

Centralization and Stability in Formal Constitutions

Yotam Gafni

Weizmann Institute of Science
yotam.gafni@gmail.com

Abstract. Consider a social-choice function (SCF) is chosen to decide votes in a formal system, including votes to replace the voting method itself. Agents vote according to their ex-ante preference between the incumbent SCF and the suggested replacement. The existing SCF then aggregates the agents' votes and arrives at a decision of whether it should itself be replaced. An SCF is self-maintaining if it can not be replaced in such fashion by any other SCF. Our focus is on the implications of self-maintenance for centralization. We present results considering optimistic, pessimistic and i.i.d. approaches w.r.t. agent beliefs, and different tie-breaking rules. To highlight two of the results, (i) for the i.i.d. unbiased case with arbitrary tie-breaking, we prove an “Arrow-Style” Theorem for Dynamics: We show that only a dictatorship is self-maintaining, and any other SCF has a path of changes that arrives at a dictatorship. (ii) If we take into account wisdom of the crowd effects, for a society with a variable size of ruling elite, we demonstrate how the stable elite size is decreasing in both how extractive the economy is, and the quality of individual decision-making. All in all we provide a basic framework and body of results for centralization dynamics and stability, applicable for institution design, especially in formal “De-Jure” systems, such as Blockchain Decentralized Autonomous Organizations (DAOs).

Keywords: Social Choice, Blockchain Governance, Institution Design

1 Introduction

Decentralized Autonomous Organizations (DAOs) [4] are smart contracts meant for the governance of a blockchain system. Importantly, DAOs have a *formal* definition of decision-making: The process is hard-coded into the smart contract, and must be followed. Typically, voting tokens are issued, and a process is specified where motions can be suggested by the community, and voted on by the token holders. Overall, DAOs currently manage over 13.6 Billion USD in assets [21].

We observe that the decision-making procedure is in itself *dynamic*. A motion can propose to change how decisions are made. This may be done either explicitly or implicitly. Explicitly, it would mean that the community brings up a motion to revise the voting method: For example, to move from simple majority of the voters on a certain proposal, to a qualified majority of at least (say) $\frac{2}{3}$ of the token holders. Implicitly, it may suggest that token holders earn interest on their tokens, increasing system centralization as larger holders will gain more.

Thus, in this work, we study the dynamics of the decision-making process itself. The natural question to ask, given such dynamics, is where do such dynamics converge? If we start at any voting method, which voting method do we end up with? An even more basic question, is whether there are *stable* voting methods, which are not susceptible to such changes?

For this purpose, we define the notion of *self-maintaining* SCFs. Voters that are faced with deciding between two SCFs, consider their expected utility under each SCF. That is, using their beliefs of how future decision motions are drawn, they calculate the probability of their preference to be the outcome of the SCF, and they prefer the SCF where their preference is more likely to indeed be the outcome. Importantly, it is not enough that (for example) a majority of the voters determine that they prefer to move to a different SCFs: If the current SCF is not the majority rule, then aggregating their preferences does not necessarily result in deciding to move to the new SCF. For a change to happen, it must be that within the current power structure, the motion of changing to a different SCF indeed passes using the current SCF. This is the main idea of our formulation.

Two remarks are in order. (i) we have yet to define the agents’ utilities, given their beliefs. As our baseline, we consider agents that want to maximize their probability of winning every decision. This type of utility has an extractive mindset: The underlying assumption is that there is some resource allocation that can be determined as part of the voting process, and each agent wants to have the most power over it. A more nuanced world model is the one underlying “Condorcet Jury Theorem”: Agents have independent signals over the true state of the world, and they wish to correctly decide in accordance with it. In our extension of Section 5, we consider a utility that is a mix of both the extractive and the common-good elements. (ii) We consider agents that are *myopic*, in the sense that they do not think ahead about the implications of changing the SCF. For example, if agents believe that any change in the SCF ultimately results in a dictatorship, they might avoid making any change at all. We discuss this assumption in more depth in the discussion.

We are inspired by very early political science work, going back to Plato, Aristotle, and Polybius [42, 43, 47]. Plato considered how all systems of governance end up in tyranny. Aristotle considered, based on different case-studies of his time, how government forms may change, without concluding in some inevitable end-state. Polybius presented a cyclical theory of constitution where dictatorship transforms to oligarchy, then to democracy, and then back to dictatorship. One interesting conclusion by Polybius is that to arrive at a stable constitution, one must mix elements of the different government forms (modern game theorists may call this a game that has no pure equilibrium, but has a mixed equilibrium). These type of dynamics were on display throughout history: In the French revolution [22], the monarchy was replaced by a full democracy, but then later Robespierre’s “Committee of Public Safety” (which may be thought of as a small oligarchy) assumed greater powers, and eventually the power struggle ended up with Napoleon as emperor (i.e., a return to monarchy). During the Communist revolution in Russia [13], the two factions of Bolsheviks and Mensheviks both fought for Communist rule against the White Army, but then later the Menshevik party was banned. Established modern work on civil resistance finds that the most important factor of successful protest is in the ability to shift the loyalties of regime insiders [16]. This paints a picture where “revolutions”, while seemingly an outside force, succeed mostly through realigning the current power structures from within, in accordance with our De-Jure approach.

1.1 Related Work

Stable Constitutions A burgeoning body of work, originating in papers by Barbera & Jackson [9] and Koray [39], studies the concept of stability and self-choice of voting methods. The main focus of [9] is on what we call qualified majority rules (ranging from Simple Majority to Unanimity). They consider only anonymous functions, and thus centralization and inequality are not part of the discussion. The interest in operational: Which voting rules are likely to emerge, given that simple majority is more effective (easier to get motions through), but qualified majority is more cautious. They find that when the voting rule and the constitution (meaning the rule that decides how to change the voting rule) are the same, then simple majority is the unique stable rule (Theorem 1 of [9]). This is analogous to our Theorem 13 of Appendix A. Our work gives a darker twist to this result: Seemingly, it is an encouraging result, as it shows democracy is stable. However, without the anonymity restriction (Section 4), the dynamics result in convergence to a dictatorship. This suggests that the convergence to simple majority comes from disregard to minority rights, which are more protected under qualified majority. Unlike us, [9] allow for voting rules and constitutions to be different (e.g., to make a constitutional amendment, you would need a $\frac{2}{3}$ majority, but for a regular vote you need a simple majority). We do not consider this distinction for two reasons: (i) It is not common in DAOs, (ii) We believe that motions passed through regular vote can effectively

change the constitution without doing so explicitly, making the actual constitution a dead letter. Nevertheless, this is an interesting direction for future work.

[39] takes a different approach than [9], and considers general SCFs (i.e., with rankings over alternatives, and not necessarily binary). They show that an axiom of self-selectivity that requires that for any preference profile, and with a freedom to decide tie-breaking, a voting rule should choose itself rather than an alternative rule, given that voters rank the voting rules according to their underlying rankings (and how the voting rules perform on these rankings). It shows that together with commonly used axioms (neutrality, unanimous), this axiom is equivalent to Independence of Irrelevant Alternative (IIA). This, of course, is very significant since it implies through Arrow’s theorem [2] that when there are at least 3 alternatives, the only universally self-selective, neutral and unanimous rules are dictatorial. Moreover, [36] show that this implication extends to the binary case as well. [39]’s notion of universal self-selectivity tightly relates to our “pessimistic” approach. Our most important divergence from [39] is that we avoid imposing the neutrality axiom (generally, we avoid any axiom other than self-maintenance itself). For constitution design, it is natural to consider the status-quo as different than the proposed motion. All qualified majority rules are not neutral: They are biased towards the status-quo. For example, an implication of [36, 39] is that within our pessimistic approach with status-quo tie-breaking, only dictatorships are self-maintaining. However, we show that a non-neutral method we call consensus-duopoly (that has better welfare and fairness performance than a dictatorship) is also self-maintaining.

Our modeling of voting is as a general Boolean function, but a more restricted (and largely accurate) model is that of a Weighted Voting Game (WVG), where voters have weights (tokens), and a motion passes if and only if voters with over a threshold T of the tokens support it. [5] study an extension of [9] (which, recall, only consider anonymous voting rules: WVGs are not necessarily anonymous) to WVGs. They characterize the stable rules as the ones with veto players (each with an equal weight), and normal players (also, each with an equal weight to each other). For a motion to pass, it must get the support of all veto players and some amount of the normal players. This characterization can not work in our model, and we believe it is for a good reason. Consider Unanimity (where all voters have to agree), which fits this characterization. In our model, both in the pessimistic and i.i.d. case, Unanimity would be replaced by simple or qualified majority, due to it being ineffective. Voters would prefer to risk losing sometimes, over having a deadlocked system. This is supported by significant precedents. In the US, the Articles of Confederation [18] required all 13 states to support a motion for it to pass. However, it was soon superseded by the US constitution [51], which only requires a qualified majority of 9 out of 13 states. More recently, Article 48 of the Treaty on European Union [25] originally required decision by Unanimity by the European Council. However, the Treaty of Lisbon [24] moved many policy areas to be decided by a qualified majority. Here is an emblematic quote, from the classic textbook on EU Law [19], page 133: “Whether Unanimity is the best protector of national sovereignty depends therefore on whether a state believes that maximizing the possibility of inaction [...] is better for the national interest than a qualified-majority voting rule which increases the possibility of action”. They give a concrete example from the UK: “The Unanimity rule [...] impeded the market liberalization desired by the Conservative Party, [...] hence the willingness to sacrifice the veto for the enhanced possibility of Community action”.

[35] considers stability of voting rules when voters aim to maximize their power indices in simple games, which are more general than WVGs. When applied to WVGs, their results greatly simplify and shows that (with their additional axioms) only dictatorships are minimally stable. We take a similar approach in Section 5, but also incorporate “Condorcet Jury Theorem” effects, which allows us to have a more nuanced outcome with (possibly) multiple stable rules, including a large committee.

More broadly, the works on *Institution Design* are highly relevant to us, most prominently Acemoglu & Robinson [1]. Our model differs in key aspects: (i) In our theoretical results of Sections 3, 4, all agents are apriori the same, and we do not make an apriori distinction between “rich” and “poor”. This changes in Section 5, which aims to connect more with [1]. (ii) All power in our setting is formal. There is no “De Facto” power in the form of rebellion, and change only comes from within the system itself, following its power structure. (iii) We consider a large variety of governance forms, rather than a binary state of oligarchy / democracy.

An important inspiration for this work is Kenneth Arrow’s renowned Impossibility Theorem [2]. This foundational theorem of social choice, shows that any SCF that satisfies efficiency (in the Pareto sense), unrestricted domain, and independence of irrelevant alternatives (IIA), must be a dictatorship. While this is clearly a gravely negative result, it emphasizes a different tension than ours: It puts forward *desired* properties of a voting system, and shows they can not all co-exist. Our approach aims to keep the clean social choice framing and use it for understanding institution dynamics.

The Condorcet Jury Theorem [20] is often cited in support of democracy as the optimal mechanism for “wisdom of the crowd”. The probabilistic argument shows that if voters have i.i.d. access to signals regarding the best decision, with $p > \frac{1}{2}$ accuracy, then the best aggregator is taking a majority over each of the voters’ signal. We make a nuanced application of it to SCF stability in Section 5: We show that this serves as a force to have more inclusive stable SCFs, somewhat analogous to the “Rebellion” [1]. A direct connection can be made through the concept of “Epistemic Proceduralism” [23], which posits that citizens contribute to the system in proportion to their belief in the quality of its decision-making. Thus, having a good “Condorcet Jury” will result in higher participative contribution by the citizenry, or in the case of DAOs, future engagement and revenue for the system.

DAO Governance and Centralization Research on Blockchain Decentralized Autonomous Organization (DAOs) takes different perspectives. One perspective is their effect of the larger blockchain ecosystem: How they may be used to facilitate strategic attacks (“Dark DAOs”) [3], affect bidding patterns, due to their group nature [6], how to protect minority opinion to maintain user retention [32], or quantifying the possible “Anscombe Paradox” effects of following the majority in many sequential decisions [30].

The focus that is most relevant to our study is that of *centralization* and security of DAOs [8, 29]. Centralization is generally a major concern in Blockchains [7, 15, 40]. A string of works [26, 31, 52] define and study centrality measures for DAOs, and generally find that DAOs are more centralized than what we would like to believe. [27] finds that in four major DAOs (Uniswap, Lido, Aave and Compound), if the top 3 or 4 token holders are in consensus, they could together decide almost every vote. They name this metric *minimal quorum*, and it is also commonly known as the Nakamoto coefficient [31, 45]. [38] find that in over 7.5% of DAOs, the contributors (who develop and maintain the DAO) have majority control, and have solely dictated a decision in over 21% of DAOs. Moreover, they find that token power balance (i.e., the voting structure) has shifted shortly before a decision in over 14% of motions. Similarly, [17] find abnormal trading activity in tokens around governance proposals. [10] find that not only are DAOs centralized, but are becoming more centralized with time. Our work offers a social-choice micro-foundation for all these phenomenons.

