

# Bridging Social Choice and Dynamic Epistemic Logic by Modeling a Strategic Voting Experiment Using Kripke Models

Group NOT

Alexander Müller, June Wallinga, Janke de Vries, and Leon Tanis

February 11, 2026

## Abstract

Strategic voting is traditionally analyzed using static preference profiles. In this report, we propose a dynamic epistemic model of strategic voting in repeated elections. By modeling polls as public announcements and utilizing the Kendall Tau distance metric as a heuristic for manipulation, we demonstrate how agents update their knowledge of others' preferences over time. We formalize this process using S5 epistemic logic and provide a case study where a strategic voter successfully manipulates the outcome, inadvertently creating Common Knowledge of the true preference profile.

## 1 Introduction

Many collective decision-making processes, such as elections or online ratings, involve repeated voting and voters receiving feedback about the voting results. For example, in elections, voters repeatedly cast ballots for public election polls; the results of these polls are announced to them. The results of these polls reflect the ballots that the voters have cast. Given a certain result, voters gain some information about what combinations of ballots can have and cannot have been cast, and thus what the true preferences of other voters might be. But the results of polls also influence what ballots voters will cast next. Using the information provided by the poll about what the true preferences of other voters might be, voters can change their ballots in order to manipulate the results of the election in their favour. As such, polls open the door to strategic voting, also called voting manipulation.

Strategic voting thus both draws on knowledge about the true preferences of voters, but also influences voting. We try to model these dynamics of strategic voting using Dynamic Epistemic Logic (DEL) and computational social choice theory. Using DEL, we attempt to model how polls change the epistemic states of voters Meyer and Van der Hoek 2004, while we use computational social choice theory to model how voters change their ballots based on their epistemic states Van Ditmarsch, Lang, and Saffidine 2013; Kumar and Vassilvitskii 2010. We allow manipulators to cast ballots that do not correspond to their true preference in order to influence the election in their favor. Our model ultimately shows how voters gain knowledge about each others' true preferences when manipulators are involved.

The central question of this report is: How can strategic voting by a single agent be represented in a multi-agent Kripke model? To answer this question, the following subquestions are answered:

- How does a public announcement of voting polls allow for single-voter strategic voting?
- How does a subsequent public announcement of strategically manipulated voting polls change knowledge of true preferences?
- What is a possible heuristic for a strategic voter when it is selecting an optimal manipulation after such a public announcement?

The first subquestion is answered in section 2, the second and third subquestion in sections 3. Sections 4 and 5 provide a examples of how subsequent public announcements change knowledge of true preferences when strategic voting had been involved.

The upshot of this project is at least twofold. Firstly, while the field of computational social choice theory has studied voting systems and strategic voting extensively, it often views strategic voting as a one-off deviation from truthful voting. This does not fit with how strategic voting proceeds in practice. In practice, voters can reason about the preferences of others using the information they gain about the true preferences of others via polls,

and can continuously change their votes based on that. Polls often only provide “partial” knowledge of the true preferences of voters, with many different actual true preferences remaining plausible to individual voters. This partial knowledge is difficult to model using social choice theory. Our framework does capture these dynamic and epistemic aspects of strategic voting. Secondly, in general, DEL only studies how the epistemic states of agents change, but not how agents *act* given certain epistemic states. As a result, DEL is rarely used for modelling strategic voting. By combining DEL and computational social choice, we *do* attempt to model how agents act given certain epistemic states.

In general, the report is structured as follows: Section 2 introduces basic concepts and notation from voting theory that are used throughout the report, including definitions for true preference profiles, the Borda voting rule, and manipulation. This section also introduces the epistemic framework used to represent voters’ knowledge, drawing on DEL and Kripke models. Building on this framework, we formalize polls as public announcements that restrict the number of possible true preference profiles. Section 3 defines strategic voting model in which a strategic agent or “manipulator” updates their knowledge after each poll and chooses a strategic ballot based on a strategy, or so-called “strategic heuristic”. Section 4 presents an illustrative case in more detail. In Section 5, the case study is modelled using multi-agent Kripke models. Lastly, section 6 contains a discussion of the project and possible directions for future work.

## 2 Preliminaries

### 2.1 Voting Definitions

Our model resembles a classic election. The agents in our model are voters. Voters are agents that cast ballots. The alternatives which they vote about are candidates. We define voters and candidates as follows, following the standard definitions from computational social choice theory (Arrow 1951; Fishburn 1973; Van Ditmarsch, Lang, and Saffidine 2013).

**Definition: Voters and Candidates** Let  $\mathcal{A} = \{1, \dots, n\}$  be a finite set of voters and  $\mathcal{C} = \{a, b, c, \dots\}$  be a finite set of candidates.

Each voter in an election has internal preferences for which candidate they would like to win. These preferences inform the votes they cast on their ballot. However, the cast ballot may differ from the genuine preference, as is the case when a voter tries to manipulate the election. To model this distinction, we separately define true preferences and ballots:

**Definition: Linear order** A *linear order* (or *total order*) (Arrow 1951; Fishburn 1973) on a set  $\mathcal{C}$  is a binary relation  $\succ \subseteq \mathcal{C} \times \mathcal{C}$  that is

1. *complete*: for all distinct  $a, b \in \mathcal{C}$ , either  $a \succ b$  or  $b \succ a$ ;
2. *transitive*: for all  $a, b, c \in \mathcal{C}$ , if  $a \succ b$  and  $b \succ c$ , then  $a \succ c$ ;
3. *antisymmetric*: for all  $a, b \in \mathcal{C}$ , if  $a \succ b$ , then not  $b \succ a$ .

In the context of voting systems, we write  $a \succ b$  to indicate that candidate  $a$  is strictly preferred to candidate  $b$ . The word “strictly” refers to the property that no indifference is allowed between  $a$  and  $b$ . This eliminates the possibility for a tied preference between two candidates. Equivalently, on a finite set of candidates, a linear order corresponds to a ranking of *all* candidates from least preferred to most preferred, excluding the possibility for ties. Note, therefore, that this linear order denoting “preference”, can mean both the pair-wise preference relations between candidates as well as the preference ranking of all candidates.

**Definition: True preference** A *true preference* is the genuine preference of voter that a voter  $i$  has with regard to the candidates. This true preference is defined as the linear order  $\succ_i$  over  $\mathcal{C}$ .

**Example** For example, suppose the set of voters  $\mathcal{A}$  consists of  $\{1, 2, 3\}$ . Suppose the set of candidates  $\mathcal{C}$  is  $\{David, Emma, Finn\}$ . If voter 1 genuinely prefers David over Emma over Finn, then voter 1’s true preference is David over Emma over Finn. We write this as  $d \succ_1 e \succ_1 f$ .

**Definition: Ballot** A *ballot* is the linear order over all candidates  $\mathcal{C}$  corresponding to the ballot that voter  $i$  casts. We write  $a \succ_i b$  if  $i$  ranks  $a$  over  $b$  on their ballot and their ballot is identical to their true preference, while we write  $a \succ'_i b$  if  $i$  ranks  $a$  over  $b$  on their ballot and their ballot differs from their true preference (e.g. in the case of strategic voting). We use  $\succ_i$  to denote agent  $i$ 's truthful ballot and  $\succ'_i$  to denote agent  $i$ 's manipulated ballot.

