# SCIENTIFIC REP⚙RTS

**OPEN**

# Modelling opinion dynamics in the age of algorithmic personalisation

Nicola Perra & Luis E. C. Rocha

Modern technology has drastically changed the way we interact and consume information. For example, online social platforms allow for seamless communication exchanges at an unprecedented scale. However, we are still bounded by cognitive and temporal constraints. Our attention is limited and extremely valuable. Algorithmic personalisation has become a standard approach to tackle the information overload problem. As result, the exposure to our friends' opinions and our perception about important issues might be distorted. However, the effects of algorithmic gatekeeping on our hyper-connected society are poorly understood. Here, we devise an opinion dynamics model where individuals are connected through a social network and adopt opinions as function of the view points they are exposed to. We apply various filtering algorithms that select the opinions shown to each user (i) at random (ii) considering time ordering or (iii) its current opinion. Furthermore, we investigate the interplay between such mechanisms and crucial features of real networks. We found that algorithmic filtering might influence opinions' share and distributions, especially in case information is biased towards the current opinion of each user. These effects are reinforced in networks featuring topological and spatial correlations where echo chambers and polarisation emerge. Conversely, heterogeneity in connectivity patterns reduces such tendency. We consider also a scenario where one opinion, through nudging, is centrally pushed to all users. Interestingly, even minimal nudging is able to change the status quo moving it towards the desired view point. Our findings suggest that simple filtering algorithms might be powerful tools to regulate opinion dynamics taking place on social networks.

Various disciplines, as for example Sociology, Psychology and Behavioral Genetics, investigate the mechanisms leading to opinion formation in groups of people[1–5]. Generally speaking, opinions are formed by a combination of self-reflection, external information sources and real-world experiences that contribute to the individual reasoning process. Furthermore, it has been argued that social interactions are fundamental in this process since they can be used to spread and exchange information (e.g. word-of-mouth and communication) as well as to shape opinion formation (e.g. social pressure and social support) through peer's social influence. In this context, both the authority and number of contacts affect the opinion of individuals[6–8]. Apart from the natural process of influencing the opinions as well as behaviours of family, friends and co-workers[9,10], social influence has been studied in marketing, politics and health interventions[11–13].

The increasing popularity of social media augments the share on daily social interactions and thus their importance in the dynamics of information exchange, potentially enhancing the power of online social influence. In fact, the low cost of interacting with multiple people, simultaneously or asynchronously, without geographical constrains, facilitates the exchange of information at a faster pace and with a diverse set of contacts. Despite such unprecedented possibilities for interactions and consumption of information, we are still bounded by cognitive and temporal constraints[14–25]. As a consequence, ideas, memes, individuals, companies and institutions compete for our limited attention which, in the current landscape, has acquired a real economic value[26,27]. In this context, social platforms use algorithms tailored to improve the online experience by ordering and filtering information judged relevant to a particular individual (or social media user). Seemingly innocuous, algorithms are designed to keep users engaged and select which information is presented to them. They act as gatekeepers and intermediaries of information, functions that traditionally have been covered by news papers and thus by hand curated editorial choices[28]. The shift towards algorithmic curation seems a natural consequence of the digital revolution. However, its short and long term societal effects are far from clear and matter of a heated debate. Before summarising some of the central points of such discussion, it is important to provide a better idea about the mechanisms driving

Centre for Business Network Analysis, Business School, University of Greenwich, SE10 9LS, London, United Kingdom. Correspondence and requests for materials should be addressed to N.P. (email: n.perra@gre.ac.uk) or L.E.C.R. (email: luis.rocha@gre.ac.uk)

algorithmic personalisation in the context of online social platforms. Indeed, algorithms are ruling a wealth of applications that go beyond social media. Examples are search engines, online shops and services.

