



Americas Branch
One Liberty Plaza, 20th Floor
New York, NY 10006

Tel : 212 337 5980
Fax : 212 337-5959
E-mail: ethrone@cambridge.org

Proofreading Instructions

Dear Contributor:

Attached please find the page proofs for your article scheduled to be published in:

Social Philosophy & Policy

Please follow these procedures:

1. **Proofreading:** Proofread your article carefully. Check especially the spellings of names and places as well as the accuracy of dates and numbers. Please answer all queries, but list the responses to the queries and other corrections separately. A revised PDF is available only upon request.
2. **Text:** Changes in the text are limited to typographical and factual errors. Rewriting or other stylistic changes are not permitted.
3. **Corrections:** Return the corrected proofs within 3 days of receipt by email to:

Pamela Phillips
Editorial Associate
Social Philosophy and Policy Foundation
E-Mail: pphillianda@yahoo.com
Phone: 419-819-7359

Please identify the corrections by page number, paragraph, and line. Please indicate the present errant copy followed by the correct copy. The corrections to the proofs should be sent within 3 days of receipt. Please note that delay in returning your proofs may require publication without your corrections.

4. **All Other Orders:** A message with a link to access a free PDF of your paper will be sent to you. To order reprints or offprints of your article or extra bound copies of the issue, please visit the Cambridge University Reprint Order Center online at: www.sheridan.com/cup/eoc

Thank you for your prompt attention to these proofs. If you have any questions about publishing, please feel free to contact the Editorial Associate (or the Cambridge Journals Department).

Best regards

Emily Throne
Cambridge University Press
One Liberty Plaza, 20th Floor
New York, NY 10006
Ph: (212) 337-5980
Fax: (212) 337-5959
E-mail: ethrone@cambridge.org

Author Queries

QA	The distinction between surnames can be ambiguous, therefore to ensure accurate tagging for indexing purposes online (eg for PubMed entries), please check that the highlighted surnames have been correctly identified, that all names are in the correct order and spelt correctly.
AQ1	Please provide the Citation for fig 4, 5, 6.

1 THREE CONCEPTS OF POLITICAL STABILITY: AN AGENT-
2 BASED MODEL*

3
4 QA BY KEVIN VALLIER
5

6
7 Abstract: *Public reason liberalism includes an ideal of political stability where justified*
8 *institutions reach a kind of self-enforcing equilibrium. Such an order must be stable for the*
9 *right reasons — where persons comply with the rules of the order for moral reasons, rather*
10 *than out of fear or self-interest. John Rawls called a society stable in this way well-ordered.*

11 *In this essay, I contend that a more sophisticated model of a well-ordered society, specif-*
12 *ically an agent-based model, yields a richer and more attractive understanding of political*
13 *stability — durability, balance, and immunity. A well-ordered society is one that possesses*
14 *a high degree of social trust and cooperative behavior among its citizens (durability) with*
15 *low short-run variability (balance). A well-ordered society also resists destabilization*
16 *caused by noncompliant agents in or entering the system (immunity).*

17 *Distinguishing between these three concepts complicates the necessary reformulation of*
18 *the idea of a well-ordered society. Going forward, public reason theorists must now distin-*
19 *guish between types of assurance, specify heretofore unknown aspects of reasonable behavior,*
20 *and reconceive of the nonideal preconditions for forming a stable, ideal social order.*

21 KEY WORDS: stability, stability for the right reasons, public reason, public
22 justification, public reason liberalism, well-ordered society, agent-based model
23

24 I. INTRODUCTION
25

26 Public reason liberalism¹ unites advocacy of liberal democratic institu-
27 tions with a constraint on the use of coercion or political decision-making.
28 This constraint holds that political action is permitted only when each
29 person, suitably idealized, has sufficient reason of her own to accept or
30 endorse the major social and economic institutions under which she lives.
31 When all suitably idealized persons have sufficient reason to endorse the
32 political actions that govern them, these activities, which typically involve
33 state coercion, are *publicly justified*.
34
35

36 * Thanks to Aaron Michelson, Joseph Bulger, and Ryan Muldoon for helping me learn to
37 build my own models. In doing so, I have drawn on the decision-making heuristic found in
38 Ryan Muldoon, Michael Borgida, and Michael Cuffaro, "The Conditions of Tolerance," *Politics,*
39 *Philosophy, and Economics* 11, no. 3 (2012): 322–44, and Ryan Muldoon, Chiara Lisciandra,
40 Cristina Bicchieri, Stephan Hartmann, and Jan Sprenger, "On the Emergence of Descriptive
41 Norms," 13, no. 1 (2014): 3–22. I have also drawn on Uri Wilensky's Iterated N-person Prisoners'
42 Dilemma code in the Netlogo models library, found here: <http://ccl.northwestern.edu/netlogo/models/PDN-PersonIterated>

¹ Sometimes described as political liberalism or justificatory liberalism.

1 The idea of public justification² includes an ideal of *political stability*
 2 where justified institutions reach a kind of self-enforcing equilibrium.³
 3 Citizens of a stable society generally recognize that all, or nearly all,
 4 people have sufficient reason to comply with directives issued by publicly
 5 justified institutions, such that unilateral deviations from those
 6 directives leads to a worse outcome from the defector's point of view.
 7 John Rawls and contemporary public reason liberals often describe an
 8 order that is stable for the right reasons as a *well-ordered society*, whose
 9 order is based on diverse persons having moral reasons to comply with
 10 the directives of just institutions, and not merely reasons to comply
 11 based on fear of punishment. The latter form of stability is typically
 12 understood as a mere *modus vivendi*.⁴ Public reason liberals favorably
 13 contrast the former with the latter.

14 In this essay, I contend that a more sophisticated model of a well-ordered
 15 society, specifically an agent-based model, yields a richer and more accurate
 16 ideal of political stability. In particular, an agent-based model helps us to distinguish
 17 between three concepts of political stability — durability, balance,
 18 and immunity. A well-ordered society is one that possesses a high degree
 19 of social trust and cooperative behavior among its citizens (*durability*) with
 20 low short-run variability (*balance*). A well-ordered society also resists destabilization
 21 caused by non-compliant agents in the system (*immunity*).

22 Distinguishing between these three concepts has two critical implications.
 23 First, previous work on political stability within public reason liberalism
 24 has depended upon a single, coherent notion of stability. Accordingly,
 25 my tripartite distinction weakens attempts to elaborate, defend, and refute
 26

27 ² Kevin Vallier and Fred D'Agostino, "Public Justification," <http://plato.stanford.edu/entries/justification-public/>.

28 ³ The literature on modeling stability within a well-ordered society is new and focuses
 29 almost exclusively on how to understand Rawls's account. For some older pieces, see Larry
 30 Krasnoff, "Consensus, Stability, and Normativity in Rawls's Political Liberalism," *Journal of*
 31 *Philosophy* 95 (1998): 269–92, and Thomas E. Hill, "The Problem of Stability in Political Liberalism," *Pacific Philosophical Quarterly* 75 (1994): 333–52. Much of the literature begins with Stephen
 32 Macedo and Gillian K. Hadfield, "Rational Reasonableness: Toward a Positive Theory of
 33 Public Reason," *Law and Ethics of Human Rights* 6, no. 1 (2012): 7–46; Paul Weithman, *Why*
 34 *Political Liberalism? On John Rawls's Political Turn* (New York: Oxford University Press, 2010);
 35 Gerald Gaus, "The Turn to a Political Liberalism," in John Mandel and David Reidy, eds.,
 36 *A Companion to Rawls* (Chichester: Wiley, 2013), 235–50; John Thrasher and Kevin Vallier,
 37 "The Fragility of Consensus: Public Reason, Diversity, and Stability," *The European Journal*
 38 *of Philosophy* 23, no. 4 (2015): 933–54; George Klosko, "Rawls, Weithman, and the Stability
 39 of Liberal Democracy," *Res Publica* 21 (2015): 235–49; Paul Weithman, "Reply to Professor
 40 Klosko," *Res Publica* 21 (2015): 251–64; George Klosko, "Stability: Political and Conception:
 41 A Response to Professor Weithman," *Res Publica* 21 (2015): 265–72; Paul Weithman, "Inclusivism,
 42 Stability, and Assurance," in Tom Bailey and Valentina Gentile, eds., *Rawls and Religion*
 43 (New York: Columbia University Press, 2015), 75–96; Paul Weithman, "Relational Equality,
 44 Inherent Stability, and the Reach of Contractualism," *Social Philosophy and Policy* 31, no. 2
 45 (2015): 92–113; and John Garthoff, "Rawlsian Stability," *Res Publica* (2015): 1–15. Two unpublished
 pieces are also helpful. See Sharon Lloyd, "Private Reasons, Public Judgments, and the
 Requirements of Reciprocity," University of Southern California, 2015, along with Andrew
 Lister, "Public Reason and Reciprocity," Queens University, 2015.

⁴ John Rawls, *Political Liberalism* (New York: Columbia University Press, 2005), pp. xl-xli.

1 public reason views that employ a single, coherent notion of stability.⁵
2 Second, distinguishing three notions of stability poses three new chal-
3 lenges in formulating the idea of a well-ordered society: (i) distinguishing
4 among types of assurance, (ii) resolving a critical ambiguity in the idea of a
5 reasonable person, and (iii) figuring out how to transition from a non-ideal
6 social order society to an ideal, well-ordered one.