**Example** Suppose on the one hand that the cast ballot of our voter 1 is *identical* to their true preference. Furthermore, their ballot ranks David over Emma over Finn. Therefore, their ballot is a truthful ballot, denoted as  $d \succ_1 e \succ_1 f$ . Suppose on the other hand that voter 1's ballot is *different* from their true preference. Suppose that their ballot ranks David over Finn over Emma. Then, her ballot is a manipulated ballot, denoted as  $d \succ'_1 f \succ'_1 e$ .

**Definition: Shorthand notation** For brevity, we mostly use string notation to denote linear orders. Thus, we denote  $a \succ_i b$ , as  $ab_i$ , and we denote  $a \succ'_i b$  as  $ab'_i$ .

**Example** Thus, coming back to our example, we write voter 1's true preference and truthful ballot  $d \succ_1 e \succ_1 f$  as  $def_1$ , while we write voter 1's manipulated ballot  $d \succ'_1 f \succ'_1 e$  as  $dfe'_1$ .

In order to execute a voting round, we must aggregate the votes and preferences of voters into profiles. Preference profiles are tuples of preferences over all candidates, one tuple for each voter. A preference profile thus represents all true preferences of all voters. Besides preference profiles, we define ballot profiles. Ballot profiles are tuples containing the ballots cast in a single voting round. A ballot profile thus represents all ballots cast in a single voting round by all voters. We distinguish between a truthful ballot profile and a manipulated ballot profile. The former does not contain manipulated ballots and is thus identical to the preference profile, while the latter does contain manipulated ballots and is thus different from the preference profile. We maintain the following definitions:

**Definition: Preference profile** A *preference profile*  $P$  is the tuple of the true preferences  $\succ_i$  of all voters  $\mathcal{A}$ , such that  $P = (\succ_1, \dots, \succ_n) \in \mathcal{L}(\mathcal{C})^n$ , where  $\mathcal{L}(\mathcal{C})$  is the set of all linear orders over  $\mathcal{C}$ .

**Example** Suppose that agent 1's true preference is  $def_1$ . Suppose that agent 2's true preference is also  $def_2$ , and agent 3's true preference is  $efd_3$ . Then, the preference profile is  $P = (def_1, def_2, efd_3)$ . For clarity,  $def_1, def_2$ , and  $efd_3$  are elements of the set of all linear orders over the set of all candidates – David, Emma and Finn – denoted as  $\mathcal{L}(\mathcal{C})^n$ . This set is the set of all possible rankings of  $n$  candidates. In the case of these three candidates:  $\mathcal{L}(\mathcal{C}) = \{def, dfe, edf, efd, fde, fed\}$ .

**Definition: Ballot profile** A *ballot profile*  $B$  (if the ballot profile is a truthful ballot profile) or  $B'$  (if the ballot profile is different from the preference profile and thus contains manipulated ballots) is the tuple of the ballots  $\succ_i$  or  $\succ'_i$  of all voters  $\mathcal{A}$ , such that  $B = (\succ_1, \dots, \succ_n) \in \mathcal{L}(\mathcal{C})^n$ , or  $B' = (\succ'_1, \dots, \succ'_n) \in \mathcal{L}(\mathcal{C})^n$ , where  $\mathcal{L}(\mathcal{C})$  is the set of all linear orders over  $\mathcal{C}$ . We denote a ballot profile as  $B$  if the ballot profile is identical to the preference profile and thus does not contain manipulated ballots. We call  $B$  a truthful ballot profile. We denote a ballot profile as  $B'$  if the ballot profile is different from the preference profile and thus contains manipulated ballots. We call  $B'$  a manipulated ballot profile.

**Example** Suppose that voter 1's, voter 2's, and voter 3's ballots are identical to their true preferences, thus  $def_1, def_2$ , and  $efd_3$ . Then, the ballot profile is a truthful ballot profile  $B = (def_1, def_2, efd_3)$ . Suppose that agent 2's and agent 3's ballots are identical to their true preferences, but voter 1's ballot different from her true preference, thus  $dfe'_1, def_2$ , and  $efd_3$ . Then, the ballot profile is a manipulated ballot profile  $B' = (dfe'_1, def_2, efd_3)$ .

We have now established the basic notation of social choice, which allows voters to have preferences over candidates, be it either truthful or manipulated. The shorthand notation makes it possible to show these linear orders in a single string. Furthermore, aggregating over preferences or ballots gives us profiles.

Now, based on a ballot profile, candidates get scores. The way in which they get a scores is via a scoring function. An example of a scoring function is a the plurality function, that assigns one point to candidates if they rank first in a ballot and zero otherwise. This is the simplest scoring function which essentially says: "Only your most preferred candidate counts".

Another example of a scoring function is the *Borda* scoring function. The Borda scoring function assigns a different number of points to candidates in a ballot, depending on their rank. To be specific: the lowest ranked candidate gets 0 points, the second to lowest candidates gets 1 point, all the way up to the highest ranked candidate, who receives  $m - 1$  points.

Borda scores contain more information about voters' preferences than the plurality function, because each candidate receives a unique number of points per ballot, depending on their rank. Later on, this will also make the polls (containing Borda scores) more informative about voters' true preferences, allowing for easier manipulation.

**Definition: Scoring function** Let  $\sigma$  be a scoring function that transforms linear orders over candidates (ballots) into a score vector containing scores for each candidate  $c \in \mathcal{C}$ .

In this project, we use the Borda scoring function as our scoring function, which we define as follows:

**Definition: Borda scoring function** Let  $\sigma$  be the Borda scoring function (Borda 1784), such that for  $m$  candidates, the Borda score of candidate  $c$  in profile  $B$ ,  $\sigma_c(B)$ , is:

$$\sigma_c(B) = \sum_{i \in \mathcal{A}} (m - \text{rank}_i(c)) \quad (1)$$

where  $\text{rank}_i(c)$  is the position of  $c$  (1st = 1, 2nd = 2, etc.) in the vote of agent  $i$ .

**Definition: Score vector** The vector  $\vec{\sigma}$  contains the Borda scores  $\sigma_c(B)$  of all candidates  $c$ , and is called the *score vector*.

