



**HAL**  
open science

# Homo moralis goes to the voting booth: coordination and information aggregation

Ingela Alger, Jean-François Laslier

► **To cite this version:**

Ingela Alger, Jean-François Laslier. Homo moralis goes to the voting booth: coordination and information aggregation. *Journal of Theoretical Politics*, 2022, 34 (2), pp.280-312. 10.1177/09516298221081811 . hal-03682814

**HAL Id: hal-03682814**

**<https://hal.science/hal-03682814>**

Submitted on 5 Jan 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

October 2021

“Homo moralis goes to the voting booth:  
coordination and information aggregation”

Ingela Alger and Jean-François Laslier

# Homo moralis goes to the voting booth: coordination and information aggregation\*

Ingela Alger<sup>†</sup>

Jean-François Laslier<sup>‡</sup>

today: October 24, 2021

## Abstract

This paper revisits two classical problems in the theory of voting—viz. the divided majority problem and the strategic revelation of information—in the light of evolutionarily founded partial Kantian morality. It is shown that, compared to electorates consisting of purely self-interested voters, such Kantian morality helps voters solve coordination problems and improves the information aggregation properties of equilibria, even for modest levels of morality.

**Keywords:** voting, *Homo moralis*, Kantian morality, social dilemmas

## 1 Introduction

The question of individual cooperation is a puzzle for social theories because cooperation should be sustained when efficient for the group but might be in contradiction with efficiency at the individual level. This puzzle appears under various disguises in different disciplines: Evolutionary Biology (Nowak and Sigmund 2005 [40]), Ethology (de Waal 1996 [50]), Economics (Moulin 1995 [36]), Political theory (Ostrom 1998 [41]) or Social Philosophy (Binmore

---

\*I.A. acknowledges funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 789111 - ERC EvolvingEconomics), as well as IAST funding from the French National Research Agency (ANR) under grant ANR-17-EURE-0010 (Investissements d’Avenir program). We are grateful for comments from Jörgen Weibull as well as audiences at the EEA 2021 Meetings, the “Political Economy: Theory Meets Empirics” workshop, Aix-Marseille School of Economics, Paris School of Economics, Toulouse School of Economics, and University of Bilkent. Olivier Lision provided excellent research assistance.

<sup>†</sup>Toulouse School of Economics, CNRS, University of Toulouse Capitole, and Institute for Advanced Study in Toulouse, France. [ingela.alger@tse-fr.eu](mailto:ingela.alger@tse-fr.eu)

<sup>‡</sup>Paris School of Economics (CNRS)

1994 [9]). In the light of recent results from the literature on the evolutionary foundations of human motivation, we use formal game theory to revisit two classic cooperation dilemmas faced by voters whose ability to communicate with each other is limited: the *divided majority problem* and the *Condorcet jury theorem*.

**The divided majority problem.** In elections with at least three candidates, it is sometimes in the interest of supporters of two of the candidates to coordinate their votes on one of them in order to block the other candidates (Cox 1997 [14]). Instrumentally motivated voters can be induced to vote strategically if their vote stands a chance of being pivotal. In large enough electorates, however, pivotality becomes irrelevant. Hence it may be necessary to resort to other explanations for why and how voters achieve such coordination.

**The Condorcet jury theorem.** In situations where individuals receive private and informative signals about the true state of nature, and where the preferences of the voters are aligned—e.g., a jury that wishes to convict a person only if she is guilty—efficient *information aggregation* is typically achieved if voters vote blindly according to the signal they receive (Condorcet, 1785 [13]). However, as is known since Austen-Smith and Banks (1996 [6]), a sophisticated voter would realize that, given that the others vote according to the signal they receive, she should condition her vote on the event that she is pivotal. In some settings this may lead her to conclude that it is better to vote against the signal she observes. This surprising result casts doubt on the efficiency of majority (and super-majority) rules as a procedure to aggregate information.

We contribute to both literatures by modeling voting behaviors based on preferences that entail a form of universalization. Specifically, we adopt the view that voters have *Homo moralis* preferences, which have been shown to be favored by evolution by natural selection (Alger and Weibull, 2013 [2]). *Homo moralis* reasons as follows: when contemplating a course of action, she evaluates what her material payoff would be if—hypothetically—each other individual of the population she belongs to were to follow the same course of action with probability  $\kappa \in [0, 1]$ . Although one might dispute whether this behavior captures the whole significance of Immanuel Kant’s morality, it incorporates a key ingredient of this construct (and, arguably, of most moral theories): the “universalization” principle (Kant 1785 [27], Roemer 2019 [43]). One obtains the standard materialistic *Homo oeconomicus* for  $\kappa = 0$  and the Kantian model of Laffont (1975 [31]) for  $\kappa = 1$ . Values of  $\kappa$  between 0 and 1 trigger partial Kantian universalization, and the parameter  $\kappa$  can be interpreted as a level of morality. Hence, *Homo moralis* preferences reconcile the theory of ethical voting with that of purely instrumental voters: it encompasses the purely instrumental motive as a special case, and spans a continuum of degrees of partial universalization, up to and including full

universalization. So the model also stands on its own feet as a model of ethical behavior, independently of its evolutionary foundations. In practice, moral judgments and actions are often guided by the universalization principle (Levine et al. 2020 [33]) and *Homo moralis* can be read as modeling a particular “ethical” behavior based on partial universalization. Moreover, there is extensive evidence that moral motivations are important for most voters (Blais 2000 [10]).

Theories with ethically motivated voters are not new. In the literature the most common formalization of an ethical voter comes in the form of the rule utilitarian (Harsanyi 1980, 1992 [25, 26]) who selects a voting strategy which, if chosen by *all* other rule-utilitarian voters, would maximize their aggregate material payoff. The partial Kantian morality captured by *Homo moralis* preferences is a less demanding ethical concept than rule utilitarianism in two respects. First, the *Homo moralis* player does not take into account the payoffs of the other players, only her own. Second, *Homo moralis* preferences induces the voter to ponder what the outcome would be, if —hypothetically— some fraction (not all) of the other voters selected the same strategy as her. It will be seen that even small values of the universalization parameter  $\kappa$  sometimes lead to results that differ significantly from the ones obtained with instrumentally motivated voters.

Our formalization of ethical voters should not be confused with group-based voting models (Coate and Conlin 2004 [12], Feddersen and Sandroni 2006 [20]) where, by assumption, strategic decisions are made at the collective level: in these models an ethical voter chooses a strategy based on the anticipation that other ethical voters will effectively also choose the same strategy. Along this line, two recent contributions have examined one of the problems studied in this paper (the divided majority problem). Li and Pique (2020 [34]) adopt the view that some voters are rule utilitarians, who select a voting strategy which, when chosen by all voters in the divided majority, maximizes their utility. Bouton and Ogden (2021 [11]) assume that voting strategies are taken at the level of the group, so that each supporter of a particular candidate acts in the interest of the group of like-minded supporters (it is as if they applied were rule utilitarians on behalf of the group). By contrast, in our model decisions are individually decided.

Since Jean-Jacques Rousseau (1755 [44]), many scholars have expressed the idea that political psychology should not be cut from its possible biological roots: see for instance Shubert (1982 [45]), Petersen (2015 [42]), Sidanus and Kurzban (2013 [47]) or Bergner and Hatemi (2017 [8]). Within this stream of research, *Homo moralis* is a theoretical model that does not link to empirical genetics but to the pure theory of evolution and stability of interactive behavior. As in Economics (Lesourne et al. 2006 [32]) the evolutionary approach

complements the now-standard but still criticized theory of rational choice (Downs 1957 [16], Green and Shapiro, 1994 [24], Stephenson et al. 2018 [48]), by refining the concept of Nash equilibrium (which does not contain by itself notions of stability or convergence<sup>1</sup>). In this work, we use a distinction between (material) payoff and (subjective) utility, which is justified by evolutionary considerations (Alger and Weibull 2019 [4]). We therefore alter the objective function. Still, we do compute the (Nash) equilibria of the games played with the altered objective. We will here not attempt a full evolutionary study of these games, but we will make the essential distinction between strict equilibria, where best responses are well defined and which are known to be robust and stable under most dynamics, and flat equilibria, in which all strategies are indifferent and that lack stability and robustness.

The paper is organized as follows. In section 2 we formally describe the *Homo moralis* model. Then we consider two classical problems in the theory of voting: the divided majority problem under plurality rule in section 3 and the question of strategic revelation of information in the Condorcet jury setting in section 4. For each of these questions we study the implications of the hypothesis of evolutionary Kantian morality. The last section is a short conclusion and proofs are in the Appendix.

## 2 Who is *Homo moralis*?

The evolutionary argument that implies precisely the behavior termed “*Homo moralis*” rests on the ideas that at least part of the fitness an individual achieves depends on the material payoff she achieves in social interactions, that her subjective utility (whose maximization drives individual behavior) is transmitted to her (biological or cultural) offspring, and that when a mutant utility function appears in the population its carriers are more exposed to interaction with other mutants than are non-mutants (because interactions are local). From these premises, Alger and Weibull 2013 [2] showed that the toolkit of evolutionary game theory (Maynard Smith 1982 [35]) can be used to prove that evolutionarily stable preferences are of the *Homo moralis* type.

Here we provide a formal definition of *Homo moralis* preferences. We start by the simple case of a two-player ( $n = 2$ ) symmetric normal form game. Let  $X$  denote the set of pure strategies and  $\pi(x, y)$  the material payoff for a player playing  $x$  in her interaction with a player playing  $y$ . Then a *Homo moralis* with degree of morality  $\kappa \in [0, 1]$  achieves the

---

<sup>1</sup>See Van Damme 1997 [15].

following *utility* from using strategy  $x$  when the opponent uses strategy  $y$ :

$$U(x, y) = (1 - \kappa) \cdot \pi(x, y) + \kappa \cdot \pi(x, x). \quad (1)$$

The first term is the individual's material payoff, given the strategies effectively used. The second term captures the Kantian moral concern: it induces the individual to ponder what her material payoff would be if, hypothetically, the other individual were to use the same strategy as her. A *Homo moralis* with degree of morality  $\kappa$  thus chooses a strategy that maximizes the weighted sum of her own material payoffs computed at the actual strategy profile and at the hypothetical universalized one, the weight attached to the latter term being  $\kappa$ . In his *Grundlegung zür Metaphysik der Sitten* (1785), Immanuel Kant wrote "Act only according to that maxim whereby you can at the same time will that it should become a universal law." In this vein, *Homo moralis* can be said to "act according to that maxim whereby you can at the same time will that others should do likewise with some probability."

The probability interpretation is particularly compatible with the logic whereby *Homo moralis* preferences have been shown to have a strong evolutionary foundation. The argument is as follows. Consider a large population where each individual inherits her preferences from an individual in the preceding generation (be it culturally or biologically), and in each generation individuals are matched at random to interact in pairs according to the material game described above. *Homo moralis* preferences are justified based on the observation that in essentially all populations, new cultural variants or genetic mutations spread locally. From this fact it follows that, even if the probability of two similar mutants being matched is very small, a mutant is still relatively more likely than non-mutants to be matched with a mutant. Taking into consideration this phenomenon, Alger and Weibull ([2]) show that, for a value of  $\kappa$  that precisely equals the probability that mutants are matched when mutants are vanishingly rare, individuals who maximize utility of the *Homo moralis* form have an evolutionary advantage over those who would behave differently, when preferences are passed on from one generation to the next and the number of individuals to whom an individual passes her preferences is determined by the material payoff she obtains.

We turn now to the more general case of an  $n$ -player interaction ( $n \geq 2$ ). The material payoff of a player  $i$  depends on her own strategy  $x_i \in X$  and on the strategies  $y_1, \dots, y_{n-1}$  used by the other players. Writing  $(y_1, \dots, y_{n-1}) \equiv \mathbf{y}_{-i}$ , the material payoff is denoted  $\pi(x_i, \mathbf{y}_{-i})$ . Let the symbol  $\mathbb{E}$  denote mathematical expectation.

**Definition 1** *In a symmetric  $n$ -player game  $\pi$ , an individual is a **Homo moralis** with degree of morality  $\kappa$  if her utility function  $U$  satisfies  $U(x, \mathbf{y}) \equiv \mathbb{E}[\pi(x, \tilde{\mathbf{y}})]$  where  $\tilde{\mathbf{y}} =$*

$(\tilde{y}_1, \dots, \tilde{y}_{n-1})$  is a random strategy profile for the other players, with each component  $\tilde{y}_i$  being the actual strategy used by opponent  $i$  ( $y_i$ ) with probability  $1 - \kappa$  and the individual's own strategy ( $x$ ) with probability  $\kappa$ .

For instance, in an interaction between three individuals the utility of a *Homo moralis* with degree of morality  $\kappa \in [0, 1]$  from using strategy  $x$  when the others use strategies  $y_1$  and  $y_2$  is:

$$U(x, y_1, y_2) = (1 - \kappa)^2 \cdot \pi(x, y_1, y_2) + \kappa(1 - \kappa) \cdot [\pi(x, x, y_2) + \pi(x, y_1, x)] + \kappa^2 \cdot \pi(x, x, x).$$

The two voting problems which will be studied in this article are cases of many-player interactions (the population of voters) who have access to the same strategies (the possible ballots to cast). A key feature of voting games is that they are *aggregative*, in the sense that (a) any individual's payoff depends only on how he/she votes and the vector of voting strategies played by the other individuals, and (b) the individual's payoff would not be affected if other individuals swapped their strategies. For instance the outcome of the vote only depends on the total number of votes obtained by each candidate. In the study of equilibrium, it will be sufficient for our purposes to state the utility that a *Homo moralis* with degree of morality  $\kappa$  achieves from playing strategy  $x$  when all the others use the same strategy, say  $y$ . In an aggregative game, this simplifies the writing of the general  $\kappa$ -moral utility function specified in Definition 1 to the following expression:

$$U(x, \mathbf{y}^{(n-1)}) = \sum_{m=1}^n \binom{n-1}{m-1} \kappa^{m-1} (1 - \kappa)^{n-m} \pi(x, \mathbf{x}^{(m-1)}, \mathbf{y}^{(n-m)}), \quad (2)$$

where  $\mathbf{y}^{(\ell)}$  is the  $\ell$ -dimensional vector whose components all equal  $y$  and  $\mathbf{x}^{(\ell)}$  the  $\ell$ -dimensional vector whose components all equal  $x$ . When all the other individuals use strategy  $y$ , the random strategy profile  $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_{n-1})$  in the general definition is such that each component other than the first one (which is the individual's own strategy) is a random variable that follows a binomial distribution, taking the value  $y$  with probability  $1 - \kappa$  and the value  $x$  with probability  $\kappa$ . Since the game is aggregative, one needs only keep track of the number of times that exactly  $m$  out of the  $n - 1$  components in  $\tilde{\mathbf{y}}$  take the value  $x$ .