1.2 Our Approach and Contribution

We start by defining agent *belief* as a distribution over possible voting vectors, corresponding to the agent’s assessment of future decisions’ distribution. We then define the agent’s ex-ante utility as the probability that a decision by the SCF, over a voting vector sampled from their belief, will

match their preference. We say that an SCF f is self-maintaining if for any other SCF f' , if voters according to their ex-ante preference between the two SCFs, and the vote is aggregated using f , then the final decision is to not replace f . Since tie-breaking over how agents vote in case f and f' have the same ex-ante utility can make a notable difference, we consider two tie-breaking rules: *Arbitrary*, that considers all possible agent votes in this case, *status-quo bias* (SQB), that always votes to keep f .

To have a robust understanding of self-maintenance, we also consider three approaches regarding beliefs:

- *Optimistic*. In the optimistic approach, we consider an SCF to be self-maintaining if there *exists* a common belief that supports it.
- *Pessimistic*. In the pessimistic approach, we only consider an SCF to be self-maintaining if *all* common beliefs support it.
- *I.i.d.* We assume either a biased or unbiased i.i.d. common belief, and consider an SCF to be self-maintaining if it is supported by it. Concretely, it means each of the agent's preference w.r.t. a decision is a Bernoulli variable with some parameter p .

These approaches correspond to standard epistemic approaches in game theory: The optimistic approach corresponds to Rationalizability [12, 46], meaning whether we can find a belief supporting a certain rule. All rules not supported are deemed extremely unlikely. The pessimistic approach corresponds to Robustness [11], meaning whether the rule can remain stable under any condition. I.i.d. is a middle ground where we assume independence, which is a natural case.

In our results, we combine the different approaches with different tie-breaking rules, to derive the following results:

- With an optimistic approach, almost every SCF of interest can be supported by some common belief to be self-maintaining. In particular, we use what we call a lexicographic belief. This holds true regardless of the tie-breaking rule.
- With a pessimistic approach, the results depend on the tie-breaking rule: With arbitrary tie-breaking, *no* SCF is self-maintaining, but with SQB tie-breaking, several SCFs are self-maintaining, including, but not limited to, dictatorships. We present necessary and sufficient conditions in this case, with some gap left between the two.
- In the unbiased i.i.d. case, we prove an “Arrow-Style” Theorem: We show that with arbitrary tie-breaking, only dictatorships are self-maintaining, and any other SCF has a path of changes that ends up in a dictatorship. With SQB tie-breaking, we find that the SCF we call “consensus-duopoly”, is self-maintaining and has better efficiency (social welfare) and fairness (Nash welfare) than a dictatorship.

Lastly, in Section 5, we extend our model to consider “Condorcet Jury Theorem” effects of information aggregation. To remain tractable, we limit the voting rules to oligarchies, ranging from a dictatorship (oligarchy of 1) to full democracy (oligarchy of n). The important parameters of this model are how much of the utility comes from taking advantage of existing resources, vs. making correct decisions (λ), and how good are agents at knowing the correct decision (p). We demonstrate that the stable oligarchy size is decreasing in both parameters, and that there may be multiple stable oligarchies with significant size differences.

2 Model

We consider a setting with $n \geq 3$ agents.

Definition 1. A social-choice function (SCF) $f : \{0, 1\}^n \rightarrow \{0, 1\}$ is a mapping from n binary preferences into a binary decision.

For any n , and any n -agent vector v , we use the notation $\eta(\mathbf{v}, x) = |\{i | v_i = x\}|$ to mean the number of votes for x in \mathbf{v} .

Definition 2. The Unanimity social-choice function decides 1 if and only if all agents vote 1. I.e.,

$$f_{\text{unanimity}}(\mathbf{v}) = \begin{cases} 1 & \forall 1 \leq i \leq n, v_i = 1 \\ 0 & \text{Otherwise.} \end{cases}$$

The simple-majority social-choice function decides 1 if and only if a majority of the agents vote 1.

$$f_{\text{sm}}(\mathbf{v}) = \begin{cases} 1 & \eta(\mathbf{v}, 1) > \eta(\mathbf{v}, 0) \\ 0 & \text{Otherwise.} \end{cases}$$

The qualified q -majority social choice function decides 1 if and only if a q fraction of the agents vote for 1 (where we assume $q > \frac{1}{2}$).

$$f_{\text{maj}}^q(\mathbf{v}) = \begin{cases} 1 & \eta(\mathbf{v}, 1) > q \cdot n \\ 0 & \text{Otherwise.} \end{cases}$$

The i^* -dictatorship social-choice function decides 1 if and only if agent i votes for 1.

$$f_{\text{dict}}^{i^*}(\mathbf{v}) = \begin{cases} 1 & v_{i^*} = 1 \\ 0 & \text{Otherwise.} \end{cases}$$

The i^* -anti-dictatorship social-choice function decides 1 if and only if agent i votes for 0.

$$f_{\text{anti-dict}}^{i^*}(\mathbf{v}) = \begin{cases} 0 & v_{i^*} = 1 \\ 1 & \text{Otherwise.} \end{cases}$$

The S -oligarchy social-choice function decides 1 if and only if a majority of agents in S vote for 1. Let \mathbf{v}_S be the vector \mathbf{v} restricted only to agents in S , then:

$$f_{\text{olig}}^S(\mathbf{v}) = \begin{cases} 1 & \eta(\mathbf{v}_S, 1) \\ 0 & \text{Otherwise.} \end{cases}$$

We now introduce a “self-referential” method of reasoning about SCFs to consider whether an SCF prevents itself from being changed. To allow this self-reference, it is necessary to define agent utilities w.r.t. different SCFs, and a pre-requisite for that is to define their beliefs. Our starting point is thus defining beliefs:

Definition 3. Agent belief. An agent belief is a probability distribution F over the discrete space $\Omega = \{0, 1\}^n$. I.e., for any $v \in \Omega$,

$$Pr_{\hat{v} \sim F}[\hat{v} = v] \geq 0,$$

and

$$\sum_{v \in \Omega} Pr_{\hat{v} \sim F}[\hat{v} = v] = 1.$$

Remark 1. We choose to talk about “beliefs” rather than “individual preferences”, but it is possible to assign the belief formulation an interpretation where it encodes an agent’s normalized ex-ante preferences over different possible decisions. For example, if the agent knows of an upcoming important decision, where it knows all the agents’ YES/NO preferences, then its belief can assign the full 1 weight over this voting vector, regardless of whether it knows of other upcoming votes. In this sense, it makes sense to talk about individual and incompatible beliefs. However, our results mostly consider common beliefs, where all agents have the same distribution, simply because for the optimistic case, this is a stricter requirement (and still, we show that all reasonable SCFs are self-maintaining). For the pessimistic case, this is a looser requirement (and still we show that almost no SCFs are self-maintaining).

Definition 4. Let F be the agent’s belief.

Agent i ’s ex-ante utility of an SCF f takes the following form:

$$u_i(f, F) = E_{v \sim F}[1[f(v) = v_i]] = Pr_{v \sim F}[f(v) = v_i].$$

Agent i ’s preferred choice among two SCFs f_1, f_2 is:

$$c_i(f_1, f_2, F) = \arg \max_{f \in \{f_1, f_2\}} u_i(f, F) \quad (1)$$

If F is clear from the context, we simply write $u_i(f)$ and $c_i(f, f')$.

It is natural to extend the notion of utility to that of welfare:

Definition 5. $SW(f, F) = \sum_{i=1}^n u_i(f, F)$.

The above rule for agent’s choice (Eq. 1) does not clarify how to deal with the case when $u_i(f_1) = u_i(f_2)$. We consider two possible approaches to it.

- *Arbitrary Tie-Breaking.* With this approach, an agent may vote either For or Against a motion when they are agnostic to the winner.
- *Status-Quo Bias (SQB).* With this approach, an agent votes Against a motion when the agent is agnostic to the winner.

Notice that having $c_i(f_1, f_2)$ for every agent naturally defines a vector function $\mathbf{c}(f_1, f_2) = (c_1(f_1, f_2), \dots, c_n(f_1, f_2))$. To be compatible with prior notation, we treat c as a function from $\{0, 1\}^n \rightarrow \{0, 1\}^n$, by encoding f_1 as 0 and f_2 as 1.

Arbitrary tie-breaking is more *restrictive* w.r.t. which functions are considered self-maintaining. That is because agents with agnostic preferences have the freedom to vote in different ways, and thus *expand* the space of preference outcomes, which in turn increases the possibility of coinciding with such preference vectors that result in changing the voting rule. As we will see, indeed there are settings with voting rules that are self-maintaining under *status-quo bias* but not with *arbitrary tie-breaking*.

Definition 6. We say an SCF f is self-maintaining for belief F if for any other SCF f' ,

$$f(\mathbf{c}(f, f', F)) = 0,$$

i.e., any motion to replace f with another function f' , where agents vote according to their ex-ante preference, and the vote is aggregated using f fails.

Our self-maintenance condition is asymmetric in nature, in that it fixes the label for maintaining f at 0 and to replace it at 1. Usually with social choice alternatives, neutrality w.r.t. the labels makes sense, since the labeling of candidates is arbitrary, and indeed notable theorems such as May's theorem use this condition. However, since our motivating setting is DAOs, or alternatively, government forms in societies, there is a natural distinction between the status quo (the current inner state of the smart contract, the current set of laws), and the suggested change. So non-neutral SCFs such as *Unanimity* make sense in this setting. To allow that, we do not require the *neutrality* axiom.

3 Optimistic and Pessimistic Approach to Stability

3.1 Optimistic Approach

In this subsection, we say an SCF f is self-maintaining (omitting the qualifier “for belief F ”) if there is some common belief F so that f is self-maintaining for belief F .

Definition 7. We fix an ordering of the voting vectors in Ω : $v^1, \dots, v^{|\Omega|}$ so that $v^{|\Omega|-1} = \mathbf{0}$, $v^{|\Omega|} = \mathbf{1}$, and define the following distribution F :

$$Pr_{\hat{v} \sim F}[\hat{v} = v^i] = \begin{cases} \frac{1}{2^i} & i < |\Omega| \\ \frac{1}{2^{|\Omega|-1}} & i = |\Omega| \end{cases}.$$

We call F “lexicographic” because it puts as much weight on a voting profile v^i as on all subsequent voting profiles $v^{i'}$ with $i' > i$ combined.

Definition 8. Vector negation. We say a vector $\neg v$ is the negation of vector v if for any $1 \leq i \leq n$, $\neg v_i = 1 - v_i$. We say that an SCF is never negation-agnostic if there is no $v \in \Omega$ so that $f(v) = f(\neg v) = 1$.

Theorem 1. In the optimistic approach with arbitrary tie-breaking, any SCF f is self-maintaining if and only if it is never negation-agnostic.

Proof. (Characterization \implies self-maintaining)

Consider some other SCF f' and let i^* be the first index i where $f'(v^i) \neq f(v^i)$. For every agent j so that $v_j^{i^*} = f(v^{i^*})$, we have $u_j(f) \geq u_j(f')$, since F is lexicographic. For the same reason, for every agent j so that $v_j^{i^*} = f'(v^{i^*})$, we have the opposite: $u_j(f') \geq u_j(f)$. By arbitrary tie-breaking, we can choose $c_j(f, f') = 0$ for the former type, and $c_j(f, f') = 1$ for the latter. This yields $\mathbf{c} \in \{v^{i^*}, \neg v^{i^*}\}$. More specifically, if $f(v^{i^*}) = 0$, it is v^{i^*} , and if $f(v^{i^*}) = 1$, it is $\neg v^{i^*}$. In the former case, we conclude that $f(\mathbf{c}(f, f')) = f(v^{i^*}) = 0$, and f is self-maintaining. In the latter case, by our theorem's assumption that $f(v) = 1 \implies f(\neg v) = 0$, and since $f(v^{i^*}) = 1$, we conclude that $f(\mathbf{c}(f, f')) = f(\neg v^{i^*}) = 0$, and f is self-maintaining.

(Does not follow characterization \implies not self-maintaining)

Consider f so that there is some v with $f(v) = f(\neg v) = 1$. Let F be any distribution over Ω .

If $v \notin \text{supp} F$, then consider f' so that $f'(v) = \neg f(v)$, and for all other v' , $f'(v') = f(v')$, then $\forall 1 \leq j \leq n, u_j(f') = u_j(f)$. Therefore, by arbitrary tie-breaking, we can choose any \mathbf{c} , and in particular $\mathbf{c}(f, f') = v$, and have

$$f(\mathbf{c}(f, f')) = f(v) = 1,$$

i.e., f is not self-maintaining for F .

If $v \in \text{supp}F$, then consider the same construction of f' . We have that $u_j(f') > u_j(f)$ whenever $v_j = 0 \neq 1 = f(v)$, and $u_j(f') < u_j(f)$ whenever $v_j = 1 = f(v)$. It thus holds that $\mathbf{c}(f, f') = \neg v$, and so $f(\mathbf{c}(f, f')) = f(\neg v) = 1$, and f is not self-maintaining for F .