7 Advances in computer science have produced modeling software that
8 allows us to develop a dramatically richer model of a well-ordered society
9 (WOS), specifically through computational *agent-based modeling* (ABM).⁶
10 An agent-based model is a class of computational models that simulate the
11 actions and interactions of *autonomous agents* (either individual or collec-
12 tive agents like groups) with a view to assessing their effects on the system
13 as a whole. ABMs encode “the behavior of individual agents in simple
14 rules so that we can observe the results of these agents’ interactions.”⁷
15 ABMs contrast with standard mathematical modeling in describing a
16 system, not by variables representing the state of the whole system, but
17 rather with a system’s individual components and their behaviors. ABMs
18 model the individual, and determine system states by the emergent prop-
19 erties of agents interacting with the environment and other agents, which
20 is why ABMs are sometimes referred to as *individual-based* models.⁸

21 The main point of building an ABM of a WOS is to distinguish between
22 types of stability, not to represent a WOS in full detail. Accordingly, many
23 of my simplifying assumptions are grounded in the goal of distinguishing
24 types of stability rather than constructing a plausible representation of
25 the most important dynamics of a WOS.⁹ My overarching aim is to make
26 the *already* agent-based elements of a well-ordered society model more
27 explicit to uncover system-level properties that emerge from a complex
28 adaptive system like a WOS.

29 I introduce my ABM in three stages. First, I develop a simple WOS model
30 that contains only reasonable agents choosing whether to comply or defect
31 from norms of cooperation. This simple model generates a distinction
32 between the capacity of a system to stabilize its constituent norms via the
33 production and maintenance of social trust, which I call *durability*, and the
34 short-run variability of cooperative behavior, which I call *balance*.

35
36 ⁵ This includes Thrasher and Vallier, “The Fragility of Consensus: Public Reason, Diversity,
37 and Stability.” Also see Weithman’s summary of Rawls’s approach in Weithman, “Inclusivism,
38 Stability, and Assurance.”

39 ⁶ I promise that “WOS” and “ABM” are the only acronyms I use in the paper. For discussion
40 of the power of these models within social science, see John H. Miller and Scott E. Page, *Complex
41 Adaptive Systems: An Introduction to Computational Models of Social Life* (Princeton, NJ: Princeton
42 University Press, 2007).

43 ⁷ Uri Wilensky and William Rand, *An Introduction to Agent-Based Modeling: Modeling Natural,
44 Social, and Engineered Complex Systems with Netlogo* (Cambridge, MA: MIT Press, 2015), 22.

45 ⁸ Steven Railsback and Volker Grimm, *Agent-Based and Individual-Based Modeling: A Practical
46 Introduce* (Princeton, NJ: Princeton University Press, 2012), 10.

⁹ I thank Steven Stich for encouraging me to make my reasons for building an ABM more
explicit.

1 In stage two, I relax the assumption of full compliance by introducing
 2 a small number of agents who maximize their expected utility in their
 3 interactions with others.¹⁰ They are *merely rational* in that they are not con-
 4 ditional cooperators, and so are not reasonable.¹¹ I show that a dynamic
 5 found among groups of reasonable agents — network reciprocity — can,
 6 under favorable conditions, enable them to maintain some durability and
 7 balance despite relentless defection from merely rational agents.¹²

8 The third feature of the model specifies conditions for the entry and exit
 9 of agents in the system. This enables merely rational agents or reasonable
 10 agents to take over the population. A WOS whose reasonable agents can
 11 resist invasion and replacement by merely rational agents is *immune*.¹³

12 This essay proceeds in nine parts. I will first describe the problem of
 13 stability and the idea of a well-ordered society, both as found in Rawls (§2)
 14 and in my refinement of a standard WOS model (§3). I will then introduce,
 15 in stages, the major features of my WOS ABM. I outline a simple version
 16 of the WOS model that contains only reasonable agents (§4) and present
 17 the results in (§5). The simple WOS model helps distinguish and define
 18 durability and balance. I then relax the compliance assumption by intro-
 19 ducing non-compliant agents into the system (§6), and discuss the out-
 20 come, which involves a significant depression of durability (§7). Finally,
 21 I allow agents to enter and exit the system under various conditions (§8).
 22 The entry-exit dynamic allows me to illustrate the idea of immunity (§9)
 23 and explore the relationship between the three concepts of stability. I con-
 24 clude by suggesting some avenues for further development in the study
 25 of political stability (§10).

26 II. RAWLS'S WELL-ORDERED SOCIETY MODEL

27
 28
 29 Throughout his work, Rawls used the idea of a well-ordered society
 30 as an account of a realistic utopia, one where a society's basic structure
 31 is regulated by principles of justice, that this fact is publicly known, and
 32 that citizens all have an effective sense of justice so as to comply with the
 33 directives of the basic structure, which they regard as just.¹⁴ This sense
 34 of justice drives people to impose the rules of their society on *themselves*,
 35 and their compliance renders a WOS stable *for the right reasons*. Further,
 36 a WOS contains stabilizing forces, such that "when infractions occur,
 37 [these] should exist that prevent further violations and tend to restore the
 38
 39

40 ¹⁰ John Rawls, *A Theory of Justice* (New York: Oxford University Press, 1971), 8.

41 ¹¹ Rawls, *Political Liberalism*, 48–54.

42 ¹² Martin Nowak, "Five Rules for the Evolution of Cooperation," *Science* 314 (2006): 1560–63,
 at 1561.

43 ¹³ Since immunity describes the ability of a system to recover from external shocks, and
 44 there are many kinds of external shocks, that there will be many types of immunity. I expand
 on this point below.

45 ¹⁴ Rawls, *Political Liberalism*, 35.

1 arrangement."¹⁵ A WOS must be able "to generate its own support" rather
 2 than have it imposed from without.¹⁶

3 The final model of a well-ordered society, then, is understood as a rep-
 4 resentation of three social facts, the first of which I will update based on
 5 Rawls's embrace of reasonable pluralism about justice.¹⁷ A well-ordered
 6 society must satisfy these three conditions:

- 8 (1) Everyone accepts, and knows that everyone else accepts, some
 9 member of a limited set of reasonable political conceptions of
 10 justice, which establish shared points of view from which citizens'
 11 claims on society can be adjudicated.¹⁸
- 12 (2) Its basic structure — that is, its main political and social institu-
 13 tions and how they fit together as one system of cooperation — is
 14 publicly acknowledged or with good reason believed, to satisfy
 15 these principles (or some mix of them).
- 16 (3) Its citizens have a normally effective sense of justice and so gen-
 17 erally comply with society's basic institutions, which they regard
 18 as just (if not fully just).

19
 20 While this description is fairly rich, we need much say more to explain the
 21 sense in which a WOS is an equilibrium. First, for Rawls, a stable society is in
 22 equilibrium on a conception of justice, not *primarily* on institutions.¹⁹ But
 23 institutional rules must also be self-stabilizing because they institution-
 24 alize a conception of justice. Since I want to allow for conceptions of justice
 25 to vary, as Rawls ultimately did, I will not presume that a WOS is in equi-
 26 librium on a conception of justice; but I will assume it is in equilibrium
 27 on the rules that lead officials and citizens to issue behavioral directives
 28 to others. These rules must be issued by social structures or organizations
 29 that citizens can, on reflection, see as institutionalizing their (presumably
 30 reasonable) conceptions of justice.

31 A WOS is a kind of Nash equilibrium, which is suggested by Rawls's claim
 32 that in a WOS, compliance with justice is each person's "best reply ... to
 33 the corresponding demands of the others."²⁰ Rawls also says that equilib-
 34 rium is reached when each person's "plan of life" is his "best reply to the
 35 similar plans of his associates."²¹ The idea, then, appears to be that no one
 36
 37

38 ¹⁵ John Rawls, *A Theory of Justice* (Cambridge: Belknap Press, 1999), 6.

39 ¹⁶ *Ibid.*, 119. Penal institutions are meant to supplement the forces maintaining stability for
 40 the right reasons. Rawls, *Theory of Justice*, 502–503. I thank Steven Stich for encouraging me
 41 to make this point explicit.

42 ¹⁷ Rawls, *Political Liberalism*, p. xlvii.

43 ¹⁸ This condition allows that different reasonable persons accept different reasonable polit-
 44 ical conceptions from one another. That is, they can converge on different conceptions at the
 45 same time and in the same society. *Ibid.*, 35.

46 ¹⁹ Weithman, "Reply to Professor Klosko," 254.

47 ²⁰ Rawls, *Theory of Justice*, 103.

48 ²¹ *Ibid.*, 497.

1 can improve her position unilaterally through defecting from compliance
 2 with justice. But this is too strong, since a society-wide equilibrium does
 3 not require perfect compliance. Instead, a society's basic rules must be
 4 "more or less regularly complied with" and that "when infractions occur"
 5 stabilizing forces should tend to restore conditions of cooperation.²² Finally,
 6 Rawls says that in a WOS "inevitable deviations from justice are effectively
 7 corrected or held within tolerable bounds by forces within the system."²³
 8 So deviations from justice are inevitable, and only need to be held within
 9 tolerable bounds. Rawls allows individual behavior to rationally deviate in
 10 at least a few cases.

11 The WOS also requires assurance if it is to be in equilibrium. For reason-
 12 able citizens only have sufficient practical reason to act justly if they are
 13 assured that others will generally do likewise. Fortunately, Rawls has a
 14 lot to say about assurance, which he understands in terms of "publicity."
 15 For Rawls, and other Rawlsians, a WOS is by definition "regulated by an
 16 effective *public* conception of justice," which means that citizens must be
 17 able to determine for themselves whether their institutions comply with jus-
 18 tice, and determine that others can do likewise from their own perspectives.²⁴
 19 Publicity has three levels, but we need only focus on the first, which is real-
 20 ized when "citizens accept and know that others likewise accept those princi-
 21 ples [of justice], and this knowledge is in turn publicly recognized," everyone
 22 sees that "the institutions of the basic structure of society are just (as defined
 23 by those principles)," and everyone sees this further fact. So equilibrium also
 24 requires assurance.