In one of the models below we will consider infinitely large populations, modeled as a continuum. To study equilibria in this setting we will take the utility of *Homo moralis* who plays strategy  $x$  in a population where all others play strategy  $y$  to be the material payoff should a share  $\kappa$  of the population (hypothetically) play  $x$  instead of  $y$ . We rely on the de Moivre-Laplace theorem to argue that this is a good approximation of the expression in (2)



when  $n$  tends to infinity. Indeed, this theorem says that the probability mass function of the random number of times that  $y$  is replaced by  $x$  in  $n$  independent trials converges to the probability density function of the normal distribution with mean  $n\kappa$  and standard deviation  $\sqrt{n\kappa(1-\kappa)}$  as  $n \rightarrow \infty$ .

Throughout the paper we use the Nash equilibrium concept: the originality of the *Homo moralis* model is to introduce a distinction between material payoff and utility, but it is amenable to application of the standard Nash notion. A Nash equilibrium strategy profile is a vector of strategies such that each player uses a strategy that maximizes her utility, given the strategies used by the others. We will say that a Nash equilibrium is *strict* if each player would obtain a strictly lower utility by any deviation from her strategy, *partially strict* if at least some players would obtain a strictly lower utility by some deviation from their strategy, and *flat* if all players are indifferent between deviating or not.

### 3 Does *Homo moralis* vote strategically in the divided majority model?

We first tackle the question of strategic voting in the divided majority setting. This is a coordination game between players that form a majority but lose the election if they split their votes among two candidates of their camp.

#### 3.1 The divided majority model

Consider an infinite population (a continuum of mass one) of voters who are to elect one candidate. There are two political parties,  $A$  and  $B$ , and three candidates: one from party  $B$ , whom we simply call  $B$  and two from party  $A$ , whom we call  $A1$  and  $A2$ . Some electors of party  $A$  rank candidate  $A1$  over  $A2$  while others rank  $A2$  over  $A1$ . Letting  $n_{A1}$ ,  $n_{A2}$ ,  $n_B$  denote the respective shares of supporters, we adopt the following assumption:

$$0 < n_{A2} < n_{A1} < n_B < n_{A1} + n_{A2} < 1. \quad (3)$$

This implies

$$1/3 < n_B < 1/2 \quad (4)$$

and  $n_B$  can take any value within these bounds. In other words, candidate  $B$  is supported by a minority of size  $n_B < 1/2$ , while a majority of size  $n_{A1} + n_{A2} > 1/2$  would be better

	$A1$	$A2$	$B$
$A1$	$1 + \varepsilon$	$1 - \varepsilon$	$0$
$A2$	$1 - \varepsilon$	$1 + \varepsilon$	$0$
$B$	$0$	$0$	$1$

Figure 1: Each cell shows the material payoff that a supporter of the candidate in the first column obtains when the candidate in the top row wins.

off with a candidate from party  $A$ . However, the majority is divided into two groups, each of them smaller than the minority. We are therefore studying the situation where, if voters vote sincerely,  $B$  is elected.<sup>2</sup>

The table in Figure 3.1 shows the (material) payoff that a voter gets depending on which candidate is elected, where  $\varepsilon \in (0, 1)$  is a parameter that measures the disagreement between  $A1$ - and  $A2$ -supporters.

We examine whether voters may be expected to vote sincerely—i.e., for the candidate they would like to win given the material payoffs—or strategically—i.e., for another candidate. We concentrate on supporters of party  $A$ , who face a coordination problem, and assume that  $B$ -supporters vote for candidate  $B$ . Under plurality voting, if  $A$ -supporters vote sincerely, the minority wins ( $B$  is elected); however, the majority can win by coordinating their votes on either  $A1$  or  $A2$ , a coordination that requires some voters to vote strategically. Myerson and co-authors [21, 37, 38] use this game to show that Approval Voting can help solve this dilemma between strategic and sincere voting, and Myerson and Weibull (2015 [39]) take a similar game as a case-study in their theory of coordination. Here we will show how Kantian morality in the form of *Homo moralis* preferences can help “solve” the dilemma under plurality voting in the sense that sincere voting is sometimes not an equilibrium, and coordination sometimes is a strict equilibrium. In some cases the resolution is partial because multiplicity of equilibria remains, and in some cases a high enough degree of morality leads to selecting the best equilibrium, a full resolution of the dilemma.

Recall that *Homo moralis* can be said to “act according to that maxim whereby you can at the same time will that others should do likewise with some probability.” Defining *Homo moralis* preferences precisely thus requires defining who the “others” are. We distinguish

---

<sup>2</sup>This contrasts with research on the divided majority problem under “aggregate uncertainty” (Li and Pique 2020 [34], Bouton and Ogden 2021 [11]) that make the assumption that if voters vote sincerely all of  $A1, A2$ , and  $B$  can be elected.

between two scenarios, the *ex post* and the *ex ante* one. In the *ex ante* scenario the voter does not know yet if her preferred candidate is  $A1$  or  $A2$ , and the reference population is the whole population of  $A$ -voters; in the *ex post* scenario, the piece of information is known and the reference population is the group of  $A1$ -supporters for an  $A1$ -supporter, and the group of  $A2$ -supporters for an  $A2$ -supporter. The model that would take as the reference group the whole population, including  $B$ -supporters would be interesting, too, but is out of the scope of this study, which concentrates on the divided majority dilemma.

In all cases, in this section we consider a large population (a continuum of voters) so that the result of the election is deterministic, defined by the fractions of the population that vote for each candidate. Our goal being to characterize symmetric equilibria, we rely on the approximation of the utility in (2) described in Section 2 for infinitely large populations. Here the material payoff of a voter is  $1 + \varepsilon$ ,  $1 - \varepsilon$ , or  $0$ , depending on whether it is her preferred  $A$ -candidate, the other  $A$ -candidate, or candidate  $B$  who wins. Hence, for a given strategy played by all others in her reference population, a *Homo moralis* voter with degree of morality  $\kappa$  evaluates each strategy by pondering what her payoff would be if, hypothetically, a share  $\kappa$  of the voters in the reference population would also use this strategy.

### 3.2 The divided majority: the *ex post* scenario

In the *ex post* scenario, all the  $A$ -supporters first learn their payoffs from candidates  $A1$  and  $A2$ , and if they have *Homo moralis* preferences with some positive degree of morality  $\kappa$  they use the group of voters who have the same payoffs as reference group. Thus, for an  $A1$ -voter (resp.  $A2$ -voter), the reference group has  $n_{A1}$  (resp.  $n_{A2}$ ) voters.

We first examine whether sincere voting is an equilibrium. By sincere voting we mean the situation in which all voters vote for their preferred candidate. Sincere voting leads to the election of  $B$ . For  $\kappa = 0$ , this is a flat equilibrium, because each voter only considers how her action impacts her material payoff, the material payoff only depends on who is elected, and, with a continuum of voters none of them is pivotal. Turning now to *Homo moralis* preferences with a positive degree of morality  $\kappa$ , this conclusion does not necessarily hold, as shown in the following proposition.<sup>3</sup>

**Proposition 1 (Ex post sincere voting)** *Suppose that all the voters are Homo moralis*

---

<sup>3</sup>For expositional simplicity, throughout we disregard any knife-edge case where  $\kappa$  exactly equals the threshold value at hand. Clearly, the set of equilibria for such knife-edge cases would depend on the assumption we would then have to make about tie-breaking rules, and this is not related to our argument.

with degree of morality  $\kappa \in [0, 1]$ . Let  $\kappa^* = \frac{n_B - n_{A1}}{n_{A2}}$ . Then  $0 < \kappa^* < 1$  and in the *ex post* scenario sincere voting (A1-supporters vote for A1 and A2-supporters for A2) is:

- a flat Nash equilibrium if  $\kappa < \kappa^*$ ;
- not a Nash equilibrium if  $\kappa > \kappa^*$ .

When the degree of morality is high enough, sincere voting is not sustainable because either some or all voters from the divided majority then prefer to vote strategically. Although such a deviation has no effect on the actual outcome, it brings satisfaction to a sufficiently moral *Homo moralis* to know that candidate *B* would be beaten, should a share  $\kappa$  of the voters in her reference group also vote strategically. Because A1-voters are more numerous than A2-voters, for any given degree of morality the deviation to strategic voting by an A2-voter is more effective than a deviation to strategic voting by an A1-voter, because *B*'s advantage is the smallest when A1-voters vote sincerely. The threshold value  $\kappa^*$  is the degree of morality above which A2-voters strictly prefer to vote strategically vote for A1 rather than voting for their preferred candidate A2, given that A1-voters cast their ballots for candidate A1.

A strong enough Kantian moral concern further makes coordination of the divided majority sustainable as a strict Nash equilibrium, as shown next. Notice that, in this *ex post* scenario, strict coordination occurs either on A1 or on A2, with the same threshold value for  $\kappa$ .

**Proposition 2 (Ex post coordination)** *Suppose that all the voters are Homo moralis with degree of morality  $\kappa \in [0, 1]$ . Let  $\kappa_1^{**} = \frac{n_A - n_B}{n_{A1}}$  and  $\kappa_2^{**} = \frac{n_A - n_B}{n_{A2}}$ . Then  $0 < \kappa_1^{**} < \kappa_2^{**} < 1$  and, in the *ex post* scenario, coordination on candidate either A1 and A2, is:*

- a flat Nash equilibrium if  $\kappa < \kappa_1^{**}$ ;
- a partially strict Nash equilibrium if  $\kappa_1^{**} < \kappa < \kappa_2^{**}$ ;
- a strict Nash equilibrium if  $\kappa > \kappa_2^{**}$ .

With a continuum population of voters, each individual voter has no real effect on the election outcome and derives no utility from expressing their opinion. Intuition thus suggests that each voter might as well vote on a random candidate. Remarkably, this is not true under *Homo moralis* preferences. As shown in the proposition, sufficiently pronounced Kantian moral concerns break the indifference and induce a strict preference to vote for one of the majority candidates, should all other majority supporters do the same. *Homo moralis* preferences thus allow the divided majority to sustain coordination on one of the *A*-candidates

and win the election against  $B$ . However, they they do not indicate on which candidate such coordination should occur; mechanisms outside of the model, e.g., focal points, would be necessary to solve this issue.

### 3.3 The divided majority: the *ex ante* scenario

Here we consider the *ex ante* situation, where each  $A$ -voter knows that she prefers party  $A$  over party  $B$ , that there are three candidates  $B, A1, A2$  with respective supports  $n_{A1}$ ,  $n_{A2}$ , and  $n_B$  as well as the associated payoffs, but does not yet know if she will rank  $A1$  above  $A2$  or *vice versa*. Unlike in the *ex post* setting, the reference population for an  $A$ -voter is the whole population of  $A$ -voters. This is not an artificial situation: it models, for instance, the reasoning of a citizen who wonders whether it is better to vote according to her ideas (voting sincerely, or naively for the candidate she will happen to prefer) or to (strategically) coordinate her vote with voters of the same camp. This is not an artificial situation: it models, for instance, the reasoning of a citizen who wonders whether it is better to vote (sincerely) for the candidate that would give her the highest material payoff or to (strategically) coordinate her vote with voters of the same camp.

For a supporter of party  $A$ , a (behavior) strategy specifies the candidate to vote for, depending on her ranking of  $A1$  and  $A2$ . We will write such a strategy  $\alpha = (\alpha_1, \alpha_2)$ ; for instance

$$\alpha^{\text{si}} = (A1, A2)$$

is the “sincere” strategy and

$$\alpha^{(1)} = (A1, A1)$$

is the strategy which dictates to vote  $A1$  in any case. In what follows we will not consider the reversed strategy  $(\alpha_1, \alpha_2)$ .<sup>4</sup> Hence, there are three strategies:  $\alpha^{(1)}$ ,  $\alpha^{(2)}$  and  $\alpha^{\text{si}}$ .

Let  $C^*(\alpha)$  denote the candidate that wins given that all  $A$ -voters use strategy  $\alpha$ . Then, the following expression shows the expected material payoff of a  $A$ -voter when all other  $A$ -voters use  $\alpha$ .

$$\pi_A(\alpha) = \begin{cases} 1 + \varepsilon \cdot (n_{A1} - n_{A2}) / (n_{A1} + n_{A2}) & \text{if } C^*(\alpha) = A1 \\ 1 - \varepsilon \cdot (n_{A1} - n_{A2}) / (n_{A1} + n_{A2}) & \text{if } C^*(\alpha) = A2 \\ 0 & \text{if } C^*(\alpha) = B. \end{cases} \quad (5)$$

---

<sup>4</sup>The reader will easily check that adding this possibility does not alter the results.

In words, from behind the veil of ignorance as to her ranking of candidates  $A1$  and  $A2$ , an  $A$ -supporter gets a material payoff that is sometimes slightly above 1 and sometimes slightly below 1 if an  $A$ -candidate wins. Since she is more likely to rank  $A1$  over  $A2$  (i.e.,  $n_{A1} > n_{A2}$ ) she would prefer  $A1$  to win from an *ex ante* perspective. Such a voter is, moreover, certain to get material payoff 0 if  $B$  wins. Note that since there is a continuum of voters the winning candidate does not depend on the strategy used by the individual voter at hand.