Since this covers every possible F , we conclude that f is not self-maintaining. \square

On the positive side, since arbitrary tie-breaking is a stricter tie-breaking form than SQB, any SCF characterized as self-maintaining is also self-maintaining with SQB. On the flip side, there may be SCFs that are not *never negation-agnostic*, and are self-maintaining with SQB tie-breaking. These SCFs seem less interesting to characterize. For example, all SCFs in Example 2 are never negation-agnostic.

3.2 Pessimistic Approach

In this subsection, we say an SCF f is self-maintaining if for *any* common belief F , f is self-maintaining for belief F .

We start by considering arbitrary tie-breaking.

Theorem 2. *In the pessimistic approach with arbitrary tie-breaking, only the constant 0 SCF is self-maintaining.*

Proof. Clearly, $f = 0$ is self-maintaining as regardless of voting profile \mathbf{c} , the alternative SCF f' , the distribution F , and the tie-breaking choice induce, we have $f(\mathbf{c}) = 0$.

Otherwise, it must be that there is some $v^* \in \Omega$ with $f(v^*) = 1$. Consider F so that $\Pr_{\hat{v} \sim F}[\hat{v} = v^*] = 1$. Then, for some $v \neq v^*$, consider f' so that $f'(v) \neq f(v)$, and for all other vectors v' , $f'(v') = f(v')$. Then, for all agents i ,

$$u_i(f') = u_i(f)$$

(as the only difference in outcome is w.r.t. a vector v which is not in the support of F). Thus by arbitrary tie-breaking we can choose $\mathbf{c}(f, f') = v^*$, and have $f(\mathbf{c}(f, f')) = f(v^*) = 1$. \square

However, this argument seems pretty arbitrary. For example, for a dictator to give up their dictator power in f_{dict}^i (which regardless of the distribution F gives them the optimal utility of 1), they must have very high trust in the new system of governance, and also favor change despite having perfect utility under f_{dict}^i .¹

We thus consider our second tie-breaking rule, status-quo bias. We provide both a necessary condition, and a sufficient condition (a class of SCFs) that partially characterize pessimistic self-maintenance in this setting. To provide intuition, we name and shortly discuss the different conditions.

Definition 9. *For a voting profile v_i , we let $S_i = \{j | v_j^i = 1\}$ be the set of agents that support the decision under v_i .*

Downward closed. *For any $v^1, v^2 \in \Omega$ with $v^1 \neq v^2$, consider v so that $v_j = 1$ if and only if $j \in S_1 \cap S_2$. Then, if $f(v^1) = f(v^2) = 0$, then so is $f(v) = 0$. In words, if two voting profiles are*

¹ Given a long enough history, this can happen. See, for example, the Roman emperor Diocletian, who is known as the only emperor to voluntarily abdicate his title [33]. Diocletian was sickly, and has been experimenting with delegating his authority throughout his rule. It stands to reason that his belief was that foregoing power would preserve the stability of the decision-process he had built, and “tie-breaking” favored it due to his sickness. A similar dynamic is presented theatrically in Shakespeare’s “King Lear” [49], where the king abdicates his rule in favor of his three beloved daughters, due to his old age. In both cases (Lear and Diocletian), the beliefs driving the abdication decision turned out to be unsubstantiated.

ruled as against under the SCF f , then the voting profile where the voters supporting the motion are the intersection of supporters on both other profiles is also ruled against.

Respects rejective consensus. If $v = \mathbf{0}$ is the voting profile where all voters are against a decision, $f(v) = 0$.

Bounded Monotonicity Violation. For any two voting profiles $v^1, v^2 \in \Omega$ so that $S_1 \subset S_2$, and $f(v^1) = 1, f(v^2) = 0$, then for any v^3 so that $(S_2 \setminus S_1) \cap S_3 = \emptyset, S_1 \subset S_3$, we have $f(v^3) = 1$. In words, if the SCF violates monotonicity for v^1, v^2 , in the sense that v^1 has a strict subset of supporting voters but is accepted while v^2 is not, then any other v^3 that extends v^1 (in terms of supporting voters) is accepted.

Bounded Downward Sensitivity. For any v^1 with $|S_1| \geq 3$ and $f(v^1) = 1$, there are at most 2 vectors v^2, v^3 with $S_2 \subset S_1, S_3 \subset S_1$ and $|S_2| = |S_3| = |S_1| - 1$ so that $f(v^2) = f(v^3) = 0$. In words, if a voting profile with k voters in favor of the motion passes under f , then removing one voter should also pass, with the exception of possibly two voters that when removed fail the vote.

Strong Duo. There is $v^1 \in \Omega$ with $|S_1| = 2$ so that $f(v^1) = 1$. In words, there are two voters so that if they are in favor of a decision and all others are against, the decision passes under f .

The list of definitions starts with more benign properties of an SCF, and ends with two properties that seem undesired. All of the properties are necessary.

Theorem 3 (Necessary Conditions for Pessimistic Self-Maintenance with SQB tie-breaking).

Any SCF f that is not the constant 0 and is self-maintaining is never negation agnostic, downward-closed, respects rejective consensus, has bounded monotonicity violation, bounded downward sensitivity, and a strong duo.

Proof. (Never negation-agnostic) Notice that this is the necessary (and sufficient) condition in Theorem 1. However, we now consider on one hand a pessimistic (rather than optimistic) approach, which is stricter, but, on the other hand, we use SQB tie-breaking rather than arbitrary tie-breaking, which is looser.

If there is some v so that $f(v) = f(\neg v) = 1$, then consider F so that $Pr_{\hat{v} \sim F}[\hat{v} = v] = 1$. Let f' so that $f'(v) = 0 \neq 1 = f(v)$, and for any other v' , $f'(v') = f(v')$. Then for any agent i with $v_i = 0$, we have $u_i(f) = 0 < 1 = u_i(f')$, and for any agent i with $v_i = 1$, we have $u_i(f) = 1 > 0 = u_i(f')$, and so $\mathbf{c}(f, f') = \neg v$. Then, $f(\mathbf{c}(f, f')) = f(\neg v) = 1$, in contradiction to self-maintenance.

(Downward-closed) Consider some v^1, v^2, v and their respective sets of indices where the agent votes 1, $S_1, S_2, S_1 \cap S_2$. We have that $f(v^1) = f(v^2) = 0$. Let F so that $Pr_{\hat{v} \sim F}[\hat{v} = v^1] = Pr_{\hat{v} \sim F}[\hat{v} = v^2] = \frac{1}{2}$. Let f' be an SCF with $f'(v^1) = f'(v^2) = 1$, and for any other v' , $f'(v') = f(v')$. Then, for any agent i in $S_1 \cap S_2$, the agent has $v_i^1 = v_i^2 = 1$, and thus $u_i(f) = 0 < 1 = u_i(f')$, and thus $c_i(f, f') = 1$. For any agent i in $S_1 \setminus S_2$ (or, symmetrically, in $S_2 \setminus S_1$), we have $v_i^1 = 1, v_i^2 = 0$, and so $u_i(f) = \frac{1}{2} = u_i(f')$, and by SQB tie-breaking, $c_i(f, f') = 0$. For any agent i in $[n] \setminus (S_1 \cup S_2)$, we have $v_i^1 = v_i^2 = 0$, and so $u_i(f) = 1 > 0 = u_i(f')$, and $c_i(f, f') = 0$. Overall, $f(v) = f(\mathbf{c}(f, f')) = 0$ by the self-maintenance condition.

(Respects rejective consensus) Consider F so that for some $v \in \Omega$, $Pr_{\hat{v} \sim F}[\hat{v} = v] = 1$, and consider f' so that for some $\tilde{v} \neg v$, $f'(\tilde{v}) = 1 - f(\tilde{v})$, and for any other v' , $f'(v') = f(v')$ (in particular, $f'(v) = f(v)$). Then, for any agent i , $u_i(f) = u_i(f')$, and by SQB tie-breaking, $\mathbf{c}(f, f') = \mathbf{0}$. We then have by self-maintenance that $f(\mathbf{0}) = 0$.

(Bounded monotonicity violation) We have $S^3 \cap S^2 = S^1$, so by downward-closed, if $f(v^3) = f(v^2) = 0$, then so must be $f(v^1) = 0$, in contradiction. This is essentially a different way to state downward-closed, that emphasizes the implications of too many instances of non-monotonicity of the outcome in the number of agents voting 1.

(Bounded downward sensitivity) Assume towards contradiction there are 3 vectors v^2, v^3, v^4 with $|S_2| = |S_3| = |S_4| = |S_1| - 1 \geq 2$ and $f(v^2) = f(v^3) = f(v^4)$. Consider F so that $Pr_{\hat{v} \sim F}[\hat{v} = v^2] = Pr_{\hat{v} \sim F}[\hat{v} = v^3] = Pr_{\hat{v} \sim F}[\hat{v} = v^4] = \frac{1}{3}$, and f' so that $f'(v^2) = f'(v^3) = f'(v^4) = 1$, and for any other v' , $f'(v') = f(v')$. Then for any agent $i \in S_1$, $u_i(f) \leq \frac{1}{3}$ (it is $\frac{1}{3}$ if for some $k \in \{2, 3, 4\}$, $S_1 \setminus S_k = \{i\}$, and 0 otherwise), and $u_i(f') \geq \frac{2}{3}$. For any other agent i , $u_i(f) = 1 > 0 = u_i(f')$. Thus, $f(c(f, f')) = f(v^1) = 1$, in contradiction to self-maintenance.

(Strong Duo) Since $f \neq 0$, there must be some v^1 so that $f(v^1) = 1$. If for this v^1 , $|S_1| \geq 2$, we show the condition holds by reverse induction, using (V). I.e., there are $|S_1|$ different possible vectors \tilde{v}^i so that $|\tilde{S}_i| = |S_1| - 1$, and at most two of them have $f(\tilde{v}^i) = 0$, so if $|S_1| > 2$, there must be a vector with one less agent voting 1 so that $f(v) = 1$, and so on.

If for v^1 we have $|S_1| < 2$, it must be that $|S_1| = 1$ by respecting rejective consensus. Then, if we assume towards contradiction there are no vectors v with $|S| = 2$ and $f(v) = 1$, then consider the bounded monotonicity violation condition. Let i be the index of the agent voting 1 in v^1 , and let j be some other index, and let v^2 so that $S_2 = \{i, j\}$. Then any v^3 that has $i \in S_3, j \notin S_3$, has $f(v^3) = 1$ (by bounded monotonicity violation), in particular, if we take for some $k \notin \{i, j\}$ the vector v^3 so that $S_3 = \{i, k\}$.

□

Next, we provide a class of SCFs that are self-maintaining under the pessimistic approach with status-quo bias tie-breaking. To provide intuition before we embark on formal definitions and lemmas, consider an SCF where two voters that always get their way when they are in agreement. Then, as a pair, this SCF is optimal for them, in the sense that there is no SCF f' they both strictly prefer. It can not benefit them when they are in agreement, so it can only benefit one on the expense of the other when they are in disagreement. Thus, if we also always need their mutual agreement to move away from f , we would not never get it (under SQB tie-breaking). This holds, for example, for an oligopoly of three, where decisions are made according to the majority of some given three voters. We identify a general class where this applies.

Definition 10. We say two agents i, j are entangled under f if for any v so that $v_i = v_j$, $f(v) = v_i$.

The class of entangled SCFs f are these where a subset $S(f)$ of the set of all agent pairs are entangled. We say a vector v is entangled under f if it has $v_i = v_j$ for some $(i, j) \in S(f)$.

We say an entangled SCF f has a flower form if there is some agent i that appears in all pairs of $S(f)$.

A 3-oligopoly SCF decides every voting profile following the majority among voters in a set S with $|S| = 3$.

Lemma 1. Every entangled SCF is either the 3-oligopoly, or has a flower form.

Proof. Consider an entangled SCF which is not a 3-oligopoly or has flower form. Since it does not have flower form, it means that for every entangled agent i there is a pair $(i', j') \in S$ so that i is not in the pair.

We claim that it must be that there is a pair (i, j) so that both agents are different from another pair (i', j') .

Assume otherwise. There is some pair that contains i , since it is an entangled agent. If the pairs (i, i') and (i, j') are in S , then this is a 3-oligopoly. If i is part of two pairs in S but not the above, then at least one of the pairs has two agents different than (i', j') (namely, i and the agent it is paired with). Otherwise, (w.l.o.g.) (i, i') is the unique pair containing i in S . Since the SCF is not flower form, there is some pair that does not contain i' , and both agents are not i, i' (since (i, i') is the only pair where i appears).