25 That said, Rawls has little to say about *how* a WOS generates assurance.
 26 In response, Paul Weithman has argued that we should interpret Rawls as
 27 arguing that the use of *public reasons* serves as an assurance mechanism.
 28 Public reasons are those derived from shared public values — they are the
 29 reasons of the public based on their conception of justice.²⁵ By using public
 30 reasons in political discussion on matters of basic justice and constitutional
 31 essentials, citizens publicly signal their allegiance to just institutions.

32 We have reason to be concerned about Rawls's WOS model. First, there is
 33 good reason to worry about whether public reasons can provide adequate
 34 assurance. John Thrasher and I have argued that they cannot. Even reason-
 35 able agents can face conditions of communicative drift, noise, and cheap
 36 talk that undermine the capacity of deliberation based on public reasons to
 37 provide adequate assurance.²⁶ Rawls also does not entertain the possibility
 38 of emergent destabilizing elements in his model. In particular, he does not
 39 acknowledge the possibility that different institutional demands may coun-
 40 teract one another. My model suggests that the degree of disorder within
 41

42 ²² Also see *ibid.*, 6.

43 ²³ *Ibid.*, 272.

44 ²⁴ Rawls, *Political Liberalism*, 66–72.

45 ²⁵ Weithman, "Inclusivism, Stability, and Assurance," 88–90.

²⁶ Thrasher and Vallier, "The Fragility of Consensus: Public Reason, Diversity, and Stability."

1 a WOS depends upon the character of reasonable agents, in particular the
2 extent to which agents can preserve stability by caring relatively little about
3 the behavior of other plays vis-à-vis their natural inclination to engage in
4 reasonable behavior. Rawls also does not specify how a WOS can stabilize
5 itself given the threat of external shocks, such as limits on resources or the
6 entry of uncooperative agents into the system. This suggests we need to
7 describe the conditions under which a WOS can and must resist invasion.

8 Fortunately, we can build a subtler model to show how elements of a
9 WOS are logically consistent in a way that yields an attractive and feasible
10 ideal. An ABM will reveal *dimensions and degrees* of political stability that
11 Rawlsian WOS models cannot. That is, the ABM allows us to distinguish
12 types of stability and to treat the factors that generate stability as continuum
13 notions, rather than as binary.

14 15 III. A SIMPLIFIED WOS MODEL

16
17 Enough Rawls. I now want to take the essential elements of his approach
18 and simplify them enough that they are subject to modeling. I will under-
19 stand a WOS as follows: (i) its citizens are generally good-willed and care
20 about engaging in reciprocal, cooperative behavior, (ii) they regard the norms
21 that govern their fundamental institutions as mostly just and legitimate, and
22 so comply with the directives of those institutions and the demands of others
23 to follow them. Finally, (iii) they believe that other members of society regard
24 their situation similarly, despite their diverse personal points of view; that is,
25 they have a high degree of assurance.

26 I employ the idea of an N-person Nash equilibrium, where agents play
27 strategies in pairs, and generate a stable, high degree of cooperation as an
28 emergent property of the system. A WOS is therefore best understood as a
29 macro-level equilibrium. Only mass defection needs to be self-correcting.
30 Macro-level stability is a function of local compliance, but it does not
31 require individual compliance in every case. The dominant mixed-strategy
32 of agents is to adopt a high probability of compliance with rules and direc-
33 tives established by just institutions, such that across the history of their
34 interactions with others, agents cannot improve their position by adopting
35 a non-cooperative strategy.²⁷

36 I forgo appeals to common knowledge. Instead, each agent merely makes
37 reliable judgments about the level of compliance within her environment.
38 My ABM instructs each agent to observe the fraction of cooperative plays in
39 the system at any one time, and partly base her choice in her next encounter
40 on that observation. Agents do not know what other players know.

41 When agents cooperate, I assume they comply with social norms of
42 justice, or regular social practices, enforced by social demands, ostracism,
43

44 ²⁷ Here I understand an agent's position and improvements upon that position as includ-
45 ing their moral commitments and personal projects.

1 and blame, and in some cases, the law.²⁸ They are norms of justice because
 2 infractions of the norms are seen not only as wrong or immoral but unjust.
 3 Importantly, these norms need not be part of a unified conception of justice.
 4 Instead, agents must merely see following them as a matter of justice
 5 and be subject to disapprobation when the norms are violated. In general,
 6 given the importance of assurance and the idea of social trust discussed
 7 below, I understand cooperative behavior as *trustworthy* behavior, where
 8 persons comply with expectations set by social conventions within that
 9 society. This means that they will not only forgo pursuing gains from defec-
 10 tion not merely because they are in public, but because they are indepen-
 11 dently motivated to cooperate.

12 We can understand defection, following David Rose, as a kind of *oppor-*
 13 *tunism*.²⁹ Rose defines opportunism as “acting to promote one’s welfare by
 14 taking advantage of a trust extended by an individual, group, or society as
 15 a whole.”³⁰ This trust is based on the expectation that everyone complies
 16 with the norms of justice present in that society. A critical feature of oppor-
 17 tunism is that it does not always cause perceptible harm, or even any harm
 18 at all. If a society is sufficiently large, small acts of opportunism are not in
 19 themselves sources of harm. For example, suppose John downloads an
 20 episode of *Game of Thrones* without an HBO subscription; it does no per-
 21 ceptible harm, but it is opportunistic. And with sufficient opportunism,
 22 some parties will be harmed. This is partly because a society’s amount of
 23 social trust in future cooperative behavior will tend to decrease, reducing
 24 the efficacy of many social institutions.³¹ So when agents interact within
 25 just institutions, defection involves breaking the mutually agreed upon terms
 26 of implicit and explicit agreements that are not unjust. In particular, I will
 27 understand defection in terms of *first-degree opportunism*, which involves
 28 taking advantage of the imperfect enforceability of contracts by renegeing
 29 on contracts.³² I so focus because first-degree opportunism is relatively
 30 simple and detectable.

31 32 IV. A SIMPLE AGENT-BASED MODEL OF A WELL-ORDERED SOCIETY

33
34 I begin introducing my ABM based on elements in the previous section.
 35 In the simple WOS model, all the agents are of a single type that I call
 36
37

38 ²⁸ I agree with Rawls that a WOS is best modeled as requiring some coercion but whose
 39 stability is driven almost entirely by the voluntary choices of citizens.

40 ²⁹ David Rose, *The Moral Foundation of Economic Behavior* (New York: Oxford University
 41 Press, 2014).

42 ³⁰ *Ibid.*, 21.

43 ³¹ For a review of the empirical literature on the benefits of social trust, see Sanjay
 44 Banerjee, Norman Bowie, and Carla Pavone, “An Ethical Analysis of the Trust Relationship,”
 45 in Reinhard Bachmann and Akbar Zaheer, eds., *Handbook of Trust Research* (Northampton:
 Elgar, 2008), 318–31.

³² Rose, *The Moral Foundation of Economic Behavior*, 30.

1 “reasonable,” following Rawlsian terminology.³³ But to avoid norma-
2 tively thick and loaded conceptions of reasonableness, let us say that
3 all reasonable agents are committed to reciprocity when they will coop-
4 erate with other agents so long as they believe others will do likewise.
5 We can therefore model reasonable agents as *conditional cooperators*,
6 who cooperate given the expectation that others will do the same. This
7 means that the main difficulty faced by a society of reasonable people is
8 assuring one another that they will respond cooperatively to cooperative
9 behavior.

10 I understand cooperation, defection, and associated game-theoretic concep-
11 tions maximally capaciously. Appealing to game theory does not require
12 representing agents as merely instrumentally rational, for instance. There is
13 no reason that game-theoretic modeling cannot model persons as having
14 utility functions that include rich moral commitments and cooperative
15 dispositions. I understand the idea of “utility” with similar capaciousness as
16 representing whatever agents regard as choiceworthy. In sum, the tools of
17 game theory do *not* require a *homo economicus* conception of the individual
18 and make no significant individualist methodological assumptions.³⁴
19 Or so I assume here forth.

20 A reasonable agent’s decision-making heuristic is a simple propensity
21 to cooperate. The agent calculates her propensity based on two factors:
22 the intrinsic utility she derives from cooperating successfully and her
23 observation of the percentage of cooperative agents in the system. The
24 first factor is the agent’s *intrinsic propensity* to cooperate, which is its gen-
25 eral liking of cooperation, or how much the agent would cooperate if
26 it were indifferent to how others treat it. The second factor involves the
27 agent calculating the ratio between the agents who cooperated in their last
28 interaction to the total number of agents in the system, yielding some ratio
29 between 0 and 1.³⁵ An agent’s *social sensitivity* is understood as the relative
30 weighting of the observed ratio of cooperation and an agent’s intrinsic
31 propensity. When an agent’s observation is combined with her intrinsic
32 propensity according to the weighting specified by its social sensitivity,
33 this determines its *effective propensity* or the probability with which she
34 will cooperate in a given interaction.

35 Social sensitivity can be set at any value between 0 and 1, such that sen-
36 sitivity functions as a weighting relative to an agent’s intrinsic propensity,
37 which is set as an input by the modeler. Suppose that an agent’s intrinsic
38 propensity to cooperate is 90% (.90). If social sensitivity is set at .5, and the
39
40
41

42 ³³ Rawls, *Political Liberalism*, 48–54. I do not include the requirement that agents recognize
43 the burdens of judgment, as it would unnecessarily complicate the model. See *ibid.*, 55–58.

44 ³⁴ Gerald Gaus, *On Philosophy, Politics, and Economics* (Belmont, CA: Wadsworth, 2007),
45 19–27.

³⁵ Agents with longer memories unnecessarily complicate the model.