Using  $C^{(\kappa)}(\alpha', \alpha)$  to denote the candidate that would win if—hypothetically—a share  $\kappa$  of the other  $A$ -voters used the same strategy as her ( $\alpha'$ ) instead of the strategy that they are actually using ( $\alpha$ ), the *utility* of a *Homo moralis* with degree of morality  $\kappa$  is:

$$U_A^{(\kappa)}(\alpha', \alpha) = \begin{cases} 1 + \varepsilon \cdot (n_{A1} - n_{A2}) / (n_{A1} + n_{A2}) & \text{if } C^{(\kappa)}(\alpha', \alpha) = A1 \\ 1 - \varepsilon \cdot (n_{A1} - n_{A2}) / (n_{A1} + n_{A2}) & \text{if } C^{(\kappa)}(\alpha', \alpha) = A2 \\ 0 & \text{if } C^{(\kappa)}(\alpha', \alpha) = B. \end{cases} \quad (6)$$

A *Homo moralis* would derive satisfaction from knowing that a candidate that gives her a high expected material payoff would win, in the hypothetical scenario that a share  $\kappa$  of the other  $A$ -supporters follow her lead. Note, however, that this satisfaction is independent of whether there actually are other voters with *Homo moralis* preferences, who would reason in a similar fashion.<sup>5</sup>

As a benchmark, consider first briefly the (standard) case in which voters care only about their material payoffs (i.e., they have *Homo moralis* preferences with  $\kappa = 0$ ). Since each individual vote has no impact on the outcome, any strategy profile—both sincere and strategic voting—is a Nash equilibrium. Indeed, given that all other voters play a certain strategy, any given voter is indifferent between all the available strategies. Turning next to *Homo moralis* preferences, we will see how these eliminate certain equilibria when the degree of morality is pronounced enough.

Turning now to *Homo moralis* preferences, we first show that sincere voting is not necessarily a Nash equilibrium.

**Proposition 3 (Ex ante sincere voting)** *Suppose that all the voters are Homo moralis with degree of morality  $\kappa \in [0, 1]$ . For the same  $\kappa^* = \frac{n_B - n_{A1}}{n_{A2}}$  as in the ex post scenario (see Proposition 1), in the ex ante scenario, sincere voting (the strategy profile  $\alpha^{\text{si}} = (A1, A2)$ )*

---

<sup>5</sup>This contrasts with group-based voting models considered by Coate and Conlin 2004 [12], Feddersen and Sandroni 2006 [20], Li and Pique (2020 [34]), and Bouton and Ogden (2021 [11]), in which each ethical voter selects a strategy based on the anticipation that other ethical voters will effectively also use the same strategy.

is:

- a flat Nash equilibrium if  $\kappa < \kappa^*$ ;
- not a Nash equilibrium if  $\kappa > \kappa^*$ .

As the reader can see, with respect to sincere voting the result in the *ex ante* case is strictly identical to the one in the *ex post* case (see Proposition 1). The reason is clear: like in the *ex post* scenario  $B$ 's advantage is the smallest when  $A1$ -voters vote sincerely; hence, the threshold value  $\kappa^*$  is (again) the degree of morality above which the  $A$ -voters can make  $A1$  win by way of voting strategically for  $A1$  independent of their ranking over candidates  $A1$  and  $A2$ .

Secondly, we show that strategic voting whereby  $A$ -voters coordinate on candidate  $A1$  is sustainable in a strict sense as a Nash equilibrium when  $A$ -voters have *Homo moralis* preferences with a sufficiently high degree of morality.

**Proposition 4 (Ex ante coordination on  $A1$ )** *Suppose that all the voters are Homo moralis with degree of morality  $\kappa \in [0, 1]$ . For the same  $\kappa_2^{**} = \frac{n_A - n_B}{n_{A2}}$  as in Proposition 2, in the ex ante scenario coordination on candidate  $A1$  (the strategy profile  $\alpha^1 = (A1, A1)$ ) is:*

- a flat Nash equilibrium if  $\kappa < \kappa_2^{**}$ ;
- a strict Nash equilibrium if  $\kappa > \kappa_2^{**}$ .

Comparing this proposition to Proposition 2, we see that coordination on  $A1$ , the strongest  $A$ -candidate, obtains for the same degrees of morality in the two scenarios. Such is not the case for coordination on  $A2$ , the weak  $A$ -candidate, as shown next.

**Proposition 5 (Ex ante coordination on  $A2$ )** *Suppose that all the voters are Homo moralis with degree of morality  $\kappa \in [0, 1]$ . Let  $\kappa^{***} = \frac{n_A - n_B}{n_A}$  and  $\kappa^{****} = \frac{n_B}{n_A}$ . Then  $0 < \kappa^{***} < 1/2 < \kappa^{****} < 1$  and, in the ex ante scenario, coordination on the candidate  $A2$  (the strategy profile  $\alpha^2 = (A2, A2)$ ) is:*

- a flat Nash equilibrium if  $\kappa < \kappa^{***}$ ;
- a strict or partially strict Nash equilibrium if  $\kappa^{***} < \kappa < \kappa^{****}$  and  $\kappa < \kappa_2^{**}$ ;
- not a Nash equilibrium if  $\kappa > \kappa^{****}$  or  $\kappa > \kappa_2^{**}$ .

It is more challenging to obtain coordination on  $A2$  than on  $A1$ , because an  $A$ -supporter is more likely to end up being an  $A1$ - than an  $A2$ -supporter. Hence, an *ex ante* commitment to vote for  $A2$  entails a sacrifice in terms of expected utility, which is not sustainable for degrees of morality so large that a deviation to a commitment to vote for  $A1$  instead entails a hypothetical victory for candidate  $A1$ .

The last two propositions show how Kantian morality in the form of *Homo moralis* preferences sometimes leads to a full resolution of the divided majority dilemma in the *ex ante* scenario, in the sense that a high enough degree of morality implies that the best equilibrium is selected: indeed, for degrees of morality  $\kappa$  above  $\max\{\kappa^*, \kappa_2^{**}, \kappa^{****}\}$  there exists a unique Nash equilibrium, in which  $A$ -supporters coordinate on the stronger candidate  $A1$ .

### 3.4 Concluding remarks on the divided majority problem

Summing up the insights generated by the analysis above, for low degrees of morality  $\kappa$  any strategy profile is a non-strict equilibrium, following the “ocean of voters” logic: what I do personally does not matter and the same remains true if I have only a small number of (hypothetical) followers. Hence, for low degrees of morality no theoretical prediction arises. By contrast, for  $\kappa$  large enough, Kantian morality solves the coordination problem in the divided majority setting: sincere voting (which means non-coordination) is no longer an equilibrium, but coordination becomes a strict equilibrium. In the *ex post* scenario, coordination is sustained in the same way on any one of the two  $A$ -candidates, while in the *ex ante* scenario, coordination on the strongest of the two  $A$ -candidates is more readily obtained, and for some values of  $\kappa$  it is the only equilibrium.

Notice that in the analysis above we assumed symmetric material payoffs, in the sense that both  $A1$ - and  $A2$ -supporters gain  $2\epsilon$  from seeing their preferred rather than their second preferred candidate win (see Table 3.1). The proofs of the results derived under this assumption, however, reveal that they are robust to variations in the values of the payoffs, in the sense that the conditions on  $\kappa$  stated in the propositions are *necessary* to sustain the relevant equilibria, whether the voters *wish* to sustain them or not. Indeed, these conditions depend solely on the sizes of the groups supporting the three candidates,  $n_{A1}$ ,  $n_{A2}$  and  $n_B$ . Varying the payoffs (for instance by having four different values of  $\epsilon$  in the table) will, however, affect the utilities attached to each strategy, and hence, voters’ *wish* to vote sincerely or strategically.

In a similar fashion, our results are robust to the introduction of an intrinsic value attached to expressive voting; the conditions for a change in the result of the election still



only depend on  $\kappa$  and on  $n_{A1}$ ,  $n_{A2}$  and  $n_B$ . Changing the payoff may only affect the fact that these deviations are desirable or not.

Having thus examined the coordination problem, we turn to the information aggregation problem.

## 4 Should Homo Moralis sit in the jury ?

### 4.1 The jury model

Consider a group or jury of  $n = 2m + 1$  members. These persons have to take a binary decision, say 0 or 1, according to a simple majority rule. There are two states of Nature, also labeled 0 and 1. All jurors agree that the right decision in each state  $\omega \in \{0, 1\}$  equals  $\omega$ , but they do not know the state of Nature. For the sake of simplicity we suppose that the material payoff of a juror is 1 if the decision is correct and 0 if not. Each juror's expected material payoff is thus the probability of a correct decision.<sup>6</sup>

The jurors share the common prior belief that the state is  $\omega$  with probability  $\mu_\omega \in (0, 1)$ , where  $\mu_0 + \mu_1 = 1$ . Each jury member  $i$  also receives a private "signal"  $s_i \in \{0, 1\}$ , a random variable that is positively correlated with  $\omega$ :

$$\begin{cases} \Pr [s_i = 0 \mid \omega = 0] = p_0 = 1 - q_0 > 1/2 \\ \Pr [s_i = 1 \mid \omega = 1] = p_1 = 1 - q_1 > 1/2. \end{cases}$$

A player's pure strategy specifies her vote as a function of her signal  $s_i$ . The set of strategies  $X$  thus consists of the following four strategies:

$$\xi^{\text{inf}} : \begin{cases} 0 \mapsto 0 \\ 1 \mapsto 1 \end{cases} ; \quad \xi^0 : \begin{cases} 0 \mapsto 0 \\ 1 \mapsto 0 \end{cases} ; \quad \xi^1 : \begin{cases} 0 \mapsto 1 \\ 1 \mapsto 1 \end{cases} ; \quad \xi^{\text{inv}} : \begin{cases} 0 \mapsto 1 \\ 1 \mapsto 0 \end{cases}$$

The classical Condorcet jury theorem concludes that the majority decision is informatively efficient in a large jury. This conclusion is reached under the assumption that all jurors truthfully report their signals, that is they all use the informative strategy  $\xi^{\text{inf}}$ . However, such a strategy profile is not necessarily a (Bayesian) Nash equilibrium. As a juror, if I believe that all the other jurors vote informatively, correct Bayesian reasoning makes me

---

<sup>6</sup>Notice incidentally that this assumption has a straightforward interpretation in term of fitness as survival probability: If the group decision is correct, all members will survive with probability one, if not they all die. Think of honeybees that "vote" to choose where to locate their new hive.

condition my vote on the event that the other jurors' vote are in a tie and makes my vote decisive. This is true even in a large jury, where the condition for the (improbable) event that the other votes are exactly in a tie becomes more informative than my own signal, which I should therefore neglect. This surprising result casts doubt on the efficiency of majority (and super-majority) rules as a procedure to aggregate information. It initiated an important literature, dealing with political elections, criminal juries, or board decisions (Feddersen and Pesendorfer, 1997, 1998 [18, 19], Gerardi and Yariv 2008 [22], Gersbach and Hann 2008 [23]).

Because the voter is conditioning on an event of low probability, the non-equilibrium conclusion is not robust to small variations of the model (Laslier and Weibull 2013 [30]). Moreover, careful analysis shows that the collective inefficiency resulting from individual rational Bayesian behavior is not that strong (Koriyama and Szentes 2009 [28]). Still, these game-theoretical analyses fall short of a justification of informative voting behavior in the jury setting. In this section we compare a jury whose jurors have *Homo moralis* preferences with a jury whose jurors have *Homo oeconomicus* preferences. We ask whether informative voting is a Bayesian Nash equilibrium for a larger set of parameter constellations in the former than in the latter jury. We will further examine whether this set coincides with that in which informative voting by all jury members is efficient in the sense that it maximizes the probability that a correct decision is made.

## 4.2 A jury with three members

For a single decision-maker (a “jury of  $n = 1$  member”) who holds a very dissymmetrical prior and/or receives signals of very low quality, it may well be rational to always vote for the most probable state  $\omega$ , without taking into account the received signal. By Bayes' law:

$$Pr[\omega = 0 | s_i = 0] = \frac{p_0 \mu_0}{p_0 \mu_0 + (1 - p_1) \mu_1} \text{ and } Pr[\omega = 1 | s_i = 0] = \frac{(1 - p_1) \mu_1}{p_0 \mu_0 + (1 - p_1) \mu_1},$$

so that deciding for  $\omega = 0$  upon receiving signal 0 is optimal if and only if  $p_0 \mu_0$  is larger than  $(1 - p_1) \mu_1$ . When this condition and the symmetric one for state 1 are met, informative voting is efficient for the single decision-maker. This happens for moderate values of the prior odds ratio:

$$\frac{1 - p_1}{p_0} < \frac{\mu_0}{\mu_1} < \frac{p_1}{1 - p_0}. \tag{7}$$

Likewise, in a jury with three members, the symmetric strategy profile according to which all jury members vote informatively,  $\boldsymbol{\xi}^{\text{inf}} \equiv (\xi^{\text{inf}}, \xi^{\text{inf}}, \xi^{\text{inf}})$ , yields a higher probability of a correct decision than the symmetric strategy profiles according to which all jury members

vote 0 or all vote 1, if and only if

$$\frac{(1-p_1)^2}{p_0^2} \cdot \frac{1+2p_1}{3-2p_0} < \frac{\mu_0}{\mu_1} < \frac{p_1^2}{(1-p_0)^2} \cdot \frac{3-2p_1}{1+2p_0}. \quad (8)$$

These inequalities are obtained by comparing the probability of a correct decision under the three alternative symmetric strategy profiles, which is equal to  $\mu_0$  if all jurors play  $\xi^0$ , to  $\mu_1$  if all jurors play  $\xi^1$ , and

$$\mu_0[p_0^3 + 3p_0^2(1-p_0)] + \mu_1[p_1^3 + 3p_1^2(1-p_1)]$$

if all jurors play  $\xi^{\text{inf}}$ .

