the second paragraph of the proof of Theorem 5.3. Let $P_0, \ldots, P_6$ refer to the same profiles as in that proof. The end result of the proof is to conclude that $F(P_0) = \emptyset$, a contradiction. Suppose that $F(P_0)$ intersects $\{a, b\} \cup X$. Then in $P_1$, where one *dcba* voter is reversed, we also have $F(P_1)$ intersecting $\{a, b\} \cup X$, by pessimistic half-way monotonicity (since the minimum of the *dcba* voter must have gone weakly down). Suppose in fact that $F(P_1)$ intersects $\{a\} \cup X$. Then in $P_2$, after reversing two *bdca* voters, we again have that $F(P_2)$ intersects $\{a\} \cup X$, by pessimistic half-way monotonicity. But since $c$ is the Condorcet winner in $P_2$, we have $F(P_2) = \{c\}$, a contradiction. Hence $F(P_1)$ must intersect $\{b\}$, i.e., $b \in F(P_1)$. But then, similarly, $F(P_3)$ must intersect $\{b, d\}$, contradicting $F(P_3) = \{a\}$ by Condorcet-consistency. So $F(P_0)$ cannot intersect $\{a, b\} \cup X$ after all, hence must intersect $\{c, d\}$. Following the proof steps about $P_4, P_5, P_6$, we see that this also leads to contradiction.

For optimistic half-way monotonicity, we work through the proofs "the other way around", starting with the profiles $P_2, P_3, P_5, P_6$, and working our way up to conclude that $F(P_0) = \emptyset$, a

contradiction. Again let us mirror parts of the proof of Theorem 5.3 to illustrate the idea. Since $d$ is the Condorcet winner at $P_6$, we have $c \notin F(P_6)$. After reversing three $cdba$ voters, we obtain $P_4$, and must have $c \notin F(P_6)$, since the optimum according to the $cdba$ voters needs to weakly get worse, by optimistic half-way monotonicity. Similarly since $b$ is the Condorcet winner at $P_5$, we have $d \notin F(P_5)$. After reversing two $dbac$ voters, we obtain $P_4$ again, and now we see that $d \notin F(P_4)$, by optimistic half-way monotonicity. So $c, d \notin F(P_4)$. After reversing two $dcba$ voters in $P_4$, we obtain $P_0$. By optimistic half-way monotonicity, we must have $c, d \notin F(P_0)$. By similarly following the proof of Theorem 5.3 for $P_1, P_2, P_3$, we can establish that $F(P_0) \cap (\{a, b\} \cup X) = \emptyset$, which gives a contradiction. □

Of course, the resolute Condorcet extensions of Proposition 5.2 induce set-valued Condorcet extensions satisfying both optimistic and pessimistic half-way monotonicity. Hence the bounds of Theorem 6.2 are also tight.

Sanver and Zwicker [2012, Section 4.2] consider a different set extension (the Gärdenfors extension, also known as Fishburn's extension). They show that a strong version of half-way monotonicity with this set extension is incompatible with Condorcet-consistency (and our impossibilities for $n = 15$ and $n = 24$ apply to this setting as well), but certain irresolute rules like the top cycle satisfy a weak version. On the other hand, Brandt and Geist [2016, Footnote 8] show that there is no tournament solution (a voting rule depending only on the majority relation) that satisfies this weak version and is also a significant refinement of the top cycle (in the sense of refining the *uncovered set*). This result depends on an interesting observation that, for tournament solutions and for all set extensions, weak half-way monotonicity is equivalent to weak strategyproofness [Brandt and Geist, 2016, Theorem 1].

## 6.2  Strong Preference Reversal Paradoxes

A (resolute) voting rule suffers from the *strong* preference reversal paradox if a voter, by reversing their preferences, can cause their *most*-preferred alternative to win. This is a rather astounding phenomenon – for example, it could be that if $i$ truthfully votes $a > b > c > d$ the outcome would be $b$, while if $i$ submits $d > c > b > a$ then the outcome would be $a$, which $i$ now ranks last! Thankfully, it is in principle possible for Condorcet extensions to avoid this paradox: maximin is an example. But, in fact, maximin is the *only* of the commonly considered Condorcet extensions which avoids this behaviour, as we will now see. The following argument is an adaptation of Pérez [2001]. For the results in this section, we have made no effort to minimise the number of voters required.