**Example** Building on the previous examples, we suppose that the truthful ballot profile  $B$  for voters 1, 2 and 3 is  $B = (def_1, def_2, efd_3)$ . If we consider the Borda score of candidate David, then our candidate  $c = d$ . Note that  $m = 3$ . Then,

$$\begin{aligned} \sigma_d(B) &= \sum_{i \in \mathcal{A}} (3 - \text{rank}_i(d)) = \\ \sigma_d(B) &= (3 - \text{rank}_1(d)) + (3 - \text{rank}_2(d)) + (3 - \text{rank}_3(d)) = \end{aligned}$$

Remember that in voter 1's ballot, David ranks first, in voter 2's ballot, David ranks third, and in voter 3's ballot, David ranks second. Hence, we have that:

$$\sigma_d(B) = (3 - 1) + (3 - 1) + (3 - 3) = 4$$

Likewise, if we consider the Borda score of Emma (who is ranked second twice and ranked first once), we have:

$$\sigma_e(B) = (3 - 2) + (3 - 2) + (3 - 1) = 4$$

And if we consider the Borda score of Finn (who is ranked last twice and ranked second only once), we have:

$$\sigma_f(B) = (3 - 3) + (3 - 3) + (3 - 2) = 1$$

The score vector  $\vec{\sigma}$  contains the Borda scores of all candidates – David, Emma, and Finn. The score vector  $\vec{\sigma}$  in our example is thus  $(\sigma_d, \sigma_e, \sigma_f) = (4, 4, 1)$

In order to determine the eventual winner of the voting round, the score vector is used by an outcome function. An outcome function uses the score vector as input and produces a ranking of all candidates. This ranking can be seen as the outcome of the voting round. The simplest outcome function simply ranks the candidates on the basis of how high their score is. The candidate with the highest score is ranked first, the candidate with the second-to-highest score is ranked second, et cetera. However, there might also be situations in which the candidate with the lowest score is ranked first – this depends on what scoring function is used. Considering the scoring function that we use, our outcome function is defined as follows:

**Definition: Outcome function** Let  $\mathcal{O}(\vec{\sigma})$  be the *outcome function* that maps a score vector  $\vec{\sigma}$  to an outcome, that is, a partial order over all candidates:

Formally, for any  $x, y \in C$ , the relation induced by  $\mathcal{O}(\vec{\sigma})$  is defined as

$$\begin{cases} x \succ y & \text{if } \sigma_x > \sigma_y, \\ x \equiv y & \text{if } \sigma_x = \sigma_y, \end{cases}$$

The outcome function essentially produces a linear order, just like the preferences/ballots. Only this order is not one denoting preferences. Rather, the outcome function orders candidates on how they are ranked in the voting round, according to the scoring function. Note, also, that this is a *partial* order, meaning there may be some candidates which are not preferred over others (tie). In our case, this happens when both candidates have the same Borda score.

**Example** Remember that in our example, the score vector  $\vec{\sigma} = (4, 4, 1)$ . Then, the outcome is a partial order over all candidates – David, Emma, and Finn – where one candidate is ranked over another if their Borda score is higher than the other, and ranked equal if their Borda score is equal. Thus, the outcome is  $\mathcal{O}(\vec{\sigma}) = (d \equiv e \succ f)$ .

The last element of social choice that requires defining is manipulation. Outcomes can sometimes be manipulated by voters to their advantage if they change their ballot.

**Definition: Manipulation** We define an outcome to be *manipulated* by a single voter  $i$  when there exist two ballot profiles

$$B = (\succ_1, \dots, \succ_i, \dots, \succ_n) \text{ and } B' = (\succ_1, \dots, \succ'_i, \dots, \succ_n) \quad (2)$$

which produce two different score vectors  $\vec{\sigma}$  and  $\vec{\sigma}'$ , respectively, such that  $\mathcal{O}(\vec{\sigma}') \succ_i \mathcal{O}(\vec{\sigma})$ , meaning that the outcome of the manipulated ballot profile  $B'$  is preferred by agent  $i$  to the truthful ballot profile  $B$  (or preference profile  $P$ ) according to  $i$ 's true preference  $\succ_i$  (Grossi 2025). We then call  $\vec{\sigma}$  an unmanipulated score vector and  $\mathcal{O}(\vec{\sigma})$  an unmanipulated outcome, while we call  $\vec{\sigma}'$  a manipulated score vector and  $\mathcal{O}(\vec{\sigma}')$  a manipulated outcome.

**Example** Recall the manipulated ballot profile from the example earlier:  $B' = (dfe'_1, def_2, efd_3)$ . Skipping a few steps, we then have the manipulated score vector  $\vec{\sigma}' = (\sigma'_d = 4, \sigma'_e = 3, \sigma'_f = 2)$ . Then, the manipulated outcome is  $\mathcal{O}(\vec{\sigma}') = (d \succ' e \succ' f)$  or  $def'$ . This manipulated outcome, containing a ballot manipulated by voter 1, corresponds exactly to voter 1's true preference  $def'$ , while the unmanipulated outcome  $\mathcal{O}(\vec{\sigma})$  does not. Hence, we have a case of manipulation by voter 1.

## 2.2 Epistemic Logic Framework

We model agent knowledge using a standard S5 Kripke model  $M = \langle W, \sim, V \rangle$ .

- $W$ : a non-empty set of possible worlds. In our models, these possible worlds are preference profiles.
- $\sim_i \subseteq W \times W$ : the indistinguishability relations for each agent  $i \in \mathcal{A}$ . We assume that these are equivalence relations (that is, that they are reflexive, symmetric, and transitive), corresponding to the S5 axioms.
- $V$ : A valuation function mapping atomic propositions to subsets of  $W$ . The subset of worlds to which an atomic proposition is mapped, are the worlds in which the atomic proposition is true.

**Logical Language** To formally model the reasoning of agents about specific candidate orderings, we define our atomic propositions as pairwise preferences. Let  $a, b \in C$  and  $i \in \mathcal{A}$ . The language  $\mathcal{L}$  is defined recursively:

$$\varphi ::= a \succ_i b \mid \neg\varphi \mid \varphi \wedge \varphi \mid K_i\varphi \mid C_G\varphi \mid [\varphi]\varphi$$

These should be read as follows:  $\neg\varphi$  is the negation of  $\varphi$ .  $\varphi \wedge \varphi$  is the conjunction of  $\varphi$  and  $\varphi$ .  $K_i$  reads “agent  $i$  knows  $\varphi$ ”.  $C_G\varphi$  reads “it is common knowledge among group  $G \subseteq \mathcal{A}$  that  $\varphi$ ”. (We use  $G$  to denote a group instead of  $B$ , which is often used in literature (e.g., Van Ditmarsch, Der Hoek, and Kooi 2008) since we already use  $B$  in the Borda scoring function.)

The atom  $a \succ_i b$  represents the fact “agent  $i$  strictly prefers candidate  $a$  over  $b$ ”. We derive this atom from Van Ditmarsch, 2013. The set of atoms is defined as  $\mathcal{L}(\mathcal{C})$ , such that  $\mathcal{L}(\mathcal{C}) = \{a \succ_i b \mid a, b \in C \text{ and } i \in \mathcal{A}\}$ . As

mentioned,  $\mathcal{A} = \{1, \dots, n\}$  is the finite set of agents, while  $\mathcal{C} = \{a, b, c, \dots\}$  is the finite set of candidates. Since both  $\mathcal{A}$  and  $\mathcal{C}$  are finite sets,  $\mathcal{L}(\mathcal{C})$  is a finite set too. Hence, The set  $L_K^m(\mathcal{L}(\mathcal{C}))$  of epistemic formulas  $\varphi, \psi, \dots$  over  $\mathcal{A}$  is also a finite set.

