


# ”Y’all are just too sensitive”:

A computational ethics approach to understanding how prejudice against marginalized communities becomes epistemic belief

Johannah Sprinz 

June 30, 2022

LMU Munich<sup>1</sup>

**Abstract** Members of marginalized communities are often accused of being “too sensitive” when subjected to supposedly harmless acts of microaggression. This paper explores a simulated society consisting of marginalized and non-marginalized agents who interact and may, based on their individually held convictions, commit acts of microaggressions. Agents witnessing a microaggression might condone, ignore or condemn such microaggressions, thus potentially influencing a perpetrator’s conviction. A prototype model has been implemented in NetLogo, and possible applications are briefly discussed.

---

<sup>1</sup>This research was originally conducted as part of the Master’s practical seminar ”Computational Ethics” at LMU Munich in the winter term of 2021, supervised by Prof. Dr. François Bry. Released as an open-access preprint licensed CC BY 4.0; (c) 2022 Johannah Sprinz. Significant contributions by Matthias Fruth are acknowledged.

# 1 Introduction

**Motivation** Members of marginalized communities are often accused of being *too sensitive* when it comes to supposedly harmless acts of microaggression. Sexist comments are regarded as *locker room talk*. Racist jokes are viewed as *harmless fun*. Misgendering of trans people is downplayed as *a mistake that can happen*.

Microaggressions have been shown to cause significant emotional harm. Nonetheless, victims are often met with incomprehension, defensiveness, or even ridicule when pointing out microaggressions. Stereotypes like the *snowflake liberal* and the *shrill feminist who can't take a joke* paint members of marginalized communities as irrational and dismiss the pain and harm inflicted upon them.

**Contribution** This paper applies computer simulation to the phenomenon of microaggressions. Observing interactions of marginalized and non-marginalized agents in simulated societies with varying prevailing convictions should enable a better understanding of how prejudice against marginalized communities becomes an epistemic belief.

**Outline** The following section highlights notable related research in the fields of computational ethics and the social sciences. Next, the methodological principles for a model are outlined. A prototype for a NetLogo simulation has been implemented and its behavior is described. Finally, possible directions for future work on this topic are discussed.

## 2 Related Work

**Computational Ethics** Computational ethics makes use of "agent-based simulation [to apply a] a computational perspective to ethics theory" [1, 76] by simulating agents capable of adopting malleable ethical principles and observing how interactions impact their ethical principles. Such simulations can provide descriptive<sup>2</sup> insights into how individual ethical principles impact societal dynamics.

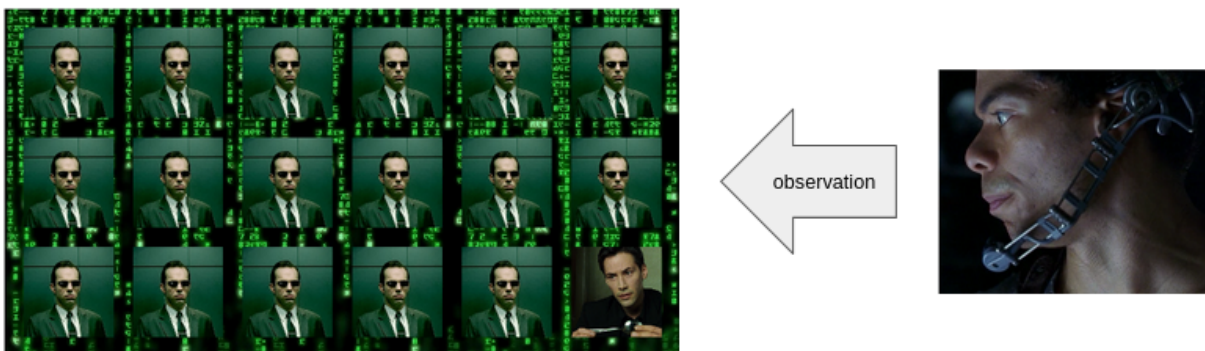


Figure 1: Computational ethics observes simulated societies consisting of interacting agents with individual ethical principles. Photos (c) 1999 Warner Bros.

---

<sup>2</sup>It can only be concluded how the input relates to the output (descriptive), not what input is required to achieve a certain output (perscriptive) [1, 76].

**Microaggressions** Microaggressions were initially conceptualized to describe unacknowledged racism experienced by African Americans [2]. Since other marginalized groups are subjected to similar types of discrimination, the definition has been extended and generalized over the years. More recent definitions literature defines microaggressions as "brief and commonplace daily verbal, behavioral, and environmental indignities, whether intentional or unintentional, that communicate hostile, derogatory, or negative racial, gender, sexual orientation, and religious slights and insults to the target person or group" [3]. Since prevailing biases influence all members of society, even well-intentioned people will unconsciously commit microaggressions. The relentlessness of these unintentional indignities can cause emotional harm comparable to more obvious intentional and outright malicious types of discrimination such as slurs and acts of physical violence [3] [4] [5].

**Colorblind Ideology** Some people, however, profess the idea of a *colorblind ideology* by which they claim to be uninfluenced by societal bias or even deny the existence thereof [6]. This kind of ideology provides individuals with numerous ways to downplay and deny the impact of racism [6]. Eisen identifies four types of colorblind racism: abstract liberalism, naturalization, cultural racism, and minimization of racism [6]. "By employing these frames, individuals divert their attention away from racism and effectively engage in racialized discourse without appearing racist" [6]. Individuals often react defensively when their beliefs are challenged. As such, reluctance to acknowledge microaggressions is a common reaction [6]. Such defensive reactions might include the argument that those affected by microaggressions are just too sensitive. However, people who proffer this argument rarely understand the nature of microaggressions and how constant exposure may affect people. For an unintentional perpetrator being criticized, it might seem that the person is overreacting because they are usually oblivious to the countless other incidents. Closer examination however reveals that "victims of microaggression need not to just develop 'thick skin' to overcome the epistemic harm of microaggression but need some form of support or epistemic resources to overcome these disadvantages" [5, 21].

**Approach** By a qualitative analysis, one can understand how the microaggressions and prejudices against people pointing them out may emerge on the individual level. Applying a computational ethics simulation might provide further insights into this dynamic on the societal level.

## 3 Methodology

This chapter introduces the conceptual ideas that constitute the model and describes the prototype implementation in the multi-agent simulation framework NetLogo [7].

### 3.1 Model

The model is constructed according to the following principles:

1. A society consists of a number of agents with individually held beliefs. We observe how much agents agree with the following two statements:
  1. "Microaggressions do not constitute a wrong."