With two distinct pairs $(i, j), (i', j')$ of entangled agents, consider a vector v with $v_i = v_j = 1, v_{i'} = v_{j'} = 0$. Then, since $(i, j) \in S$ and they are in agreement over v , it must be that $f(v) = 1$. But since $(i', j') \in S$ and they are in agreement over v , it must be that $f(v) = 0$, leading to a contradiction. \square

Consider a set \hat{S} of pairs that has flower form, i.e., there is some single agent i so that i is a member of each pair. Notice that \hat{S} does not uniquely determine an SCF f , since if i is in disagreement with all the other agents in the flower form, then the outcome is not restricted. Therefore, there is a class of entangled SCFs f so that $S(f) = \hat{S}$. We say that f is the *conservative* member of this class if for any $v \in \Omega$ so that $\forall j \in (\bigcup \hat{S}) \setminus \{i\}, v_i \neq v_j$, we have $f(v) = 0$, i.e., in the case of disagreement we specified, the SCF always decides against. It follows from the definitions that this uniquely defines an SCF:

Lemma 2. *For a set S of pairs that has flower form, there is a unique conservative entangled SCF.*

Theorem 4. *The 3-oligopoly, and every conservative flower-form entangled SCF, are self-maintaining.*

Proof. We start by showing that 3-oligopoly is self-maintaining in the pessimistic approach with SQB tie-breaking. Let S be the oligopoly agents.

The key property is that for any F and $i, j \in S$, $u_i(f) + u_j(f) = \max_{f'} u_i(f') + u_j(f')$. This holds because the outcome follows i and j consensus whenever such exists. Thus, for any f' and $i \in S$ so that $u_i(f) > u_i(f')$, it must be that for any $j \in S, j \neq i$, $u_j(f) < u_j(f')$. Thus, the choice vector $\mathbf{c}(f, f')$ will always have a majority 0 among the agents in S , and so we must have $f(\mathbf{c}(f, f')) = 0$, and so self-maintenance is satisfied.

The argument to show that a conservative flower form entangled SCF f is self-maintaining is very similar: If we let i be the unique agent that appears in all pairs $S(f)$ (due to the flower form), then the only vectors v so that $f(v) = 1$ are ones where $v_i = 1$ and $v_j = 1$ for some $(i, j) \in S(f)$, by the conservative form requirement. However, for any F to have $u_i(f) > u_i(f')$ (in order for $c_i(f, f') = 1$) must imply $u_j(f) < u_j(f')$ and $c_j(f, f') = 0$, and thus by self-maintenance, $f(\mathbf{c}(f, f')) = 0$. \square

We show two concrete examples for the conservative flower form entangled SCFs (For $n = 3$, these are all such SCFs):

Definition 11. *The consensus-duopoly SCF is, for some two agents $i \neq j$,*

$$f_{\text{cons-duo}}^{i,j}(v) = \begin{cases} 1 & v_i = v_j = 1 \\ 0 & v_i = v_j = 0 \\ 0 & v_i \neq v_j. \end{cases}$$

I.e., if the two agents are in agreement that the motion should pass, it passes, and otherwise no.

Definition 12. *The 3-oligopoly with veto player SCF is, for some three agents i, j, k , and where i is the special veto player,*

$$f_{\text{olig-veto}}^{i,j,k}(v) = \begin{cases} 1 & v_i = 1 \wedge (v_j = 1 \vee v_k = 1) \\ 0 & \text{Otherwise.} \end{cases}$$

I.e., the SCF follows the majority among the agents i, j, k , unless $v_i = 0$, in which case it outputs 0 regardless of the other agents' values.

Consensus-Duopoly corresponds to $S = \{(i, j)\}$. 3-oligopoly with veto player corresponds to $S = \{(i, j), (i, k)\}$.

To conclude our discussion, we show a gap between our necessary and sufficient conditions, i.e., an SCF outside the class of Theorem 4 that is self-maintaining.

Example 1. With $n = 3$, consider the following SCF:

$$f(v) = \begin{cases} 1 & v \in \left\{ \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right\} \\ 0 & \text{Otherwise.} \end{cases}$$

It is clear that f is not an entangled SCF since $f(\mathbf{1}) = 0$, so no two agents can be entangled. However, it is self-maintaining, since both vectors so that $f(v) \neq v_1$ has $v_1 = v_3$. Thus, for any $f' \neq f$, it holds that $u_1(f') - u_1(f) \leq u_3(f') - u_3(f)$ (changing any of the two vectors with $f(v) = 1$ would add an equal amount to both u_1 and u_3 , and a change in any other vector will subtract at least as much from u_1 as it does from u_3). To have $c_1(f, f') = 1$, we must have $u_1(f') - u_1(f) > 0$, and so we must have $u_3(f') - u_3(f) > 0$ and $c_3(f, f') = 1$. However, for any vector with $v_1 = v_3 = 1$ we have $f(v) = 0$, satisfying the self-maintenance condition.

4 The I.i.d Approach

In this section, we consider any belief F to be a joint distribution of an i.i.d. distribution over $\{0, 1\}$. I.e., for some $F_p = \text{Bernoulli}(p)$, $F = F_p \times \dots \times F_p$. If $p > \frac{1}{2}$, i.e., the probability of each voter to prefer 1 in a randomly drawn decision exceeds $\frac{1}{2}$, we say the voters are *change-inclined*. If $p = \frac{1}{2}$, we say the voters are *unbiased*. If $p < \frac{1}{2}$, we say the voters are *change-averse*.

Consider i.i.d. preference voters with parameter p . First, we can identify the following simple self-maintaining rules:

Theorem 5. *With i.i.d. preference voters and arbitrary tie-breaking, any i -dictatorship, i -anti-dictatorship, and the constant 0 are self-maintaining.*

Proof. The constant SCF f is self-maintaining, as regardless of the agents' preferences of replacing it with another SCF f' , it outputs 0, satisfying the self-maintenance constraint.

We show that i -dictatorship is self-maintaining.

We can directly calculate:

$$u_i(f_{dict}^i) = \Pr_{v_j \sim F_p}[f_{dict}^i(v) = v_i] = 1.$$

Since for any f , 1 is a trivial upper bound over $u_i(f)$ (as it is a probability of some event), and, moreover, i -dictatorship is the unique SCF that has $u_i(f) = 1$, we have $c_i(f_{dict}^i, f') = 0$, for any other f' . Then, regardless of $\mathbf{c}_{-i}(f_{dict}^i, f')$ (since the i -dictatorship output only depends on v_i):

$$f_{dict}^i(\mathbf{c}(f_{dict}^i, f')) = 0.$$

The proof that i -anti-dictatorship is self-maintaining is similar: Since it is the unique SCF that has $u_i(f) = 0$, we have $c_i(f_{anti-dict}^i, f') = 1$, and thus

$$f_{anti-dict}^i(\mathbf{c}(f_{anti-dict}^i, f')) = 0.$$

□

Let us now partition the agents into three subsets, depending on their utility level under f :

$$S_1 = \{i | u_i(f) < 2p(1-p)\}, S_2 = \{i | 2p(1-p) \leq u_i(f) < p^2 + (1-p)^2\}, S_3 = \{i | u_i(f) \geq p^2 + (1-p)^2\}.$$

Then, we can characterize a necessary condition for any SCF other than the ones in Theorem 5.

Lemma 3 (Main Structural Lemma). *With i.i.d. preference voters with parameter p , any self-maintaining SCF f that is not a dictatorship, anti-dictatorship, or the constant 0, must satisfy the following property: For every vector v that has $v_i = 1$ for any $i \in S_1 \cup S_2$, $f(v) = 0$. For any vector v that has $v_i = 0$ for any $i \in S_2 \cup S_3$, $f(v) = 0$.*

Proof. First, notice that for any $0 \leq p \leq 1$, $2p(1-p) \leq p^2 + (1-p)^2$, as this inequality is a rearrangement of $0 \leq (1-2p)^2$. Thus, the three sets constitute a full partition of the agents.

Consider v so that $\forall i \in S_1 \cup S_2, v_i = 1$, and $\forall i \in S_3, v_i = 0$. Then, it must be that $f(v) = 0$. Otherwise, if $f(v) = 1$, consider f' to be the i^* -dictatorship for some agent in $S_1 \cup S_2$. It has $u_{i^*}(f') = 1$, and $\forall i \in S_1 \cup S_2 \setminus \{i^*\}, u_i(f') = Pr_{v \sim F}[v_i = f(v)] = Pr_{v_i \sim F_p, v_{i^*} \sim F_p}[v_i = v_{i^*}] = \sum_{x=0}^1 Pr_{v_i \sim F_p}[v_i = v_{i^*} | v_i = x] \cdot Pr_{v_i \sim F_p}[v_i = x] = p^2 + (1-p)^2$. Therefore, $\mathbf{c}_i(f, f') = 1$ for any $i \in S_1 \cup S_2$, as they have $u_i(f) < u_i(f')$, and for any $i \in S_3$, $\mathbf{c}_i(f, f') = 0$, as $u_i(f) \geq u_i(f')$. We then have $f(\mathbf{c}(f, f')) = f(v) = 1$, violating self-maintenance.

Consider v so that $\forall i \in S_1 \cup S_2, v_i = 1$, for some $i^* \in S_3, v_{i^*} = 1$, and for all other $i \neq i^*$ so that $i \in S_3, v_i = 0$. Then, it must be that $f(v) = 0$. Otherwise, if $f(v) = 1$, consider f' to be the i^* -dictatorship. It has $u_{i^*}(f') = 1$, and $\forall i \neq i^*, u_i(f') = p^2 + (1-p)^2$. Therefore, $\mathbf{c}_i(f, f') = 1$ for any $i \in S_1 \cup S_2$, $\mathbf{c}_{i^*}(f, f') = 1$, and $\mathbf{c}_i(f, f') = 0$ for any $i \in S_3 \setminus \{i^*\}$. I.e., $f(\mathbf{c}(f, f')) = f(v) = 1$, violating self-maintenance.

We now show by induction that for any v so that $\forall i \in S_1 \cup S_2, v_i = 1, f(v) = 0$. The induction is on the number k of agents $i \in S_3$ with $v_i = 1$. The $k = 0, 1$ base cases were directly shown before. Now consider that $f(v) = 0$ for all vectors v with $k' \leq k$ for some $k \geq 1$, and we wish to show that this holds for any vector v with $k' = k + 1$. Consider v so that $v_i = 1$ for all $i \in S_1 \cup S_2$, and there is $S' \subseteq S_3$ of size $k + 1$ so that $\forall i \in S', v_i = 1$, and $\forall i \in S_3 \setminus S', v_i = 0$. Let $i_1, i_2 \in S'$ be two different agents in S' . Define v' so that $v'_i = 1$ for all $i \in S_1 \cup S_2 \cup (S' \setminus i_1)$, and $v'_i = 0$ otherwise, and v'' so that $v''_i = 1$ for all $i \in S_1 \cup S_2 \cup (S' \setminus i_2)$, and $v''_i = 0$ otherwise. By the induction assumption, $f(v') = f(v'') = 0$. Assume towards contradiction that $f(v) = 1$. Consider f' which is identical to f , besides $f'(v') = f'(v'') = 1$. Then, the utility of all agents in $S_1 \cup S_2 \cup (S' \setminus \{i_1, i_2\})$ improves (as they have $v'_i = v''_i = 1$, and $f'(v') = f'(v'') = 1$). The utility of agents i_1, i_2 does not change:

$$\begin{aligned} & u_{i_1}(f') \\ &= u_{i_1}(f) - Pr_{v \sim F}[v = v'] + Pr_{v \sim F}[v = v''] \\ &= u_{i_1}(f) - p^{|S_1|+|S_2|-1+k}(1-p)^{|S_3|+1-k} + p^{|S_1|+|S_2|-1+k}(1-p)^{|S_3|+1-k} \\ &= u_{i_1}(f), \end{aligned}$$

and similarly for $u_{i_2}(f')$. The utility of agents $S_3 \setminus S'$ declines (as they have $v'_i = v''_i = 0$, and $f'(v') = f'(v'') = 1$). Overall, by arbitrary tie-breaking, we can choose $\mathbf{c}(f, f') = v$, and have $f(\mathbf{c}(f, f')) = f(v) = 1$, in contradiction to self-maintenance. We thus conclude that $f(v) = 0$.

To prove the second part of the lemma statement, we apply a very similar argument, now using *anti-dictatorship* as the base for our induction.

Consider v so that $\forall i \in S_2 \cup S_3, v_i = 0$, and $\forall i \in S_1, v_i = 1$. Then, it must be that $f(v) = 0$. Otherwise, if $f(v) = 1$, consider f' to be the i^* -anti-dictatorship for some agent in $S_2 \cup S_3$. It has $u_{i^*}(f') = 0$, and $\forall i \in S_2 \cup S_3 \setminus \{i^*\}, u_i(f') = Pr_{v \sim F}[v_i = f(v)] = Pr_{v_i \sim F_p, v_{i^*} \sim F_p}[v_i \neq v_{i^*}] =$

$\sum_{x=0}^1 Pr_{v_{i^*} \sim F_p}[v_i \neq v_{i^*} | v_i = x] \cdot Pr_{v_i \sim F_p}[v_i = x] = 2p(1-p)$. Therefore, $\mathbf{c}_i(f, f') = 1$ for any $i \in S_1$, as they have $u_i(f) < u_i(f')$, and for any $i \in S_2 \cup S_3$, $\mathbf{c}_i(f, f') = 0$, as $u_i(f) \geq u_i(f')$. We then have $f(\mathbf{c}(f, f')) = f(v) = 1$, violating self-maintenance.