Consider five alternatives, $A = \{x, y, z, u, t\}$, and the following 41-voter profile $P_\dagger$:

| 5 | 7 | 3 | 6 | 1 | 2 | 3 | 5 | 8 | 1 |
|---|---|---|---|---|---|---|---|---|---|
| $x$ | $x$ | $y$ | $y$ | $y$ | $y$ | $z$ | $z$ | $u$ | $t$ |
| $z$ | $t$ | $x$ | $x$ | $u$ | $t$ | $y$ | $y$ | $z$ | $y$ |
| $y$ | $u$ | $u$ | $t$ | $z$ | $u$ | $u$ | $t$ | $t$ | $z$ |
| $t$ | $z$ | $z$ | $u$ | $x$ | $x$ | $x$ | $x$ | $y$ | $u$ |
| $u$ | $y$ | $t$ | $z$ | $t$ | $z$ | $t$ | $u$ | $x$ | $x$ |

THEOREM 6.3. *If $f$ is a Condorcet extension that avoids the strong preference reversal paradox, then* $f(P_\dagger) = t$.

The punchline is that most popular Condorcet extensions do *not* choose $t$ when faced with $P_\dagger$: Black chooses $y$ (the Borda winner), the unique Kemeny ranking is $zyxtu$ with $z$ on top, Baldwin and Nanson choose $z$, Dodgson chooses $y$ (in 9 swaps, $t$ takes 15), and the uncovered set is $\{x, y, z\}$, so tie-broken versions of all the common tournament solutions (those that are refinements of the

uncovered set, like Copeland, Slater, TEQ, or the bipartisan set) will not select $t$ either. The "correct choice" of $t$ is made by maximin, as well as Young's rule, Schulze's method and Ranked Pairs. For the latter three, one can construct other examples where they, too, suffer from the strong preference reversal paradox.

PROOF OF THEOREM 6.3. Suppose, for a contradiction, that $f$ is a Condorcet extension that avoids the strong preference reversal paradox, but $f(P_\dagger) \neq t$.

- If $f(P_\dagger) = x$, then $x$ is also selected after 8 $uztyx$ voters reverse their preferences (to avoid paradox), but in the resulting profile $y$ is Condorcet winner, a contradiction.
- If $f(P_\dagger) = y$, then $y$ is also selected after 7 $xtuzy$ voters reverse their preferences, but in the resulting profile $z$ is Condorcet winner, a contradiction.
- If $f(P_\dagger) = z$, then $z$ is also selected after 6 $yxtuz$ voters reverse their preferences, but in the resulting profile $u$ is Condorcet winner, a contradiction.
- If $f(P_\dagger) = u$, then $u$ is also selected after 5 $xzytu$ voters reverse their preferences, but in the resulting profile $t$ is Condorcet winner, a contradiction.

Since each case leads to a contradiction, we conclude that $f(P_\dagger) = t$.                  □

An alternative approach to the strong preference reversal paradox can be taken by following Duddy's [2014] interpretation of the strong no-show paradox. He considers *weak* preferences (allowing indifferences). In this setting, we say that a voting rule suffers from the strong preference reversal paradox if a voter, by reversing their preferences, can cause *one of* their most-preferred alternatives to win. Duddy considered the analogous strong no-show paradox, and showed that all Condorcet extensions suffer from it. One can prove an analogue of this result for the strong preference reversal paradox by 'doubling' all voters in Duddy's proof, following the remark after the proof of Theorem 5.4.

## 7  A VERSION OF THE GIBBARD–SATTERTHWAITE THEOREM

The famous Gibbard–Satterthwaite Theorem states that every non-dictatorial voting rule that has full range must be manipulable, when there are at least 3 alternatives. The theorem is somewhat opaque in that it does not tell us *what kinds* of manipulations will be successful. Here we will combine two results to show that every non-trivial voting rule is either *needlessly* or *egregiously* manipulable.