2. "Marginalized agents are overly sensitive."
2. An agent might identify as part of the marginalized community within the society.
3. Agents interact with each other randomly.
4. A microaggression does not need to be targeted towards an individual agent.
5. The decision to commit a microaggression does not have to be conscious, nor does the perpetrator have to be aware that their action constituted a microaggression. Whether or not an agent will commit a microaggression depends on their conviction that microaggressions do not constitute a moral wrong. As such, they may commit a microaggression if and only if their agreement with statement 1 is sufficiently high.
6. An Agent witnessing a microaggression will react positively, neutrally, or negatively based on whether or not they consider microaggressions a moral wrong, i.e., their agreement with statement 1. A perpetrator's conviction may be influenced by the reaction they received, and the reacting agent's convictions might themselves be influenced by the reaction they have given. Negative feedback might be accepted or rejected by the perpetrator. If a perpetrator rejects negative feedback from a marginalized agent, the conviction that marginalized agents are overly sensitive may be reinforced.
7. Consumed media and other background noise might subconsciously influence agents' convictions positively or negatively independent from interactions.
8. For simplicity, only one type of marginalization is observed at a time. Intersectionality<sup>3</sup> is considered out of scope as of now.

## 3.2 Simulation

This section explains the prototype implementation<sup>4</sup> of the previously described model in as an agent-based simulation in NetLogo.

**Fundamental Principles** Marginalized and Non-Marginalized agents can be initialized with normal-distributed agreement percentages for two convictions:

- **c1:** "microaggressions do not constitute a wrong"
- **c2:** "members of the marginalized group are overly sensitive"

The simulation runs until one of three possible end conditions is reached:

1. **equilibrium:** no potential perpetrators: Microaggressions have effectively been eradicated; there are no agents with `c1 >= action_threshold`.

---

<sup>3</sup>The concept of intersectionality plays an important role when discussing topics of marginalization, as it enables the analysis of situations where a person is affected by multiple types of marginalization at once. The term was coined by Kimberle Crenshaw in 1989, highlighting that the marginalization experienced by black women "cannot be understood through an analysis of patriarchy rooted in white experience" [8, 156f].

<sup>4</sup>The model is available under <https://neothethird.gitlab.io/ceth-seminar/model.nlogo>.

2. **equilibrium: no negative reactors:** Microaggressions have become so normalized that no-one speaks up against them; there are no agents with  $c1 \leq \text{negative\_threshold}$ .
3. **deadlock: society is too polarized for change:** All agents either have very strong agreement or disagreement with  $c1$ . Change is improbable due to the high polarization.

Agents can interact with one another one-to-one and may commit microaggressions in these interactions. An agent witnessing a microaggression may react positively, negatively, or neutrally. This behavior is configured through thresholds that control Poisson-distributed probability values calculated for every interaction. Example: Suppose the `action_threshold` is set to 75%. All agents with  $c1 \geq 75$  may now theoretically become perpetrators. Whether or not an agent actually commits a microaggression is determined by  $c1 \geq \text{individual\_action\_threshold}$ , where `individual\_action\_threshold` is a poisson-distributed probability variable with a mean of  $(\text{action\_threshold} + (100 - \text{action\_threshold}) / 2) = 75 + 12.5 = 87.5$ . Thus, the higher the  $c1$  of an agent, the higher their chance of becoming a perpetrator.

Reactions are Poisson-distributed as well. Positive reactions are controlled by the `individual\_positive\_threshold`, a Poisson-distributed probability variable with a mean of  $(\text{positive\_threshold} + (100 - \text{positive\_threshold}) / 2)$ ; halfway between 100 and `positive\_threshold`. Negative reactions are controlled by `individual\_negative\_threshold`, a poisson-distributed probability variable with a mean halfway between 0 and `negative\_threshold`, i.e.,  $\text{negative\_threshold} / 2$ .

**Visualization** Agents are rendered on a two-dimensional plane (see Figure 2) based on their convictions with the  $x$ -axis representing  $c1$  and the  $y$ -axis representing  $c2$ .

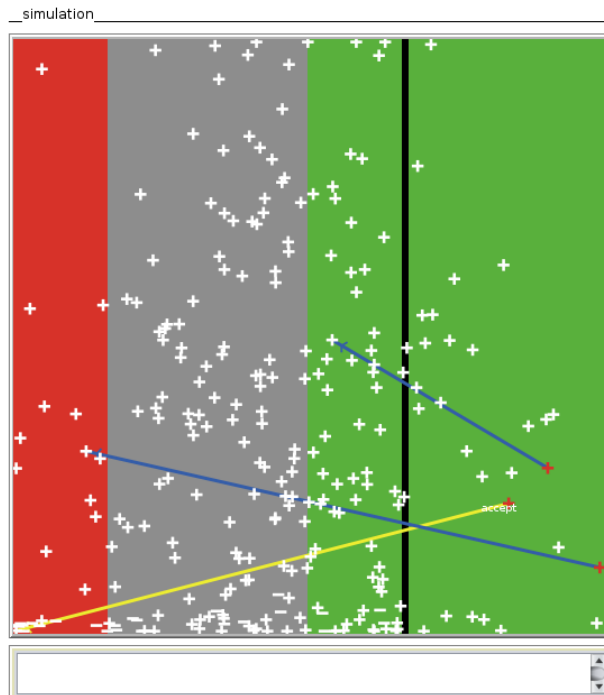


Figure 2: Screenshot of the NetLogo model's visualized society.

The `action_threshold` is represented by a vertical black line. Agents right of this line are potential perpetrators. The green area represents the `positive_threshold`; Agents in this area might react positively to microaggressions. The red area represents the `negative_threshold`; Agents in this area might react negatively to microaggressions. The grey area contains agents that will always react neutrally to microaggressions, as their `c1` is below the `positive_threshold` and above the `negative_threshold`. Marginalized agents are represented by a `-`, non-marginalized agents by a `+`. If an agent is tinted red, it will commit a microaggression in this tick. Otherwise, agents are tinted white. An arrow points from the perpetrator of a microaggression to the reacting agent. A pink arrow indicates a positive reaction to the microaggression; blue indicates neutral, and yellow indicates negative. If a perpetrator receives a negative reaction and accepts the criticism, they will be marked with the label "accept". Otherwise, with the label "reject". The output field beneath the map displays end conditions for the simulation. Monitors reporting values in the form of `[non_marginalized + marginalized = total]` (see Figures 3, 5, and 6) are used to show how many agents a specific claim applies to at any given time; overall and broken down by whether or not they are marginalized. The user interface also includes various plots (see Section 4) that are updated in real time.