1 agent observes only 60% (.6) of agents cooperating, then she calculates
 2 her effective propensity like so:

$$3 \quad \text{Effective propensity} = \text{social sensitivity} * \text{percent cooperating} + (1 - \text{social} \\ 4 \quad \text{sensitivity}) * \text{intrinsic propensity.}$$

6
 7 In this case, then, we have $.5(.6) + (1-.5)(.90)$ or .75. This means that the
 8 next time the agent plays a game with another agent her effective propen-
 9 sity to cooperate is 75%. She rolls a four-sided die with three directives
 10 to cooperate and one directive to defect, and follows the directive rolled.
 11 Notice that the agent begins with a very high propensity to cooperate
 12 in the absence of information about cooperation in the system. If we set
 13 social sensitivity to 0, the agent will cooperate 90% of the time, since she
 14 entirely discounts her information about what others are doing. Once
 15 social sensitivity is positive, however, the agent adjusts her propensity to
 16 cooperate based on her observation.

17 In this way, social sensitivity is an assurance parameter because the
 18 agent's estimation of the ratio of cooperative agents to defecting agents can
 19 be understood as a kind of assurance that others will behave cooperatively.
 20 I believe that social sensitivity allows the model to capture the essence of
 21 reasonable Rawlsian agents and a thin notion of publicity. Unlike Rawlsian
 22 agents, however, reasonable agents make finer-grained judgments about
 23 the likelihood of cooperation and they will defect in proportion to their
 24 observation of defection.

25 Our final piece of the simple WOS model is the output variable — social
 26 trust. Social trust is calculated by taking the average of the effective propen-
 27 sities of all agents at a single time. Social trust, therefore, represents
 28 the degree to which the system as a whole is prepared to cooperate with
 29 others. This technical definition of social trust, then, bears some resemblance
 30 to more common usages of the term.³⁶

31 Reasonable agents also have a general desire to engage with other
 32 reasonable agents rather than unreasonable agents. In the model, an agent
 33 calculates the average position of cooperative agents and turns towards
 34 it after playing a game. She also turns away from the average position
 35 of those who defected (regardless of their hard-wired strategy). She
 36 does *not* also move towards cooperators or away from defectors. Since
 37 she only turns, she takes a somewhat more random walk tilted towards
 38 recently cooperative populations and away from recently somewhat less
 39 cooperative populations. Reasonable agents recalculate this location after
 40 each play.

41
 42
 43 ³⁶ In another work, I appeal to the notion of social trust defined in Christiano Castelfranchi
 44 and Rino Falcone, *Trust Theory: A Socio-Cognitive and Computational Model* (West Sussex: Wiley,
 45 2010). For a survey of different views, see Rose, *The Moral Foundation of Economic Behavior*,
 19–38.

1 In the simple model with only reasonable agents, this form of correla-
2 tion has a significant effect, leading players to congregate in a central
3 hub. The reason for this is intriguing. With a simple turn towards cooper-
4 ators, and 90 degree range of movement left or right, reasonable agents
5 will start to encounter one another more often, but the more often they
6 interact, the more often they turn towards one another, which lead them
7 to cooperate more often, increasingly centralizing the cooperators.
8 The formation of the cooperative hub is an emergent feature of a very
9 simple dynamic: after each play, face those who have just cooperated,
10 and take a somewhat random walk.

11 The important point here is that reasonable agents have the ability to
12 correlate their behavior by congregating so as to play many games with
13 one another. In this way, I draw on the idea of “network reciprocity”
14 in game theory, where previous cooperation leads agents to interact with
15 one another more often, forming a pro-social network.³⁷ The dynamic that
16 drives network reciprocity in the WOS model is rudimentary. Agents do
17 not learn the agent-strategies of other players, so they do not distinguish
18 between agent types. Nor do they record the effective propensities of other
19 agents at any one time. Instead, *all they know* is each player’s present play,
20 or its last play, if it is not partnered with another agent when the observa-
21 tion occurs. So there is no complex reputation effect. Each player merely
22 has a drive to face players who have cooperated and face away from those
23 who have defected.

24 Even if all reasonable agents are naturally disposed to cooperate in every
25 interaction with others, they will cooperate less if they do not realize that
26 others are similarly disposed. Cooperation among reasonable agents can
27 quickly break down if the agents lack assurance. Consequently, we should
28 not model a WOS by assuming that all reasonable persons always cooper-
29 ate with one another.³⁸ This critical alteration to Rawls’s model helps us
30 to more accurately represent the dilemmas faced by cooperative agents in
31 a well-ordered society.³⁹

32 Reasonable agents should have two further features, both of which
33 I omit for presentation purposes. First, agents do not make observational
34 errors, whereas the most accurate modeling assumptions would allow
35 for mistakes. Further, reasonable agents are prepared to defect when the
36 cost of complying with a rule is too great, even if others are prepared to
37
38

39 ³⁷ Nowak, “Five Rules for the Evolution of Cooperation,” 1561.

40 ³⁸ It also explains why public reason needs an assurance mechanism to make sense of sta-
41 bility for the right reasons. See Weithman, *Why Political Liberalism?* 327–35.

42 ³⁹ Though, clearly Rawls thought the generation of publicity was critical for maintaining
43 stability; but he was not at all clear about how facts about cooperation are made public
44 knowledge. It appears that he believed that public reasoning functions as an assurance
45 mechanism. As noted, John Thrasher and I have argued that public reasoning is not an
effective form of assurance. In light of that paper, my model bases assurance on observed
cooperation with others.

1 cooperate. Even a reasonable agent should defect if she expects cooperation
2 will kill her! I have omitted this “cost caveat” from the model because
3 my aim is to model a well-ordered society operating under favorable
4 conditions, such that cooperation should almost always prove beneficial,
5 if less beneficial than getting away with defection.

6 We can understand the simple WOS ABM as describing citizen-agents
7 interacting in a legal environment that applies payoffs to agents based
8 on cooperative or non-cooperative behavior, such that cooperative behavior
9 is rewarded and uncooperative behavior is punished. Importantly,
10 however, reasonable agents in the simple model do not care about their
11 payoffs, and so focus exclusively on reciprocity. I will only partly relax
12 this assumption in the more complex model, and only for what I will call
13 merely rational agents, not reasonable agents.

14 V. RESULTS OF THE SIMPLE WOS ABM

15
16
17 In the simple WOS ABM, each iteration of the model — a “tick” in
18 Netlogo terminology — brings reasonable agents closer to a system-
19 wide cooperative equilibrium whose social trust is equal to the average
20 *intrinsic* propensity of agents to cooperate. In other words, given that
21 all agents are reasonable, the fact that these agents are socially sensitive
22 to what other reasonable agents are doing does not discourage them from
23 cooperation in the long-run. Second, the agents quickly cluster into a tight
24 network, with only some agents moving around outside of the core cluster.
25 The clustering effect is robust across the number of agents in the system and
26 reasonable agents’ intrinsic propensity to cooperate.

27 The significant feature of the model is that social sensitivity has a
28 substantial short-term effect on the range of social trust found across
29 each run (which I define as 500 ticks). While the average equilibrium
30 level is set early in each run, typically in the first 100 ticks, the variability
31 of short-run social trust increases as the social sensitivity of the agents
32 increases. Compare the average level social trust in a game where fifty
33 agents each have an intrinsic propensity of .9. In Figure 1, social sensi-
34 tivity is set to .3 and the second to .9:

35 Recall that a system’s level of social trust is the average of the effective
36 propensities of all agents at a time; the graphs show the variation of social
37 trust over time. Display 1 only represents social trust during one run of
38 the model. While the system is non-deterministic, such that different runs
39 yield different curves, Figure 1 nonetheless represents general system
40 behavior at the two levels of social sensitivity.

41 The model shows that the more agents care about what others agents
42 are doing, the more social trust will vary in the short run. However,
43 the average level of social trust will remain constant over the long-run.
44 Average social trust is determined by the intrinsic propensities of the
45 agent set by the modeler, such that the higher the intrinsic propensity,

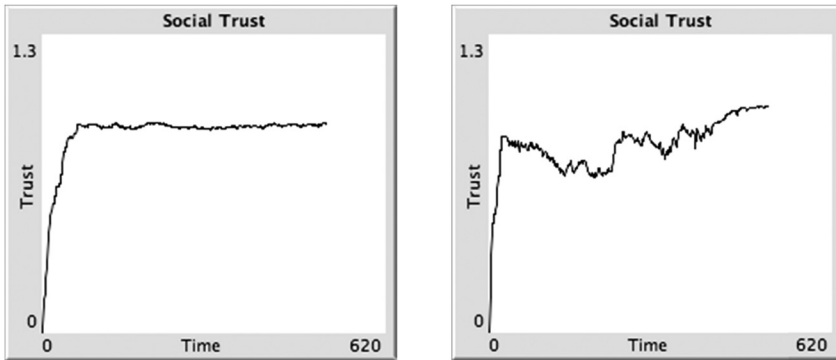


FIGURE 1. Social Sensitivity and Trust Variability – Steady and Shaky

the higher average level of social trust.⁴⁰ The variance of social trust, however, is a function of social sensitivity. As social sensitivity increases, the volatility of the system increases, though this has almost no effect on the average level of social trust, as demonstrated in Figure 1.

Consider now the two societies represented in Figure 1. Steady and Shaky have the same high average level of social trust. In this sense, then, both are stable for the right reasons. But in Steady, the variation in social trust across time is quite small (as it is at .3 social sensitivity), whereas in Shaky, the variation in social trust is quite large (at .9 social sensitivity). Which society is better, given Rawlsian aspirations? I believe Steady should come out ahead, as I argue below. It not only has the capacity to stabilize its constituent norms by maintaining a high level of social trust and cooperation, but exhibits low variance in that capacity.