First, let us define some axioms. Fix some electorate $N$, and let $\mathcal{D} \subseteq A!^N$ be a *domain*, i.e., a subcollection of profiles. A voting rule *on $\mathcal{D}$* is a map $f : \mathcal{D} \rightarrow A$. We say that

- $f$ is *non-imposed* (or is *onto*, or has *full range*) if for all $a \in A$, there is $P \in \mathcal{D}$ with $f(P) = a$;
- $f$ is *non-dictatorial* if there is no $i \in N$ such that $f(P) = \max_{>_i} A$ for all profiles $P \in \mathcal{D}$;
- $f$ is *unanimous* if for all $a \in A$, we have that $f(P) = a$ for all $P \in \mathcal{D}$ such that every voter $i \in N$ ranks $a$ in top position;
- $f$ is *anonymous* if $f(\sigma P) = f(P)$ for all permutations $\sigma$ of $N$;
- $f$ is *manipulable* if there exists a profile $P$, a voter $i$, and a linear order $>_i'$ such that both $(P_{-i}, >_i') \in \mathcal{D}$ and $(P_{-i}, >_i) \in \mathcal{D}$, and $f(P_{-i}, >_i') >_i f(P_{-i}, >_i)$. Thus, voter $i$ strictly prefers misrepresenting their preferences.

Note that unanimity is stronger than non-imposition, and that anonymity is stronger than non-dictatorship.

Besides the incompatibility between Condorcet-consistency and half-way monotonicity that has been the topic of this paper, the other tool we need is the Campbell–Kelly Theorem. We let $\mathcal{D}_{\text{Condorcet}}$ denote the domain of all profiles $P \in A!^N$ that admit a Condorcet winner.

THEOREM 7.1 (CAMPBELL AND KELLY, 2003, 2016). *Suppose $N$ contains an odd number of voters and $|A| \geqslant 3$. Let $f : \mathcal{D}_{Condorcet} \to A$ be an onto and non-dictatorial voting rule. Then $f$ is not manipulable if and only if $f$ is identical to the* Condorcet rule *, i.e., $f(P)$ is the Condorcet winner of $P$ for every $P \in \mathcal{D}_{Condorcet}$.*

We will be interested in an implication of the Campbell–Kelly Theorem for voting rules $f : A!^N \to A$ defined for *all* profiles.

COROLLARY 7.2. *Suppose $N$ contains an odd number of voters and $|A| \geqslant 3$. Let $f : A!^N \to A$ be a voting rule on the full domain. Suppose that $f$ is anonymous and unanimous. If $f$ is not Condorcet-consistent, then $f$ is manipulable on $\mathcal{D}_{Condorcet}$.*

PROOF. If $f$ is anonymous, then $f|_{\mathcal{D}_{Condorcet}}$ is also anonymous and thus non-dictatorial. Similarly, if $f$ is unanimous, then $f|_{\mathcal{D}_{Condorcet}}$ is non-imposed, since all profiles in which an alternative is ranked top by everyone is contained in $\mathcal{D}_{Condorcet}$. Thus, by Theorem 7.1, $f|_{\mathcal{D}_{Condorcet}}$ is manipulable. □

This implies that if $f$ is *not* a Condorcet extension, then $f$ admits a successful manipulation that occurs between two profiles that have Condorcet winners. Thus, $f$ is, in a sense, *needlessly* manipulable since it could avoid this manipulation if only it selected the Condorcet winners of those profiles.

On the other hand, we have seen in this paper that Condorcet extensions are *also* guaranteed to be manipulable (when $n$ and $m$ are large enough) by Theorem 5.3. Precisely, they all suffer from the preference reversal paradox, instances of which we might describe as *egregious* manipulation instances. Putting these two results together, we obtain a "disjunctive Gibbard–Satterthwaite theorem" that uses slightly stronger assumptions,[5] but more explicitly identifies the nature of manipulability. This disjunctive version was first proposed by Zwicker [2016, Corollary 2.8] using slightly different assumptions.

COROLLARY 7.3 (ZWICKER'S COROLLARY). *Suppose there are at least 4 alternatives and an odd number[6] of at least 15 voters. Let $f$ be an anonymous and unanimous voting rule. Either $f$ is manipulable on $\mathcal{D}_{Condorcet}$, or $f$ is manipulable by preference reversal.*

PROOF. If $f$ is a Condorcet extension then this follows from Theorem 5.3; otherwise it follows from Corollary 7.2. □

## 8 CONCLUSIONS AND FUTURE WORK

In this paper, we have undertaken a detailed study of the preference reversal paradox. We have seen that many known results about the no-show paradox transfer to our setting, but that these impossibilities require additional voters to go through. An interesting contrast appeared for set-valued rules: while imposing optimistic or pessimistic *participation* allows for stronger positive results (in terms of number of voters supported), we saw that optimistic or pessimistic half-way monotonicity is essentially as strong as its resolute version, not allowing for additional voters. Our computer-aided approach leaves some questions unresolved. In particular, our positive results are somewhat frustrating, since the voting rules produced are only available as lookup tables, and it

---

[5]Namely, higher bounds on $n$ and $m$, the requirement that $n$ is odd, anonymity instead of non-dictatorship, and unanimity instead of non-imposition. The latter two assumptions can evidently be weakened to non-dictatorship and non-imposition on $\mathcal{D}_{Condorcet}$, but this makes the statement seem less appealing.