**General Controls** Sliders `population` and `margin_size` (see Figure 3) control the number of agents and the percentage of marginalized agents. The `setup` button will initialize a society based on the current configuration, `go` will start the simulation. the `go once` and `go 5x` buttons can be used for step-by-step execution. Pre-configured scenarios can be loaded using the `load scenario...` button. Every scenario is explained after selection.

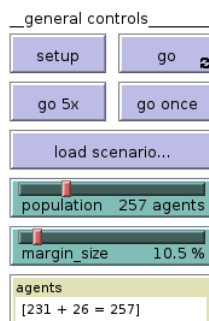


Figure 3: Screenshot of the NetLogo model's general controls. The society consists of 257 agents, 10.5% (=26) of which are marginalized, leaving 231 non-marginalized agents.

**Conviction Initialization** Convictions one (`c1`) and two (`c2`) for marginalized (`m`) and non-marginalized (`p`) agents are initialized as normal-distributed values with a mean and a deviation using sliders (see Figure 4) in the form of `<p|m>_<c1|c2>_<deviation|mean>`.

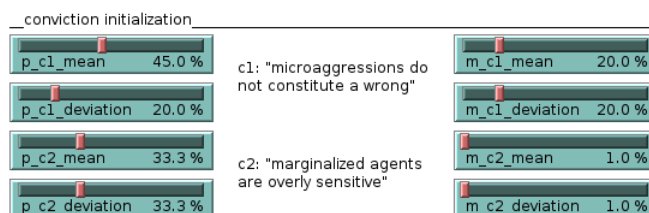


Figure 4: Screenshot of the NetLogo model's conviction sliders.

**Action Behavior** The `action_threshold` slider (see Figure 5) controls the minimum required `c1` for an agent to commit a microaggression.

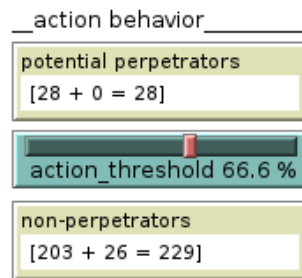


Figure 5: Screenshot of the NetLogo model's action controls.

**Reaction Behavior** The `positive_threshold` (see Figure 6) controls the minimum required `c1` for an agent to react positively to a microaggression. Similarly, the `negative_threshold` specifies the maximum allowed `c1` for an agent to react negatively to a microaggression.

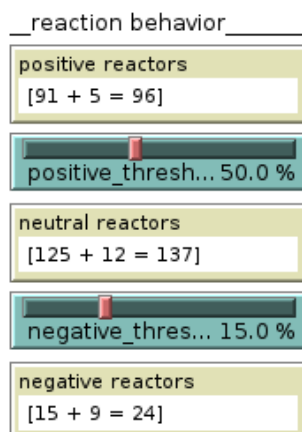


Figure 6: Screenshot of the NetLogo model's reaction controls.

**Miscellaneous Controls** The `critical_faculty` slider (see Figure 7) allows setting the likelihood of a perpetrator accepting criticism when made aware of a microaggression they committed. The `stealth` slider enables the simulation of types of marginalization that are less recognizable to other agents by specifying the likelihood of an agent not recognizing another agent as marginalized.

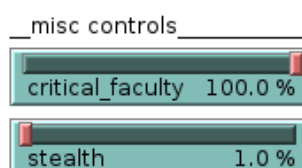


Figure 7: Screenshot of the NetLogo model's miscellaneous controls.

**Reactions** Changes in convictions one (c1) and two (c2) due to interactions can be set individually for marginalized (m) and non-marginalized (p) agents. If no microaggression occurs in an interaction, the `idle` value is applied. Otherwise, a `positive`, `neutral`, or `negative` reaction might be given to or received from a marginalized (m) or non-marginalized (p) agent. A received `negative` reaction might be `accepted` or `rejected`. Sliders (see Figure 8) take the form `<p|m>.<c1|c2>_on.<idle|<<positive.<to|from>|neutral.<to|from>|negative.<to|accepted_from|rejected_from>>.<p|m>>`.

Example: A privileged agent, Bob, commits a microaggression in an interaction with a marginalized agent, Alice, who reacts neutrally. Bob’s convictions change based on `p_c1_on_neutral_from_m` and `p_c2_on_neutral_from_m` respectively, while Alice’s convictions change based on `m_c1_on_neutral_to_p` and `m_c2_on_neutral_to_p`.

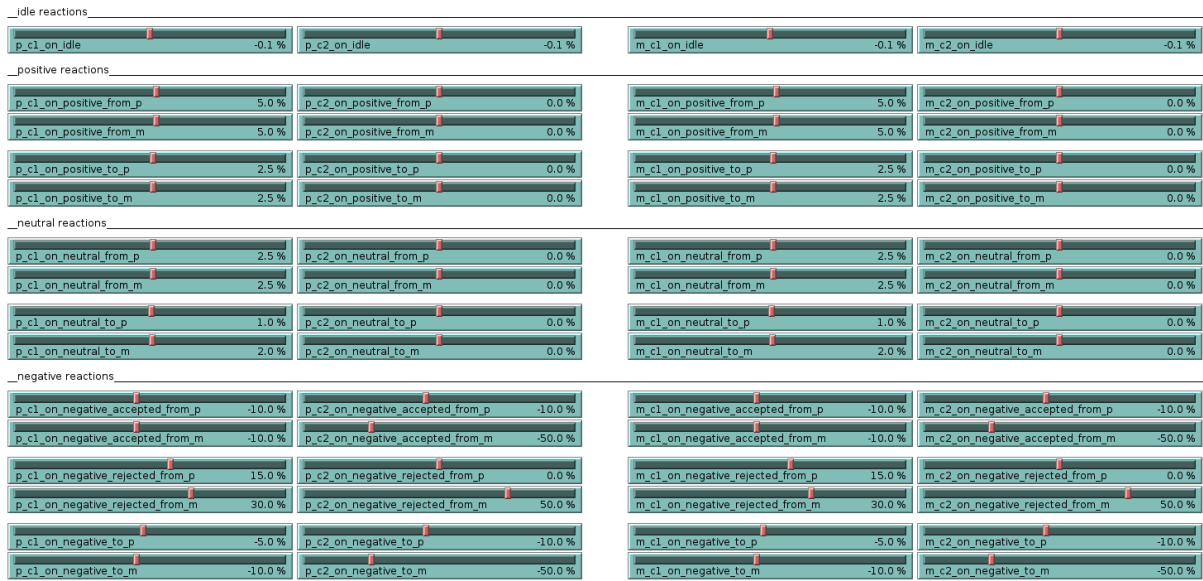


Figure 8: Screenshot of the NetLogo model’s reactions sliders.