**Definition: Further Shorthand Notation** For brevity, we mostly use string notation to denote a linear order for an agent, as described in Section 2. However, in this report, almost all focus will be on the case of three agents and three candidates. An example case could be  $abc_i$ , which denotes that agent  $i$  has the preference  $a \succ b \succ c$ . Formally, this is defined as the conjunction of the pairwise atoms:

$$abc_i \equiv (a \succ_i b) \wedge (b \succ_i c) \wedge (a \succ_i c)$$

Because the indistinguishability relations  $\succ_i$  for each agent  $i \in \mathcal{A}$  are transitive, we can write this more succinctly as:

$$abc_i \equiv (a \succ_i b) \wedge (b \succ_i c)$$

**Definition: Semantics** The interpretation of formulas in a model  $M = \langle W, \sim, V \rangle$  is defined as follows (Van Ditmarsch, Lang, and Saffidine 2013; Meyer and Van der Hoek 2004):

$$\begin{aligned} M, w \models a \succ_i b &\iff w \in V_{a \succ_i b} \\ M, w \models \neg \varphi &\iff M, w \not\models \varphi \\ M, w \models \varphi \wedge \psi &\iff M, w \models \varphi \text{ and } M, w \models \psi \\ M, w \models K_i \varphi &\iff \forall t \in W : (w \sim_i t \implies M, t \models \varphi) \\ M, w \models C_G \varphi &\iff \forall t \in W : (w \sim_G^* t \implies M, t \models \varphi) \end{aligned}$$

**Definition: Common knowledge:** The relation  $\sim_G^*$  is defined as the reflexive, transitive closure of the union of individual accessibility relations:  $\sim_G^* = (\bigcup_{i \in G} \sim_i)^*$ .

**Definition: Public Announcement Update:** The semantics of the dynamic operator are defined by a model restriction.  $M, w \models [\varphi]\psi \iff M, w \models \varphi \implies M|\varphi, w \models \psi$ . For  $[\varphi]\psi$ , read: ‘‘after every (public and truthful) announcement of  $\varphi$ , it holds that  $\psi$ ’’. The updated model  $M|\varphi = \langle W', \sim', V' \rangle$  is defined as:

- $W' = \{w \in W \mid M, w \models \varphi\}$  (The set of worlds where the announcement is true).
- $\sim'_i = \sim_i \cap (W' \times W')$  (The restriction of relations to the new domain).
- $V'(a \succ_i b) = V(a \succ_i b) \cap W'$  (The restriction of valuation).

## 2.3 Model Assumptions

Our framework relies on several strong epistemic and game-theoretic assumptions that allow for the precise logical modeling of the voting process.

1. **Truthful Initial State:** We assume that in the initial state  $M_0$ , prior to any polls, agents are introspective regarding their own true preferences  $\varphi (K_i(\varphi))$  but are ignorant of other agents’ preferences. This implies the initial model consists of all logically possible linear orders  $(\mathcal{L}(\mathcal{C})^n)$ .
2. **Truthfulness of the First Poll:** We assume the first poll is conducted on *sincere* ballots. This follows most logically from the truthful initial state in which agents only know their own true preferences and are ignorant of other agents’ preferences. In this case, it makes the most sense for their ballot to reflect their preferences. Thus, the first poll is a reflection of the agents’ true preferences.
3. **Common Knowledge of Protocol:** The scoring and outcome function, and the fact that polls are public announcements, are common knowledge among all agents.
4. **Perfect Recall & Logic (S5):** Agents are modelled as ideal reasoners with perfect memory of previous announcements.
5. **Single Manipulator:** The manipulator models all other agents as truthful reasoners. They thus assume that the ballots of the other agents in all rounds will be identical to their true preferences. This allows the strategic heuristic in Section 3.2 to work.

## 2.4 Distance Metric: Kendall Tau

We utilize the Kendall Tau Distance ( $d_\tau$ ), which serves as a metric for the dissimilarity between two rankings (Kumar and Vassilvitskii 2010). In our context, this metric will quantify how far an election outcome deviates from a voter’s true preference. This metric will then be used to decide the optimal vote  $\succ^*$  for the manipulator, defined in the Strategic Heuristic in Section 3.2. We originally thought of using the Spearman’s Footrule Distance (Diaconis and Graham 1977), but decided to use the Kendall Tau Distance as it felt as a more intuitive metric in the case where relative ordering feels more important than absolute ordering. See Appendix A for a discussion on this.

**Definition: Kendall Tau distance metric** Formally, let  $\succ_1$  and  $\succ_2$  be two linear orderings (rankings) over the set of candidates  $\mathcal{C}$ . The distance is defined as:

$$d_\tau(\succ_1, \succ_2) = \sum_{x,y \in \mathcal{C}} \bar{K}_{x,y}(\succ_1, \succ_2) \quad (3)$$

where the penalty term  $\bar{K}_{x,y}$  is defined based on the relative ordering of the pair  $\{x, y\}$ :

- $\bar{K}_{x,y} = 1$  if the pair is **discordant** (the ordering of this pair is reversed between  $\succ_1$  and  $\succ_2$ ).
- $\bar{K}_{x,y} = 0$  if the pair is **concordant** (the ordering is preserved).
- $\bar{K}_{x,y} = 0.5$  if exactly one of the rankings contains a **tie** for the pair.

**Example** To illustrate the calculation of the Kendall Tau Distance ( $d_\tau$ ), consider a scenario where an agent’s true preference is  $abc$  and the election outcome is  $acb$ . To compute the distance, we evaluate the relative ordering of all possible pairs  $\{a, b\}$ ,  $\{b, c\}$ , and  $\{c, a\}$ :

- **Pair  $\{a, b\}$ :** In the true preference,  $a \succ b$ ; in the outcome,  $a \succ b$ . The ordering is concordant, resulting in a penalty of 0.
- **Pair  $\{b, c\}$ :** In the true preference,  $b \succ c$ ; in the outcome,  $c \succ b$ . The ordering is discordant, resulting in a penalty of 1.
- **Pair  $\{c, a\}$ :** In the true preference,  $a \succ c$ ; in the outcome,  $a \succ c$ . The ordering is concordant, resulting in a penalty of 0.

Summing these individual penalties, the total Kendall Tau Metric is  $d_\tau(abc, acb) = 0 + 1 + 0 = 1$ .

## 3 The Model: Polls and Strategy

### 3.1 Polls as Public Announcements

We treat a poll as a truthful public announcement of the aggregate Borda scores  $\vec{\sigma}$ . Note that although this poll is a *truthful* public announcement, the aggregate Borda scores resulting from ballots may be manipulated. It is therefore important to note that “truthful” refers only to the faithfulness of the semantic formula which is the public announcement, although the underlying ballots may represent manipulated preference.

**The Truthfulness of Strategic Polls** A crucial distinction must be made between the *content* of the announcement and the *ballots* that generated it. While a strategic voter may cast a “deceptive” ballot (differing from their true preference), the election authority truthfully announces the resulting score. Thus, the announcement that would reflect the possibility that “Candidate A has 5 points” is factually true in the world where those ballots were cast. The epistemic update eliminates all worlds where the aggregate score could not possibly be 5, regardless of whether those scores came from sincere or strategic votes.

**Translation to Logic** Let  $\Pi$  be the set of all possible preference profiles. We define the translation  $\tau(\vec{\sigma})$  which maps a score vector to a proposition in our logical language.

The announcement formula is the disjunction of all profiles  $P \in \Pi$  that result in the score vector  $\vec{\sigma}$ . Each profile  $P$  is logically defined as the conjunction of the individual preference orderings of the agents in that profile. Using the shorthand notation defined in Section 2.2:

$$\tau(\vec{\sigma}) = \bigvee_{\{P \in \Pi \mid \sigma(P) = \vec{\sigma}\}} \left( \bigwedge_{i \in \mathcal{A}} \succ_i^P \right) \quad (4)$$

where  $\succ_i^P$  denotes the specific preference ordering (e.g.,  $abc_i$ ) held by agent  $i$  in profile  $P$ .