In this way, the ABM allows us to distinguish two types of stability: *durability* and *balance*. *Durability* is defined as a system's average level of social trust over time. A durable system is one with high average social trust among agents. It may initially seem strange to describe an *average* level of *anything* as a form of stability. Stability as a concept seems to denote *variability* in a level of some feature of a complex system, not the level itself.⁴¹ But remember what sort of level we're talking about — it is a level of social trust, which is a degree of cooperation with the constitutive moral and legal conventions in a WOS. Justified moral and legal conventions are maintained by cooperation, which involves stabilizing them. Stabilization requires a lot of cooperation, then, and mass defection leads conventions to collapse. To illustrate, imagine a social system with an invariant, but *low* level of cooperation. That society is unstable in the sense that its justified conventions are fragile equilibria if they are equilibria *at all*, given the high

⁴⁰ For a review of the model, and my data sets, see <http://www.kevinvallier.com/stability>.

⁴¹ I'm grateful to Alan Hamlin for helping me to see this distinction.

1 concentration of defection. So to understand durability as a kind of sta-
2 bility, we must remember that cooperation in a WOS maintains justified
3 moral and legal conventions, such that large amounts of defection lead to
4 unstable conventions.

5 We can understand durability as a kind of first-order stability, the
6 capacity of a system to stabilize its constituent norms. *Balance* refers to
7 the variability of social trust over some time period. That is, balance is
8 a kind of *second-order* stability, or the stability of a measure of stability.
9 Thus, durability and balance are different answers to the question,
10 “Stability of what?” Durability measures the stability of a society’s con-
11 stituent social norms by measuring its degree of social trust. Balance
12 measures the stability of a society’s level of social trust.⁴²

13 We can distinguish forms of balance in two ways: first, by the length
14 of the relevant time series over which variability is measured, say over
15 100 or 500 ticks, and second, by distinguishing between frequency mea-
16 surement and amplitude measurement. Consider the following two series
17 in Figure 2, each of which has an average of 100 units and cover the same
18 length of time.

19 The two series have the same amplitude but different frequencies and
20 there is a good argument that Series B is more stable than Series A, as it
21 transitions more slowly through more fine-grained states. For our purposes,
22 I restrict balance to *short-run frequency*. A social system is balanced when,
23 like Steady, it has a low variability in social trust over the short-run. If the
24 system resembles Shaky, we can call it volatile.

25 I argue that a WOS should be understood as both balanced and durable.
26 Durability enables the system to maintain a high level of social trust,
27 which is obviously critical. The case for balance is less obvious, but still
28 strong. Low variance is a social good because high variance systems
29 contain highly undesirable periods, and most people will prefer steady,
30 reliable expectations even if it means fewer high points. This is likely
31 true simply in virtue of human risk aversion. Big highs aren’t worth big
32 lows. Second, Shaky might yield negative social consequences based
33 on the fact that some agents will recognize the volatility and act less
34 cooperatively as a result. That is, variability in social trust in the short-
35 run may reduce the *level* of social trust in the long-run. The simple
36 model lacks a complex learning algorithm required to test this claim,
37 but we could program reasonable agents to periodically measure the
38 maximum and minimum level of social trust over a lengthy period of
39 time, and then adjust their effective propensities up or down based
40 on the size of the range. As complexity increases, generating the log-
41 ical consequences of the model becomes a much more computationally

42
43
44 ⁴² We will see below that immunity is a kind of first-order stability, though we could also
45 measure the variability of immunity to create a fourth concept of stability.

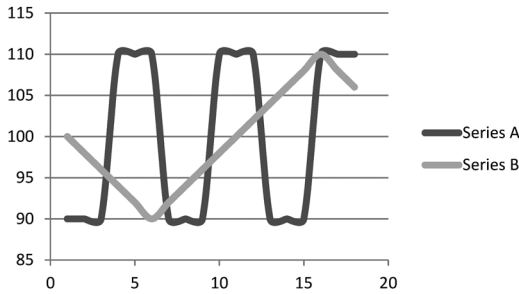


FIGURE 2. Frequency and Amplitude

demanding task; so I do not present that model here. But the point should be clear enough. We not only care about having a high average level of social trust, but that it not vary too much.

An unbalanced system also varies randomly, such that it is impossible to predict the system's capacity to stabilize its constituent norms at any one time. This creates challenges for actors who wish to alter and improve their legal order, since predicting the effect of a legal change depends upon the predictor's understanding of how the new law will be received by those subject to it. A balanced order will be predictable: either people will be relatively more likely to comply with the new law, or relatively less likely, but laws can be based on a high degree of certitude about a society's level of social trust. An unbalanced order makes these predictions much harder, and at the limit, impossible.

A third problem with an unbalanced order derives from the fact that a highly variable level of social trust will increase transactions costs between agents at the local level. Agents cannot form expectations about how likely others are to exhibit certain behavior, for cooperation or defection. And this makes engaging in any social activity risky, since at least agents in a low-trust society will know not to stick their necks out. This will likely lead to a less cooperative and trusting social order over the long-run, given that fewer people will be likely to take the risks necessary to create a productive, high-trust order.

The potential upside of an unbalanced order is that it is harder to take advantage of cooperative agents. Since balance is a function of social sensitivity, agents in an unbalanced order will take the behavior of others into account more than they would in a balanced order. This means that they will generally be more responsive to non-cooperative behavior from other agents. That said, the complex ABM discussed below suggests that an unbalanced order is *less* immune from invasion by merely rational agents than otherwise. Increases in social sensitivity depress immunity somewhat. So we have at least some evidence that the purported downside balance fails to materialize in the model.

1 An unbalanced order may also seem to realize the good of having a
2 dynamic and disruptive order, which will be more effective at discovering
3 new ways for persons to coordinate and cooperate. However, while I think
4 discovery is an important part of a WOS that has been almost entirely
5 ignored in Rawlsian models, discovery orientation is compatible with
6 balance. Balance obtains when the level of compliance with norms is
7 low; if these norms permit and encourage social experimentation, then
8 balance promotes discovery rather than discouraging it.⁴³

9 We can already start to see the challenges the ABM poses to traditional
10 political liberal approaches to stability. First, public reason liberals, including
11 me, have not distinguished between durability and balance, or even
12 between forms of balance. Reaching durability and balance may require
13 different social mechanisms, a fact which undermines the case for appealing
14 to a single dynamic to establish stability. The model also uncovers an
15 ambiguity in the idea of a reasonable person, particularly in specifying the
16 extent to which reasonable people should cooperate with others based on
17 their intrinsic character traits or their observations of how others behave.
18 Thus, we must decide how much reasonable people care about the actions
19 of others. And since a WOS should be both durable and balanced, we have
20 a richer ideal of stability.

21 To ensure that we understand the results of the model, let's take an
22 informal tour through the experience of a single agent in an environment
23 with some similarity to the one in my model. Call her Reba. Reba is good-
24 willed in general, and eagerly cooperates when she sees others doing
25 likewise. Much like us, Reba underestimates the extent to which others
26 affect her evaluations, thinking that her intrinsic propensity to cooperate
27 drives her behavior. But in reality, her behavior is determined much more
28 by social expectations than stable features of her character. Let's assume
29 other agents are similarly spirited. For this reason, we assign each agent
30 in that system a 90% intrinsic propensity to cooperate with .9 social
31 sensitivity – they care nine times as much about what they experience
32 with others than drawing on their own character. That is a lot, I grant,
33 but it will prove illuminating.

34 Now imagine that Reba has moved to a new neighborhood; she has
35 heard good things about it, but she knows no one there. She is optimistic
36 that other people will be as good-willed as she is. So long as others are
37 kind to her, Reba thinks, she will be kind in return. Next imagine that
38 once Reba has moved into a neighborhood, she begins to interact with
39 her neighbors on a regular basis. Let's understand her interactions com-
40 mercially: she is engaged in economic exchanges with others, say through
41
42

43 For discussion, see Gerald Gaus, *The Tyranny of the Ideal: Justice in a Diverse Society*
44 (Princeton, NJ: Princeton University Press, 2016), and Ryan Muldoon, *Diversity and the*
45 *Social Contract* (New York: Routledge, 2017).

1 yard sales, maintaining a neighborhood watch, keeping the neighborhood
2 park clean, contributing to the yearly block party, and so on.

3 Most of Reba's neighbors are kind. They almost always offer fair prices
4 at yard sales, participate in the neighborhood watch, and throw away
5 trash at the park. They even bring food to the block party every year. But
6 other neighbors are not as kind: John sometimes shirks; he is tempted to
7 opportunism. So John charges too much at yard sales, or offers too little
8 for valuable items; he sometimes participates in the neighborhood watch,
9 but he gets lazy on Mondays, and this led to a successful burglary last
10 month. John picks up trash when others are around, but he occasionally
11 litters in their absence. And he sometimes forgets to bring anything to the
12 block party. Reba notices that her other neighbors start to shirk a bit more
13 after observing John's behavior, and she resents investing so heavily in
14 the good of the neighborhood. Her resentment leads her to reduce her
15 participation. When others observe Reba, community pillar, contributing
16 less, they draw back as well. And this lowers the level of social trust in the
17 neighborhood.

18 But Reba doesn't give up. Buoyed by her mostly positive interactions
19 with others, she tends to spend more time at their homes, focuses more on
20 their yard sales, and doesn't go to John's anymore, and so on. John sees
21 this behavior, feels guilty, and starts to clean up his act a bit. The result
22 of all these actions is a neighborhood with a high level of social trust but
23 with some variability. When John starts to shirk, Reba withdraws, which
24 reduces shirking and maintains social trust. In our terms, the neighbor-
25 hood is durable, but less than fully balanced. The average level of social
26 trust is high, such that beneficial social norms can stabilize and generate
27 desirably high levels of compliance. However, Reba and other neighbors
28 care a great deal about what others are doing, and so generate periods of
29 variable neighborly spirit.