**Noise** Background noise changes to convictions one (c1) and two (c2) for marginalized (m) and non-marginalized (p) agents are modeled as normal-distributed values with a mean and a `deviation` using sliders in the form of `<p|m>.<c1|c2>_noise.<deviation|mean>`.

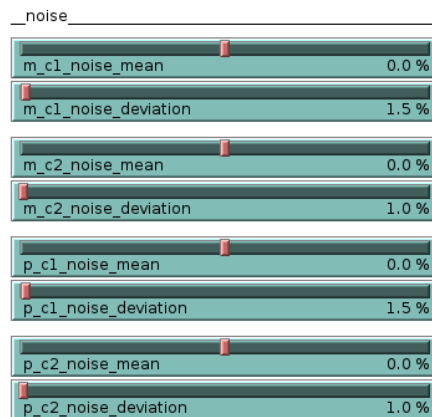


Figure 9: Screenshot of the NetLogo model’s noise sliders.



## 4 Basic Validation

This section provides some experimental results generated by the prototype. Note that this is not yet an attempt at full validation of the model and only serves to show how the model behaves. The two simulation runs with different conviction initialization values and otherwise identical configurations are compared to observe potential behavioral differences. The next subsection describes the configuration, followed by an outline of the results and noteworthy observations.

### 4.1 Configuration

The model is configured as follows; the full configuration is provided in Appendix A. The population consists of 500 agents, with a `margin_size` of 10.5%. The `stealth` likelihood is negligibly low at 1%, and `critical_faculty` is set to 50%. Ergo, marginalized agents are almost always read as such, and perpetrators of microaggressions accept criticism in half of the cases. A relatively low `positive_threshold` has been chosen at 50%, allowing agents to commit microaggressions even if they only slightly agree with `c1`. This is meant to model the unintentional nature of many microaggressions. The `action_threshold` is only slightly higher at 66.6%. Contrarily, the `negative_threshold` is fairly low at 15%, requiring high disagreement with `c1` for agents to react negatively. Non-marginalized agents are initialized to have generally higher agreements with both convictions, although there is significant variance. Marginalized agents generally have lower agreement with `c1`, and very low agreement with `c2`. For trial 1, `p_c1_mean` is set to 45%; for trial 2, `p_c1_mean` is set to 66.6%. This models a society where microaggressions are less and more common, respectively.

### 4.2 Results

**Trial 1:** `p_c1_mean = 45` Trial 1 starts out with a fairly balanced `c1` and `c2` among non-marginalized agents, and balanced `c1` and low `c2` among marginalized agents (see top-left image in Figure 10). Both convictions show a mostly steady decline (see top image in 11), resulting in a steadily growing number of negative reactors (see top image in Figure 12). The simulation results in the stop condition `equilibrium: no potential perpetrators` after 710 ticks (see top-right image in Figure 10). At this point, overall average `c1` is just below the `negative_threshold` of 15%.

**Trial 2:** `p_c1_mean = 66.6` Trial 2 starts out with generally higher `c1` and balanced `c2` among non-marginalized agents, and balanced `c1` and low `c2` among marginalized agents (see bottom-left image in Figure 10). At the start, about 80% of non-marginalized agents are potential positive reactors. The value quickly rises and stalls at about 99% (see top-right image in Figure 15). Marginalized agents start out at only about 9.5% positive reactors, but follow this trend with a small delay and stalling at about 85% (see bottom-right image in Figure 15). Overall average `c2` conviction falls and saturates at about 10% (see bottom image in 11). The simulation results in the stop condition `deadlock: society is too polarized for change` after 1508 ticks (see bottom-right image in Figure 10). At this point, overall average `c1` is about 97.5% (see bottom image in Figure 11) with the remaining agents maintaining low `c1` mostly belonging to the marginalized community (see left column in Figure 14).