Having dealt with the normative properties of informative voting, we turn to the following positive question: under what condition is  $\xi^{\text{inf}}$  a Bayesian Nash equilibrium in a jury with three jurors? For comparison, and to prepare the ground for analysis of a jury composed of individuals with *Homo moralis*, we examine this question under the assumption that all the jurors are *Homo oeconomicus*, i.e., their utility coincides with their material payoff.

As is known since Austen-Smith and Banks (1996 [6]), to check whether  $\xi^{\text{inf}}$  is a best response for a *Homo oeconomicus* individual  $i$  to  $(\xi^{\text{inf}}, \xi^{\text{inf}})$ , it is necessary and sufficient to examine how a deviation to another strategy would affect  $i$ 's material utility in states of the world where  $i$  is pivotal, since these are the only states where the deviation would affect the outcome of the vote. Since the other two jury members—call them  $j$  and  $k$ —play  $\xi^{\text{inf}}$ , it is thus sufficient to compare the expected costs and expected benefits of deviating in states where  $j$  and  $k$  received different signals,  $s_j \neq s_k$ . In this subsection we will without loss of generality assume that  $\omega = 0$  is the least likely state, i.e., that  $\mu_0 < 1/2$ .

Consider first a deviation by  $i$  from  $\xi^{\text{inf}}$  to  $\xi^1$ . Such a deviation alters the vote outcome from 0 to 1 if  $s_j \neq s_k$  and  $s_i = 0$ . The change in the vote outcome raises the material utility if the state of Nature is  $\omega = 1$  but lowers it if the state of Nature is  $\omega = 0$ . Hence,  $i$  strictly prefers not to deviate to  $\xi^1$  if and only if the probability that  $s_j \neq s_k$ ,  $s_i = 0$ , and  $\omega = 0$ , exceeds the probability that  $s_j \neq s_k$ ,  $s_i = 0$ , and  $\omega = 1$ :<sup>7</sup>

$$\mu_0 \cdot p_0 \cdot 2p_0(1-p_0) > \mu_1 \cdot (1-p_1) \cdot 2p_1(1-p_1).$$

Likewise, a deviation by  $i$  from  $\xi^{\text{inf}}$  to  $\xi^0$  alters the outcome from 1 to 0 if  $s_j \neq s_k$  and  $s_i = 1$ ,

---

<sup>7</sup>Like Austen-Smith and Banks (1996 [6]) we disregard parameter constellations where jurors are indifferent between strategies and thus focus on strict inequalities.

and this raises the material utility if the state of Nature is  $\omega = 0$  but lowers it if the state of Nature is  $\omega = 1$ . Hence, the condition for the deviation to  $\xi^0$  to be unviable boils down to:

$$\mu_1 \cdot p_1 \cdot 2p_1(1 - p_1) > \mu_0 \cdot (1 - p_0) \cdot 2p_0(1 - p_0).$$

In sum,  $\xi^{\text{inf}}$  is a strict Bayesian Nash equilibrium in a jury composed of three *Homo oeconomicus*, if and only if

$$\frac{(1 - p_1)^2}{p_0^2} \cdot \frac{p_1}{1 - p_0} < \frac{\mu_0}{\mu_1} < \frac{p_1^2}{(1 - p_0)^2} \cdot \frac{1 - p_1}{p_0}. \quad (9)$$

In view of this equation, we see that the individual rationality of informative voting in a group of three is neither implied by rationality of informative voting for a single individual (compare with equation (7)) nor by its efficiency for the group (compare with equation (8)).

We turn now to a jury composed of jurors with *Homo moralis* preferences with some degree of morality  $\kappa \in [0, 1]$ . The reasoning is similar to the one above, in the sense that it is necessary and sufficient to consider situations in which the deviation affects the outcome of the vote. Thus, like above, a juror who ponders a certain deviation evaluates her expected material payoff if the deviation affects the outcome of the vote because she is pivotal. However, in addition, a *Homo moralis* also evaluates how the deviating strategy would affect her expected material payoff if, hypothetically, with probability  $\kappa$  each other juror were also to play this strategy instead of the one she is actually playing. We will say that jury member  $i$  is  $\kappa$ -pivotal if a deviation by  $i$  from strategy  $\xi$  to some strategy  $\xi'$  would affect the outcome of the vote in the hypothetical scenarios envisaged by *Homo moralis* in which at least one of the other jury members also play  $\xi'$ .

Like above, consider again juror  $i$ , and assume that  $j$  and  $k$  play  $\xi^{\text{inf}}$ . Table 1 lists the signal combinations  $(s_i, s_j, s_k)$  for which  $i$  is either pivotal or  $\kappa$ -pivotal. The first three columns show the signals received by the three jury members. The fourth column shows the vote outcome if  $i$  plays  $\xi^{\text{inf}}$ . The last four columns display the outcome of the vote if  $i$  were to deviate to  $\xi^1$ , and each of these columns corresponds to a different scenario that *Homo moralis* envisages. Thus, the column labeled (a) is the outcome when  $i$  deviates to  $\xi^1$  while  $j$  and  $k$  play  $\xi^{\text{inf}}$ . In columns (b) (resp. (c)),  $i$  ponders what the outcome would be if, hypothetically,  $j$  but not  $k$  (resp.  $k$  but not  $j$ ) played  $\xi^1$  instead of  $\xi^{\text{inf}}$ , a scenario to which a *Homo moralis* with degree of morality  $\kappa$  attaches weight  $\kappa(1 - \kappa)$ . Finally, in the last column  $i$  ponders what the outcome would be if, hypothetically, both of the other jurors played  $\xi^1$  instead of  $\xi^{\text{inf}}$ , a scenario to which a *Homo moralis* with degree of morality  $\kappa$  attaches weight  $\kappa^2$ . The signal realizations  $(s_i, s_j, s_k)$  not listed in this table are irrelevant,

because for all of them the outcome of the vote is 1, whether or not  $i$  deviates to  $\xi^1$ .

Table 1: Signal realizations  $(s_i, s_j, s_k)$  for which a deviation by  $i$  from  $\xi^{\text{inf}}$  to  $\xi^1$  would alter the vote outcome, if: (a)  $j$  and  $k$  play  $\xi^{\text{inf}}$ , (b) hypothetically,  $j$  were also to play  $\xi^1$  instead of  $\xi^{\text{inf}}$ , (c) hypothetically,  $k$  were also to play  $\xi^1$  instead of  $\xi^{\text{inf}}$ , (d) hypothetically, both  $j$  and  $k$  were also to play  $\xi^1$  instead of  $\xi^{\text{inf}}$

$s_i$	$s_j$	$s_k$	$(\xi^{\text{inf}}, \xi^{\text{inf}}, \xi^{\text{inf}})$	(a) $(\xi^1, \xi^{\text{inf}}, \xi^{\text{inf}})$	(b) $(\xi^1, \xi^1, \xi^{\text{inf}})$	(c) $(\xi^1, \xi^{\text{inf}}, \xi^1)$	(d) $(\xi^1, \xi^1, \xi^1)$
0	0	0	0	0	1	1	1
0	0	1	0	1	1	1	1
0	1	0	0	1	1	1	1
1	0	0	0	0	1	1	1

Since a switch in the vote outcome from 0 to 1 is costly in state  $\omega = 0$  and beneficial in state  $\omega = 1$ , a juror  $i$  with *Homo moralis* preferences and degree of morality  $\kappa$  strictly prefers strategy  $\xi^{\text{inf}}$  to  $\xi^1$  if and only if the probability that  $i$  would either be pivotal or  $\kappa$ -pivotal when the state is  $\omega = 0$  exceeds the probability of the same event when the state is  $\omega = 1$ . Counting in Table 1 the pivotal cells with their associated probabilities in states  $\omega = 0$  and 1, as well as the weight attached to them by *Homo moralis*, one finds the following necessary and sufficient condition for the deviation to  $\xi^1$  to be unappealing:

$$\begin{aligned} & \mu_0 \cdot [ 2p_0^2(1-p_0) + p_0^2 [2\kappa(1-\kappa) + \kappa^2] ] \\ > \mu_1 \cdot [ 2p_1(1-p_1)^2 + (1-p_1)^2 [2\kappa(1-\kappa) + \kappa^2] ]. \end{aligned} \quad (10)$$

In a similar manner, Table 2 lists the signal combinations  $(s_i, s_j, s_k)$  for which  $i$  is either pivotal or  $\kappa$ -pivotal if she deviates from  $\xi^{\text{inf}}$  to  $\xi^0$ . Since a switch in the vote outcome from 1 to 0 is costly in state  $\omega = 1$  and beneficial in state  $\omega = 0$ ,  $i$  strictly prefers strategy  $\xi^{\text{inf}}$  to  $\xi^0$  if and only if the probability that  $i$  would either be pivotal or  $\kappa$ -pivotal when the state is  $\omega = 1$  exceeds the probability of the same event when the state is  $\omega = 0$ :

$$\begin{aligned} & \mu_1 \cdot [ 2p_1^2(1-p_1) + p_1^2 [2\kappa(1-\kappa) + \kappa^2] ] \\ > \mu_0 \cdot [ 2p_0(1-p_0)^2 + (1-p_0)^2 [2\kappa(1-\kappa) + \kappa^2] ]. \end{aligned} \quad (11)$$

Table 2: Signal realizations  $(s_i, s_j, s_k)$  for which a deviation by  $i$  from  $\xi^{\text{inf}}$  to  $\xi^0$  would alter the vote outcome, if: (a)  $j$  and  $k$  play  $\xi^{\text{inf}}$ , (b) hypothetically,  $j$  were also to play  $\xi^0$  instead of  $\xi^{\text{inf}}$ , (c) hypothetically,  $k$  were also to play  $\xi^0$  instead of  $\xi^{\text{inf}}$ , (d) hypothetically, both  $j$  and  $k$  were also to play  $\xi^0$  instead of  $\xi^{\text{inf}}$

$s_i$	$s_j$	$s_k$	$(\xi^{\text{inf}}, \xi^{\text{inf}}, \xi^{\text{inf}})$	(a) $(\xi^0, \xi^{\text{inf}}, \xi^{\text{inf}})$	(b) $(\xi^0, \xi^0, \xi^{\text{inf}})$	(c) $(\xi^0, \xi^{\text{inf}}, \xi^0)$	(d) $(\xi^0, \xi^0, \xi^0)$
1	1	1	1	1	0	0	0
1	1	0	1	0	0	0	0
1	0	1	1	0	0	0	0
0	1	1	1	1	0	0	0

Defining the two threshold values

$$\bar{\lambda}^{(\kappa)} \equiv \frac{p_1^2}{(1-p_0)^2} \cdot \frac{\kappa(2-\kappa) + 2(1-p_1)}{\kappa(2-\kappa) + 2p_0} \quad (12)$$

and

$$\underline{\lambda}^{(\kappa)} \equiv \frac{(1-p_1)^2}{p_0^2} \cdot \frac{\kappa(2-\kappa) + 2p_1}{\kappa(2-\kappa) + 2(1-p_0)}, \quad (13)$$

we note that the condition derived above for  $\xi^{\text{inf}}$  to be a Nash equilibrium in a jury of *Homo oeconomicus* (see (9)) can be written  $\underline{\lambda}^{(0)} < \mu_0/\mu_1 < \bar{\lambda}^{(0)}$ , and that the condition for  $\xi^{\text{inf}}$  to be efficient (see (8)) can be written  $\underline{\lambda}^{(1)} < \mu_0/\mu_1 < \bar{\lambda}^{(1)}$ .

The reasoning above, together with the observation that a deviation to the reverse strategy  $\xi^{\text{inv}}$  is clearly dominated by a deviation to both  $\xi^1$  and  $\xi^2$ , allows us to state necessary and sufficient conditions for informative voting to be an equilibrium.

**Proposition 6** *In a jury consisting of three jurors with Homo moralis preferences with degree of morality  $\kappa$ , for any  $\mu_0 \in (0, 1/2)$ :*

1.  $\xi^{\text{inf}}$  is a strict Bayesian Nash equilibrium if and only if  $\underline{\lambda}^{(\kappa)} < \frac{\mu_0}{1-\mu_0} < \bar{\lambda}^{(\kappa)}$ ;
2.  $\underline{\lambda}^{(\kappa)}$  is strictly decreasing and  $\bar{\lambda}^{(\kappa)}$  is strictly increasing in  $\kappa$ , for all  $\kappa \in [0, 1]$ ;
3. if  $\kappa = 1$ ,  $\xi^{\text{inf}}$  is a strict Bayesian Nash equilibrium if and only if  $\xi^{\text{inf}}$  is efficient.

This result is illustrated in Figure 2, which shows the values of  $\mu_0/\mu_1$  for which  $\xi^{\text{inf}}$  is a strict Bayesian Nash equilibrium, both for a jury composed of three *Homo oeconomicus*,

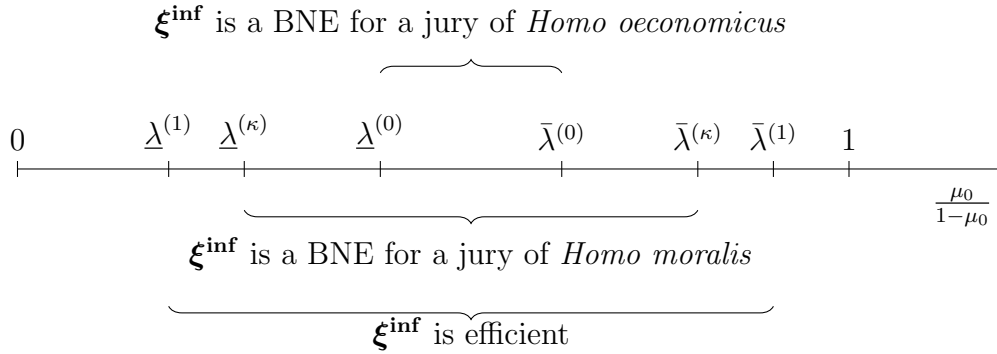


Figure 2: Values of  $\mu_0$  for which  $\xi^{\text{inf}}$  is a strict Bayesian Nash equilibrium

and for a jury composed of three *Homo moralis* with a positive degree of morality  $\kappa \in [0, 1]$ . It also shows the values of  $\mu_0/\mu_1$  for which  $\xi^{\text{inf}}$  is efficient.