30 Notice also that neighborhood stability is based on the moral motiva-
31 tions associated with being a good neighbor, and how those motives inter-
32 act with a desire for personal gain. This is clear from John and Reba's
33 behavior. The neighborhood, then, is not merely durable, but durable for
34 the right reasons. While it lacks some balance, it is still largely balanced
35 for the right reasons. Similarly, if Reba and other neighbors begin to care
36 less about what others are doing, and commit themselves to contributing
37 simply because it is the right thing to do, the neighborhood will remain
38 durable for the right reasons, and become more balanced for the right
39 reasons. Consequently, Reba and John's neighborhood becomes more
40 well-ordered.

41

42

43

VI. RELAXING COMPLIANCE

44

45 I would now like to model a well-ordered society that relaxes the assump-
tion that agents fully comply with its rules. While reasonable agents may

1 decide not to cooperate *only if* they lack assurance, some agents are now
 2 allowed to defect *even if* they have assurance. I relax the compliance
 3 assumption for two reasons. First, relaxing compliance allows us to
 4 isolate dynamics that allow a less well-ordered society to develop into
 5 a more well-ordered society, which is a critical part of nonideal theory
 6 understood as the theory of transition from present circumstances to
 7 the ideal.⁴⁴ Second, relaxing compliance allows us to identify a third attrac-
 8 tive conception of stability — *immunity*. Immunity in general specifies the
 9 degree to which a WOS can recover from external shocks. Immunity resem-
 10 bles durability in this way because both are forms of first-order stability,
 11 though they differ in what they measure. Durability measures social trust,
 12 whereas immunity measures the survival rates of reasonable agents in com-
 13 petition with rational agents.

14 Critically, there are as many conceptions of immunity as there are
 15 kinds of external shocks to a social system. A society might be immune
 16 because it can repel an invasion of defectors. Alternatively, a society
 17 might be immune from unexpected events like an economic supply
 18 shock on the grounds that it can quickly recover from a change in eco-
 19 nomic circumstances. In this paper, I will focus exclusively on immu-
 20 nity against the invasion of small numbers of defectors who can enter
 21 at variable rates depending upon how well defectors in the system per-
 22 form vis-à-vis reasonable agents.

23 I call the second WOS model a complex, *real* well-ordered society model
 24 to connote both that the model represents a wide range of social phe-
 25 nomena and that it is a closer representation of our real social challenges.
 26 The most important addition to the simple WOS model is the second type
 27 of agent that simply acts to maximize her utility in each play. I call these
 28 “merely rational” agents, or “rational” agents for short.⁴⁵ Their expected
 29 utility is determined by the game they believe they are playing with other
 30 agents. To demonstrate the robustness of political stability despite the
 31 presence of agents willing to defect, merely rational players are given pris-
 32 oner’s dilemma payoffs, such that their dominant strategy is defection.

33 Notice, then, that I have relaxed the compliance assumption by adding
 34 merely rational agents, and not by altering the strategies of reasonable
 35 agents. This assumption employs simpler decision-making algorithms,
 36 and so simplifies the model. It also allows us to see effectively isolate the
 37 damage that merely rational players can wreak on political stability. We will
 38 see that they hit durability hard.

39 To properly flesh out the model, we must recall the mechanism by
 40 which information about cooperation and defection is transmitted. In the
 41

43 ⁴⁴ For discussions of Rawls’s approach to ideal theory, see A. John Simmons, “Ideal and
 44 Nonideal Theory,” *Philosophy and Public Affairs* 38, no. 1 (2010): 5–36 and Amartya Sen,
 45 *The Idea of Justice* (Cambridge, MA: Harvard University Press, 2009), 52–74.

⁴⁵ The name is not meant to imply that the reasonable agents act irrationally in any sense.

1 real WOS, reasonable agents do not distinguish between themselves and
2 rational agents. *All they know* is which agents are presently cooperating or
3 defecting, or who cooperated or defected in their last play with an agent.
4 *All they do* with this information is face the general direction of agents who
5 cooperated and turn away from those who have defected. The emergent
6 effect is that reasonable agents cluster with one another and flee rational
7 agents. This creates what is, in effect, a form of network reciprocity, where
8 reasonable agents play games with reasonable players more often than
9 with rational players, though they do so without recognizing this fact. In
10 other versions of the complex model, I have given reasonable agents more
11 information, such as the effective propensities of each agent. In that case,
12 reasonable agents are much more effective at building cooperative net-
13 works, but I wanted to see if the same effect would hold under conditions
14 of extremely limited information. And, in fact, it does hold.

15 Rational agents also face agents who have cooperated, but they do not
16 *also* turn away from defectors. This difference is vital. Reasonable agents
17 can effectively run from defecting agents and towards one another, and
18 since rational agents always or nearly always defect, over time reason-
19 able agents evade rational agents, even though they cannot act based on
20 identifying a rational agent as a rational agent. What justifies this differ-
21 ential treatment? The main justification is that rational agents are not as
22 enthusiastic about cooperating *in comparison to* reasonable agents. They
23 want to interact with reasonable agents, but they do not dislike one
24 another. In contrast, reasonable agents dislike and flee merely rational
25 agents because they recognize that they, the reasonable agents, are sub-
26 ject to exploitation. Merely rational agents like cooperating, but they will
27 defect in order to benefit; reasonable agents like cooperating so much that
28 they prefer not to defect even when it benefits them. So reasonable agents
29 are driven into cooperative interactions with one another, despite having
30 very little information about how each other will behave. Note that the
31 rational agents use the same information that reasonable agents do, no
32 more and no less, but they are less desperate to escape defectors since
33 they have no intrinsic liking of cooperation. Again, an important, emer-
34 gent result of these different motivations is that reasonable agents find one
35 another quickly, such that they can increase the frequency of the games
36 they play with one another, allowing them to accumulate more resources.
37 This is because in doing so they can outscore the merely rational players
38 through a simple higher frequency of interactions. Even if they sometimes
39 interact with rational players and lose, the reasonable agents do better over
40 time on the whole. We will see that the capacity to build networks will have
41 important effects on the entry-exit dynamic introduced below by helping
42 reasonable agents signal for new agents to enter the system at a rate faster
43 than rational agents, which enables reasonable agents to resist invasion.

44 I think my assumption is backed up by some recent work in social evo-
45 lutionary theory, which shows that cooperative agents can out-compete

1 uncooperative agents by correlating their behavior. The property assigned
 2 to cooperative agents is similar to the notion of network reciprocity, which
 3 again involves cooperative agents finding ways to interact primarily with
 4 other cooperative agents.⁴⁶ Defector agents have a generally harder time
 5 forming cooperative relations because they are uncooperative in general,
 6 and network reciprocity requires consistent cooperation if it is to form and
 7 grow.

9 VII. RESULTS OF THE REAL WOS ABM

11 The introduction of merely rational players into the environment has
 12 significant negative effects. Suppose we have fifty reasonable agents and
 13 five merely rational agents. So long as the merely rational players play
 14 Prisoner's Dilemmas, under a broad range of payoffs, they will depress
 15 durability. Intrinsic propensities and durability are tightly correlated in
 16 the simple WOS model, as an intrinsic propensity of .9 will always yield
 17 a long-run average social trust level of .9. The introduction of rational
 18 agents drags durability down from that level. Further, as intrinsic pro-
 19 pensity increases, the degree of durability depression increases, from
 20 around 11.2% when reasonable agents have an intrinsic propensity of .7 to
 21 16% when reasonable agents have an intrinsic propensity of 1. The reason
 22 for this increasing degree of durability depression is that merely rational
 23 agents can more easily exploit players with very high intrinsic propen-
 24 sities to cooperate, as reasonable agents' love of cooperation will make
 25 them vulnerable to defection more and more often.

26 Durability is *dramatically* depressed as social sensitivity increases,
 27 especially at an intrinsic propensity of 1, as represented in Figure 3.⁴⁷ At a
 28 social sensitivity level of .2, durability is only depressed 2 percent by
 29 the presence of a few rational agents. But depression increases to 12% at
 30 a social sensitivity of .6 and a whopping 49% at a social sensitivity of .9.⁴⁸
 31 The reason for the increased depression is that, as social sensitivity increases,
 32 an agent's intrinsic propensity to cooperate counts for less and less in deter-
 33 mining agent choices, from 80% at a social sensitivity of .2 to 10% at a social
 34 sensitivity of .9. So the dynamics of observation and correlation take over at
 35 high social sensitivities, such that hard-wired rational agents can dramatically
 36 disrupt the system.⁴⁹

37 Let's now extend our walkthrough to the real WOS model. Assume
 38 that John and Reba's neighborhood has achieved durability and balance.

40
 41 ⁴⁶ I worry that network reciprocity effects only work for small groups. To adjust for this
 42 possibility, I have instructed each reasonable agent not to keep tabs on its own record with
 43 each agent, but rather instead to calculate a general location that is the average x and y coor-
 44 dinates of all agents who have cooperated in the previous turn or are presently cooperating.

⁴⁷ I omit social sensitivity levels of 0 and 0.1 since they give little or no weight to observation.

⁴⁸ I ignored a social sensitivity of 1, as this eliminates the influence of intrinsic propensity.

⁴⁹ Merely rational players have a small effect on balance, so I set it aside here.