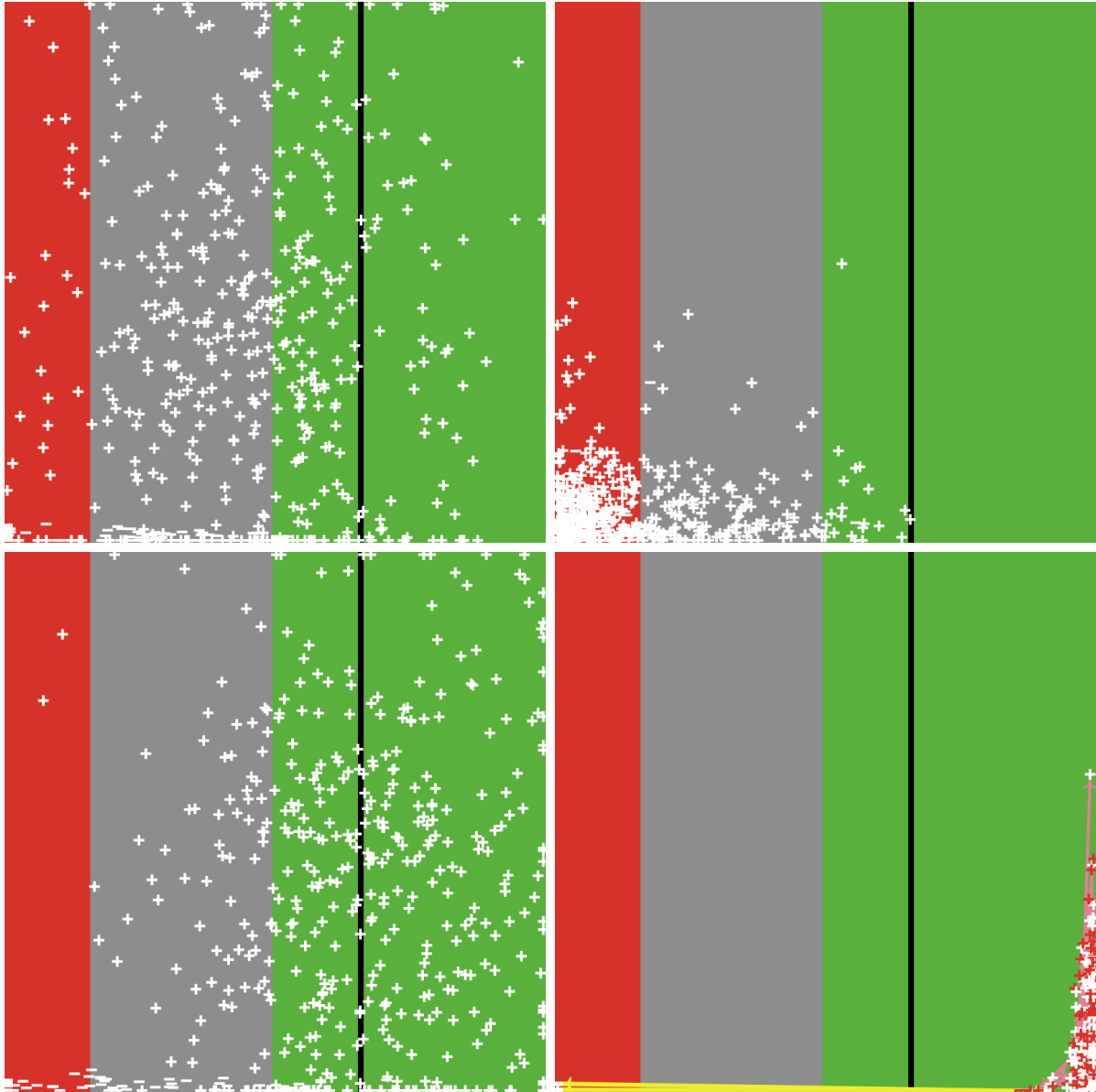


Figure 10: Visualization of the simulated society for trial 1 with  $p_{c1\_mean} = 45$  on top and trial 2 with  $p_{c1\_mean} = 66.6$  below at initialization on the left and after reaching the stop condition on the right.

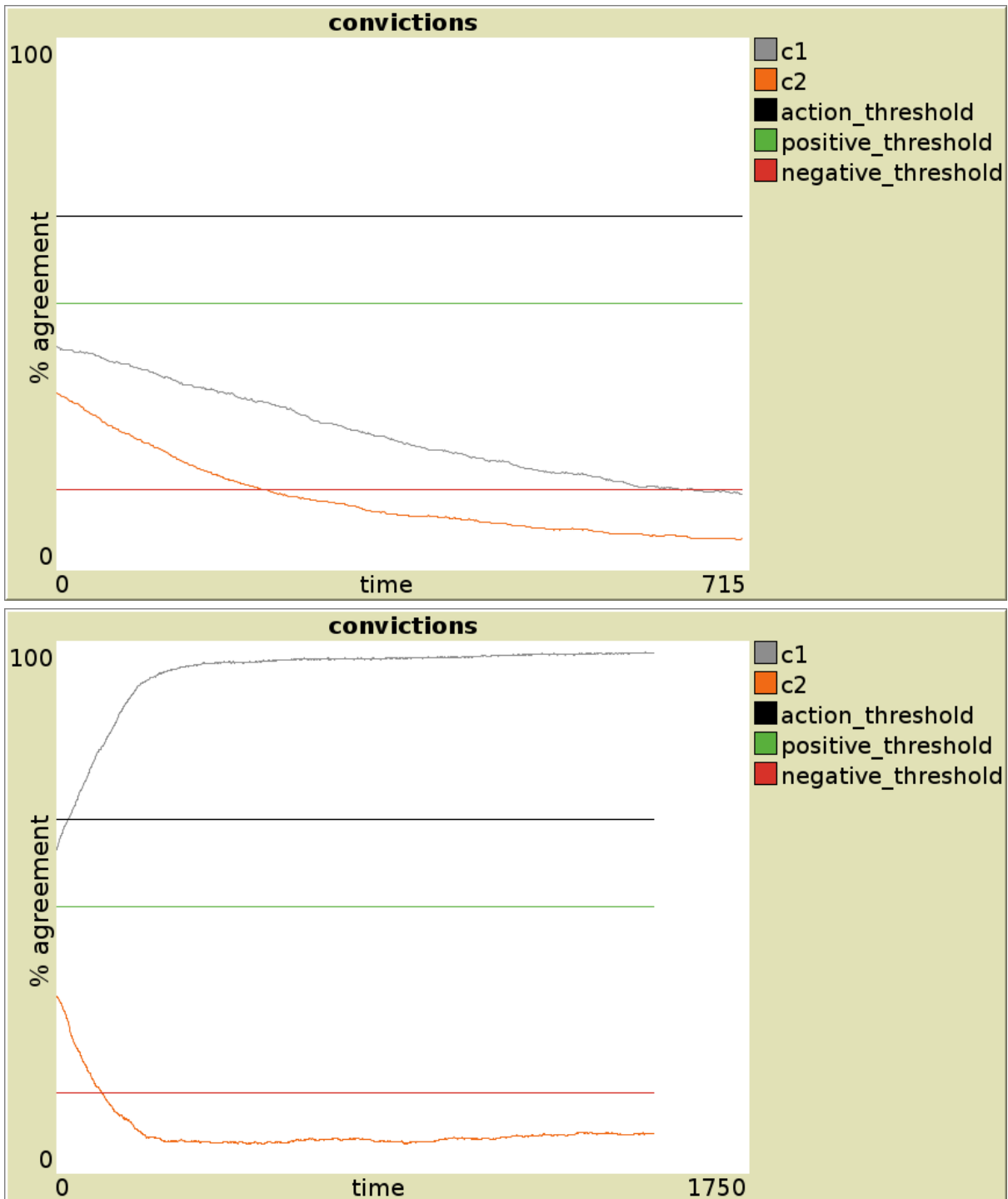


Figure 11: The convictions plot depicts the average agreement with c1 (grey) and c2 (orange) in the simulated society. For orientation, the `action_threshold` (black), `positive_threshold` (green), and `negative_threshold` (red) are visualized as well. The upper plot represents trial 1, the lower plot represents trial 2.