In sum, then, *Homo moralis* preferences unambiguously render a small jury more efficient, in the sense that they render informative voting individually rational in settings where such voting is efficient but not individually rational for *Homo oeconomicus*.

The intuition for this result is clear. *Homo oeconomicus* deviates from the informative strategy if this improves the odds that the correct decision is reached, conditional on her vote being pivotal. *Homo moralis* follows a similar reasoning, but conditional on her vote being  $\kappa$ -pivotal, meaning that the outcome would change if not only she but also each other juror were to deviate with probability  $\kappa$ : this exaggerates the impact of the deviation, thus making it less appealing.

It will now be shown that this intuition takes its full force when the size of the jury is large, for in that case the probability  $\kappa$  approximates a fraction  $\kappa$  of the population.

### 4.3 A large jury

We restrict attention to the case of an odd number  $n = 2m + 1$  of voters and check whether the strategy profile  $\xi^{\text{inf}}$  whereby all  $n$  voters use the informative strategy  $\xi^{\text{inf}}$  is a (Bayesian) Nash equilibrium for large  $n$ .

**Proposition 7** *For any fixed values of the parameters  $\mu_0, \mu_1, p_0, p_1, \kappa \in (0, 1)$  such that  $p_1 \geq p_0 > 1/2$ , there exists an integer  $N$  such that for any jury of size  $n = 2m + 1 \geq N$ :*

- if  $\kappa > 1 - p_0/p_1$ , the informative voting strategy profile  $\xi^{\text{inf}}$  is a strict Nash equilibrium;

- if  $\kappa < 1 - p_0/p_1$  the informative voting strategy profile  $\xi^{\text{inf}}$  is not a Nash equilibrium;
- if  $\kappa = 1 - p_0/p_1$  the informative voting strategy profile  $\xi^{\text{inf}}$  is a strict Nash equilibrium if  $p_0\mu_0 < \mu_1$  and is not a Nash equilibrium if  $p_0\mu_0 > \mu_1$ .

Having thus shown that *Homo moralis* preferences pave the way for informative voting to be sustained as a Nash equilibrium in large juries, as (implicitly) posited by Condorcet, we next ask whether such preferences would also help large juries escape from the use of non-informative voting strategies.

Thus, consider now the situation in which all jurors use a non-informative strategy, for instance  $\xi^0$ . Then the decision is 0 independently of the signals received by the jurors. When the jury is large this is clearly a (flat) Nash equilibrium in a jury composed of jurors with *Homo oeconomicus* preferences, since the outcome would then be unaffected by the deviation of a single juror. The expected payoff at this strategy profile is simply  $\mu_0$ , the prior probability that the state of nature is 0. We now prove that in a jury composed of jurors with *Homo moralis* preferences, for any strictly positive value of the morality parameter  $\kappa$  each juror strictly prefers to deviate to the informative strategy. While it is still the case that such a deviation by a single juror has no effect on the actual outcome, *Homo moralis* evaluates the outcome under the informative strategy, should a share  $\kappa$  of the other jurors also use it. For a large enough jury, this makes the individual voter  $\kappa$ -pivotal, thus inducing a strict preference ranking over the strategies.

**Proposition 8** *For any fixed values of the parameters  $\mu_0, \mu_1, p_0, p_1, \kappa \in (0, 1)$  there exists an integer  $N$  such that if  $n = 2m + 1 \geq N$ , a strategy profile in which each juror always votes for the same option (0 or 1) independently of her signal, is not a Nash equilibrium.*

Note that the result holds for any  $\kappa$ : even if the value of  $\kappa$  is small, for  $n$  large enough, uninformative voting is not an equilibrium. This contrasts with Proposition 7, in which the constraint  $\kappa > 1 - p_0/p_1$  appears.

#### 4.4 Conclusion on the large jury problem

It follows from the results reported in Propositions 7 and 8 that the Condorcet Jury Theorem (asymptotic efficient revelation of information) holds for large juries composed of jurors with *Homo moralis* preferences, when the degree of morality  $\kappa$  is large enough. The following theorem sums up these results by characterizing the equilibrium properties of all symmetric



strategy profiles (the first two points are taken from the two previous propositions, and we leave to the reader the proof of the third point, which can follow the same lines).

**Theorem 1** *For any fixed values of the parameters  $\mu_0 \in [0, 1]$  and  $p_1 \geq p_0 > 1/2$ , for any  $\kappa \in (0, 1]$ , if the size of the jury  $n = 2m + 1$  is large enough:*

- *informative voting, whereby each juror casts a vote for the decision that was suggested to him/her through the private signal (i.e., all jurors use the informative strategy  $\xi^{\text{inf}}$ ), is a strict Nash equilibrium if  $\kappa > 1 - p_0/p_1$ , and is not a Nash equilibrium if  $\kappa < 1 - p_0/p_1$ ;*
- *uninformative voting, whereby all jurors always vote for the same decision (i.e., all jurors either use strategy  $\xi^0$  or strategy  $\xi^1$ ), is not a Nash equilibrium;*
- *all players using the reversed strategy  $\xi^{\text{inv}}$  is not a Nash equilibrium.*

## 5 Conclusion

In this paper we take the evolutionary foundations of *Homo moralis* preferences as our starting point and study the consequences of these preferences in two distinct settings: first when voters face a pure coordination problem, and second when voters face an information aggregation problem.<sup>8</sup> Interest in these questions is warranted for many reasons, in particular because answers to these questions are necessary to properly evaluate voting rules from a normative point of view. But the theory of voting, on top of being of political relevance *per se*, contains the study of several archetypal situations of interaction, or “games” that are of broad interest. For instance the political game of coordination of a divided majority can be seen as a toy model of social coordination in general.<sup>9</sup> The two problems we studied can both be seen as instances of “social dilemmas” but they are different.

The first example—the divided majority problem—is a pure multi-person coordination problem. It provides a simple but non-trivial exercise to illustrate how *Homo moralis* preferences can help solve, at least partially, coordination problems (see also Alger and Weibull 2017 [3]). Our study of the divided majority problem highlighted the fact that the *Homo*

---

<sup>8</sup>The question of participation when the electorate is large and voting is costly is tackled in a companion paper [1].

<sup>9</sup>Interestingly, the literature on animal behavior describes several collective phenomena similar to voting, from Honeybees to African wild dogs, relying on the interpretation of various techniques of social communication. See Walker et al. 2017 [51], Seeley 2010 [46], Sumpter 2010 [49].

*moralis* template can be used with various reference groups (akin to various places where the “veil of ignorance” is placed). Whoever wants to use the *Homo moralis* model for an application should wonder what is the appropriate reference identification group. The evolutionary justification provides a hint on this point: the reference group is the one among which interactions have taken place over evolutionary time. However, this is ultimately an empirical question.

The second example—the rational approach to the Condorcet jury theorem—adds a question of information processing to the coordination issue. This raises the same kind of difficulties that appear in the evolutionary theory of language (Laslier 2003 [29], Demichelis and Weibull 2008 [17], Benz et al. 2011 [7]), and analyzing these issues require the use of Bayesian equilibrium. We find that *Homo moralis* preferences help improve the information aggregation that the jurors can achieve by voting based solely on their private information, without communicating with each other.

In the two problems we studied, the partial morality built into *Homo moralis* preferences impacts the predictions. Importantly, these predictions remove or lessen the phenomena often described as “paradoxes” or “curses” that typically appear in the standard model with materially self-interested individual. This observation is an invitation to add the *Homo moralis* model to the toolkit of political economy for descriptive purposes, and also to deepen the evolutionary analysis of political games.

## Appendix

### Proof of Proposition 1

If everyone votes sincerely, candidate  $B$  wins, and  $A$ -supporters achieve utility 0.

1. Consider now an  $A_2$ -supporter who ponders deviating to a vote for  $A_1$ . As a *Homo moralis* with the population of  $A_2$ -supporters as her reference group, when computing her utility for this deviation she evaluates what the outcome would be if, hypothetically, a fraction  $\kappa$  of the  $n_{A_2}$  supporters of candidate  $A_2$  would also vote for  $A_1$ . She thus considers what the outcome would be if the number of votes in favor of  $A_2$  went down from  $n_{A_2}$  to  $(1 - \kappa)n_{A_2}$ , the number of votes for  $A_1$  went up from  $n_{A_1}$  to  $n_{A_1} + \kappa n_{A_2}$ , and the number of votes for  $B$  remained at  $n_B$ . This voter would benefit in utility terms from this deviation if and only if the candidate that would win was  $A_1$  instead of  $B$ . The condition for this deviation to be favorable is thus  $n_{A_1} + \kappa n_{A_2} > n_B$ , or

$\kappa > \kappa^*$ .

2. In a similar manner, an  $A1$ -supporter would strictly prefer to deviate and vote for  $A2$  rather than voting sincerely for  $A1$  if and only if  $n_{A2} + \kappa n_{A1} > n_B$ . Since  $n_{A2} < n_{A1}$ , we have  $n_{A2} + \kappa n_{A1} < n_{A1} + \kappa n_{A2}$ ; in other words, an  $A1$ -supporter strictly prefers to deviate only if an  $A2$ -supporter also prefers to deviate.
3. In sum, if all other  $A$ -supporters vote sincerely, both  $A1$ - and  $A2$ -supporters are indifferent between voting sincerely and deviating to some other strategy if  $\kappa < \kappa^*$ , while  $A2$ -supporters strictly prefer to deviate if  $\kappa > \kappa^*$ .

## Proof of Proposition 2

Suppose that all  $A$ -supporters vote for the same candidate  $Ak$ ,  $k \in \{1, 2\}$ . Then candidate  $Ak$  gets the score  $n_{A1} + n_{A2} > n_B$  and thus wins. Now note that:

1. In the *ex post* scenario an  $A1$ -supporter who considers some deviation, evaluates what the outcome of the vote would be should a proportion  $\kappa$  of her fellow  $A1$ -supporters also play the deviating strategy. In this hypothetical scenario candidate  $Ak$  would get the score  $(1 - \kappa)n_{A1} + n_{A2}$ . The deviation entails a drop in utility if in this hypothetical scenario candidate  $B$  wins, i.e., if the score  $(1 - \kappa)n_{A1} + n_{A2}$  falls short of  $n_B$ , the score of candidate  $B$ , i.e., if  $\kappa > \frac{n_A - n_B}{n_{A1}}$ . The deviation has no effect on the deviator's utility if  $Ak$  still wins in this hypothetical scenario, i.e., if  $\kappa < \frac{n_A - n_B}{n_{A1}}$ .
2. Likewise: in the *ex post* scenario an  $A2$ -supporter who considers some deviation, evaluates what the outcome of the vote would be should a proportion  $\kappa$  of her fellow  $A2$ -supporters also play the deviating strategy. In this hypothetical scenario candidate  $Ak$  would get the score  $n_{A1} + (1 - \kappa)n_{A2}$ . The deviation entails a drop in utility if in this hypothetical scenario candidate  $B$  wins, i.e., if the score  $n_{A1} + (1 - \kappa)n_{A2}$  falls short of  $n_B$ , the score of candidate  $B$ , i.e., if  $\kappa > \frac{n_A - n_B}{n_{A2}}$ . The deviation has no effect on the deviator's utility if  $Ak$  still wins in this hypothetical scenario, i.e., if  $\kappa < \frac{n_A - n_B}{n_{A2}}$ .
3. These observations imply the statement in the proposition, since  $\frac{n_A - n_B}{n_{A2}} > \frac{n_A - n_B}{n_{A1}}$ .

## Proof of Proposition 3

If everyone votes sincerely, candidate  $B$  wins, and  $A$ -supporters achieve utility 0.

1. Consider now an  $A$ -supporter who ponders deviating to strategy  $\alpha^{(1)}$ . In the *ex ante* scenario, this voter evaluates what her expected utility would be if, hypothetically, a share  $\kappa$  of other  $A$ -supporters would also use strategy  $\alpha^{(1)}$ . The score of candidate  $A1$  would in this hypothetical scenario be  $n_{A1} + \kappa n_{A2}$ , that of candidate  $A2$  would be  $(1 - \kappa)n_{A2}$ , and that of candidate  $B$  would be  $n_B$ . Since  $n_{A2} < n_B$ , it follows that  $C^{(\kappa)}(\alpha^{(1)}, \alpha^{si}) \neq B$  if and only if  $C^{(\kappa)}(\alpha^{(1)}, \alpha^{si}) = A1$ , which is true iff  $\kappa > (n_B - n_{A1})/n_{A2}$ . The expected utility of the deviating  $A$ -voter is thus (see (6))

$$U_A^{(\kappa)}(\alpha^{(1)}, \alpha^{si}) = \begin{cases} 1 + \varepsilon \cdot (n_{A1} - n_{A2})/(n_{A1} + n_{A2}) & \text{if } \kappa > (n_B - n_{A1})/n_{A2} \\ 0 & \text{otherwise.} \end{cases} \quad (14)$$

Hence,  $U^{(\kappa)}(\alpha^{(1)}, \alpha^{si}) > U^{(\kappa)}(\alpha^{si}, \alpha^{si})$  if  $\kappa > (n_B - n_{A1})/n_{A2}$ , while  $U^{(\kappa)}(\alpha^{(1)}, \alpha^{si}) = U^{(\kappa)}(\alpha^{si}, \alpha^{si})$  otherwise.