**Examples:** If  $\vec{\sigma} = (6, 3, 0)$ , only one profile is possible (everyone votes  $abc$ ). The announcement is:

$$\tau(6, 3, 0) = (abc_1 \wedge abc_2 \wedge abc_3)$$

If  $\vec{\sigma} = (6, 2, 1)$ , the announcement becomes a disjunction of the profiles compatible with these scores:

$$\tau(6, 2, 1) = (abc_1 \wedge abc_2 \wedge acb_3) \vee (abc_1 \wedge acb_2 \wedge abc_3) \vee (acb_1 \wedge abc_2 \wedge abc_3)$$

**Update Semantics** The public announcement of  $\tau(\vec{\sigma})$  updates the model  $M$  to  $M \mid \tau(\vec{\sigma}) = \langle W', \sim', V' \rangle$ , where  $W' = \{w \in W \mid M, w \models \tau(\vec{\sigma})\}$ . This effectively deletes all worlds incompatible with the poll results.

## 3.2 The Strategic Heuristic

In order to enhance our single epistemic agent with the capacity for strategising, we define a strategic heuristic that helps the agent to choose its optimal strategic ballot  $\succ_k^*$ . The strategic heuristic of our agent is defined as follows: A strategic agent  $k$  seeks to bring the final election outcome as close as possible to their true preference. The other agents' votes  $\{B\} \setminus \succ_k^*$  are assumed to be fixed because they can not be individually derived from the public announcement poll.

The agent simulates the election result for every possible ballot they could cast. To apply the distance metric, we must first convert the resulting Borda scores into a ranking.

Let  $\vec{\sigma}$  be the aggregated vector of Borda scores resulting from agent  $k$  casting ballot  $\succ'$  while the set of remaining votes  $\{B\} \setminus \succ'_k$  remains unchanged (denoted as  $\vec{\sigma}(\succ', \{B\} \setminus \succ'_k)$ ).

Agent  $k$  chooses the optimal strategic ballot  $\succ_k^*$  by minimizing the distance between the *induced outcome ranking*  $\mathcal{O}(\vec{\sigma}(\succ', \{B\} \setminus \succ'_k))$  and their *true preference*  $\succ_k$ :

$$\succ_k^* = \arg \min_{\succ' \in \mathcal{L}(\mathcal{C})} d_\tau \left( \succ_k, \mathcal{O}(\vec{\sigma}(\succ', \{B\} \setminus \succ'_k)) \right) \quad (5)$$

This minimization iterates over all possible linear orders  $\succ' \in \mathcal{L}(\mathcal{C})$ . The agent calculates the scores for each potential ballot, converts those scores into a ranking, and measures how “far” that ranking is from their true preference. It chooses the ballot  $\succ_k^*$  that leads to the voting outcome that is closest to its true preference  $\succ_k$ . Importantly, no epistemic knowledge is currently used when the agent calculates its next decisions. This will happen when the number of manipulators  $n \geq 2$ , as there is then uncertainty in what the other manipulator did. It is unclear yet what the strategic heuristic will look like in that scenario.

## 4 Case Study: Successful Manipulation

To demonstrate that minimizing the Kendall Tau distance effectively guides an agent toward a successful manipulation, we consider a voting scenario with 3 agents  $\mathcal{A} = \{1, 2, 3\}$  and 3 candidates  $\mathcal{C} = \{a, b, c\}$ . We use the standard Borda scoring vector for 3 candidates:  $(2, 1, 0)$ . This means that if agent  $i$  turns in a ballot of  $(a \succ b \succ c)$ , then  $a$  will get a Borda score of 2 from voter  $i$ ,  $b$  will get a score of 1 and  $c$  will get a score of 0.

### 4.1 Setup and Initial State

Let the true preference profile  $P = (\succ_1, \succ_2, \succ_3)$  of the agents be defined as follows: Agents 1 and 2 both hold the preference  $\succ_1 = (\succ_2 = a \succ b \succ c)$ , while Agent 3 holds the preference  $\succ_3 = (c \succ b \succ a)$ . Let the manipulator in this example be agent 3.

Suppose the election begins with a sincere poll. If all agents vote truthfully, the resulting Borda scores are calculated as follows according to the Borda scoring function in Eq. 1: Candidate  $a$  receives 4 points (2 from Agent 1, 2 from Agent 2, 0 from Agent 3); Candidate  $b$  receives 3 points (1 from each agent); and Candidate  $c$  receives 2 points. The outcome vector is  $\vec{\sigma} = (4, 3, 2)$ . Out of the initially possible  $(3!)^3 = 216$  profiles, only 15 possible profiles  $P$  remain. These 15 profiles are visualized in Table 1, where the columns represent the ballot for voter 1, the rows represent the ballot for voter 2. The cells in the table represent the only possible ballot for voter 3, given the ballots of voter 1 and 2 and scoring vector  $\vec{\sigma} = (4, 3, 2)$ . If a cell is left empty, this means that the combination of ballots of voter 1 and 2 is not possible, given scoring vector  $\vec{\sigma} = (4, 3, 2)$ . In our example, the upper-left cell in Table 1 represents the real world.

Voter 1 \ Voter 2	abc	acb	bac	bca	cab	cba
abc	[cba]	bca	cab	acb	bac	abc
bac	cab				abc	
acb	bca			abc		
cab	bac		abc			
bca	acb	abc				
cba	abc					

Table 1: Possible worlds consistent with score vector  $\vec{\sigma} = (4, 3, 2)$ . The cells represent the required preference for Voter 3. The bracketed entry is the true world.

## 4.2 Scenario A: Sincere Voting Outcome

Under sincere voting, the outcome ranking is  $\mathcal{O} = a \succ b \succ c$ . Agent 3 compares this outcome to their true preference  $c \succ b \succ a$ . Since the outcome is the exact inverse of their preference, every pair of candidates is discordant. Specifically, for the pairs  $\{c, b\}$ ,  $\{c, a\}$ , and  $\{b, a\}$ , the ordering in the outcome opposes the agent's desire. Consequently, the Kendall Tau distance (Eq. 3) is maximized:

$$d_\tau(\text{true preference, outcome}) = 1 + 1 + 1 = 3$$

This high distance signals to Agent 3 that the current path leads to their worst-case scenario.

## 4.3 Scenario B: Strategic Manipulation

We have taken agent 3 to be our strategic voter. Which means that this agent will now apply the strategic heuristic from Section 3.2. This heuristic returns the best alternative ballot  $\succ_3^*$  for agent 3, which minimizes the distance between the induced outcome ranking  $\mathcal{O}(\vec{\sigma}(\succ_3', \{B\} \setminus \succ_3))$  with alternative ballot  $\succ_3'$  and their true preference  $\succ_3$ . By applying Eq. 5 we find  $\succ_3^* = (b \succ c \succ a)$ . As only voter 3 is allowed to change their ballot, both voter 1 and 2 keep their initial ballot:  $\succ_1 = \succ_2 = (a \succ b \succ c)$ . We can now calculate the new Borda scores for each candidate  $c \in C$  with Eq. 1. The new scores become  $\sigma_a = 4$  (unchanged),  $\sigma_b = 4$  (increased by 1), and  $\sigma_c = 1$  (decreased by 1). The result is a tie between  $a$  and  $b$ . According to our distance metric definition in Section 2.4, a tie results in a penalty of 0.5 for the relevant pair. The induced ranking is  $b \equiv a \succ c$ .