Every time we access Facebook, or similar platforms, we experience algorithmic curation. In fact our *News Feed* is the complex result of what people and companies connected to the Internet do. In 2013 Facebook reported that every time we visit our News Feed we could be exposed on average to 1,500 stories from our friends, people and pages we follow and other content that is sponsored[29,30]. Out of this pool, they prioritise around 300 stories. While the details of the algorithm behind the personalisation are a corporate secret, we know that they are based on a combination of explicit information we provide such as demographics, likes, friendship relations, interests and implicit information derived from previous interactions with the platform, IP addresses, locations, cookies, posts age, page relationships, and content quality etc[28,31]. All this information is then fused with three main principles of content curation: popularity, semantic and collaborative filtering[32]. The first refers to the practice of promoting content that is popular across the platform. The second recommends posts similar to those previously consumed. Similarity is measured according to a range of features such as topics, words, authors, among others. The third refers to the tendency to suggest what people similar or close to us consume.

The critics of algorithmic personalisation consider such practices possibly very dangerous. Recommendations might increase the divide, exacerbating the phenomena of filter bubbles and echo chambers[33–35]. They might remove opposing view points and homogenise the type of information we are exposed to. The phenomenon is often described as the *threat of invisibility* because, in a society where more and more people use social platforms to gather information and form opinions, also what is not shown to us might be very important[36–38]. The homogenisation of information based on past preferences and actions might also reduce the possibility of serendipitous discoveries of new interests, sources, people, etc[31]. Algorithmic gatekeeping might suffer from technical biases induced by the limitations of the databases they use and the codes they are based on[39,40]. Algorithms are data hungry. This raises ethical concerns about the access and sharing of potentially sensitive information[41,42]. Finally, social media platforms are large for-profit corporations driven by shareholders' and stakeholders' interests. There is clearly, at least in principle, a friction between societal and commercial objectives that might result in the so called *corrupt personalisation*[41,43]. The proponents of algorithmic personalisation paint a far less scarier picture. First, without algorithmic filtering social platforms would be much more demanding and probably less appealing to users[41,44]. Indeed, as mentioned above, the data produced nowadays is simply too much for any type of manual filtering. Personalisation might be positive with respect to alternatives simply based on popularity metrics which could induce even stronger homogenisation of information[45]. Some recent results suggest that the possible biases created by different personalisation algorithms are very similar to those resulting from hand curation[32]. Others show how the homogeneity and popularity bias are not the same across platforms and services[46]. Finally, algorithmic gatekeeping can be developed having in mind the perils of filter bubbles and the importance of serendipity and heterogeneity in information[33]. To this point however, it is important to stress how the effects of algorithmic gatekeeping are largely not well understood. Thus the design of *optimal* solutions are far from trivial and possibly counter intuitive[47]. Indeed, the complexity of the processes involved makes predicting the effects of new features very hard and the very idea of optimal is relative. In fact, at its roots, it depends also on the specific definition of democracy considered[33].
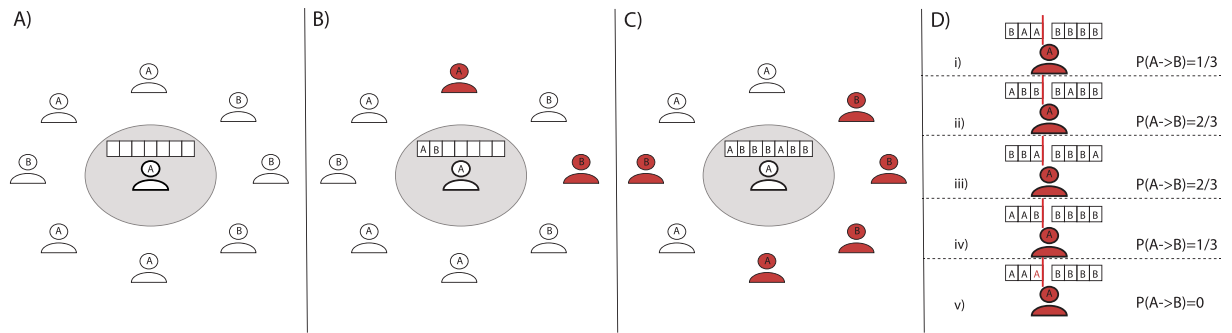
In this background, we study, from a theoretical point of view, the effects of algorithmic personalisation on a population of networked individuals. We model online interactions where users both share their opinions in their timeline and update their own following the information exposed to them. In these settings, we consider the impact of different personalisation algorithms that prioritise the information shown to each user according to time ordering or their opinions expressed in previous interactions with the system. We then study a scenario in which a particular opinion is centrally posted in all users' timelines. We refer to this as *nudge* and we assume that the social platform introduces, promotes, and pushes a particular opinion to all users thus nudging the population towards it. In all cases, we study the effect of each filtering method as a function of different (i) initial conditions, and (ii) network topologies where the interactions take place.