[6]Because Theorem 7.1 only holds for odd electorates, this corollary also requires this assumption. Campbell and Kelly [2015] show that their theorem also holds for even electorates if we require the stronger axioms of anonymity and neutrality (on $\mathcal{D}_{Condorcet}$). However, on the full domain, anonymity and neutrality are usually incompatible with resoluteness. It would be interesting to find appealing conditions that allow an analogue of Corollary 7.3 to go through for even electorates as well.

would be desirable to find concise descriptions. Moreover, the approach does not extend well for $m \geqslant 5$ alternatives, since searching over profiles of larger size quickly becomes infeasible; thus, we do not know if our impossibilities are proveable with fewer voters as $m$ increases. As we have seen, half-way monotonicity is weaker than participation; it would be interesting to find natural voting rules that satisfy half-way monotonicity but fail participation. It appears that the only known natural rules satisfying half-way monotonicity are scoring rules; some other artificial rules are discussed by Moulin [1988, p. 63], Campbell and Kelly [2002], and Núñez and Sanver [2017].

Our disjunctive Gibbard–Satterthwaite theorem gives a more explicit account of the types of manipulations that are needed to prove it. A related result appears in the literature: the Gibbard–Satterthwaite Theorem holds even if we allow voters to report only preferences that are obtainable by at most one (adjacent) swap from their honest vote [Caragiannis et al., 2012, Sato, 2013]. It is plausible that there are further interesting results of this sort; gaining a better understanding of them could complement the literature on *dictatorial domains* [Aswal et al., 2003] which studies restricted preference domains that still lead to a Gibbard–Satterthwaite-style impossibility.

Given our results about set-valued voting rules in Section 6.1, one may hope to get an analogue of the Duggan–Schwartz Theorem [Duggan and Schwartz, 2000] in the style of Corollary 7.3. The Duggan-Schwartz Theorem states that any (non-trivial) set-valued voting rule that is not manipulable by optimists or pessimists must have a *nominator*, that is, a voter whose top choice is always part of the returned choice set. For set-valued Condorcet extensions (which cannot have nominators), we know from Theorem 6.2 that they cannot satisfy optimistic or pessimistic half-way monotonicity. However, the Campbell–Kelly Theorem does not admit a direct analogue to the set-valued context, because rules with nominators are also strategyproof on $\mathcal{D}_{\text{Condorcet}}$. It would be interesting to obtain a justification for Condorcet-consistency in the style of the Campbell–Kelly Theorem for set-valued voting rules, perhaps by adding additional axioms.

## ACKNOWLEDGMENTS

## REFERENCES

N. Aswal, S. Chatterji, and A. Sen. 2003. Dictatorial domains. *Economic Theory* 22, 1 (2003), 45–62.

Gilles Audemard and Laurent Simon. 2009. Predicting Learnt Clauses Quality in Modern SAT Solvers. In *Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI)*. 399–404.

A. Belov and J. Marques-Silva. 2012. MUSer2: An efficient MUS extractor. *Journal on Satisfiability, Boolean Modeling and Computation* 8 (2012), 123–128.

A. Biere. 2013. Lingeling, Plingeling and Treengeling entering the SAT competition 2013. In *Proceedings of the SAT Competition 2013*. 51–52.

F. Brandl, F. Brandt, and C. Geist. 2016. Proving the Incompatibility of Efficiency and Strategyproofness via SMT Solving. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI)*. AAAI Press, 116–122.

F. Brandl, F. Brandt, C. Geist, and J. Hofbauer. 2015. Strategic Abstention based on Preference Extensions: Positive Results and Computer-Generated Impossibilities. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI)*. AAAI Press, 18–24.

F. Brandt and C. Geist. 2016. Finding Strategyproof Social Choice Functions via SAT Solving. *Journal of Artificial Intelligence Research* 55 (2016), 565–602.