Consider v so that $\forall i \in S_2 \cup S_3, v_i = 0$, for some $i^* \in S_1, v_{i^*} = 0$, and $\forall i \in S_1 \setminus \{i^*\}, v_i = 1$. Then, it must be that $f(v) = 0$. Otherwise, if $f(v) = 1$, consider f' to be the i^* -anti-dictatorship. It has $u_{i^*}(f') = 0$, and $\forall i \in S_2 \cup S_3, u_i(f') = 2p(1-p)$. Therefore, $\mathbf{c}_i(f, f') = 1$ for any $i \in S_1 \setminus \{i^*\}$, as they have $u_i(f) < u_i(f')$, and for any $i \in S_2 \cup S_3 \cup \{i^*\}$, $\mathbf{c}_i(f, f') = 0$, as $u_i(f) \geq u_i(f')$. We then have $f(\mathbf{c}(f, f')) = f(v) = 1$, violating self-maintenance.

We now show by induction that *every* vector v so that $\forall i \in S_2 \cup S_3, v_i = 0$, must have $f(v) = 0$. We do so by induction over k , where k is the number of agents in S_1 with $v_i = 0$. What we proved so far constitutes the base cases of $k = 0, 1$, showing that any v with up to 1 agent in S_1 that has $v_i = 0$, has $f(v) = 0$. We now show that if this holds for some k , then it also holds for $k + 1$. Take a vector v that has for some $S' \subseteq S_1$ of size $k + 1$, $\forall i \in S_1 \setminus S', v_i = 1$, and $\forall i \in S' \cup S_2 \cup S_3, v_i = 0$. Now choose two different $i_1, i_2 \in S'$, and consider v' so that $v'_{i_1} = 1$, and all other elements are equal to v . Consider v'' so that $v''_{i_2} = 1$, and all other elements are equal to v . Then, v', v'' both have k agents in S_1 with $v'_i = 0$. Thus, by the induction assumption $f(v') = f(v'') = 0$. Consider f' so that it is identical to f , but has $f'(v') = f'(v'') = 1$. The utility of all agents in $S_1 \setminus S' \cup \{i_1, i_2\}$ improves (as they have $v'_i = v''_i = 1$, and $f'(v') = f'(v'') = 1$). The utility of agents i_1, i_2 does not change:

$$\begin{aligned} u_{i_1}(f') &= u_{i_1}(f) - Pr_{v \sim F}[v = v'] + Pr_{v \sim F}[v = v''] \\ &= u_{i_1}(f) - p^{|S_1|-k}(1-p)^{|S_2|+|S_3|+k} + p^{|S_1|-k}(1-p)^{|S_2|+|S_3|+k} \\ &= u_{i_1}(f), \end{aligned}$$

and similarly for $u_{i_2}(f')$.

The utility of agents $i \in S_2 \cup S_3 \cup S' \setminus \{i_1, i_2\}$ declines (as they have $v'_i = v''_i = 0$, and $f'(v') = f'(v'') = 1$). Overall, by arbitrary tie-breaking, we can choose $\mathbf{c}(f, f') = v$, and have $f(\mathbf{c}(f, f')) = f(v) = 1$, in contradiction to self-maintenance. We thus conclude that $f(v) = 0$. \square

Theorem 6. *For unbiased voters, an SCF f is self-maintaining if and only if it is either a dictatorship, anti-dictatorship, or the constant 0.*

Proof. We recall that *unbiased* voters are such that have $p = \frac{1}{2}$. In this case, $2p(1-p) = p^2 + (1-p)^2$, and so $S_2 = \emptyset$.

Theorem 5 shows that f being one of the specified rule types results in it being self-maintaining. We thus consider, going forward, that f is some other SCF.

Let v be any vector in $\{0, 1\}^n$. We use the characterization of Lemma 3 to show that $f(v) = 0$. Consider the following two vectors: Let v' be equal to v for all agents in S_1 , but have $v'_i = 0$ for all $i \in S_2 \cup S_3$. Let v'' be equal to v for all agents in $S_2 \cup S_3$, but have $v''_i = 1$ for all $i \in S_1$. Consider f' so that it is identical to f , but has $f'(v') = f'(v'') = 1$.

The utility of all agents in S_1 that have $v_i = 1$ increases, as $v'_i = v''_i = 1$, and we have $c_i(f, f') = 1$. The utility of all agents in $S_2 \cup S_3$ that have $v_i = 0$ decreases, as $v'_i = v''_i = 0$, and we have $c_i(f, f') = 0$.

For an agent in S_1 that has $v_i = 0$, the utility does not change. That is since it has $v'_i = 0, v''_i = 1$, and so:

$$u_i(f') = u_i(f) + Pr_{\hat{v} \sim F}[\hat{v} = v'] - Pr_{\hat{v} \sim F}[\hat{v} = v''] = u_i(f) + \frac{1}{2^n} - \frac{1}{2^n} = u_i(f).$$

For an agent in $S_2 \cup S_3$ that has $v_i = 1$, the utility does not change. That is since it has $v'_i = 0, v''_i = 1$, and so:

$$u_i(f') = u_i(f) - Pr_{\hat{v} \sim F}[\hat{v} = v'] + Pr_{\hat{v} \sim F}[\hat{v} = v''] = u_i(f) - \frac{1}{2^n} + \frac{1}{2^n} = u_i(f).$$

Overall, we conclude that by arbitrary tie-breaking we can set $\mathbf{c}(f, f') = v$, and so by the self-maintenance constraint, must have $f(v) = f(\mathbf{c}(f, f')) = 0$. □

We see the emergence of anti-dictatorship in the characterization as a quirk of our formalism. In the discussion section, we show why they are essentially equivalent to dictatorships: In essence, if the anti-dictator knows the rule is to always go against their preference, they can always misstate their preference and arrive at a dictatorship.

We next go beyond the symmetric assumption that voters are *unbiased*, i.e., $p = \frac{1}{2}$. We show that more self-maintaining SCFs emerge once we relax this assumption. We fully characterize the case of $n = 3$. Interestingly, all self-maintaining SCFs with $n = 3$ satisfy that their *welfare*, as well as *Nash welfare*, is at most that of a dictatorship. It is plausible that this property holds with general n , but for now we leave it as a conjecture.

Definition 13. *The social welfare of an SCF f with distribution F is $SW(f) = \sum_{i=1}^n u_i(f)$.*

The Nash welfare of an SCF f with distribution F is $NW(f) = \pi_{i=1}^n u_i(f)$.

Social welfare is a standard way to measure the *efficiency* of an outcome, while the Nash welfare [37] has proven to be a good measurement for the *fairness* of the outcome [14].

Claim. In a single-item auction, for any \mathbf{v}, \mathbf{b} , the joint utility is at most equal to the highest valuation: $u_{joint}(\mathbf{b}; \mathbf{v}) \leq \max \mathbf{v}$.

Theorem 7. *With $n = 3$, and i.i.d. voters with $0 \leq p \leq 1$, there are self-maintaining SCFs outside the characterization of Theorem 6. However, there is no non-constant self-maintaining SCF with a higher welfare, or Nash welfare, than a dictatorship.*

We now consider our second tie-breaking rule, that of *status-quo bias*. We can show a construction of an SCF for general n , in particular $n = 3$ as well, that is self-maintaining and has higher social and Nash welfare than a dictatorship.

Recall Definition 11 of the consensus-duopoly SCF.

Theorem 8. *Consensus-duopoly is self-maintaining with SQB tie-breaking for any $0 \leq p \leq 1$. Moreover, it has better social and Nash welfare than a dictatorship with change-averse and unbiased voters.*

We remark that as a function of n , there are some *change-inclined* values of p where the social and Nash welfare of consensus-duopoly surpass that of a dictatorship, as well.

5 Incorporating Common Value

We now present a simplified model, where the class of voting rules are a sequence of oligarchies O_1, \dots, O_n . An oligarchy O_i is composed of a set of S_i voters with $|S_i| = i$, and $S_{i-1} \subset S_i$. The vote is decided by a majority of the set S_i , where a motion needs a strict majority to pass. This model is meant to study whether a dictatorship (O_1) would prefer to open up to an oligarchy (O_i

with $1 < i < n$) or a full democracy (O_n), and vice versa. We now specify the voters' utilities, which is the main point of divergence from the model presented in Section 2. The voter utility has two components, *extractive* E and *participative* P , weighted by a parameter $0 \leq \lambda \leq 1$ so that $u_j(O_i) = \lambda \cdot E_j(O_i) + (1 - \lambda) \cdot P_j(O_i)$. The extractive component simply considers whether the voter is part of the oligarchy:

$$E_j(O_i) = \frac{1}{i} \cdot 1[j \leq i].$$

For the participative part, we draw from the insight of the Condorcet Jury Theorem. We consider that there is a true state s , and voters have an i.i.d. signal with $p > \frac{1}{2}$ w.r.t. it. This means the voters' signal vector v is drawn from the binomial distribution with parameter p .

Thus,

$$P_j(O_i) = \Pr_{v \sim \text{Binomial}(p, i)}[|v_j = s|_j > \frac{i}{2}],$$

is the probability that the right decision is made under oligarchy O_i . The self-maintenance condition for oligarchy i given parameters λ, p is that for any $i' \neq i$, a (weak) majority of the voters have $u_j(O_i) \geq u_j(O_{i'})$.

We make a few observations. With $\lambda = 1$ (a fully extractive economy), and $i \geq 2$, if we let $i' = \lceil \frac{i+1}{2} \rceil$, then all voters $j \leq i'$ prefer $O_{i'}$ over O_i , and they constitute a strict majority in O_i . Thus, the only self-maintaining oligarchies are O_1 (a dictator) and O_2 (a consensus-duopoly). With $\lambda = 0$ (a fully participative economy), the probability that the majority of a Binomial distribution matches the true state is monotone increasing in i , and therefore the only self-maintaining rule is the full democracy O_n . We thus consider a mixed case of $0 < \lambda < 1$ going forward.

To explicitly calculate the full set of self-maintaining oligarchies for given values of λ, p , we do the following. We start with a necessary condition for self-maintenance, based on the idea we applied with $\lambda = 1$. I.e., we require $u_j(O_i) \geq u_j(O_{i'})$ for the voters $j \leq i'$ with $i' = \lceil \frac{i+1}{2} \rceil$ and with $i' = 2i - 1$. This yields

$$\begin{aligned} \frac{\lambda}{i} + (1 - \lambda) &\geq \frac{\lambda}{i} + (1 - \lambda) \Pr_{v \sim \text{Binomial}(p, i)}[|v_j = s|_j > \frac{i}{2}] \\ &\geq \frac{\lambda}{i'} + (1 - \lambda) \Pr_{v \sim \text{Binomial}(p, i')}[|v_j = s|_j > \frac{i'}{2}] \stackrel{\text{Hoeffding}}{\geq} \frac{\lambda}{i'} + (1 - \lambda)(1 - e^{-2i'(p-0.5)^2}). \end{aligned}$$

It turns out these simple bounds are already quite good at filtering out large values of i . With the remaining candidates for self-maintaining oligarchies, we require:

Definition 14. *Right-Maximal.* For a function f and a value i^* , we say it is right-maximal if $\forall i' > i^*, f(i^*) \geq f(i')$.

Halving-Safe. For a function f and a value i^* , we say it is halving-safe if $\forall \frac{i^*}{2} < i' < i^*, f(i^*) \geq f(i')$.

It is clear from our discussion so far that it is necessary and sufficient for a self-maintaining oligarchy i to satisfy both properties. This allows us to precisely characterize the outcomes for given parameters, as the below figures show.