We recalculate the distance between Agent 3's true preference ( $c \succ b \succ a$ ) and this new strategic outcome ( $b \equiv a \succ c$ ):

- The pair  $\{c, b\}$  remains discordant (True:  $c \succ b$ , Outcome:  $b \succ c$ ), adding **1.0** to the distance.
- The pair  $\{c, a\}$  remains discordant (True:  $c \succ a$ , Outcome:  $a \succ c$ ), adding **1.0** to the distance.
- The pair  $\{b, a\}$  is now a tie (True:  $b \succ a$ , Outcome:  $b \equiv a$ ), adding only **0.5** to the distance.

$$d_\tau = 1.0 + 1.0 + 0.5 = 2.5$$

By voting strategically, Agent 3 reduces the epistemic distance from 3 to 2.5. This reduction correctly identifies that the manipulated outcome (a tie between their 2nd and 3rd choice) is preferable to the sincere outcome (a strict win for their 3rd choice).

Intuitively this can also be described as a successful manipulation. Initially agent 3 observes from the poll that  $a$  is winning and  $c$  (their favourite) has no chance. However, they strictly prefer  $b$  over  $a$ . To minimize the distance to the outcome, Agent 3 calculates the effect of changing their ballot to  $b \succ c \succ a$ . This strategy effectively “boosts”  $b$  while keeping  $a$  at the bottom. The outcome is a tie of  $b$  and  $a$ , while nothing changes in the position of  $c$ . Strictly speaking, this is an improvement. The only change that occurs is that  $b$  is now tied with  $a$ , instead of ranked below  $a$ , as was the case before. This would be a preferable outcome if you prefer  $b$  over  $a$ .

## 5 Case Study: Epistemic changes

### 5.1 First epistemic situation

Three epistemic situations can be distinguished in the case study. The first epistemic situation is the situation after the agents have determined their own true preferences. A representation of this situation as a multi agent Kripke model is not included in the report, yet should be obvious: the multi agent Kripke model contains 216 worlds, and there are accessibility relations for an agent between worlds in which the true preference of that agent is identical. For instance, there are accessibility relations for agent 1 between all worlds in which the true preference of agent 1 is  $abc_1$ , between all worlds in which the true preference of agent 1 is  $cba_1$ , but not between worlds in which the true preference of agent 1 is  $abc_1$  and worlds in which it is  $cba_1$ . In this first epistemic situation, there is already knowledge: in all worlds, each agent knows its own true preference, but not that of other agents.

### 5.2 Second epistemic situation

The second epistemic situation is the situation after the public announcement of the first, unmanipulated score vector. This is thus the announcement that candidate  $a$ 's score is 4, that of  $b$  is 3 votes, and of  $c$  is 1. As mentioned, the public announcement restricts the number of possible worlds from 216 to 15. Firstly, the translation  $\tau(\vec{\sigma})$  translates the score vector  $(4, 3, 2)$  to a proposition in our logical language:

$$\tau(4, 3, 2) = (abc_1 \wedge abc_2 \wedge cba_3) \vee (acb_1 \wedge abc_2 \wedge bca_3) \vee (bac_1 \wedge abc_2 \wedge cab_3) \vee (bca_1 \wedge abc_2 \wedge acb_3) \vee (cab_1 \wedge abc_2 \wedge bac_3) \vee (cba_1 \wedge abc_2 \wedge abc_3) \vee (abc_1 \wedge acb_2 \wedge bca_3) \vee (bca_1 \wedge acb_2 \wedge abc_3) \vee (abc_1 \wedge bac_2 \wedge cab_3) \vee (cab_1 \wedge bac_2 \wedge abc_3) \vee (abc_1 \wedge bca_2 \wedge acb_3) \vee (acb_1 \wedge bca_2 \wedge abc_3) \vee (abc_1 \wedge cab_2 \wedge bac_3) \vee (bac_1 \wedge cab_2 \wedge abc_3) \vee (abc_1 \wedge cba_2 \wedge abc_3)$$

Secondly, the public announcement of  $\tau(4, 3, 2)$  updates the model  $M$  of the previous, first epistemic situation containing 216 worlds to the restricted model  $M \mid \tau(4, 3, 2)$ :

$$M \mid \tau(4, 3, 2) = \langle W', \sim', V' \rangle$$

The set  $W'$  contains all worlds of the previous, first epistemic situation in which it holds that  $M, w \models \tau(4, 3, 2)$ . These are the fifteen worlds depicted in Table 1 and Figure 1. Thus,  $W' = \{w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8, w_9, w_{10}, w_{11}, w_{12}, w_{13}, w_{14}, w_{15}\}$ . The set  $\sim'_i$  contains all relations that held between the worlds in  $W'$  that held in the previous, first epistemic situation containing 216 worlds. The set  $V'_p$  contains the valuations of all propositions for the worlds in  $W'$  that held in the previous, first epistemic situations.

Like in the first epistemic situation, there are only accessibility relations for an agent between worlds in which the true preference of that agent is identical. Figure 1 is a representation of the Kripke model of this situation. Reflexive relations are not represented in the figure; nevertheless, these should be assumed to hold in each world for each agent, since in each world, each agent considers that world possible. The accessibility relations for agent 1 are represented as dotted lines, those for agent 2 as single lines, and for agent 3 as double lines.

In the second epistemic situation, after the public announcement of the first voting poll, there are only three worlds in which an agent knows the true preferences of all agents. In world  $w_1$ , it holds that:

$$M, w_1 \models K_3(abc_1 \wedge abc_2 \wedge cba_3)$$

as it this is equivalent to:

$$\text{for all } t \in W : w_1 \sim_3 t \text{ implies } M, t \models (abc_1 \wedge abc_2 \wedge cba_3)$$

and:

all  $t$  that satisfy  $t \in W : w_1 \sim_3 t$ , are  $w_1$ , and  $M, w_1 \models (abc_1 \wedge abc_2 \wedge cba_3)$

Similarly, in world  $w_6$ , it holds that:

$$M, w_6 \models K_1(cba_1 \wedge abc_2 \wedge abc_3)$$

as this is equivalent to:

$$\text{for all } t \in W : w_6 \sim_1 t \text{ implies } M, t \models (cba_1 \wedge abc_2 \wedge abc_3)$$

and:

all  $t$  that satisfy  $t \in W : w_6 \sim_1 t$ , are  $w_6$ , and  $M, w_6 \models (cba_1 \wedge abc_2 \wedge abc_3)$

And finally, in world  $w_{15}$ , it holds that:

$$M, w_{15} \models K_2(abc_1, cba_2, abc_3)$$

for this is equivalent to:

$$\text{for all } t \in W : w_{15} \sim_2 t \text{ implies } M, t \models (abc_1 \wedge cba_2 \wedge abc_3)$$

and:

all  $t$  that satisfy  $t \in W : w_{15} \sim_2 t$ , are  $w_{15}$ , and  $M, w_{15} \models (abc_1 \wedge cba_2 \wedge abc_3)$

Thus, if an agent has voted  $cba$ , then after the epistemic update – the public announcement of the first, unmanipulated score vector – that agent knows the true preferences of all agents. Yet, after the epistemic update, there are no propositions that everybody knows, and also no propositions that are common knowledge.