F. Brandt, C. Geist, and D. Peters. 2016. Optimal Bounds for the No-Show Paradox via SAT Solving. In *Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*. 314–322.

D. E. Campbell and J. S. Kelly. 2002. Non-monotonicity does not imply the no-show paradox. *Social Choice and Welfare* 19, 3 (2002), 513–515.

D. E. Campbell and J. S. Kelly. 2003. A strategy-proofness characterization of majority rule. *Economic Theory* 22, 3 (2003), 557–568.

D. E. Campbell and J. S. Kelly. 2015. Anonymous, neutral, and strategy-proof rules on the Condorcet domain. *Economics Letters* 128 (2015), 79–82.

D. E. Campbell and J. S. Kelly. 2016. Correction to "A Strategy-proofness Characterization of Majority Rule". *Economic Theory Bulletin* 4, 1 (2016), 121–124.

I. Caragiannis, E. Elkind, M. Szegedy, and L. Yu. 2012. Mechanism design: from partial to probabilistic verification. In *Proceedings of the 13th ACM Conference on Electronic Commerce (ACM EC)*. ACM, 266–283.

C. Duddy. 2014. Condorcet's principle and the strong no-show paradoxes. *Theory and Decision* 77, 2 (2014), 275–285.

J. Duggan and T. Schwartz. 2000. Strategic Manipulability without Resoluteness or Shared Beliefs: Gibbard-Satterthwaite Generalized. *Social Choice and Welfare* 17, 1 (2000), 85–93.

P. C. Fishburn and S. J. Brams. 1983. Paradoxes of Preferential Voting. *Mathematics Magazine* 56, 4 (1983), 207–214.

C. Geist and U. Endriss. 2011. Automated Search for Impossibility Theorems in Social Choice Theory: Ranking Sets of Objects. *Journal of Artificial Intelligence Research* 40 (2011), 143–174.

C. Geist and D. Peters. 2017. Computer-aided Methods for Social Choice Theory. In *Trends in Computational Social Choice*, U. Endriss (Ed.). Chapter 13. Forthcoming.

M. J. H. Heule, O. Kullmann, and V. W. Marek. 2016. Solving and Verifying the Boolean Pythagorean Triples Problem via Cube-and-Conquer. In *Proceedings of the 19th International Conference on Theory and Applications of Satisfiability Testing (Lecture Notes in Computer Science (LNCS))*, Vol. 9710. Springer-Verlag, 228–245.

J. L. Jimeno, J. Pérez, and E. García. 2009. An extension of the Moulin No Show Paradox for voting correspondences. *Social Choice and Welfare* 33, 3 (2009), 343–459.

M. H. Liffiton, A. Previti, A. Malik, and J. Marques-Silva. 2015. Fast, flexible MUS enumeration. *Constraints* (2015), 1–28.

H. Moulin. 1988. Condorcet's Principle implies the No Show Paradox. *Journal of Economic Theory* 45 (1988), 53–64.

M. Núñez and M. R. Sanver. 2017. Revisiting the connection between the no-show paradox and monotonicity. *Mathematical Social Sciences* (2017). Forthcoming.

J. Pérez. 2001. The Strong No Show Paradoxes are a common flaw in Condorcet voting correspondences. *Social Choice and Welfare* 18, 3 (2001), 601–616.

M. R. Sanver and W. S. Zwicker. 2009. One-way monotonicity as a form of strategy-proofness. *International Journal of Game Theory* 38, 4 (2009), 553–574.

M. R. Sanver and W. S. Zwicker. 2012. Monotonicity properties and their adaption to irresolute social choice rules. *Social Choice and Welfare* 39, 2–3 (2012), 371–398.

S. Sato. 2013. A sufficient condition for the equivalence of strategy-proofness and non-manipulability by preferences adjacent to the sincere one. *Journal of Economic Theory* 148 (2013), 259–278.

P. Tang and F. Lin. 2009. Computer-aided proofs of Arrow's and other impossibility theorems. *Artificial Intelligence* 173, 11 (2009), 1041–1053.

A. D. Taylor. 2005. *Social Choice and the Mathematics of Manipulation*. Cambridge University Press.

W. S. Zwicker. 2016. Introduction to the Theory of Voting. In *Handbook of Computational Social Choice*, F. Brandt, V. Conitzer, U. Endriss, J. Lang, and A. D. Procaccia (Eds.). Cambridge University Press, Chapter 2.