## Materials and Methods

We devise a mechanistic model to analyze the impact of social interactions and algorithmic exposure on the individual and group opinion dynamics. In particular, we consider a simple bipartisan model of opinion dynamics in which users can select between two opinions: A or B. Users might decide to change their status according to the fraction of posts shown in their timelines promoting each opinion. Our model assumes that only part of the information shared by friends is shown to the users. This is the fundamental idea beyond algorithmic curation. We then study how different methods used to create the personalisation affects the evolution of opinions in the system. We first look at the effect of sorting information randomly, according to time or previous ideological leanings expressed by each user in previous posts (semantic filtering). Then we consider a scenario in which one specific opinion is pushed in the timeline of each user.

The opinion formation model contains three parts (Fig. 1). The first is the underlying social network structure connecting users through friendship ties. The second is the activation mechanism that defines the timings in which information is exchanged between users through the social network. The third part is the algorithmic filtering mechanism that selects which information is presented in the user's timeline (i.e. the posts actually seen by the user).

The model has some similarities and is inspired by previous research. Wang *et al.*[16] proposed a model to mimic the spreading of memes on Twitter. In their approach, users' timelines are populated by posts from other users. With the aim to study the mechanisms behind the popularity of memes, they considered time ordered timelines in which posts survive only for a limited period of time. However, they did not considered opinion dynamics, but the spreading of memes. Also, they did not investigate the effects of different sorting algorithms. More recently, Rossi *et al.*[48] shed some light, as we aim here, on the feedback between opinions and algorithmic personalisation.

**Figure 1.** Schematic representation of the opinion dynamics model. In the plot we focus on the user $i$ in the center of panels A–C. It is connected with eight other users that are depicted around. For simplicity, we assume that these users are connected only to the central user $i$ thus their opinions or timelines will not change or be updated until user $i$ will be active. Active users are represented in red and the list (i.e. $L_i$) of posts that user $i$ could possibly see is described by the boxes. Panel A shows the initial state. Opinions are drawn randomly and the $L_i$ is initially empty. At the next time step panel B two friends of $i$ are active and post their opinions thus $L_i$ starts to be filled. At the next time step (panel C) other friends become active and so on. In panel D we show what would happen in case user $i$ becomes active. We assume $Q = 3$ (red vertical line) and show the effect of different algorithmic personalisation schemes: (i) REF, the list is random shuffled; (ii) OLD, does not change the order, the posts are shown in the order they were posted; (iii) REC, recent posts are prioritized; (iv) PR, shows the effects of a personalisation that prioritise posts according the preference/opinion of $i$; (v) NU, shows the nudge case in which one post (in red), $z = 1/3$, shown to the user does not come from its friends but it is imposed by the social platform. For each scenario we show the probability that the user $i$ would have to change opinion. In all scenarios only the first $Q$ posts after the personalisation affect the opinion of $i$. The rest are deleted.

To this end, they proposed a model in which a single individual interacts with a news aggregator which recommends news items to the user. In their model, they considered rather realistic scenarios in which the user might be subject to confirmation biases and the aggregator tries to optimise the number of clicks. However, the model neglects the network effects induced by the interaction of multiple users connected via social ties. Geschke and colleagues[49] studied the emergence of filter bubbles considering the interplay between individual, social, and algorithmic filtering mechanisms. In doing so, they developed an agent based model where users have a limited attention and are exposed to bits of information from a range of sources such as mass media, friends, personal discoveries, and algorithmic recommendations. However, they did not investigate in details the effects of different networks' features on opinion dynamics as we do here. Finally, one of the areas that received more attention, probably due to data availability, focuses on the effects of algorithms on customer choices[50–53]. Although, some of these papers study the interplay between algorithmic personalisation and social interactions, they consider purchase or more in general items selection (i.e. books) which are driven by different processes than those we are considering here. For example, purchasing power, popularity of items, reviews of products, among others, are critical ingredients in those cases.