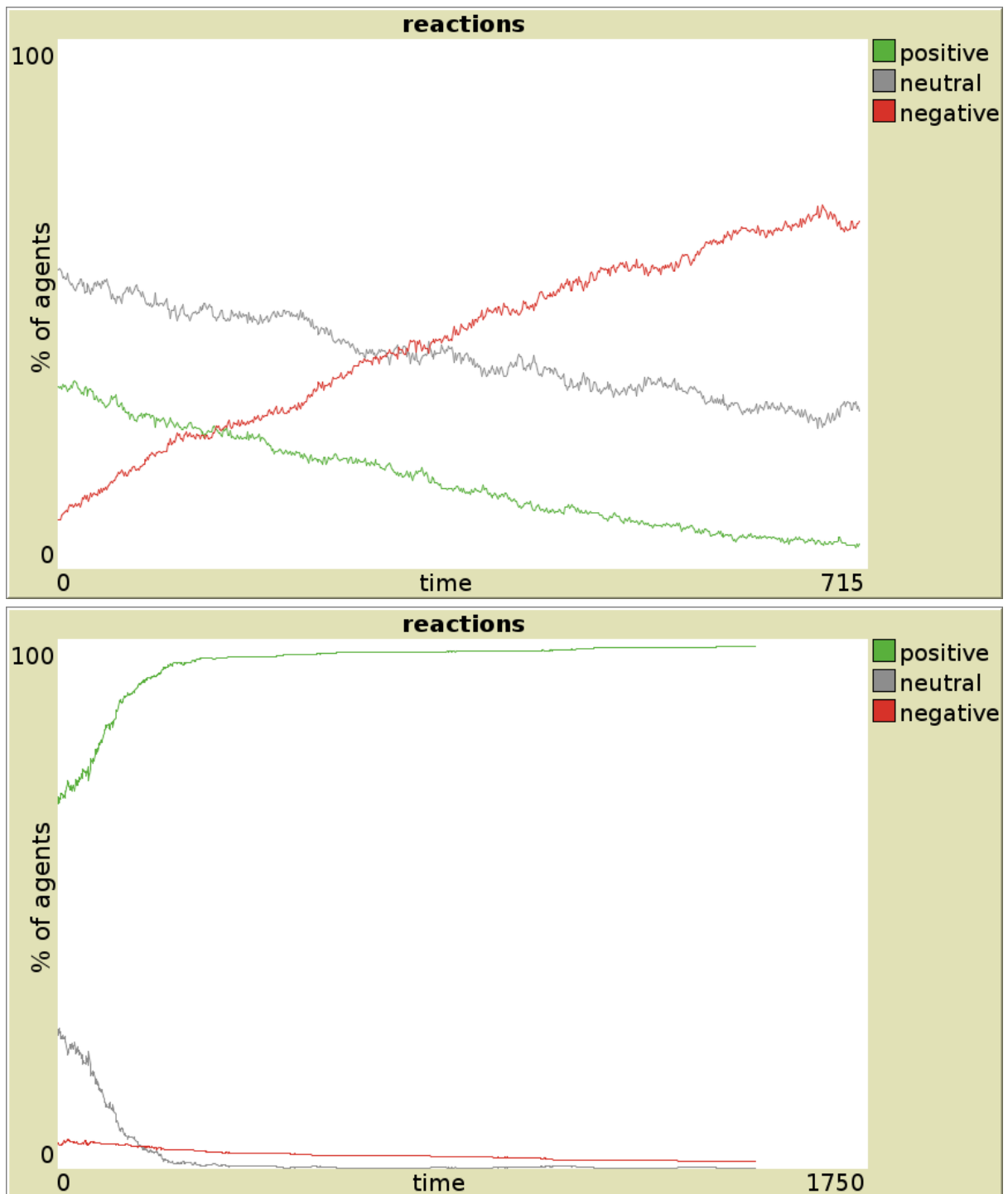


Figure 12: The reactions plot depicts the percentage of agents above the positive threshold (potential positive reactors, green), below the negative threshold (potential negative reactors, red), and in between the thresholds (neutral reactors, grey). The upper plot represents trial 1, the lower plot represents trial 2.

## 5 Discussion

Intuitively, the behavior seen in the trials seems plausible. Lower initial `c1`, as seen in trial 1, allows the agents to slowly be convinced of the harmful nature of microaggressions, while the higher initial `c1` of trial 2 results in a critical mass of agents being so radicalized that the society ends up becoming too polarized for change. This is an acceptable result for a first prototype, but the behavior should be analyzed and refined in future research.

The impact of the `noise` parameter requires further examination. With the low `negative_threshold` and `positive_threshold` used in this paper, the potential negative reactors have a much higher chance of being pushed above `negative_threshold` than a potential positive reactor has of being pulled below `positive_threshold`. Implementing different noise parameters for agents in different threshold ranges would be possible.

With the configuration used in this paper (see Section 4), overall `c2` still falls even in frequent perpetrators with high `c1`. The relevant parameters should be re-evaluated to allow more insights into the dynamic between `c1` and `c2`.

Due to NetLogo's simple design, complex constraints like intersectionality (see Section 3.2) can become challenging. Re-implementing the model in a simulation framework that supports an object-oriented class-based approach to dynamically model different types of marginalization might be advised.

## 6 Conclusion and Outlook

### 6.1 Future Work

The model should be improved and possibly re-implemented in a simulation framework that allows for more flexibility than NetLogo. This would make implementing additional configuration options for more granular control easier.

The model's default parameters should be refined and, if possible, validated against real-world empirical data. Different "scenarios" to explore the situation of specific real-world marginalized communities could be researched and pre-configured to be loaded in the model's user interface. Introducing a `false_positive_clock` configuration parameter analog to `stealth` could be added to model cases where agents might falsely read a non-marginalized agent as marginalized.

There are many possibilities for extending the background logic, such as allowing an agent's `c1` to influence their `critical_faculty`, and introducing interdependency between an agent's `c2` and `c1`.

### 6.2 Final Remarks

This paper set out to apply a computational ethics approach to the phenomenon of microaggressions and the prejudices against marginalized communities that arise from it. A basic model has been theorized, and a prototype for a simulation has been implemented in NetLogo. The fundamental viability of the approach has been demonstrated, and possible expansions for future research have been discussed. Readers are encouraged to conduct their own experiments with the model, which has been made available under <https://neothethird.gitlab.io/ceth-seminar/model.nlogo>.