2. Consider now an  $A$ -supporter who ponders deviating to strategy  $\alpha^{(2)}$ . In the *ex ante* scenario, this voter evaluates what her expected utility would be if, hypothetically, a share  $\kappa$  of other  $A$ -supporters would also use strategy  $\alpha^{(2)}$ . The score of candidate  $A1$  would in this hypothetical scenario be  $(1 - \kappa)n_{A1}$ , that of candidate  $A2$  would be  $\kappa n_{A1} + n_{A2}$ , and that of candidate  $B$  would be  $n_B$ . Since  $n_{A1} < n_B$ , it follows that  $C^{(\kappa)}(\alpha^{(2)}, \alpha^{si}) \neq B$  if and only if  $C^{(\kappa)}(\alpha^{(2)}, \alpha^{si}) = A2$ , which is true iff  $\kappa > (n_B - n_{A2})/n_{A1}$ . The expected utility of the deviating  $A$ -voter is thus (see (6))

$$U_A^{(\kappa)}(\alpha^{(2)}, \alpha^{si}) = \begin{cases} 1 - \varepsilon \cdot (n_{A1} - n_{A2})/(n_{A1} + n_{A2}) & \text{if } \kappa > (n_B - n_{A2})/n_{A1} \\ 0 & \text{otherwise.} \end{cases} \quad (15)$$

Hence,  $U^{(\kappa)}(\alpha^{(2)}, \alpha^{si}) > U^{(\kappa)}(\alpha^{si}, \alpha^{si})$  if  $\kappa > (n_B - n_{A2})/n_{A1}$ , while  $U^{(\kappa)}(\alpha^{(2)}, \alpha^{si}) = U^{(\kappa)}(\alpha^{si}, \alpha^{si})$  otherwise. But the assumptions in (3) imply  $(n_B - n_{A2})/n_{A1} > (n_B - n_{A1})/n_{A2}$ , and hence the deviation to strategy  $\alpha^{(2)}$  is beneficial only if the deviation to strategy  $\alpha^{(1)}$  is beneficial.

3. Combining the two previous paragraphs, this proves that there exists a utility-enhancing deviation if and only if  $\kappa > (n_B - n_{A1})/n_{A2}$ .

## Proof of Proposition 4

Suppose that all  $A$ -supporters use strategy  $\alpha^{(1)}$ . Then candidate  $A1$  wins and  $A$ -voters have expected utility  $1 + \varepsilon \cdot (n_{A1} - n_{A2})/(n_{A1} + n_{A2})$ . In the *ex ante* scenario an  $A$ -supporter who considers some deviation, evaluates what the outcome of the vote would be should a

proportion  $\kappa$  of all  $A$ -supporters also play the deviating strategy.

1. Consider first a deviation to strategy  $\alpha^{\text{si}}$ . The hypothetical scores that would obtain should a share  $\kappa$  of other  $A$ -voters also use strategy  $\alpha^{\text{si}}$  instead of  $\alpha^{(1)}$  are  $n_{A1} + (1 - \kappa) \cdot n_{A2}$  for candidate  $A1$ ,  $\kappa \cdot n_{A2}$  for candidate  $A2$ , and  $n_B$  for candidate  $B$ . Since  $\kappa n_{A2} < n_{A1} + (1 - \kappa)n_{A2}$ , the deviation leads to the election of  $B$  if  $n_B > n_{A1} + (1 - \kappa) \cdot n_{A2}$  or of  $A1$  in the opposite case. That is a drop in expected utility if  $n_B > n_{A1} + (1 - \kappa) \cdot n_{A2}$ , i.e., if  $\kappa > \frac{n_A - n_B}{n_{A2}}$ , and no change in expected utility otherwise.
2. Consider now a deviation from strategy  $\alpha^{(1)}$  to strategy  $\alpha^{(2)}$ . This deviation leads to the hypothetical scores  $(1 - \kappa)n_A$  for candidate  $A1$ ,  $\kappa n_A$  for candidate  $A2$ , and  $n_B$  for candidate  $B$ . The effect on the voter's utility depends on the value of  $\kappa$ :
  - if  $\kappa > 1/2$ , the deviation implies  $C^{(\kappa)}(\alpha^{(2)}, \alpha^{(1)}) \in \{A2, B\}$  and thus a drop in utility (see (6));
  - if  $\kappa < 1/2$ , there are two cases:
    - if  $\kappa > \frac{n_A - n_B}{n_A}$ , the deviation implies  $C^{(\kappa)}(\alpha^{(2)}, \alpha^{(1)}) = B$  and thus entails a drop in utility;
    - if  $\kappa < \frac{n_A - n_B}{n_A}$ , the deviation implies  $C^{(\kappa)}(\alpha^{(2)}, \alpha^{(1)}) = A1$ , leaving the utility unaffected.
3. Noting that  $\frac{n_A - n_B}{n_{A2}} > \frac{n_A - n_B}{n_A}$ , and that  $\frac{n_A - n_B}{n_A} < 1/2$  is equivalent to  $n_A < 2n_B$  (which is true by assumption (see (3))), we conclude that, if all other  $A$ -supporters use strategy  $\alpha^{(1)}$ , an  $A$ -supporter:
  - is indifferent between the strategies  $\alpha^{(1)}$ ,  $\alpha^{\text{si}}$ , and  $\alpha^{(2)}$  if  $\kappa < \frac{n_A - n_B}{n_A}$ ;
  - is indifferent between the strategies  $\alpha^{(1)}$  and  $\alpha^{\text{si}}$ , but strictly prefers not to deviate to  $\alpha^{(2)}$  if  $\frac{n_A - n_B}{n_A} < \kappa < \frac{n_A - n_B}{n_2}$ ;
  - strictly prefers not to deviate from  $\alpha^{(1)}$  if  $\kappa > \frac{n_A - n_B}{n_{A2}}$ .

## Proof of Proposition 5

Suppose that all  $A$ -supporters use strategy  $\alpha^{(2)}$ . Then candidate  $A2$  wins and  $A$ -supporters have expected utility  $1 - \varepsilon \cdot (n_{A1} - n_{A2}) / (n_{A1} + n_{A2})$ . In the *ex ante* scenario an  $A$ -supporter who considers some deviation, evaluates what the outcome of the vote would be should a proportion  $\kappa$  of all  $A$ -supporters also play the deviating strategy. A deviation is strictly profitable if it leads to the election of  $A1$ , indifferent if it leads to the election of  $A2$ , and strictly unfavorable if it leads to the election of  $B$ .

1. Consider first a deviation to strategy  $\alpha^{\text{si}}$ . The hypothetical scores that would obtain should a share  $\kappa$  of other  $A$ -voters also use strategy  $\alpha^{\text{si}}$  instead of  $\alpha^{(2)}$  are  $\kappa \cdot n_{A1}$  for candidate  $A1$ ,  $(1 - \kappa) \cdot n_{A1} + n_{A2}$  for candidate  $A2$ , and  $n_B$  for candidate  $B$ . Note that, by assumption,  $n_{A1} < n_B$ , hence the winner is either  $B$  if  $n_B > (1 - \kappa) \cdot n_{A1} + n_{A2}$  or  $A2$  in the opposite case. Thus, the deviation leads to a drop in expected utility if  $n_B > n_{A1} + (1 - \kappa) \cdot n_{A2}$ , i.e., if  $\kappa > \frac{n_A - n_B}{n_{A2}}$ , and to no change in expected utility otherwise.
2. Consider now a deviation from strategy  $\alpha^{(2)}$  to strategy  $\alpha^{(1)}$ . The hypothetical scores that would obtain with this deviation are  $\kappa n_A$  for candidate  $A1$ ,  $(1 - \kappa)n_A$  for candidate  $A2$ , and  $n_B$  for candidate  $B$ . The effect on the voter's expected utility depends on the value of  $\kappa$  as follows.
  - If  $\kappa < \min\{1/2, \frac{n_A - n_B}{n_A}\}$ , the deviation implies  $C^{(\kappa)}(\alpha^{(1)}, \alpha^{(2)}) = A2$  and thus no change in expected utility. Because  $n_B > n_A/2$ , the condition  $\kappa < \min\{1/2, \frac{n_A - n_B}{n_A}\}$  simplifies to:  $\kappa < \frac{n_A - n_B}{n_A}$ .
  - If  $\min\{1/2, \frac{n_A - n_B}{n_A}\} < \kappa < \max\{1/2, \frac{n_B}{n_A}\}$ , the deviation implies  $C^{(\kappa)}(\alpha^{(1)}, \alpha^{(2)}) = B$  and thus a drop in expected utility. Since  $n_B > 1/3$ , and therefore  $2n_B > n_A$ , the condition for this case simply writes:  $\frac{n_A - n_B}{n_A} < \kappa < \frac{n_B}{n_A}$ .
  - If  $\kappa > \max\{1/2, \frac{n_B}{n_A}\}$ , i.e., if  $\kappa > \frac{n_B}{n_A}$ , the deviation implies  $C^{(\kappa)}(\alpha^{(1)}, \alpha^{(2)}) = A1$  and thus an increase in expected utility.
3. Note that  $\frac{n_A - n_B}{n_{A2}}$  is larger than  $\frac{n_A - n_B}{n_A}$ , but can be either larger or smaller than  $\frac{n_B}{n_A}$ . Thus the two previous analyses of the two deviations from  $\alpha^{(2)}$  lead to the following conclusion. From the  $\alpha^{(2)}$ -situation,
  - if  $\kappa < \frac{n_A - n_B}{n_A}$ , the voter is indifferent between  $\alpha^{(2)}$ ,  $\alpha^{(1)}$ , and  $\alpha^{\text{si}}$ ;
  - if  $\frac{n_A - n_B}{n_A} < \kappa < \frac{n_B}{n_A}$  and  $\kappa < \frac{n_A - n_B}{n_{A2}}$ , the voter strictly prefers not to deviate to  $\alpha^{(1)}$ , but is indifferent between  $\alpha^{(2)}$  and  $\alpha^{\text{si}}$ ;
  - if  $\frac{n_B}{n_A} < \kappa < \frac{n_A - n_B}{n_{A2}}$ , the voter strictly prefers to deviate to  $\alpha^{(1)}$ , but is indifferent between  $\alpha^{(2)}$  and  $\alpha^{\text{si}}$ ;
  - if  $\kappa > \frac{n_B}{n_A}$  and  $\kappa > \frac{n_A - n_B}{n_{A2}}$ , the voter strictly prefers to deviate to  $\alpha^{(1)}$  and to  $\alpha^{\text{si}}$ .

The proposition follows directly from these cases.

## Proof of Proposition 7

Let  $\pi(\xi_i, \boldsymbol{\xi}_{-i})$  denote  $i$ 's expected material payoff if  $i$  plays strategy  $\xi_i$  and the other jurors play strategy profile  $\boldsymbol{\xi}_{-i} \in X^{2m}$ . Writing  $A_\omega(\xi_i, \boldsymbol{\xi}_{-i})$  for the probability that the decision is correct in state  $\omega$ , we thus have:

$$\pi(\xi_i, \boldsymbol{\xi}_{-i}) = \mu_0 A_0(\xi_i, \boldsymbol{\xi}_{-i}) + \mu_1 A_1(\xi_i, \boldsymbol{\xi}_{-i}). \quad (16)$$

Now, under majority rule the probability that the right decision is taken in state  $\omega$  equals the probability that at least  $m + 1$  jurors vote  $\omega$  when the state is  $\omega$ . Using the following notation for the binomial probability of  $t$  or more successes out of  $T$ :

$$B^+(p, t, T) = \sum_{k=t}^T C_T^k p^k (1-p)^{T-k}, \quad (17)$$

we immediately obtain that if all the jurors vote informatively, i.e., if  $(\xi_i, \boldsymbol{\xi}_{-i}) = (\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}})$ , the probability that the decision is correct in state  $\omega$  equals the probability that at least  $m + 1$  jurors receive the signal  $\omega$ . In other words,

$$\pi(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) = \mu_0 A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) + \mu_1 A_1(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}), \quad (18)$$

where

$$A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) = B^+(p_0, m + 1, 2m + 1) \quad (19)$$

and

$$A_1(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) = B^+(p_1, m + 1, 2m + 1). \quad (20)$$

We now derive conditions required for a juror, say  $i$ , to prefer playing  $\xi^{\text{inf}}$  to deviating to strategy  $\xi^0$ . If this juror has *Homo moralis* preferences with degree of morality  $\kappa \in [0, 1]$ , she evaluates the consequences of the deviation on her expected material payoff, should each other juror play  $\xi^{\text{inf}}$  with probability  $1 - \kappa$  and  $\xi^0$  with probability  $\kappa$ . Taking into account this reasoning, let  $v_0^{(\kappa)}$  be the number of votes 0 and  $v_1^{(\kappa)}$  be the number of votes 1 that *Homo moralis* envisages. Because  $i$  votes 0, the probability of a correct decision, based on the reasoning of *Homo moralis* with degree of morality  $\kappa$ , is:

$$\begin{aligned} \pi^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) &= \Pr[v_0^{(\kappa)} \geq m \text{ and } \omega = 0] + \Pr[v_1^{(\kappa)} \geq m + 1 \text{ and } \omega = 1] \\ &= \mu_0 A_0^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) + \mu_1 A_1^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}), \end{aligned} \quad (21)$$

where  $A_\omega^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  is the probability that at least  $m + 1$  jurors vote  $\omega$  in state  $\omega$ , taking into account the hypothetical scenario that each other juror play  $\xi^{\text{inf}}$  with probability  $1 - \kappa$  and  $\xi^0$  with probability  $\kappa$ . Formally:

$$A_0^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) = \Pr[v_0 \geq m \mid \omega = 0] = B^+(p'_0, m, 2m) \quad (22)$$

$$A_1^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) = \Pr[v_1 \geq m + 1 \mid \omega = 1] = B^+(p'_1, m + 1, 2m), \quad (23)$$

where

$$p'_0 = (1 - \kappa)p_0 + \kappa \quad (24)$$

and

$$p'_1 = (1 - \kappa)p_1. \quad (25)$$

Now, note that as  $m \rightarrow \infty$ , both  $B^+(p, m + 1, 2m + 1)$  and  $B^+(p, m, 2m + 1)$  tend either to 0 or to 1 depending on whether  $p$  is smaller or greater than  $1/2$ . Specifically:

- both  $A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}})$  and  $A_1(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}})$  are increasing in  $m$  and tend to 1 as  $m \rightarrow \infty$ , since (by assumption) both  $p_0 > 1/2$  and  $p_1 > 1/2$ ;
- $A_0^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  is increasing in  $m$  and tends to 1 as  $m \rightarrow \infty$ , since  $p'_0 = (1 - \kappa)p_0 + \kappa > p_0 > 1/2$ ;
- $A_1^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  is increasing in  $m$  and tends to 1 as  $m \rightarrow \infty$  if  $p'_1 = (1 - \kappa)p_1 > 1/2$ , but is decreasing in  $m$  and tends to 0 as  $m \rightarrow \infty$  if  $p'_1 = (1 - \kappa)p_1 < 1/2$ .