We run a Cartesian grid of parameters, with $\lambda \in \{0.1, 0.2, \dots, 0.9\}$ and $p \in \{0.6, 0.65, \dots, 0.95\}$. Let us describe the results. At the corner of low λ and low p values, we see that *the unique* solution is a large oligarchy (let us call it a committee). The size of committee is monotone-decreasing in both λ and p . Then, when we move away from the low-values corner, we get some parameter values where there are multiple self-maintaining outcomes: A large committee, a dictatorship, and also possibly a small oligarchy (with $i \in \{2, 3\}$). When λ and p values are high enough, no large committee is no

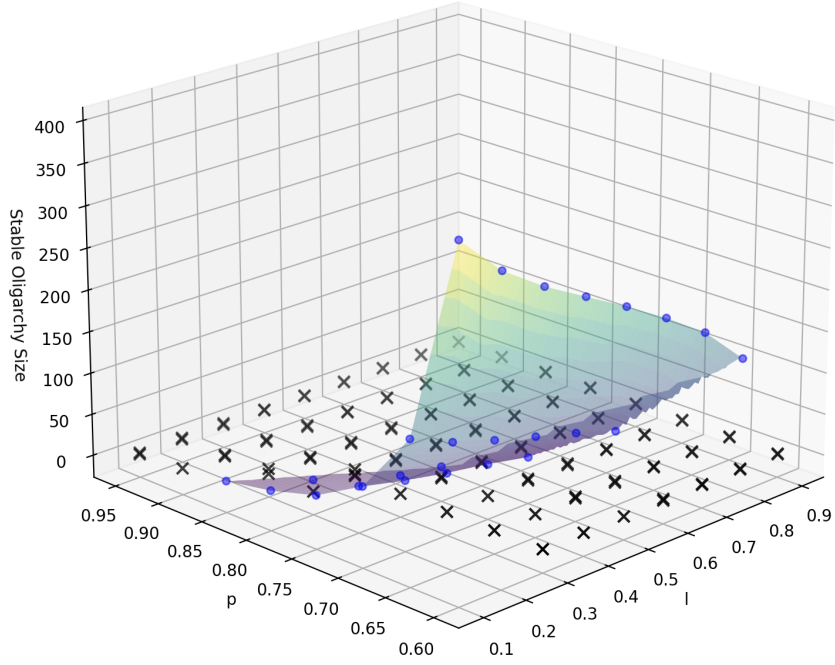


Fig. 1: Stable oligarchies as a function of λ and p . Black x marks are used for small oligarchies (such as dictatorship, consensus-duopoly and 3-oligopoly), and blue circles for larger-value oligarchies (upwards of $i = 10$).

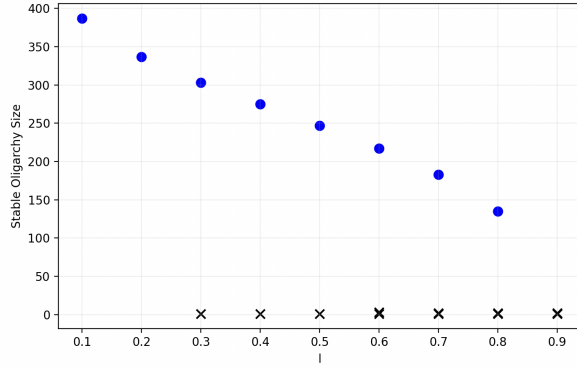
longer self-maintaining, and we remain only with the dictatorship and (possibly) a small oligarchy. As we reach the high λ , high p corner, a dictatorship is the only self-maintaining rule.

To see our description more clearly, we “slice” the graph at both $\lambda = 0.5$, and $p = 0.6$, presented in Figure 2.

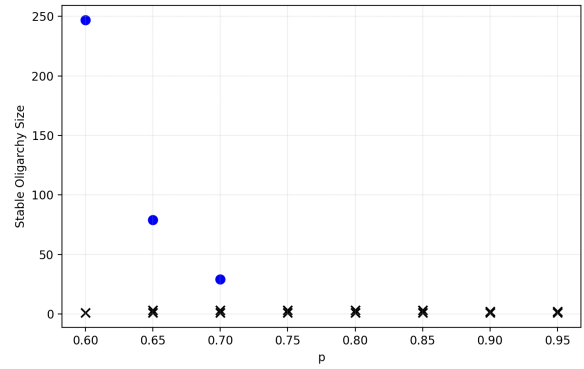
Beyond the self-maintaining points, we are also interested in dynamics. We consider $p = 0.6$ with $\lambda \in \{0.6, 0.7, 0.8, 0.9\}$, where as in Figure 2a, there are 3 equilibria (a dictatorship, a small oligarchy, and a large committee) with $\lambda \in \{0.6, 0.7, 0.9\}$, and only a dictatorship and a small oligarchy with $\lambda = 0.9$. We wish to understand where do we end up if we start in an arbitrary oligarchy. Figure 3 shows the answer.

We note the following. As λ increases, the tendency for centralization shows up not only in terms of the self-maintaining committee being of smaller i value (and disappearing as a self-maintaining solution with $\lambda = 0.9$), but also in terms of the dynamics: We have more red lines instead of purple (leading only to a small oligarchy rather also possibly to the large committee) and more purple lines instead of blue. The dictatorship is isolated and only leads to itself. We can see why for $\lambda = 0.6$, high values of i lead only to the large committee: Since the function is overall increasing after reaching the global minimum and up to the self-maintaining point at $i = 217$ ², then high enough values of $i \leq 217$ are Halving-Safe but not Right-Maximal, and so they converge to the large committee. Values of $i > 217$ are not Halving-Safe (the function is slightly decreasing in that range).

² The increase has some “waviness” depending on whether i is odd or even, due to the way tie-breaking is handled.



(a) Stable oligarchies as a function of λ with $p = 0.6$



(b) Stable oligarchies as a function of p with $\lambda = 0.5$

Fig. 2: Two-dimensional Slices of the three-dimensional graph of stable-oligarchies as a function of λ and p . Black x marks are used for small oligarchies (such as dictatorship, consensus-duopoly and 3-oligopoly), and blue circles for larger-value oligarchies.

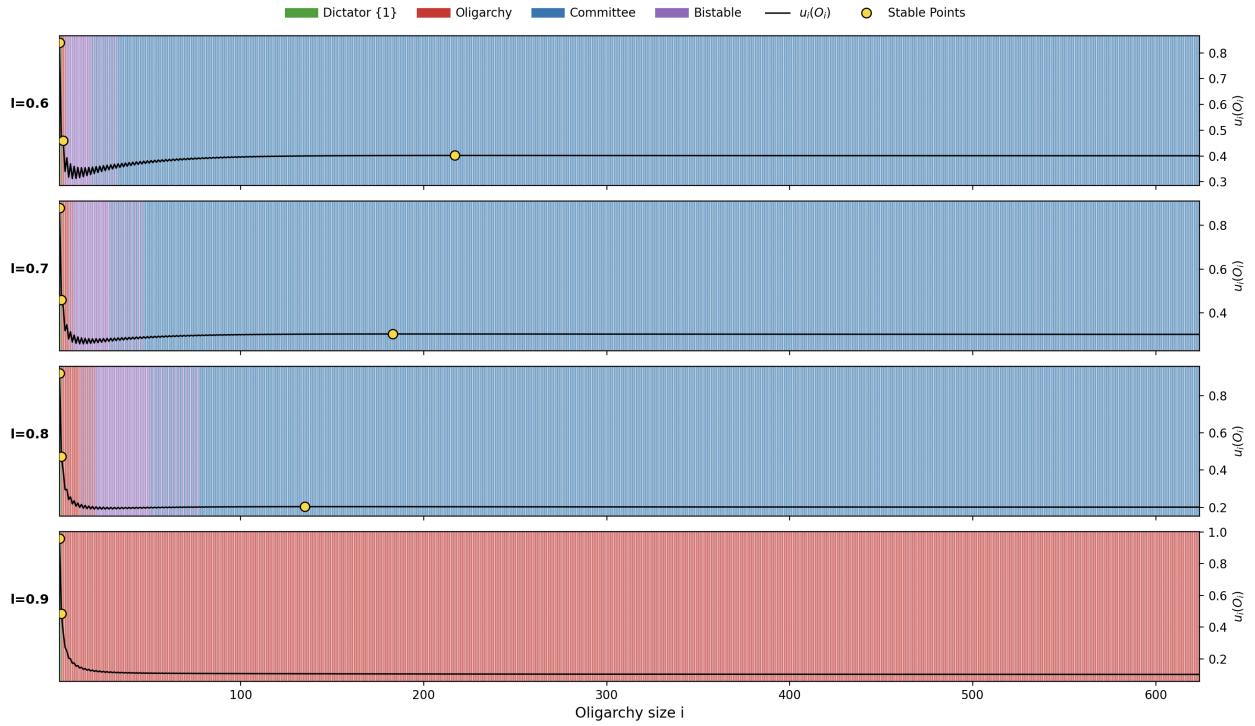


Fig. 3: Centralization dynamics for different extraction Levels ($l \in \{0.6, 0.7, 0.8, 0.9\}$). The colored ribbons represent basins of attraction for the Dictatorship (green), Oligarchy (red), and Committee (blue), with purple indicating that both Committee and Oligarchy can be reached. The black line represents the utility function $u_i(O_i)$, with stable rules marked as gold points.

6 Discussion

6.1 Anti-dictatorships and Dictatorships

In some of our results, anti-dictatorship emerges alongside dictatorship as a self-maintaining SCF. We make a small distinction that clarifies why, in essence, there is no difference between dictatorial and anti-dictatorial SCFs:

Definition 15. *Voter Best Response.* Given an SCF f , and other voters' preferences v_{-i} , a voter best response $BR(v)$ is such that:

$$\forall x \in \{0, 1\}, 1[f(BR(v), v_{-i}) = v_i] \geq 1[f(x, v_{-i}) = v_i].$$

Equilibrium SCF. Given an SCF f , we say f' is an equilibrium SCF for it, if for voter combined strategies profile e , it holds that $(e(v))_i = BR(v_i, (e(v))_{-i})$ for all agents i , and $f'(v) = f \circ e(v)$.

Theorem 9. *i -dictatorship is the unique equilibrium SCF for i -anti-dictatorship. It is also the unique equilibrium SCF for i -dictatorship.*

Proof. Fix any vector $v \in \{0, 1\}^n$. Then, $f_{anti-dict}^i(v_i, v_{-i}) = \neg v_i$, $f_{anti-dict}^i(\neg v_i, v_{-i}) = \neg(\neg v_i) = v_i$, and so $BR(v_i) = \neg v_i$. For all other agents, their votes do not influence the results, and so in any equilibrium $f'(v) = f_{anti-dict}^i \circ e(v) = f_{anti-dict}^i(\neg v_i, v_{-i}) = v_i$. Thus, $f' = f_{dict}^i$.

For the second part, again fix any vector $v \in \{0, 1\}^n$. Then, $f_{dict}^i(v_i, v_{-i}) = v_i$, and so $BR(v_i) = v_i$. For all other agents, their votes do not influence the results, and so in any equilibrium $f'(v) = f_{dict}^i \circ e(v) = f_{dict}^i(v_i, v_{-i}) = v_i$. Thus, $f' = f_{dict}^i$. \square

We conclude that if the prescribed rule is an anti-dictatorship, then effectively it becomes a dictatorship by the i anti-dictator voting opposite to their true preference. On the other hand, if it prescribed as a dictatorship, then it is incentive-compatible for all agents to report truthfully.

6.2 Dynamics vs. Statics and Voters that Plan Ahead

In our formulation, voters face a choice between two SCFs, and decide for the one where they expect a higher utility. This separates the SCFs into self-maintaining and non-self-maintaining ones. However, this is a static view, and the dynamics are also interesting. For example, it is conceivable that there is a (Polybius-style) cycle of SCFs that lead to each other. In this case, none of the SCFs would be self-maintaining, but we also do not converge to an end-state. If no such cycles exist and convergence is guaranteed, then the rate of convergence is of interest (how quickly does a government form descend into dictatorship? Are there certain SCFs that slow this convergence?)

Moreover, we have considered voters that vote to move from SCF f to another SCF f' if their utility is better under f' . However, voters may apply a more sophisticated reasoning, that is informed by our analysis. Voters may consider the following: "While f' has better utility for me than f , it is not in itself self-maintaining, and so eventually we will lead to another SCF f'' , which may not be desirable for me. Voters should then apply some form of backward induction to decide whether to switch to f' . However, this requires many modeling decisions regarding how motions to switch the SCF are generated, and we leave it as an open direction.

6.3 Abstentions, Multi-vote and Turing-complete SCFs

It is interesting to consider abstentions, as they are very common in practice, including in DAOs, where users are not required to participate in every vote, see e.g. the NounsDAO data in [32]. Abstentions may affect the results, especially when voters aim to aggregate information [28], and are important but under-considered element of voting [44]. Technically, they could help by avoiding discussion of different tie-breaking rules. On the other hand, it complicates the binary definition of the SCF.

Applying the framework to multi-vote decisions is a promising direction as well, as it could be interesting to see how famous methods like the Borda rule, Plurality, and so on fare in terms of stability. However, given our mostly negative conclusions in the binary case, it does not espouse much hope for better results there. The binary case is usually the “well-behaving” case in social choice (e.g., Gibbard-Satterthwaite’s negative conclusion applies only when there are more than two alternatives).

One thing that is unique to DAOs is that they run on top of smart contracts that are Turing-Complete. Thus, we can think of SCFs that have a very general form, much more general than matching rankings to a decision as it is commonly treated in social choice. Whether there is a way to create meaningful self-maintaining rules in this way is very interesting. One path is by restricting the set of SCFs that we may transition to: our Section 5 and Appendix A are examples for the implications of such restrictions.