## References

- [1] A. Ruvinsky, “Computational Ethics,” in *Encyclopedia of Information Ethics and Security*, M. Quigley, Ed. IGI Global, 2007, pp. 76–82, ISBN: 978-1-59140-987-8.
- [2] C. Pierce and F. B. Barbour, “Offensive mechanisms,” in *The Black Seventies*. Porter Sargent Publisher, 1970, pp. pp. 265–282.
- [3] D. W. Sue, *Microaggressions in everyday life: race, gender, and sexual orientation*. Hoboken, N.J.: Wiley, 2010, ISBN: 978-0-470-59415-5.
- [4] S. J. Kapusta, “Misgendering and Its Moral Contestability,” *Hypatia*, vol. 31, no. 3, pp. 502–519, 2016, doi: 10.1111/hypa.12259.
- [5] B.-R. Kelly, “How do Microaggressions Cause Epistemic Harm? On the Conceptual Difficulties Between Epistemic Injustice and Microaggression,” 2021, Manuscript, [https://bella-rosekelly.weebly.com/uploads/1/3/9/9/139952594/bellarose\\_how\\_do\\_microaggressions\\_cause\\_epistemic\\_harm.pdf](https://bella-rosekelly.weebly.com/uploads/1/3/9/9/139952594/bellarose_how_do_microaggressions_cause_epistemic_harm.pdf).
- [6] D. B. Eisen, “Combating the “Too Sensitive” Argument: A Demonstration That Captures the Complexity of Microaggressions,” *Teaching Sociology*, vol. 48, no. 3, pp. 231–243, Jul. 2020, doi: 10.1177/0092055X20930338.
- [7] U. Wilensky, “NetLogo,” 1999, at the Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL: <http://ccl.northwestern.edu/netlogo/>.
- [8] K. Crenshaw, “Demarginalizing the Intersection of Race and Sex: A Black Feminist Critique of Antidiscrimination Doctrine, Feminist Theory, and Antiracist Politics [1989],” in *Feminist Legal Theory*, 1st ed., K. T. Bartlett and R. Kennedy, Eds. Routledge, Feb. 2018, pp. 57–80, doi: 10.4324/9780429500480-5.

## Appendix A Model Configuration

Table 1: General configuration.

population	500
margin_size	10.5
stealth	1
critical_faculty	50

Table 2: Thresholds.

action_threshold	66.6
positive_threshold	50
negative_threshold	15

Table 3: Conviction initialization.

p.c1.mean	45 ∨ 66.6
p.c1.deviation	20
p.c2.mean	33.3
p.c2.deviation	33.3
m.c1.mean	20
m.c1.deviation	20
m.c2.mean	1
m.c2.deviation	1

Table 4: Conviction changes due to noise.

p.c1.noise.mean	0
p.c1.noise.deviation	1.5
p.c2.noise.mean	0
p.c2.noise.deviation	1
m.c1.noise.mean	0
m.c1.noise.deviation	1.5
m.c2.noise.mean	0
m.c2.noise.deviation	1

Table 5: Conviction changes due to Idleness.

p.c1.on.idle	-0.1
p.c2.on.idle	-0.1
m.c1.on.idle	-0.1
m.c2.on.idle	-0.1

Table 6: Conviction changes due to positive reactions.

p.c1.on.positive.to.p	2.5
p.c1.on.positive.from.p	5
p.c1.on.positive.to.m	2.5
p.c1.on.positive.from.m	5
p.c2.on.positive.to.p	0
p.c2.on.positive.from.p	0
p.c2.on.positive.to.m	0
p.c2.on.positive.from.m	0
m.c1.on.positive.to.p	2.5
m.c1.on.positive.from.p	5
m.c1.on.positive.to.m	2.5
m.c1.on.positive.from.m	5
m.c2.on.positive.to.p	0
m.c2.on.positive.from.p	0
m.c2.on.positive.to.m	0
m.c2.on.positive.from.m	0

Table 7: Conviction changes due to neutral reactions.

p.c1.on.neutral.to.p	1
p.c1.on.neutral.from.p	2.5
p.c1.on.neutral.to.m	2
p.c1.on.neutral.from.m	2.5
p.c2.on.neutral.to.p	0
p.c2.on.neutral.from.p	0
p.c2.on.neutral.to.m	0
p.c2.on.neutral.from.m	0
m.c1.on.neutral.to.p	1
m.c1.on.neutral.from.p	2.5
m.c1.on.neutral.to.m	2
m.c1.on.neutral.from.m	2.5
m.c2.on.neutral.to.p	0
m.c2.on.neutral.from.p	0
m.c2.on.neutral.to.m	0
m.c2.on.neutral.from.m	0

Table 8: Conviction changes due to negative reactions.

p.c1.on.negative.to.p	-5
p.c1.on.negative.accepted.from.p	-10
p.c1.on.negative.rejected.from.p	15
p.c1.on.negative.to.m	-10
p.c1.on.negative.accepted.from.m	-10
p.c1.on.negative.rejected.from.m	30
p.c2.on.negative.to.p	-10
p.c2.on.negative.accepted.from.p	-10
p.c2.on.negative.rejected.from.p	0
p.c2.on.negative.to.m	-50
p.c2.on.negative.accepted.from.m	-50
p.c2.on.negative.rejected.from.m	50
m.c1.on.negative.to.p	-5
m.c1.on.negative.accepted.from.p	-10
m.c1.on.negative.rejected.from.p	15
m.c1.on.negative.to.m	-10
m.c1.on.negative.accepted.from.m	-10
m.c1.on.negative.rejected.from.m	30
m.c2.on.negative.to.p	-10
m.c2.on.negative.accepted.from.p	-10
m.c2.on.negative.rejected.from.p	0
m.c2.on.negative.to.m	-50
m.c2.on.negative.accepted.from.m	-50
m.c2.on.negative.rejected.from.m	50



## Appendix B Plots

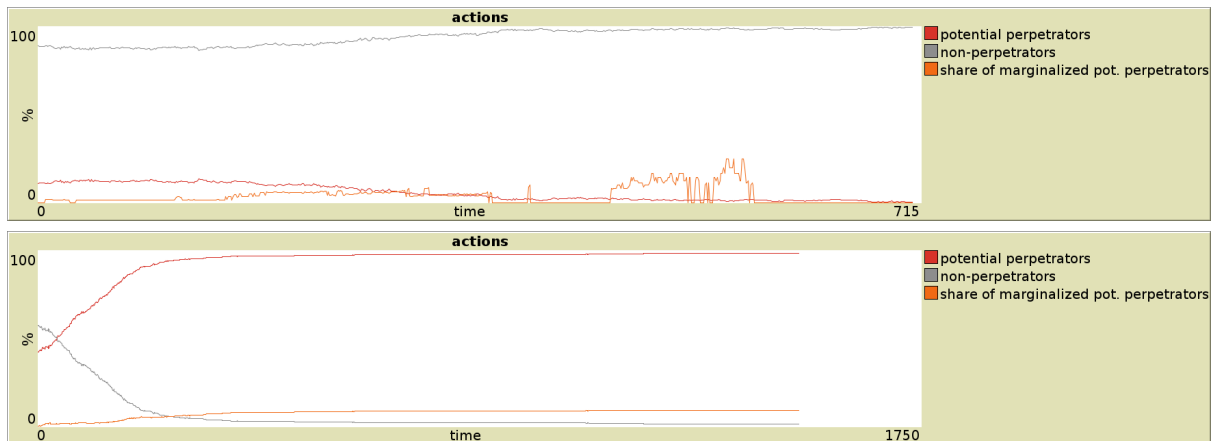


Figure 13: The actions plots depicts the percentage of potential perpetrators (red), the percentage of non-perpetrators (grey), and the share of marginalized agents among the potential perpetrators (orange) for trial 1 on top and trial 2 below.