Taken together, these observations imply the following:

- $\pi(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}})$  (see (18)) is increasing in  $m$  and tends to  $\mu_0 + \mu_1$  as  $m \rightarrow \infty$ ;
- $\pi(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  is increasing in  $m$  and tends to  $\mu_0 + \mu_1$  as  $m \rightarrow \infty$  if  $(1 - \kappa)p_1 > 1/2$ ;
- $\pi(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  tends to  $\mu_0$  as  $m \rightarrow \infty$  if  $(1 - \kappa)p_1 < 1/2$ .

Clearly, the deviation to  $\xi^0$  is thus not viable if  $(1 - \kappa)p_1 < 1/2$  and  $m$  is large enough. We take note of this result:

**Lemma 1** *Suppose that  $\kappa > 1 - 1/(2p_1)$ . If  $m$  is large enough, each juror strictly prefers to play  $\xi^{\text{inf}}$  than to deviate to  $\xi^0$ .*



By contrast, if  $\kappa < 1 - 1/(2p_1)$ , both  $\pi(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}})$  and  $\pi(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  tend to 1 as  $m \rightarrow \infty$ , hence it is a priori impossible to determine whether a deviation to  $\xi^0$  is viable. In order to go further, we remark that the kind of binomial sums that appear in the calculations can be approximated using the large deviation theory for binomial laws; see for instance Arratia and Gordon (1989 [5]).

Let  $x_n$  and  $y_n$ ,  $n \in \mathbb{N}$  be two sequences of real numbers. We say that  $x$  and  $y$  are equivalent, and we write  $x \sim y$ , when  $\lim_{n \rightarrow \infty} x_n/y_n = 1$ . Likewise, we use the notation  $x \ll y$  when  $\lim_{n \rightarrow \infty} x_n/y_n = 0$ .

**Lemma 2** *Let  $a$  and  $q$  be two real numbers such that  $0 < q < a < 1$ , and*

$$\begin{aligned} H = H(q, a) &= a \log \frac{a}{q} + (1 - a) \log \frac{1 - a}{1 - q}, \\ r = r(q, a) &= \frac{q}{1 - q} / \frac{a}{1 - a} = \frac{q(1 - a)}{a(1 - q)}, \\ \rho = \rho(q, a) &= \frac{1}{(1 - r)\sqrt{2\pi a(1 - a)}}. \end{aligned}$$

*Then, when  $n$  tends to infinity:*

$$B^+(q, an, n) \sim \rho e^{-Hn}.$$

We apply these formula to the expressions above, in order to check whether the deviation is profitable. Recall that we defined  $A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) = B^+(p_0, m + 1, 2m + 1)$ , so that  $1 - A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) = B^+(q_0, m + 1, 2m + 1)$ , where  $q_0 = 1 - p_0$ . We likewise re-write:

$$\begin{aligned} 1 - A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) &= B^+(q_0, m + 1, 2m + 1) \\ 1 - A_1(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) &= B^+(q_1, m + 1, 2m + 1) \\ 1 - A_0^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) &= B^+(q'_0, m, 2m) \\ 1 - A_1^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) &= B^+(q'_1, m + 1, 2m), \end{aligned}$$

where

$$\begin{aligned} q_0 &= 1 - p_0 \\ q_1 &= 1 - p_1 \\ q'_0 &= (1 - \kappa)q_0 \\ q'_1 &= (1 - \kappa)q_1 + \kappa. \end{aligned}$$

Writing the expected utility gain from deviating in the form:

$$\Delta^0 = \pi^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) - \pi(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}),$$

we have

$$\Delta^0 = \mu_0 A_0^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) + \mu_1 A_1^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}}) - \mu_0 A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}) - \mu_1 A_1(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}}), \quad (26)$$

which can also be written

$$\begin{aligned} \Delta^0 &= \mu_0 \cdot [B^+(q_0, m+1, 2m+1) - B^+(q'_0, m, 2m)] \\ &\quad + \mu_1 \cdot [B^+(q_1, m+1, 2m+1) - B^+(q'_1, m+1, 2m)]. \end{aligned} \quad (27)$$

Let us now apply Lemma 2 to the first of these four binomial sums. Let  $n = 2m + 1$  and  $a = 1/2$ , then  $an = m + 1/2$  so that having  $an$  or more Bernoulli successes means having  $m + 1$  or more. Because  $q_0 < 1/2$  the lemma applies and writes:  $B^+(q_0, m+1, 2m+1) \sim \rho(q_0, 1/2)e^{-H(q_0, 1/2)n}$ . Note that  $H(q_0, 1/2) = (1/2)\log[4q_0(1-q_0)]$  so that  $e^{-H(q_0, 1/2)n} = [4q_0(1-q_0)]^{-n/2}$  and we obtain:

$$B^+(q_0, m+1, 2m+1) \sim \rho(q_0, 1/2)[4q_0(1-q_0)]^{-(2m+1)/2}.$$

The same reasoning works for the other sums in (27), because  $q_0, q'_0, q_1$ , and  $q'_1$  are all strictly smaller than  $1/2$ . In each case, i.e., for  $q \in \{q_0, q_1, q'_0, q'_1\}$ , the same short computation gives the same form for  $H$ :

$$H(q, 1/2) = -(1/2)\log[4(1-q)q]$$

and we obtain:

$$B^+(q_0, m+1, 2m+1) \sim \rho(q_0, 1/2)[4q_0(1-q_0)]^{(2m+1)/2}$$

$$B^+(q_1, m+1, 2m+1) \sim \rho(q_1, 1/2)[4q_1(1-q_1)]^{(2m+1)/2}.$$

$$B^+(q'_0, m, 2m) \sim \rho(q'_0, 1/2)[4q'_0(1-q'_0)]^m$$

$$B^+(q'_1, m+1, 2m) \sim \rho(q'_1, 1/2)[4q'_1(1-q'_1)]^m.$$

Having in mind the shape of the function  $x \mapsto x(1-x)$ , one can see that, because  $q'_0 = (1-\kappa)q_0 < q_0 < 1/2$ ,  $q'_0(1-q'_0) < q_0(1-q_0)$ , so that the term  $1 - A_0^{(\kappa)}(\xi^0, \boldsymbol{\xi}^{\text{inf}})$  tends

to zero faster than the term  $1 - A_0(\xi^{\text{inf}}, \boldsymbol{\xi}^{\text{inf}})$ . Hence, when  $m$  tends to infinity:

$$B^+(q'_0, m, 2m) \ll B^+(q_0, m + 1, 2m + 1).$$

We thus reach the following conclusion for the first term in (27): for  $m$  large enough,

$$\mu_0 \cdot [B^+(q_0, m + 1, 2m + 1) - B^+(q'_0, m, 2m)] \sim \mu_0 \cdot B^+(q_0, m + 1, 2m + 1).$$

Turning now to the second term in (27), note that  $q_1 < q'_1 < 1/2$ , so that  $q'_1(1 - q'_1) > q_1(1 - q_1)$ . It follows that

$$B^+(q_1, m + 1, 2m + 1) \ll B^+(q'_1, m + 1, 2m).$$

We thus reach the following conclusion for the second term in (27):

$$\mu_1 \cdot [B^+(q_1, m + 1, 2m + 1) - B^+(q'_1, m + 1, 2m)] \sim -\mu_1 \cdot B^+(q'_1, m + 1, 2m).$$

Taken together, these observations imply that when  $m$  tends to infinity,

$$\Delta^0 \sim \mu_0 \cdot B^+(q_0, m + 1, 2m + 1) - \mu_1 \cdot B^+(q'_1, m + 1, 2m). \quad (28)$$

Recalling that  $q_0 < 1/2$  and  $q'_1 < 1/2$ , we have

$$q_0(1 - q_0) < q'_1(1 - q'_1) \iff q'_1 > q_0 \iff p_0 > p_1(1 - \kappa).$$

Hence:

- $p_0 > p_1(1 - \kappa) \implies B^+(q_0, m + 1, 2m + 1) \ll B^+(q'_1, m + 1, 2m)$ ,
- $p_0 < p_1(1 - \kappa) \implies B^+(q'_1, m + 1, 2m) \ll B^+(q_0, m + 1, 2m + 1)$ ,

and it follows that:

- If  $\kappa > 1 - p_0/p_1$ , then  $\Delta^0 \sim -\mu_1 \cdot B^+(q'_1, m + 1, 2m)$ .
- If  $\kappa < 1 - p_0/p_1$ , then  $\Delta^0 \sim \mu_0 \cdot B^+(q_0, m + 1, 2m + 1)$ .

From this we can infer the sign of  $\Delta^0$  when  $m$  is large: for instance in the first case, because the ratio  $\Delta^0/(\mu_1 \cdot B^+(q'_1, m + 1, 2m))$  tends to  $-1$  when  $m$  tends to infinity, there exists  $\bar{m}$  such that for all  $m \geq \bar{m}$ ,  $\Delta^0 < 0$ . With this reasoning we complement in a unique statement Lemma 1 and conclude:

- If  $\kappa > 1 - p_0/p_1$ , then  $\Delta^0 < 0$  for  $m$  large, in which case deviating to  $\xi^0$  is not profitable;
- If  $\kappa < 1 - p_0/p_1$ , then  $\Delta^0 > 0$  for  $m$  large, in which case deviating to  $\xi^0$  is profitable.

In the knife-edge case  $\kappa = 1 - p_0/p_1$ , that is  $q_0 = q'_1$ , equation (28) writes:

$$\Delta^0 \sim \mu_0 B^+(q_0, m+1, 2m+1) - \mu_1 B^+(q_0, m+1, 2m).$$

Decomposing the binomial sums and re-arranging the terms, one obtains:

$$\begin{aligned} \Delta^0 &\sim \mu_0 \sum_{k=m+1}^{2m+1} C_{2m+1}^k q_0^k (1-q_0)^{2m+1-k} - \mu_1 \sum_{k=m+1}^{2m} C_{2m+1}^k q_0^k (1-q_0)^{2m-k} \\ &= \sum_{k=m+1}^{2m} [\mu_0(1-q_0) - \mu_1] C_{2m+1}^k q_0^k (1-q_0)^{2m-k} + \mu_0 q_0^{2m+1} \\ &= [\mu_0(1-q_0) - \mu_1] B^+(q_0, m+1, 2m) + \mu_0 q_0^{2m+1}. \end{aligned}$$

Because  $B^+(q_0, m+1, 2m)$  decreases at rate  $[4q_0(1-q_0)]^{-m}$ , the last term ( $\mu_0 q_0^{2m+1}$ ) is negligible, so that

$$\Delta^0 \sim [\mu_0(1-q_0) - \mu_1] B^+(q_0, m+1, 2m),$$

implying that  $\Delta^0$  has the same sign as  $\mu_0(1-q_0) - \mu_1$ . We conclude that in the knife-edge case,  $\xi^0$  destabilizes informative voting if and only if  $\mu_0(1-q_0) > \mu_1$  (when  $m$  tends to infinity).

The symmetric conclusions are reached for a deviation from  $\xi^{\text{inf}}$  to  $\xi^1$ , the threshold value for  $\kappa$  being equal to  $1 - p_1/p_0$  instead of  $1 - p_0/p_1$ . Since  $\kappa$  is positive only one of these threshold values is relevant, however:

- if  $p_0 \geq p_1$ , then  $1 - p_0/p_1 \leq 0$ , and a deviation to  $\xi^0$  is not profitable for any  $\kappa \in [0, 1]$ ; but if  $p_0 > p_1$  the threshold value  $1 - p_1/p_0$  is positive, implying that a deviation to  $\xi^1$  is profitable if  $\kappa \in [0, 1 - p_1/p_0]$ ;
- as shown above, the opposite conclusion holds if  $p_1 \geq p_0$ .

Putting all this together we find the statement in the proposition, where for simplicity and without loss of generality, we only treat the case  $p_1 \geq p_0$ .

## Proof of Proposition 8

Suppose that all jurors use strategy  $\xi^0$ . The probability that the decision is correct is then:

$$\pi(\xi^0, \xi^0) = \mu_0. \quad (29)$$

Note that this is also the utility of a *Homo moralis* with any degree of morality  $\kappa$  who uses strategy  $\xi^0$  given that all other jurors do so as well.