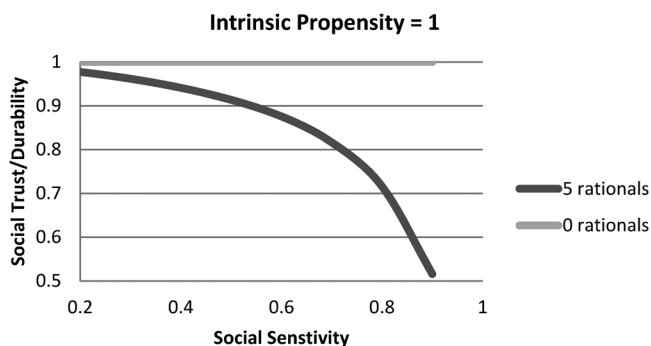


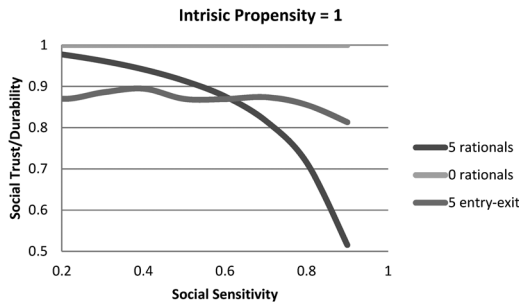
FIGURE 3. Fall in Durability as Social Sensitivity Increases

The neighborhood is nearly well-ordered. Then Sarah moves in. Sarah is a pure opportunist. She doesn't care about her neighbors at all, or what they think of her. She will engage in cooperative activities when she thinks it will benefit her and perhaps her family, but otherwise not. And in cases where she can drain social resources without consequence, she will. Sarah is the person who doesn't bring food to the block party, but drinks all the beer. She is the person whose house is safe but who never bothers to look out for dangers to the community. Sarah doesn't even recycle! If Sarah is caught engaged in uncooperative behavior, she will comply to divert attention away from her behavior, say by bringing food to the block party the next year, or bothering to peak out her windows when she knows her neighbors will see her. She picks up some trash here and there and charges below market prices at her yard sales. Since she's being watched, Sarah doesn't steal. Not for a while, at least.

John, Reba, and the other neighbors are generally able to determine that Sarah or someone else is a drain on the community. Recognizing that someone is deliberately violating social norms leads to somewhat more noncompliance by John and Reba. John thinks that, with everyone so mad at the shirker, people won't notice if he shirks a bit. And perhaps Reba gets so angry at the shirkers that she reciprocates with non-cooperative behavior. As a result, neighborhood trust falls, more defection ensues, and attempts to reestablish order are punctuated and brief, and only occasionally successful.

VIII. ENTRY AND EXIT

I now introduce the final major feature of the model — entry and exit. The real WOS model has a carrying capacity parameter that sets how many agents the system can hold. Once carrying capacity is reached, an exit algorithm kicks in by removing agents at random at a rate sufficient to keep the total number of agents under the carrying capacity. Further, all agents can



AQ1 FIGURE 4. Durability with rational agents and the entry-exit dynamic

“reproduce” up to the carrying capacity, which I represent as introducing a new agent into the system with the same agent-strategy as the “parent” agent-strategy—either reasonable or merely rational.⁵⁰

In the real WOS ABM, entry is cued by the “profit” margins that one agent strategy gains over the other. If the profit margin is 0, then whenever one *agent-strategy’s* score exceeds the others’ score, that strategy receives a new agent. If the margin is set at .02, then if one agent hard-wired to play that strategy will enter. The reasonable agents receive payoffs, just as rational agents do, even though they do not make decisions based on payoffs but rather merely on their expectation of reciprocal cooperation. So when I index entry rates to relative payoffs, I am appealing to payoff information for both rational agents and reasonable agents that reasonable agents do not consult. I find this is a helpful simplifying assumption for distinguishing conceptions of stability.

I should say a bit about the scoring system. All agents receive set payoffs depending upon the combination of plays by the agent and its partner, ranked in line with the formal structure of the Prisoner’s dilemma, where the payoffs are ranked as follows for the agent whose strategy is represented by the first variable: $DC > CC > DD > CD$. The agent gets the most when she defects and the other cooperates, less when both cooperate, still less when both defect, and even less when she cooperates and the other defects. The system then stockpiles the payoffs received in each game. An agent-strategy’s score is the sum of the stockpiles of each agent playing that strategy.

When the exit and entry algorithms are combined, one agent strategy can quickly replace another. If one strategy has a higher score than the other strategy more often than the reverse, it will reproduce faster. Since the exit algorithm removes agents from the environment at random, it will tend to remove more agents with the losing strategy than agents

⁵⁰ Agents enter with the average score of those who share its strategy.

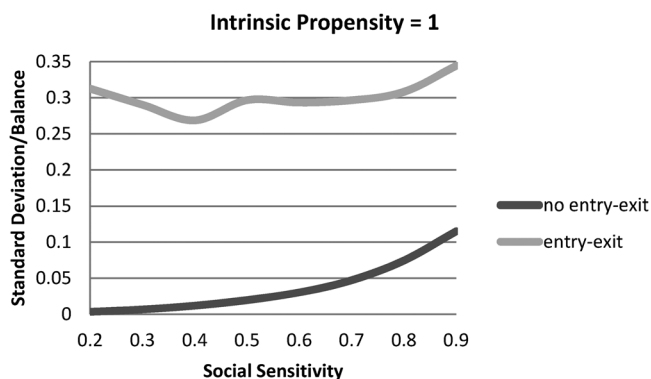


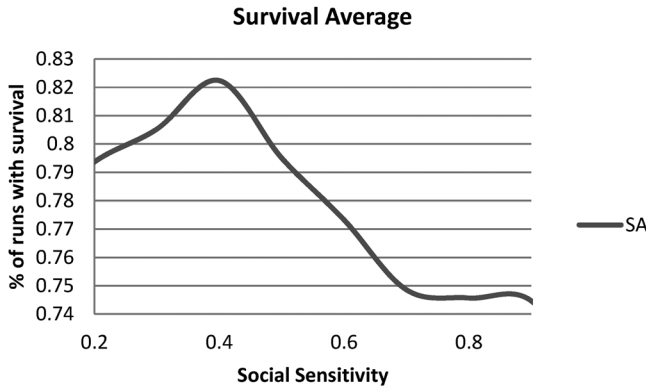
FIGURE 5. Balance variation due to the entry-exit dynamic.

with the winning strategy. Entry and exit algorithms, then, create a new kind of equilibrium condition that is similar to an *evolutionarily stable strategy* in evolutionary game theory—one that can resist invasion by alternative strategies in a particular environment when the new strategy is relatively rare.⁵¹

The primary upshot of the complex well-ordered society model is identifying our conception of immunity, which is realized when reasonable agents can repel invasion by rational agents. So a society has significant immunity when the majority population of reasonable agents resists replacement by rational agents. A society is perfectly immune when it not only resists invasion by rational players but also actively out-reproduces them, replacing them up to the full carrying capacity of the environment. If, for instance, we start with 10 rational agents and 90 reasonable agents, with a carrying capacity of 200, then a society has full immunity when it always reaches equilibrium at 200 reasonable agents and 0 rational agents at some point in the not too distant future. It has a high, but not perfect degree of immunity when its equilibrium state is, say, composed of 90% reasonable agents and 10% rational agents.

To be clear, immunity does not increase during a model's run, and can only be attributed to a system once the carrying capacity has been reached and the competition between agent-strategies begin. So the immunity of a system is a property described by its final equilibrium state. The society is immune if and only if its end state retains a much larger number of reasonable agents than rational agents; and it is fully immune if it can repel a sizeable number of rational agents entirely. So a WOS is fully immune when reasonable players resist replacement under a very robust range of conditions.

⁵¹ Gaus, *On Philosophy, Politics, and Economics*, 135–42 has a concise and clear discussion of the idea of an evolutionary stable strategy for the uninitiated.



16
17
18
19
20
21
22
23
24
25
26

FIGURE 6. Survival Rates

27
28
29
30
31
32
33
34
35
36
37
38
39

Raising immunity is a difficult social achievement, due to merely rational agents who believe they face PD payoffs with all other agents. Because of the randomness in the model, rational agents can sometimes achieve a higher average score long enough to significantly increase their proportion of the population. In a few cases, rational players take over. But there are more cases where reasonable agents repel rational agents long enough for the rational agents to die out. Once reasonable agents are victorious, *a real WOS is logically indistinguishable from a simple WOS with an entry-exit dynamic*. The simple reason is that all the rational agents have left, and the carrying capacity of the society has been reached, so reasonable agents are only interacting with one another.⁵²

40
41
42
43
44
45

The significant result is that the real WOS has a dynamic that creates a transitional path to a fully well-ordered society—a simple WOS with an entry-exit dynamic. This means that we can connect non-ideal theory and ideal theory. For Rawls, a WOS is populated exclusively by reasonable agents; Rawls classifies problems with non-compliant, unreasonable agents as non-ideal theory.⁵³ A social order with sufficient immunity can resist invasion by a large number of non-compliant agents, which means that it can transition from a non-ideal, or less-than-well-ordered society, to the ideal of a well-ordered society.⁵⁴ The WOS model developed here, then, can play an important role both in understanding how a non-ideal order can transition into a WOS.

⁵² The data I have compiled demonstrates that a simple WOS with an entry-exit dynamic is both durable and balanced.

⁵³ Rawls, *Theory of Justice*, 214. Also see Gaus, *The Tyranny of the Ideal*.

⁵⁴ For a detailed discussion of nonideal theory as a form of exploration of ways to realize certain social and political ideals, see Gerald Gaus and Keith Hankins, "Searching for the Ideal: The Fundamental Diversity Dilemma," in Kevin Vallier and Michael Weber, eds., *Political Utopias: Contemporary Debates* (New York: Oxford University Press, 2016), forthcoming.