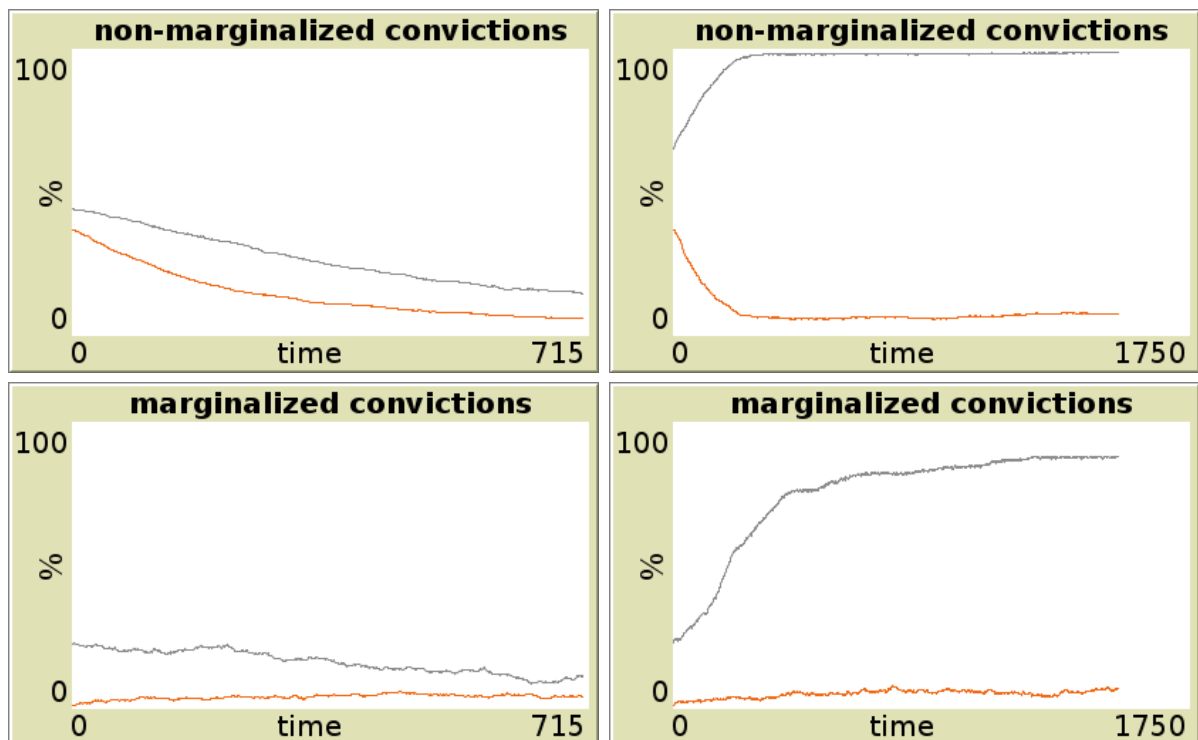


Figure 14: The (non-)marginalized convictions plots depict the average agreement with  $c_1$  (grey) and  $c_2$  (orange) among non-marginalized (top row) and marginalized (bottom row) agents in the simulated society for trial 1 on the left and trial 2 on the right.

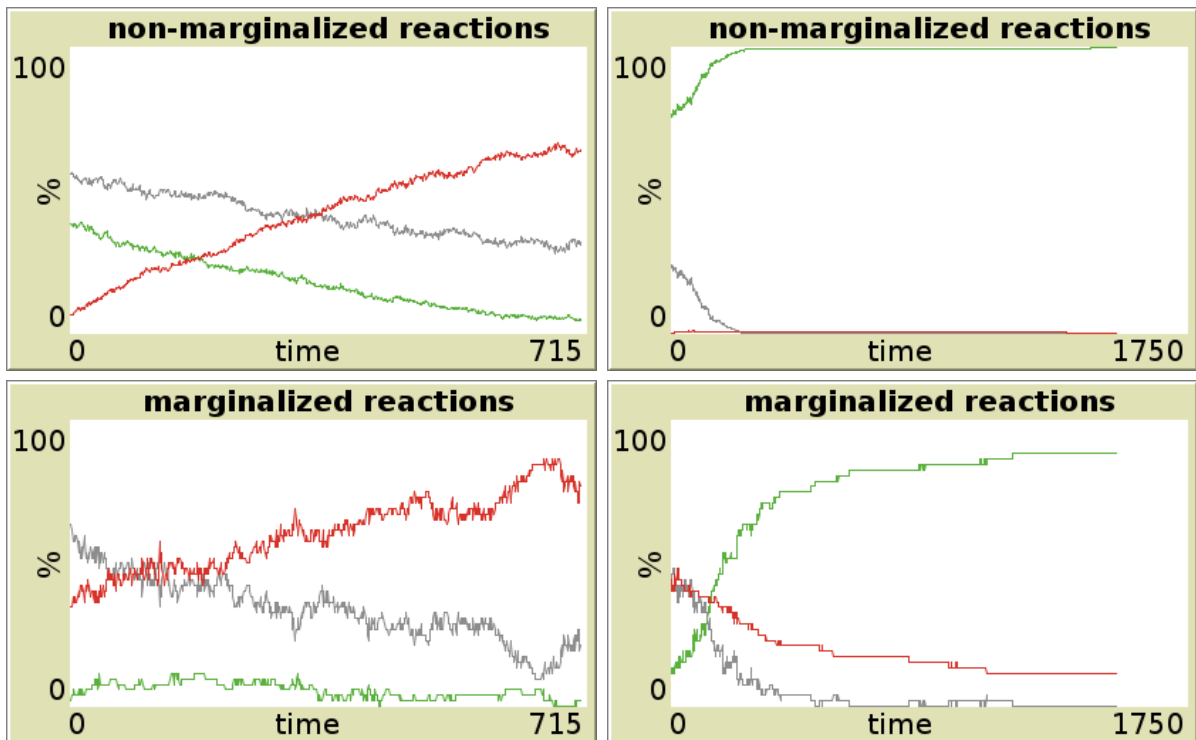


Figure 15: The (non-)marginalized reactions plots depict the percentage of non-marginalized (top row) and marginalized (bottom row) agents above the positive threshold (potential positive reactors, green), below the negative threshold (potential negative reactors, red), and in between the thresholds (neutral reactors, grey) for trial 1 on the left and trial 2 on the right.