We now derive conditions required for a *Homo moralis* juror  $i$  to prefer to deviate from  $\xi^0$  to the informative strategy  $\xi^{\text{inf}}$ . Such a juror evaluates the consequences of the deviation on her expected material payoff, should each other juror play  $\xi^0$  with probability  $1 - \kappa$  and  $\xi^{\text{inf}}$  with probability  $\kappa$ . As above, let  $v_0^{(\kappa)}$  be the number of votes 0 and  $v_1^{(\kappa)}$  be the number of votes 1 that *Homo moralis* envisages. Because  $i$  votes informatively, the probability of a correct decision, based on the reasoning of *Homo moralis* with degree of morality  $\kappa$ , is:

$$\pi^{(\kappa)}(\xi^{\text{inf}}, \xi^0) = \mu_0 \cdot [p_0 B^+(1 - \kappa q_0, m, 2m) + q_0 B^+(1 - \kappa q_0, m + 1, 2m)] \quad (30)$$

$$+ \mu_1 \cdot [p_1 B^+(\kappa p_1, m, 2m) + q_1 B^+(\kappa p_1, m + 1, 2m)]. \quad (31)$$

In this expression line (30) corresponds to the probability of a correct decision in state of Nature  $\omega = 0$ . The decision is correct in this state if either  $i$  receives signal 0 and at least  $m$  other voters vote 0, or  $i$  receives the wrong signal 1 and at least  $m + 1$  others vote 0. Any other voter votes 0 either if she uses strategy  $\xi^0$  (a scenario to which *Homo moralis* attaches weight  $1 - \kappa$ ), or when she uses strategy  $\xi^{\text{inf}}$  (a scenario to which *Homo moralis* attaches weight  $\kappa$ ), and receives signal 0; so in state  $\omega = 0$  *Homo moralis* ponders a hypothetical world in which each other voter has probability  $1 - \kappa + \kappa p_0 = 1 - \kappa q_0$  to vote 0. Line (31) likewise counts the votes 1 in state  $\omega_1$ , given that *Homo moralis* ponders a hypothetical world in which each other voter uses strategy  $\xi^0$  with probability  $1 - \kappa$  and strategy  $\xi^{\text{inf}}$  with probability  $\kappa$ ), and thus votes 1 with probability  $\kappa p_1$ .

Since  $q_0 < 1/2$  (by assumption), we have  $1 - \kappa q_0 > 1/2$ , which implies that both binomial sums in line (30) tend to 1. Hence, line (30) tends to  $\mu_0$ .

In line (31), as  $m \rightarrow \infty$  both binomial sums tend to 1 if  $\kappa p_1 > 1/2$  and to 0 if  $\kappa p_1 < 1/2$ . There are thus two cases:

- If  $\kappa p_1 > 1/2$ : the expected utility of the deviation to  $\xi^{\text{inf}}$  tends to  $\mu_0 + \mu_1$ , and since this strictly exceeds  $\mu_0$  the deviation is strictly profitable (for  $m$  large enough).
- If  $\kappa p_1 < 1/2$ : the expected utility of the deviation tends to  $\mu_0$ , so the benefit from

deviating tends to 0. In order to check its sign we write the expected benefit as follows:

$$\begin{aligned}\Delta &= \pi^{(\kappa)}(\xi^{\text{inf}}, \boldsymbol{\xi}^0) - \pi^{(\kappa)}(\xi^0, \boldsymbol{\xi}^0) \\ &= \mu_0 \cdot [p_0(B^+(1 - \kappa q_0, m, 2m) - 1) + q_0(B^+(1 - \kappa q_0, m + 1, 2m) - 1)] \\ &\quad + \mu_1 \cdot [p_1 B^+(\kappa p_1, m, 2m) + q_1 B^+(\kappa p_1, m + 1, 2m)].\end{aligned}$$

Since  $B^+(p, t, T) - 1 = -B^+(1 - p, t, T)$ , this can in turn be written as follows:

$$\begin{aligned}\Delta &= -\mu_0 \cdot [p_0 B^+(\kappa q_0, m + 1, 2m) + q_0 B^+(\kappa q_0, m, 2m)] \\ &\quad + \mu_1 \cdot [p_1 B^+(\kappa p_1, m, 2m) + q_1 B^+(\kappa p_1, m + 1, 2m)].\end{aligned}$$

Recalling that  $q_0 < 1/2 < p_1$ , so that  $\kappa q_0 < \kappa p_1 < 1/2$ , we note that the term in the first square brackets goes to 0 faster than the second one (see Lemma 2 in the proof of Proposition 7). Hence, for  $m$  large enough,  $\Delta$  is positive and deviating is strictly profitable.

The knife-edge case  $\kappa p_1 = 1/2$  is easily solved by writing the term in square brackets in line (31) as follows:

$$\begin{aligned}& p_1 B^+(\kappa p_1, m, 2m) + q_1 B^+(\kappa p_1, m + 1, 2m) \\ &= p_1 B^+(1/2, m, 2m) + q_1 B^+(1/2, m + 1, 2m) \\ &> q_1 [B^+(1/2, m, 2m) + B^+(1/2, m + 1, 2m)] \\ &= q_1,\end{aligned}$$

which proves that, in this case,  $\pi^{(\kappa)}(\xi^{\text{inf}}, \boldsymbol{\xi}^0) > \mu_0 + q_1 \mu_1 > \mu_0$ , making the deviation profitable.

The same reasoning can be applied to show that, for  $m$  large enough, a deviation from  $\xi^1$  to  $\xi^{\text{inf}}$  is profitable for any  $\kappa > 0$ . This completes the proof.

## References

- [1] Ingela Alger and Jean-François Laslier (2020) ‘‘Homo moralis goes to the voting booth: a new theory of voter turnout’’ Working Paper, Paris and Toulouse Schools of Economics.
- [2] Ingela Alger and Jorgens Weibull (2013) ‘‘Homo moralis: preference evolution under incomplete information and assortative matching’’ *Econometrica* 81: 2269—2302.

- [3] Ingela Alger and Jörgen Weibull (2017) “Strategic behavior of moralists and altruists” *Games* 8(3): 38.
- [4] Ingela Alger and Jörgen Weibull (2019) “Evolutionary models of preference formation” *Annual Review of Economics*, 11: 329–354.
- [5] Richard Arratia and Louis Gordon (1989) “Tutorial on large deviations for the binomial distribution” *Bulletin of Mathematical Biology* 51: 15—131.
- [6] David Austen-Smith and Jeffrey S. Banks (1996) “Information aggregation, rationality, and the Condorcet jury theorem” *American Political Science Review* 90: 34—45.
- [7] Anton Benz, Christian Ebert, Gerhard Jger, and Robert van Rooil (eds.) (2011) *Language, Games, and Evolution: Trends in Current Research on Language and Game Theory*. Berlin, Heidelberg: Springer.
- [8] Carisa L. Bergner and Peter K. Hatemi (2017) “Integrating genetics into the study of electoral behavior” pp.367—405 in: Kai Arzheimer, Jocelyn Evans and Michael Lewis-Beck (eds.) *The Sage Handbook of Electoral Behavior. Volume 1*. London: Sage.
- [9] Ken Binmore (1994) *Playing Fair: Game Theory and the Social Contract*. Cambridge, Mass. MIT Press.
- [10] André Blais (2000) *To Vote or Not to Vote? The Limits and Merits of Rational Choice Theory*. Pittsburgh, PA: University of Pittsburgh Press.
- [11] Laurent Bouton and Benjamin Ogden (2021) “Group-based voting in multicandidate elections” *Journal of Politics* 83: 468—482
- [12] Stephen Coate and Michael Conlin (2004) “A group rule utilitarian approach to voter turnout: theory and evidence” *American Economic Review* 94: 1476—1504.
- [13] Condorcet (1785) *Essai sur l’application de l’analyse à la probabilité des décisions rendues à la pluralité des voix*. Paris: Imprimerie Royale.
- [14] Gary W. Cox (1997) *Making Votes Count: Strategic Coordination in the World’s Electoral Systems*. Cambridge: Cambridge University Press.
- [15] Eric van Damme (1997) *Stability and Perfection of Nash Equilibria*. Heidelberg: Springer-Verlag.
- [16] Anthony Downs (1957) *An Economic Theory of Democracy*. New York: Harper and Row.
- [17] Stefano Demichelis and Jörgen Weibull (2008) “Language, meaning, and games: a model of communication, coordination, and evolution” *American Economic Review*, 98(4): 1292—1311.
- [18] Timothy Feddersen and Wolfgang Pesendorfer (1997) “Voting behavior and information aggregation in elections with private and common values” *Econometrica* 65: 1029—1058.

- [19] Timothy Feddersen and Wolfgang Pesendorfer (1998) “Convicting the innocent: the inferiority of unanimous jury verdicts under strategic voting” *American Political Science Review* 92: 23—35.
- [20] Timothy Feddersen and Alvaro Sandroni (2006) “A theory of participation in elections” *American Economic Review* 96(4): 1271—1282.
- [21] Robert Forsythe, Roger Myerson, Thomas Rietz, and Robert Weber (1993) “An experiment on coordination in multi-candidate elections: The importance of polls and election histories” *Social Choice and Welfare*, 10(3): 223—247.
- [22] Dino Gerardi and Leeat Yariv (2008) “Information acquisition in committees” *Games and Economic Behavior* 62: 436—459.
- [23] Hans Gersbach and Volker Hahn (2008) “Should the individual voting records of central bankers be published?” *Social Choice and Welfare* 30: 655—683.
- [24] Donald P. Green and Ian Shapiro (1994) *Pathologies of Rational Choice Theory. A Critique of Applications in Political Science*. New Haven, Conn.: Yale University Press.
- [25] John Harsanyi (1980) “Rule utilitarianism, rights, obligations and the theory of rational behavior” *Theory and Decision* 12: 115—133.
- [26] John Harsanyi (1992) “Game and decision theoretic models in ethics” in: *Handbook of Game Theory* edited by R.J. Aumann and S. Hart, Elsevier, volume 1, pp. 669—707.
- [27] Immanuel Kant (1785) *Grundlegung zur Metaphysik der Sitten*. Trad: Mary Gregor and Jens Timmermann (2011) *Groundwork of the Metaphysics of Morals: A German-English Edition*. Cambridge, Mass.: Cambridge University Press.
- [28] Yukio Koriyama and Balázs Szentes (2009) “A resurrection of the Condorcet jury theorem” *Theoretical Economics* 4: 227—252.
- [29] Jean-François Laslier (2003) “The evolutionary analysis of signal games” in *Cognitive Economics*, edited by P. Bourguin and J.-P. Nadal, Heidelberg : Springer, pp. 281—291.
- [30] Jean-François Laslier and Jörgen Weibull (2013) “An incentive-compatible Condorcet jury theorem” *The Scandinavian Journal of Economics* 115(1): 84—108.
- [31] Jean-Jacques Laffont (1975) “Macroeconomic constraints, economic efficiency, and ethics: An introduction to Kantian economics”, *Economica* 42: 430-437.
- [32] Jacques Lesourne, André Orléan, and Bernard Walliser (eds.) (2006) *Evolutionary Microeconomics*. Berlin: Springer.
- [33] Sydney Levine, Max Kleiman-Weinera, Laura Schulz, Joshua Tenenbaum, and Fiery Cushman (2020) “The logic of universalization guides moral judgment” *PNAS* 117 (42): 26158—26169.



- [34] Christofer Li and Ricardo Pique (2020) “A theory of strategic voting with non-instrumental motives” *Social Choice and Welfare* 55: 369—398.
- [35] John Maynard Smith (1982) *Evolution and the Theory of Games*. Cambridge, Mass.: Cambridge University Press.
- [36] Hervé Moulin (1995) *Cooperative Microeconomics: A Game-Theoretical Introduction*. Cambridge, Mass.: Princeton University Press.
- [37] Roger Myerson and Robert Weber (1993) “A theory of voting equilibria” *American Political Science Review* 87: 102—114.
- [38] Roger Myerson (2002) “Comparison of scoring rules in Poisson voting games” *Journal of Economic Theory* 103: 219—251.
- [39] Roger Myerson and Jörgen Weibull (2015) “Tenable strategies and settled equilibria” *Econometrica* 83(3): 943—976.
- [40] Martin A. Nowak and Karl Sigmund (2005) “Evolution of indirect reciprocity” *Nature* 437: 1293—1295.
- [41] Elinor Ostrom (1998) “A behavioral approach to the rational choice theory of collective action: Presidential Address, American Political Science Association, 1997” *American Political Science Review* 92(1): 1—22.
- [42] Michael B. Petersen (2015) “Evolutionary political psychology” in: D. Buss (Ed.) *Handbook of Evolutionary Psychology*. Vol. 2. p. 1084-1102. John Wiley.
- [43] John Roemer (2019) *How We Cooperate*. Yale University Press.
- [44] Jean-Jacques Rousseau (1755) *Discours sur l'origine et les fondements de l'inégalité parmi les hommes*. Reprinted in : *Ecrits politiques*, 1992, Le livre de Poche.
- [45] Glendon Schubert (1982) “Evolutionary Politics” *Western Political Quarterly* 175—193.
- [46] Thomas D. Seeley (2010) *Honeybee Democracy*. Princeton, NJ: Princeton University Press.
- [47] Jim Sidanus and Robert Kurzban (2013) “Toward an evolutionary informed political psychology” in: L. Huddy, D. O. Sears, and J. S. Levy *The Oxford Handbook of Political Psychology* pp. 205—236. Oxford University Press.
- [48] Laura B. Stephenson, John H. Aldrich, and Andr Blais (2018) *The Many Faces of Strategic Voting: Tactical behavior in electoral systems around the world*. University of Michigan Press.
- [49] David J. T. Sumpter (2010) *Collective Animal Behavior*. Princeton, NJ: Princeton University Press.

- [50] Franz de Waal (1996) *Good Natured: The origins of right and wrong in humans and other animals*. Cambridge, Mass.: Cambridge University Press.
- [51] Reena H. Walker, Andrew J. King, J. Weldon McNutt, and Neil R. Jordan (2017) Sneeze to leave: African wild dogs (*Lycaon pictus*) use variable quorum thresholds facilitated by sneezes in collective decisions *Proceedings of the Royal Society of London B: Biological Sciences*, vol. 284. DOI: 10.1098/rspb.2017.0347.