6.4 Practical Takeaways.

We attempt to draw a few informal conclusions from our results, particularly of Section 5. We see that centralization is increasing with how extractive the economy is. For DAOs, this means that if the balance between *current locked-up funds* and *expected future revenue* shifts towards the locked-up funds, we could expect centralization of power. Shocks in the decision-making quality of individuals also relates to centralization. For example, if Artificial Intelligence (AI) enhances voters ability to decide, the stable outcome may be a much smaller elite size. Our results also provide an intuitive explanation why in some countries decisions are made by unelected bodies in the size-range of hundreds to thousands members. For example, in China, the National People’s Congress has about 3000 members, and the Central Committee has about 200 members. Of particular interest are our dynamics results, where we identify a critical sensitivity to initial conditions: different starting states converge to distinct stable outcomes. This suggests that achieving an inclusive steady-state requires a highly distributive initial token allocation. This can be achieved through “airdrops” (free distribution of tokens to prospective holders) or Initial Coin Offerings (ICOs). Existing economic works indeed thinks of these methods as a commitment device by platforms to not exploit users [34, 48, 50]. While in our model it is not guaranteed that the platform stays as decentralized as after the airdrop due to the social-choice dynamics, it helps with the equilibrium choice of the more decentralized stable large committee. Another justification for ICOs is solving the coordination problem of users in platform [41], where their participation costs are independent of others, but their utility depends on others engaging with the platform, and so they might end up in the no-trade equilibrium. ICOs and airdrops can be used to bootstrap and select the participative equilibrium. This argument is orthogonal to our considerations.

Bibliography

- [1] D. Acemoglu and J. A. Robinson. A theory of political transitions. *American Economic Review*, 91(4):938–963, 2001.

- [2] K. J. Arrow. A difficulty in the concept of social welfare. *Journal of political economy*, 58(4): 328–346, 1950.
- [3] J. Austgen, A. Fábrega, S. Allen, K. Babel, M. Kelkar, and A. Juels. Dao decentralization: Voting-bloc entropy, bribery, and dark daos, 2023. URL <https://arxiv.org/abs/2311.03530>.
- [4] J. Austgen, A. Fábrega, S. Allen, K. Babel, M. Kelkar, and A. Juels. Dao decentralization: Voting-bloc entropy, bribery, and dark daos, 2023. URL <https://arxiv.org/abs/2311.03530>.
- [5] Y. Azrieli and S. Kim. On the self-(in)stability of weighted majority rules. *Games and Economic Behavior*, 100:376–389, 2016. ISSN 0899-8256. <https://doi.org/https://doi.org/10.1016/j.geb.2016.10.010>. URL <https://www.sciencedirect.com/science/article/pii/S0899825616301257>.
- [6] M. Bahrani, P. Garimidi, and T. Roughgarden. When Bidders Are DAOs. In J. Bonneau and S. M. Weinberg, editors, *5th Conference on Advances in Financial Technologies (AFT 2023)*, volume 282 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 21:1–21:21, Dagstuhl, Germany, 2023. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. ISBN 978-3-95977-303-4. <https://doi.org/10.4230/LIPIcs.AFT.2023.21>. URL <https://drops.dagstuhl.de/entities/document/10.4230/LIPIcs.AFT.2023.21>.
- [7] M. Bahrani, P. Garimidi, and T. Roughgarden. Centralization in block-building and proposer-builder separation. In J. Clark and E. Shi, editors, *Financial Cryptography and Data Security*, pages 331–349, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-78676-1.
- [8] S. Ballesti, P. Saggese, S. Kitzler, and B. Haslhofer. Slaying the dragon: The quest for democracy in decentralized autonomous organizations (daos), 2025. URL <https://arxiv.org/abs/2511.09263>.
- [9] S. Barbera and M. O. Jackson. Choosing how to choose: Self-stable majority rules and constitutions. *The Quarterly Journal of Economics*, 119(3):1011–1048, 2004.
- [10] T. Barbereau, R. Smethurst, O. Papageorgiou, J. Sedlmeir, and G. Fridgen. Decentralised finance’s timocratic governance: The distribution and exercise of tokenised voting rights. *Technology in Society*, 73:102251, 2023.
- [11] D. Bergemann and S. Morris. Robust mechanism design. *Econometrica*, 73(6):1771–1813, 2005. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/3598751>.
- [12] B. D. Bernheim. Rationalizable strategic behavior. *Econometrica*, 52(4):1007–1028, 1984. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/1911196>.
- [13] V. N. Brovkin. *The Mensheviks after October: Socialist Opposition and the Rise of the Bolshevik Dictatorship*. Cornell University Press, 1987.
- [14] I. Caragiannis, D. Kurokawa, H. Moulin, A. D. Procaccia, N. Shah, and J. Wang. The unreasonable fairness of maximum nash welfare. *ACM Trans. Econ. Comput.*, 7(3), Sept. 2019. ISSN 2167-8375. <https://doi.org/10.1145/3355902>. URL <https://doi.org/10.1145/3355902>.
- [15] N. Chemaya, A. Yaish, S. Yacouel, D. Malkhi, and L. W. Cong. Quantifying inequality in blockchain networks. *Available at SSRN 5540258*, 2025.
- [16] E. CHENOWETH and M. J. STEPHAN. *Why Civil Resistance Works: The Strategic Logic of Nonviolent Conflict*. Columbia University Press, 2011. URL <http://www.jstor.org/stable/10.7312/chen15682>.
- [17] L. W. Cong, D. Rabetti, C. C. Wang, and Y. Yan. Centralized governance in decentralized organizations. *Available at SSRN 5168660*, 2025.
- [18] Continental Congress. Articles of confederation and perpetual union. National Archives and Records Administration, 1777. URL <https://www.archives.gov/milestone-documents/articles-of-confederation>. Adopted November 15, 1777; ratified March 1, 1781.
- [19] P. Craig and G. De Búrca. *EU law: text, cases, and materials*. Oxford University Press, USA, 2011.

- [20] M. de Condorcet. *Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix*. Imprimerie Royale, Paris, 1785. URL https://archive.org/details/bub_gb_RzAVAAAAQAAJ.
- [21] DeepDAO. Deepdao page for oceandao, 2025. Accessed: 2025-12-25.
- [22] W. Doyle. *The Oxford history of the French revolution*. Oxford University Press, 2018.
- [23] D. Estlund. Epistemic proceduralism and democratic authority. In *Does Truth Matter? Democracy and Public Space*, pages 15–27. Springer, 2009.
- [24] European Union. Treaty of Lisbon amending the Treaty on European Union and the Treaty establishing the European Community — Article 1, point 56 (Article 48 TEU). Official Journal of the European Union, C 306, 2007. URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A12007L%2FTXT>.
- [25] European Union. Consolidated version of the treaty on european union — Article 48, 2016. URL http://data.europa.eu/eli/treaty/teu_2016/art_48/oj. The article detailing the ordinary and simplified revision procedures.
- [26] A. Fabrega, A. Zhao, J. Yu, J. Austgen, S. Allen, K. Babel, M. Kelkar, and A. Juels. Voting-Bloc entropy: A new metric for DAO decentralization. In *34th USENIX Security Symposium (USENIX Security 25)*, pages 1299–1318, Seattle, WA, 2025. USENIX Association. ISBN 978-1-939133-52-6. URL <https://www.usenix.org/conference/usenixsecurity25/presentation/fabrega-entropy>.
- [27] B. Falk, T. Pathan, A. Rigas, and G. Tsoukalas. Blockchain governance: an empirical analysis of user engagement on daos. *arXiv preprint arXiv:2407.10945*, 2024.
- [28] T. J. Feddersen and W. Pesendorfer. The swing voter’s curse. *The American Economic Review*, 86(3):408–424, 1996. ISSN 00028282. URL <http://www.jstor.org/stable/2118204>.
- [29] R. Feichtinger, R. Fritsch, L. Heimbach, Y. Vonlanthen, and R. Wattenhofer. SoK: Attacks on DAOs. In R. Böhme and L. Kiffer, editors, *6th Conference on Advances in Financial Technologies (AFT 2024)*, volume 316 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 28:1–28:27, Dagstuhl, Germany, 2024. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. ISBN 978-3-95977-345-4. <https://doi.org/10.4230/LIPIcs.AFT.2024.28>. URL <https://drops.dagstuhl.de/entities/document/10.4230/LIPIcs.AFT.2024.28>.
- [30] R. Fritsch and R. Wattenhofer. The price of majority support. *arXiv preprint arXiv:2201.12303*, 2022.
- [31] R. Fritsch, M. Müller, and R. Wattenhofer. Analyzing voting power in decentralized governance: Who controls daos? *Blockchain: Research and Applications*, 5(3):100208, 2024. ISSN 2096-7209. <https://doi.org/https://doi.org/10.1016/j.bcr.2024.100208>. URL <https://www.sciencedirect.com/science/article/pii/S2096720924000216>.
- [32] Y. Gafni and B. Golan. Beyond proportional individual guarantees for binary perpetual voting, 2024. URL <https://arxiv.org/abs/2408.08767>.
- [33] E. Gibbon. *The history of the decline and fall of the Roman Empire*, volume 1. Harper, 1833.
- [34] I. Goldstein, D. Gupta, and R. Sverchov. Utility tokens as a commitment to competition. *The Journal of Finance*, 79(6):4197–4246, 2024. <https://doi.org/https://doi.org/10.1111/jofi.13389>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/jofi.13389>.
- [35] H. Hermida-Rivera. Minimal stable voting rules. *Games and Economic Behavior*, 153:541–553, 2025. ISSN 0899-8256. <https://doi.org/https://doi.org/10.1016/j.geb.2025.07.006>. URL <https://www.sciencedirect.com/science/article/pii/S0899825625000995>.
- [36] H. Hermida-Rivera and T. T. Kerman. Binary self-selective voting rules. *Journal of Public Economic Theory*, 27(3):e70039, 2025. <https://doi.org/https://doi.org/10.1111/jpet.70039>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/jpet.70039>.
- [37] M. Kaneko and K. Nakamura. The nash social welfare function. *Econometrica*, 47(2):423–435, 1979. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/1914191>.

- [38] S. Kitzler, S. Ballester, P. Saggese, B. Haslhofer, and M. Strohmaier. The governance of decentralized autonomous organizations: A study of contributors’ influence, networks, and shifts in voting power. In J. Clark and E. Shi, editors, *Financial Cryptography and Data Security*, pages 313–330, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-78679-2.
- [39] S. Koray. Self-selective social choice functions verify arrow and gibbard-satterthwaite theorems. *Econometrica*, 68(4):981–996, 2000.
- [40] A. Levy, S. M. Weinberg, and C. Zhou. Analyzing the economic impact of decentralization on users. *arXiv preprint arXiv:2512.11739*, 2025.
- [41] J. Li and W. Mann. Digital tokens and platform building. *The Review of Financial Studies*, 38(7):1921–1954, 2025.
- [42] C. Lord et al. *Aristotle’s politics*. University of Chicago Press, 2013.
- [43] B. C. McGing. *Polybius’ Histories*. Oxford University Press, 2010.
- [44] R. Meir. Tyranny of the minority in social choice: a call to arms. In *Proceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*, pages 2865–2869, 2025.
- [45] C. Ovezik, D. Karakostas, M. Milad, D. W. Woods, and A. Kiayias. Sok: Measuring blockchain decentralization. In *Applied Cryptography and Network Security: 23rd International Conference, ACNS 2025, Munich, Germany, June 23–26, 2025, Proceedings, Part I*, page 184–214, Berlin, Heidelberg, 2025. Springer-Verlag. ISBN 978-3-031-95760-4. https://doi.org/10.1007/978-3-031-95761-1_7. URL https://doi.org/10.1007/978-3-031-95761-1_7.
- [46] D. G. Pearce. Rationalizable strategic behavior and the problem of perfection. *Econometrica*, 52(4):1029–1050, 1984. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/1911197>.
- [47] C. D. Reeve et al. Plato: Republic. *Hackett, Indianapolis*, 2004.
- [48] M. Reuter. Platform Precommitment via Decentralization. *IMF Working Papers*, 2024(028):1, 2 2024. ISSN 1018-5941. <https://doi.org/10.5089/9798400267284.001>. URL <http://dx.doi.org/10.5089/9798400267284.001>.
- [49] W. Shakespeare. *King lear*, volume 5. Classic Books Company, 2001.
- [50] M. Socking and W. Xiong. Decentralization through tokenization. *The Journal of Finance*, 78(1):247–299, 2023. <https://doi.org/https://doi.org/10.1111/jofi.13192>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/jofi.13192>.
- [51] U.S. Constitution. Constitution of the united states. National Archives and Records Administration, 1787. URL <https://www.archives.gov/founding-docs/constitution>. Signed September 17, 1787; effective March 4, 1789.
- [52] A. Yaish, N. Chemaya, D. Malkhi, and L. W. Cong. Inequality in the age of pseudonymity. In *Proceedings of the Fortieth AAAI Conference on Artificial Intelligence and Fortieth Conference on Innovative Applications of Artificial Intelligence and Eighteenth Symposium on Educational Advances in Artificial Intelligence, AAAI’26/IAAI’26/EAAI’26*. AAAI Press, 2026.