1 Another nice feature of the entry-exit dynamic is that it allows us to rep-
2 resent a number of social factors that influence political stability. Recall that
3 Rawls only allows agents to enter his model by birth and exit by death.⁵⁵
4 But with a sufficiently high entry and exit rate, we can model a relaxa-
5 tion of this requirement. For instance, the real WOS ABM can potentially
6 account for emigration and immigration effects on stability, along with
7 generational shifts. Alternatively, we can use the entry and exit of rational
8 and reasonable agents to represent strategy changes among people in a
9 society. Specifically, we can represent this event by removing a reasonable
10 agent from the system and replacing it with a rational agent.
11

12 IX. RESULTS OF THE ENTRY-EXIT DYNAMIC IN A REAL WOS ABM 13

14 The full real WOS model illuminates the notion of immunity based on two
15 simplifying assumptions. First, rational players always defect, since they
16 simply maximize expected utility in each game they play. They respond to
17 no incentive to become more cooperative. Second, players only know what
18 other agents are presently doing or how they played in their very last play.
19 They do not know player reputations, or each agent's strategy-type.

20 In light of these simplifying assumptions, what factors cause immu-
21 nity? One obvious factor is the ratio of reasonable to rational players in the
22 model's starting conditions. With many rational agents in the system from
23 the outset, reasonable players seldom replace rational players, and rarely
24 do so entirely. This is not a problem for the model, since it is fair to assume
25 that most people are not pure defectors, but instead are conditional coop-
26 erators who care about reciprocity and fairness. The question is whether
27 a few bad apples can spoil the bunch. A second factor is that reasonable
28 agents can flee rational agents, and so form tight networks of reciprocal
29 benefit. Reasonable agents play many games together, increasing their
30 average scores faster than merely rational agents *who deliberately gravitate*
31 towards those cooperative hubs. With a higher average score, reasonable
32 agents enter more quickly the rational agents, and usually replace them.

33 Readers can explore the parameters and results of the model through the
34 outboard appendix. But we can see from Display 4 how much the entry-
35 exit dynamic matters. Here I stick to an intrinsic propensity of 1 and graph
36 durability as a function of social sensitivity (excluding a social sensitivity of
37 1, since this makes intrinsic propensity irrelevant in agent calculations, and
38 social sensitivities of 0 and .1, which are too low).

39 Introducing five merely rational agents into the simple WOS significantly
40 decreases durability, and more so as social sensitivity increases. But with the
41 entry-exit dynamic, the depression of durability is considerably reduced.
42 As the agents become more socially sensitive, they are *much* better at resisting
43
44

45 ⁵⁵ Rawls, *Political Liberalism*, p. xliii.

1 invasion by rational agents vis-à-vis their counterparts in the system with-
2 out the entry-exit dynamic; the slope of falling durability is much lower.
3 Over the course of a thousand runs of the model, to 500 ticks, the intro-
4 duction of the entry-exit dynamic increases the durability of the system by
5 30%, with a 20% reduction in durability from the simple WOS model with
6 only reasonable agents. Introducing merely rational agents still depresses
7 durability, but its effects are limited.

8 The entry-exit dynamic has a remarkable effect on balance. If we com-
9 pare complex WOS models, one with an entry-exit dynamic, and the other a
10 closed system, we generate sharply divergent results. Without the entry-exit
11 dynamic, the standard deviation of the system is quite low, rising from 0 at
12 a social sensitivity of .2 to .1 at a social sensitivity of 1.0. The system with the
13 entry-exit dynamic is far more volatile.

14 The introduction of the entry-exit dynamic generates a very fat tailed
15 curve, where all levels of social sensitivity generate a standard deviation
16 between .25 and .35, rather than 0 and .11. What this suggests is that
17 the full WOS model has low balance and that restricting entry and exit
18 conditions could increase balance. But notice that, just like the simple
19 WOS model, a social system can have a high degree of durability despite
20 depressed balance.

21 Survival is more complicated. At an intrinsic propensity of 1, agents
22 prefer to cooperate all of the time, but they discount this liking the
23 more socially sensitive they become. Without the entry-exit dynamic,
24 survival is irrelevant, since the number of rational and reasonable
25 agents does not change. But with the entry-exit dynamic, destruction is
26 a live option. Consider Display 6, which represents the average degree
27 of survival of 1000 runs at each tenth of social sensitivity (excluding 0, .1,
28 and 1).

29 Here we can see that survival is by no means assured in the full WOS
30 model. However, the large majority of the time (74%–84% of the time)
31 reasonable agents take over the system. And many of these runs, were
32 they allowed to exceed 500 ticks, would yield a higher level of survival.
33 Once the reasonable agents get to a sufficient size, they can resist invasion
34 permanently, and so will ultimately survive rational agents. So Display 6
35 *understates* the degree of immunity in the system.

36 Immunity, balance, and durability bear interesting and sometimes unex-
37 pected relationships to each other. A system with low immunity will be
38 destroyed, such that balance and durability will disappear. So immunity
39 is a precondition for balance and durability. However, once a society has
40 a sufficient degree of immunity, increases in balance and durability have
41 little further effect on immunity. We have already seen that durability and
42 balance come apart as well. The version of the complex model in Display
43 6 represents a society with an intrinsic propensity to cooperate of 1. If we
44 focus on versions of the system with a social sensitivity of .5, we find a
45 durability of 86%, balance at .296, and immunity of 79%. The model also

1 shows that systems with high degrees of balance (low standard devia-
2 tions) can also have high durability and immunity.

3 Let's bring these types of stability into concrete terms via our walk-
4 through. We can represent the ideas of entry and exit by movement in and
5 out of Reba and John's neighborhood. Suppose that Sarah is a successful
6 defector — a career opportunist — and so frustrates Reba that Reba decides
7 to move across town. Or perhaps Sarah convinces her like-minded friends
8 that Reba's community is ripe for the taking. In both cases, Sarah and her
9 friends will become a larger portion of the population relative to Reba's
10 group of cooperators. Alternatively, Reba and John's strategy of avoiding
11 interactions with Sarah's defector cohort might so deprive her of the rel-
12 evant community goods that she moves to another neighborhood where
13 she can take better advantage of others. Sarah will still pursue interactions
14 with John and Reba, but they are busy building one another up through
15 cooperation, and so only interact with defectors periodically.

16 Such an order will possess a high degree of immunity, due to Reba and
17 John's ability to resist the invasion of Sarah and her pack of opportunists.
18 John and Reba gradually, and perhaps without realizing it, build a coop-
19 erative neighborhood together, ignoring Sarah's occasional incursions.
20 In doing so, John and Reba establish a neighborhood with a high degree
21 of long-run social trust. However, to represent low balance, we should
22 allow that John and Reba are observant and care a great deal about what
23 others are doing. Thus, the prospect of Sarah and her pack invading, and
24 their occasional appearance can cause social trust to break down quickly,
25 but also with the ability to return to high levels of social trust in the future.

26 27 X. A WELL-ORDERED SOCIETY: DURABLE, BALANCED, AND IMMUNE

28
29 My results are tentative because my model has a number of draw-
30 backs. First, the agents in the WOS are extremely cognitively limited.
31 They cannot track different agent strategies, they make extremely coarse-
32 grained observations of cooperation, rational agents always defect, rea-
33 sonable agents don't pay attention to their payoffs, players cannot change
34 agent-strategies, and reasonable agents flee non-cooperators more effec-
35 tively than rational agents. But simple models have virtues. The weak
36 assumptions of the ABM, and ABM modeling in general, allow us to
37 enhance our conception of political stability by illuminating distinct
38 social processes that can be stable or unstable.

39 The ABM has three implications beyond showing that a well-ordered
40 society must be durable, balance, and immune for the right reasons.
41 First, even setting immunity aside, distinguishing durability and balance
42 reveals that we may not be able to establish stability via a single social
43 mechanism. Consequently, the arguments made on behalf of various assur-
44 ance mechanisms in a WOS are threatened. For example, some political lib-
45 erals have argued that complying with the requirements of public reason

1 can help citizens of a well-ordered society assure one another that they are
2 committed to one or a small set of political conceptions of justice. The point
3 of assurance is to generate a political order that is stable for the right reasons.
4 But if the ideal of political stability is deeply ambiguous, it is no longer
5 clear what standard assurance mechanisms accomplish.

6 Second, the ABM identifies three new lines of research within the public
7 reason project: (i) identifying different assurance mechanisms for different
8 types of stability, (ii) determining how socially sensitive reasonable persons
9 must be, or the extent to which we should allow their social sensitivity to vary,
10 and (iii) uncovering varied transitional paths from a society that is not well-ordered
11 to one that is. On this final point, we can see from the model that a real WOS
12 can transition into a simple WOS so long as cooperative agents can form cooperative
13 networks that allow them to systematically out-produce and replace rational agents.⁵⁶
14 Immunity connects ideal and non-ideal theory by helping us understand how a non-ideal
15 order can transition into a WOS.
16

17 Finally, I believe the ABM allows us to develop a more sophisticated
18 public reason liberal approach to constitutional choice. This paper is part
19 of a broader project that attempts to specify how contractarian parties can
20 choose a constitutional arrangement for themselves. I have developed
21 a generic model with three stages, the last of which concerns political stability
22 and that is specified by the model developed here. Within public reason liberalism,
23 constitutional arrangements must be stable for the right reasons, so whichever
24 rules are chosen must be subjected to a stability test. This essay specifies that test.
25 A constitutional rule is publicly justified only when it generates an adequate
26 degree of durability, balance, and immunity. The last factor, immunity, will
27 prove especially important if we believe that constitutional rules should be
28 selected under non-ideal conditions, such that constitutional rules will contain
29 provisions for discouraging defection. Thus, the model developed here can
30 specify the choice of constitutional rules in both ideal and non-ideal
31 circumstances.
32

33 *Philosophy, Bowling Green State University*
34
35
36
37
38
39
40
41
42
43

44 ⁵⁶ An underanalyzed assumption in the essay is that agents interacting over time can
45 accumulate resources, building on previous cooperative effort. Thus, immunity might be a function
of a growing economy. I hope to explore this effect in a future paper.