### 5.3 Third epistemic situation

The third epistemic situation is the situation after the public announcement of the second, manipulated score vector, which is manipulated by agent 3. This is thus the announcement that candidate  $a$ 's score is 4, that of  $b$  is 4 votes, and of  $c$  is 1. As it will turn out, after this public announcement, all agents know what the preference profile is – that is, what the actual true preferences of all agents are.

Since we are working in with  $S5$ , and  $S5$  takes all agents to be perfect logicians, and since all agents know which agent is the strategic voter (which is a crucial assumption), all agents can argue as follows. According to the first voting poll, candidate  $a$  got 4 points,  $b$  got 3 points, and  $c$  got 2 point. After the announcement of these scores, agent 3 changed its ballot so that  $a$  got 4 points,  $b$  got 4 points, and  $c$  got 1 point. From this, agents 1 and 2 can deduce one thing about agent 3's true preference, namely that agent 3 prefers candidate  $b$  over  $a$ . Thus, the possible worlds in figure 1 get restricted to only those worlds in which agent 3 prefers candidate  $b$  over  $a$ . The rest are eliminated. The worlds that are eliminated are  $w_3, w_4, w_6, w_8, w_9, w_{10}, w_{11}, w_{12}, w_{14}$  and  $w_{15}$ . This leaves  $w_1, w_2, w_5, w_7$  and  $w_{13}$ .

In addition, agent 1 and 2 can deduce that agent 3 does not prefer candidate  $b$  over both  $a$  and  $c$  (that is, does not rank him first in his preference), since if agent 3 did, it would not have been possible for him manipulate the score of  $b$  in such a way that  $b$ 's score increases. In other words, he would not have been able to 'boost' candidate  $b$ . Thus, the possible worlds in figure 1 get further restricted to only those worlds in which agent 3 does not rank candidate  $b$  first. The rest are eliminated. The worlds that are eliminated are  $w_2, w_5, w_7$  and  $w_{13}$ . This leaves only  $w_1$  – the actual state.

Thus, after public announcement of the second, manipulated score vector, all agents know the true preferences of all agents, and know that these are  $abc_1, abc_2$  and  $cba_3$ . Everyone knows this, and it is common knowledge:

$$M, w_1 \models C_{123}(abc_1 \wedge abc_2 \wedge cba_3)$$

as it holds that:

$$\text{for all } t \in W : w_1 \sim_{123} t \text{ implies } M, t \models (abc_1 \wedge abc_2 \wedge cba_3)$$

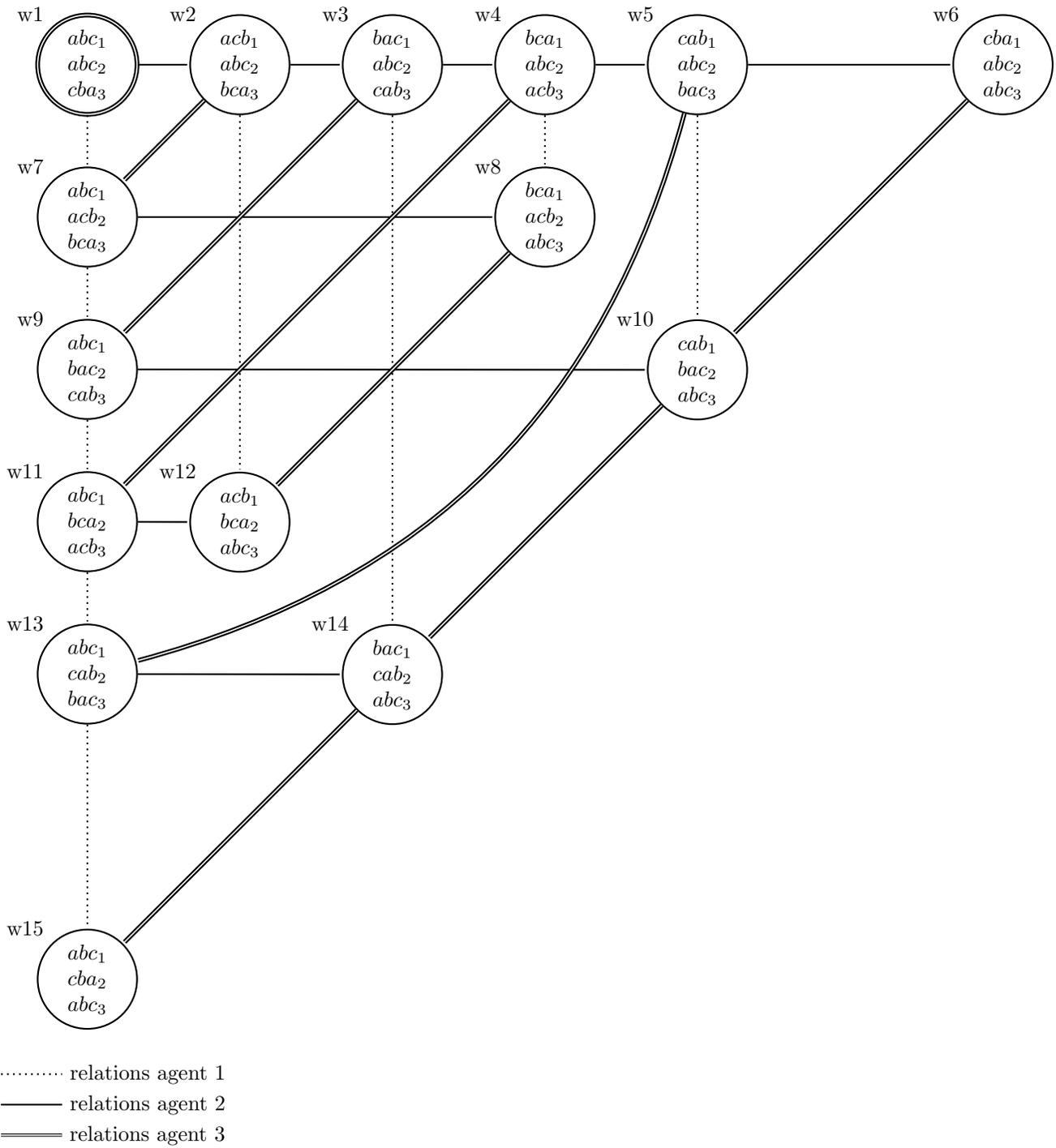
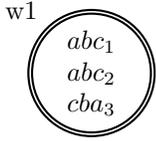


Figure 1: Kripke model of situation after the public announcement of the first voting poll



..... relations agent 1  
 ——— relations agent 2  
 = = = relations agent 3

Figure 2: Kripke model of situation after the public announcement of the second voting poll

## 6 Discussion & Future Work

Our research question, “How can strategic voting by a single agent be represented in a multi-agent Kripke model?”, was divided into three subquestions: “How does a public announcement of voting polls allow for single-voter strategic voting?”, “How do subsequent public announcements of strategically manipulated voting polls change knowledge of true preferences?” and “What is a possible heuristic for a strategic voter when it is selecting an optimal manipulation after such a public announcement?”.