A Restricting to Anonymous SCFs

In many natural settings, such as Parliaments, only *anonymous* SCFs are considered. I.e., we may require that a decision reaches consensus, or a qualified majority, or a simple majority, but we do not give preferential treatment to votes by a certain party or another ³. We give an overview of how our results translate to this case. For convenience, we only consider odd n (so that we never have ties). Overall, the conclusion is that simple majority plays a similar role to dictatorship, once

³ This is only true up to a “first approximation”. There are many cases with special roles for different voters, such as the veto countries in the UN security council, or the US vice-president having a tie-breaking vote in the senate.

anonymity is introduced, in the sense that voters would move away from “checks and balances” that they may view as prohibitive.

Definition 16. *Anonymous SCF.* We say f is anonymous if for any permutation π over the agents, and any $v \in \Omega$, $f(\pi(v)) = f(v)$.

Monotone SCF. We say f is monotone if for any two $v^1, v^2 \in \Omega$ so that $S_1 \subseteq S_2$, if $f(v^1) = 1$ then $f(v^2) = 1$. I.e., a decision may not be rejected because more agents vote for it.

There are exactly n anonymous and monotone SCFs: Anonymous implies that we can encode the SCF as a function $\tilde{f}: [n] \rightarrow \{0, 1\}$, where the input is the number of agents who vote 1. Then, adding monotone on top of it implies that the SCF takes the form of a threshold function: Some $1 \leq k \leq n$ so that for a vector v with a set S of agents with $v_i = 1$, we have $f(v) = 1[|S| \geq k]$.

It is straight-forward to redefine our different approaches in this setting. Instead of considering the optimistic and pessimistic approach w.r.t. to *all* SCFs, we restrict the space only to anonymous and monotone SCFs.

Theorem 10. *Limited to anonymous and monotone SCFs, in the optimistic approach with arbitrary tie-breaking, any SCF f is self-maintaining if and only if it is a qualified q -majority with $q \geq \frac{1}{2}$.*

This is not a direct corollary of Theorem 1 (the general optimistic case), because we are restricted only to anonymous and monotone functions. However, using the same proof idea, the same characterization (self-maintaining if and only if it is never negation-agnostic) applies. Our characterization of anonymous and monotone SCFs implies that any such SCF is a q -majority with some $0 \leq q \leq 1$, and being never negation-agnostic further means $q \geq \frac{1}{2}$.

Theorem 11. *Limited to anonymous and monotone SCFs, in the pessimistic approach with arbitrary tie-breaking, no SCF is self-maintaining besides $f = 0$.*

Proof. Proof Sketch. This repeats the same construction as in the general non-anonymous case: We fix F to be some v w.p. 1, and consider f' that has the same outcome as f over v , but not elsewhere. \square

Theorem 12. *Limited to anonymous and monotone SCFs, in the pessimistic approach with SQB tie-breaking, an SCF $f \neq 0$ is self-maintaining if and only if it is simple-majority in the case $n = 3$.*

Proof. As we later see in the i.i.d. case, simple-majority rule is preferred over all other rules for (e.g.) the unbiased i.i.d. belief, even with SQB tie-breaking. Thus, it is left to show that simple-majority is self-maintaining if and only if $n = 3$.

If $n = 3$, then simple-majority is the same as 3-oligopoly, and Theorem 4 shows 3-oligopoly is self-maintaining in the pessimistic approach with SQB tie-breaking in the general case, so in particular this holds for all anonymous and monotone SCFs.

Otherwise, consider k so that $n = 2k - 1$. k corresponds to the minimal amount of votes that constitute a majority. Thus, for any v with k agents who vote 1, $f_{sm}(v) = 1$. Consider the following three vectors: v^1 so that $v_i^1 = 0$ if and only if $i \leq k - 1$, v^2 so that $v_i^2 = 0$ for all $i \leq k - 2$, and for $i = k$ (and no other indices). v^3 so that $v_i^3 = 0$ if and only if $i \geq k + 1$. It holds that $f_{sm}(v^1) = f_{sm}(v^2) = f_{sm}(v^3) = 1$. It also holds that for any $1 \leq i \leq k$, they have $v_i = 0$ for two of the three vectors. Thus, if we choose the belief F so that $\forall 1 \leq i \leq 3, Pr_{\hat{v} \sim F}[\hat{v} = v^i] = \frac{1}{3}$, and f' be the qualified majority that requires at least $k + 1$ votes for a motion to pass, then $f'(v^1) = f'(v^2) = f'(v^3) = 0$, and $\forall 1 \leq i \leq k, u_i(f_{sm}) = \frac{1}{3} < \frac{2}{3} = u_i(f')$, and so $c_i(f_{sm}, f') = 1$, and there is a majority that supports moving from f to f' , so $f_{sm}(c(f_{sm}, f')) = 1$, contradicting self-maintenance. \square

Theorem 13. *Limited to anonymous and monotone SCFs, with i.i.d. beliefs and SQB tie-breaking, for any p (unbiased, change-averse and change-inclined) an SCF is self-maintaining if and only if it is simple-majority.*

Proof. Proof-sketch. With an anonymous SCF f and i.i.d. belief F , the overall symmetry implies that $\forall 1 \leq i \leq n, u_i(f) = \frac{1}{n}SW(f, F)$. Thus, an SCF is self-maintaining if and only if it maximizes welfare for F . Since the majority vote maximizes welfare for any F (since for any $v \in \Omega$, it benefits the most agents, as it goes with the majority), and with i.i.d. F , every $v \in \Omega$ has a positive probability to occur, and so any other SCF has strictly lower welfare, we conclude that simple-majority is the unique self-maintaining SCF. \square

B Missing Proofs for Section 4

Theorem 7. *With $n = 3$, and i.i.d. voters with $0 \leq p \leq 1$, there are self-maintaining SCFs outside the characterization of Theorem 6. However, there is no non-constant self-maintaining SCF with a higher welfare, or Nash welfare, than a dictatorship.*

Proof. Consider our partition into 3 sets based on agent utilities. Then, any agent in S_1, S_2 has lower utility than an agent $j \neq i$ under i -dictatorship. Therefore, if there is more than a single agent in S_1, S_2 , we have a 1-to-1 matching from each agent to an agent that has (weakly) higher utility under dictatorship (namely the agents *not* in S_3 match with a non-dictator, and the remaining agent matches with the dictator). This then implies the theorem both for the sum (social welfare) and product (Nash welfare) of these utilities.

We thus focus our attention on the case where $|S_3| \geq 2$. If $|S_3| = 3$, then the main structural lemma (Lemma 3) implies $f(v) = 0$. We thus may assume going forward that $|S_3| = 2$. W.l.o.g., we let agent 1 be in $S_1 \cup S_2$, and agents 2 and 3 are in S_3 .

By the main structural lemma, if $v_1 = 1$, then $f(v) = 0$. If $v_2 = v_3 = 0$, then also $f(v) = 0$. This leaves three vectors with unspecified outcome under F , and thus $2^3 = 8$ possible values for f . We thus continue with a case analysis.

- If $f(001) = f(010) = 0$. Then, it must be that $f(011) = 0$. This is since we may consider $f'(001) = f'(010) = 1$, and for any other v , $f'(v) = f(v)$. We have $Pr_{v \sim F}[v = 001] = Pr_{v \sim F}[v = 010] = p(1-p)^2$, and by arbitrary tie-breaking, $\mathbf{c}(f, f') = 011$. Then, by the self-maintenance constraint, it must be that $f(\mathbf{c}(f, f')) = f(011) = 0$. We conclude that in this case $f = 0$ is constant.
- If $f(011) = 0$, it must be that $f(001) = f(010) = 0$ (and we are thus back in the previous case). Otherwise, consider if $f(001) = 1$ (the case where $f(010) = 1$ is handled symmetrically). Then, consider f' so that $f'(101) = 1, f'(011) = 1$, and for any other v , $f'(v) = f(v)$. We have by arbitrary tie-breaking $\mathbf{c}(f, f') = 001$, and thus $f(\mathbf{c}(f, f')) = f(001) = 1$, in contradiction to self-maintenance. The two cases imply we must have $f(011) = 1$, and (w.l.o.g.) $f(001) = 1$ (we may choose $f(001) = 0, f(010) = 1$, but since our characterization of f is symmetric w.r.t. agent 2 and 3 up to this point, we can w.l.o.g. choose the former).
- $f(011) = 1, f(001) = 1, f(010) = 0$. The agent utilities in this case are $(1-p)^2, (1-p)(p+p^2 + (1-p)^2), (1-p)(1+p)$, respectively. It can be directly shown that there is no p for which the sum of these utilities is higher than the dictatorship welfare, and no p for which the product of these utilities is higher than the dictatorship Nash welfare. We show that f is self-maintaining: Notice that to violate self-maintenance, it must be that f' is (weakly) preferred

by agent 3, and (weakly) not preferred by agent 1. However, the vectors that do not follow agent 3 preference are $f(111) = 0, f(101) = 0$, where it agrees with agent 1. Thus, if f' chooses 1 for this vectors, it improves agent 1 utilities. Similarly, the vectors that *do* follow agent 1 preferences are $f(000) = 0, f(010) = 0$, where agent 1 agrees with agent 3, and so worsening agent 1 utility also worsens agent 3 utility. We conclude that the only way to achieve \mathbf{c} with $c_1 = 0, c_3 = 1$ is if $u_1(f) = u_1(f')$, and $u_3(f) = u_3(f')$, but this is impossible since (other than with $p = \frac{1}{2}$), $Pr_{v \sim F}[v = 111], Pr_{v \sim F}[v = 101], Pr_{v \sim F}[v = 010], Pr_{v \sim F}[v = 000]$ are each different.

- $f(011) = 1, f(001) = 1, f(010) = 1$. In this case, f is equal to the agent 1-anti-dictatorship, besides the fact that $f(000) = 0$ while $f_{anti-dict}^1(000) = 1$. Thus, $\forall 1 \leq i \leq 3, u_i(f) = u_i(f_{anti-dict}^1) + (1 - p)^3$. It can be shown directly that both the social welfare and Nash welfare of these utilities underperform these of the dictatorship for any $0 \leq p \leq 1$. Regarding self-maintenance, we note that any v with $f(v) = 1$ has $v_1 = 0$. Thus, agent 1 needs to weakly prefer f over f' . However, $f(000) = 0$ is the only vector so that agent 1 is satisfied with f . Changing the outcome of this vector is also bad for agents 2 and 3, but we need at least one of them to weakly improve under f' since we can not have $\mathbf{c} = 000$, as then $f(\mathbf{c}) = f(000) = 0$. Also note that $Pr_{v \sim F}[v = 000]$ is unique when $p \neq \frac{1}{2}$.

□

Theorem 8. *Consensus-duopoly is self-maintaining with SQB tie-breaking for any $0 \leq p \leq 1$. Moreover, it has better social and Nash welfare than a dictatorship with change-averse and unbiased voters.*

Proof. We start by showing consensus-duopoly is self-maintaining. Let i, j be the duopoly voters. As we know from the proof of Theorem 4, for general beliefs, and in particular in our case, it holds that $u_i(f_{cons-duo}^{i,j}) + u_j(f_{cons-duo}^{i,j}) \geq u_i(f') + u_j(f')$ for any SCF f' . It is thus impossible to have both

$$u_i(f') > u_i(f_{cons-duo}^{i,j}), \quad u_j(f') > u_j(f_{cons-duo}^{i,j}). \quad (2)$$

However, to have $f_{cons-duo}^{i,j}(\mathbf{c}) = 1$ we must have $c_i(f_{cons-duo}^{i,j}, f') = c_j(f_{cons-duo}^{i,j}, f') = 1$, which under SQB tie-breaking requires Eq. 2. We conclude that $f_{cons-duo}^{i,j}$ is self-maintaining.

We now consider social and Nash welfare. For any agent k other than i, j , their utility under $f_{cons-duo}^{i,j}$ has

$$\begin{aligned} u_k(f_{cons-duo}^{i,j}) &= (1 - p) \cdot 2p(1 - p) + p^3 + (1 - p)^3 \\ &\stackrel{\text{change-averse and unbiased}}{\geq} \frac{1}{2} \cdot 2p(1 - p) + p^3 + (1 - p)^3 \\ &= p(1 - p) + p^3 + (1 - p)^3 \\ &\geq p^2 + (1 - p)^2 = u_k(f_{dict}^i). \end{aligned}$$

It is thus a sufficient condition for the Nash welfare and social welfare of $f_{cons-duo}^{i,j}$ to surpass f_{dict}^i to have that the product and the sum of u_i, u_j is better under the former. We can directly write these conditions:

$$\begin{aligned} (p^2 + (1 - p))^2 &\geq p^2 + (1 - p)^2, \\ 2(p^2 + (1 - p)) &\geq p^2 + (1 - p)^2 + 1, \end{aligned}$$

and verify they hold for any p .

□