Our model fully answers the first question for a situation in which there are three voters, one manipulator, and common knowledge about who the one manipulator is; it shows which voting polls allow for strategic voting. Our model also answers the second question by showing how exactly knowledge of true preferences changes with each voting poll. The third subquestion is answered by showing that minimizing Kendall Tau Distance leads to the optimal manipulation.

However, if we would want to know when subsequent public announcements of voting polls allow *human voters* specifically to vote strategically, and how subsequent public announcements change the knowledge of true preferences that *human voters* specifically have, then we would have to combine our model with a theory of mind. This would place restrictions on the order of knowledge that voters can have. We could also have to combine our model with a theory of memory. This would place limits on the number of for instance votes of other candidates or poll results in other voting rounds voters can remember and thus take into account when deciding to vote strategically.

If we would want to know when subsequent public announcements of voting polls allow for strategic voting, and how subsequent public announcements change knowledge of true preferences in actual elections, then we would have to drop the assumption that agents know which and how many agents are strategic voters, since this is often not the case in *actual elections* (in which case, for instance, it is possible for all voters to be strategic voters). Since without this assumption, voters would be less sure about whether their manipulation will change the outcome of the elections, we then would have to give the voters a strategy, for instance, being risk-averse or prone to taking risks.

Furthermore, if we would want to model actual past elections, then we would have to switch to action models, since in actual elections, voters often change their preferences. Our logic does not allow for such changes.

A striking result of our case study is the emergence of Common Knowledge regarding the true preference profile ( $C_{\mathcal{A}}P$ ) through the very act of manipulation. While the strategic voter successfully shifts the outcome  $\mathcal{O}(\vec{\sigma})$  to a more desirable  $\mathcal{O}(\vec{\sigma}')$ , the transition between these score vectors provides a traceable signal that truthful agents can use to invert the strategic heuristic. By observing the specific delta in Borda scores, Agents 1 and 2 are able to prune their indistinguishability relations until only the true world  $w_1$  remains.

This leads to a notable epistemic asymmetry: the act of manipulation is informative only for the victims, not the perpetrator. In our model, Agent 3 attains full knowledge of the profile  $P$  immediately after the first sincere poll. The subsequent manipulated poll serves only to actually achieve the desired outcome, providing Agent 3 with zero additional bits of information about the other voters. Conversely, this second poll is the exact mechanism that resolves the remaining uncertainty for the truthful agents. We are left with an ironic result: by successfully manipulating the system to their advantage, the manipulator becomes the only agent who hasn’t learned anything.

## References

Arrow, Kenneth Joseph (1951). *Social Choice and Individual Values*. New York, NY, USA: Wiley: New York.

- Borda, Jean-Charles de (1784). “Mémoire sur les Élections au scrutin.”S. 657–665 in”. In: *Mémoires de L’Académie Royale des Sciences, Année 1781*.
- Diaconis, Persi and Ronald L Graham (1977). “Spearman’s footrule as a measure of disarray”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 39.2, pp. 262–268.
- Fishburn, Peter C. (1973). *The Theory of Social Choice*. Princeton, NJ: Princeton University Press.
- Grossi, Davide (2025). *Lecture Notes on Algorithmic Voting Theory*. Version 3, Lecture notes for the course Computational Social Choice. URL: <http://www.davidegrossi.me>.
- Kumar, Ravi and Sergei Vassilvitskii (2010). “Generalized distances between rankings”. In: *Proceedings of the 19th international conference on World wide web*, pp. 571–580.
- Meyer, J.J.Ch. and W. Van der Hoek (2004). *Epistemic logic for AI and computer science*. 41. Cambridge University Press.
- Van Ditmarsch, Hans, Wiebe van Der Hoek, and Barteld Kooi (2008). *Dynamic epistemic logic*. Springer.
- Van Ditmarsch, Hans, Jérôme Lang, and Abdallah Saffidine (2013). “Strategic voting and the logic of knowledge”. In: *arXiv preprint arXiv:1310.6436*.

# A Appendix: Comparison with Spearman’s Footrule Distance

In this report, we utilized the Kendall Tau distance metric to define the strategic heuristic. To justify this choice, we analyze why a different metric, the Spearman’s Footrule Distance, fails to capture the strategic incentives in our case study (Section 4).

## A.1 Definition

Spearman’s Footrule Distance measures the dissimilarity between two rankings by summing the absolute differences in the rank positions of each candidate. Let  $\text{rank}(c, \succ)$  denote the position of candidate  $c \in \mathcal{C}$  in ranking  $\succ$  (where 1st place = 1, 2nd = 2, etc.). The distance between true preference  $\succ_{true}$  and outcome  $\succ_{out}$  is:

$$D_S(\succ_{true}, \succ_{out}) = \sum_{c \in \mathcal{C}} |\text{rank}(c, \succ_{true}) - \text{rank}(c, \succ_{out})| \tag{6}$$

## A.2 Application to Case Study

Recall the setup from Section 4:

- **Agent 3’s True Preference:**  $c \succ b \succ a$ .
- **Rankings:**  $\text{rank}(c) = 1, \text{rank}(b) = 2, \text{rank}(a) = 3$ .

**Scenario A: Sincere Voting Outcome** The sincere outcome was  $a \succ b \succ c$ .

- Ranks:  $\text{rank}(a) = 1, \text{rank}(b) = 2, \text{rank}(c) = 3$ .

We calculate the Spearman distance:

$$\begin{aligned} D_S &= |\text{rank}(c)_{true} - \text{rank}(c)_{out}| + |\text{rank}(b)_{true} - \text{rank}(b)_{out}| + |\text{rank}(a)_{true} - \text{rank}(a)_{out}| \\ &= |1 - 3| + |2 - 2| + |3 - 1| \\ &= 2 + 0 + 2 = 4 \end{aligned}$$

**Scenario B: Manipulated Outcome** The manipulated outcome resulted in a tie between  $a$  and  $b$ , with  $c$  last:  $a \equiv b \succ c$ . Thus,  $a$  and  $b$  share ranks 1 and 2 (with an average of 1.5), and  $c$  is rank 3.

- Ranks:  $\text{rank}(a) = 1.5, \text{rank}(b) = 1.5, \text{rank}(c) = 3$ .

We calculate the new Spearman distance:

$$\begin{aligned} D'_S &= |1 - 3| + |2 - 1.5| + |3 - 1.5| \\ &= 2 + 0.5 + 1.5 = 4 \end{aligned}$$

Using Spearman’s Footrule,  $D_S = 4$  and  $D'_S = 4$ . The distance remains identical. Because the metric relies on absolute position, the penalty for  $c$  being last ( $|1 - 3| = 2$ ) dominates the metric in both cases. The subtle improvement of  $b$  tying with  $a$  is arithmetically cancelled out by  $a$  dropping slightly. Therefore, an agent minimizing Spearman’s distance would fail to manipulate in this scenario, even though intuitively this feels like an improvement. For this report, this felt like it validates our choice of Kendall Tau, which does identify the pairwise improvement between  $a$  and  $b$ .