**Network structure.** The social network structure is fixed and defined as a set of $N$ users (i.e. network nodes) connected by undirected links $(i, j)$. The distribution of contacts, clustering, and small world phenomena are well known features of networks that might have a strong effects on opinion dynamics or more in general complex contagion processes[8,54–58]. To isolate the role of topology and its interplay with the various algorithmic personalisation strategies, we consider three types of networks models. The first is a random network following the configuration model[59] with power-law (heterogeneous) degree distribution ($P(k) \propto k^{-2.5}$) with $k_{\min} = 2$. These parameters produce small-world networks with an average degree $\langle k \rangle = 5.5$ and an average clustering coefficient $\langle C \rangle = 10^{-3}$. The second is the Watts-Strogatz model[60] where each node is initially connected to six (i.e. $\langle k \rangle = 6$) friends in a ring configuration; with probability $g$ a link is rewired to a uniformly chosen random node. In the following, we consider three configurations of the model by setting $g = 0$, $g = 10^{-2}$, and $g = 1$. The three values of $g$ all produce networks with homogeneous degree distributions but allow us to explore very different combination of clustering and average path length. The first network is characterised by high values of clustering $\langle C \rangle = \frac{3(\langle k \rangle - 2)}{4(\langle k \rangle - 1)} = \langle C(0) \rangle = 0.6$[61] and a very large average path length which is known to scale as $N/(2\langle k \rangle)$. The second instead, by a large value of the clustering $\langle C(g) \rangle = \langle C(0) \rangle (1 - g)^3 = 0.58$[61]. However, in contrast with the first case, this network is characterised by the small-world phenomenon. The third, is instead a random network with negligible clustering $\left( \langle C \rangle = \frac{\langle k \rangle}{N - 1} \right)$ and characterised by the small-world phenomenon. The final type of network we consider is a regular $2D$ lattice with periodic boundary conditions. Each node has degree 4, the clustering coefficient is zero, and the average path length is very large due to spatial clustering between nodes. This network model is used as a theoretical baseline because of the absence of both high clustering and small-world features. In our simulations, we use $N = 10^5$ and run the simulations taking averages over $\tau = 10$ independent realisations of the networks with the same fixed parameters. We use $\tau = 1$ for the lattice because it is a deterministic model but also in this case we make stochastic simulations of the opinion dynamics. In summary (see Table 1), the first type of networks, CM for short, is characterised by a heterogeneous degree distribution,

| Network Type | Degree Distribution | High Clustering | Small World |
|---|---|---|---|
| CM | Heterogeneous | x | ✓ |
| WS, $g=0$ | Homogeneous | ✓ | x |
| WS, $g=10^{-2}$ | Homogeneous | ✓ | ✓ |
| WS, $g=1$ | Homogeneous | x | ✓ |
| LA | Homogeneous | x | x |

**Table 1.** Main networks features summary. In the table CM stands for Configuration model, WS for Watts and Strogatz model, and LA for regular lattice.

negligible clustering and short average path length. The second type of networks, WS for short, features homogeneous degree distributions as well as (i) high clustering and large average path length (for $g=0$) (ii) high clustering and short average path length (for $g=10^{-2}$) (iii) negligible clustering and short average path length (for $g=1$). Finally, the regular lattice, LA for short, is characterised by a homogeneous degree distribution, zero clustering, and large average path length. Altogether, these networks allow to isolate and contrast the interplay and interactions of the three topological features with the sorting algorithms.

**Nodes activation.**    With probability $p_i$, each user $i$ becomes active during one time-step $t$ (time is discrete). When the user becomes active, it first broadcasts its current opinion to all its friends, for example by writing a post, then updates its own looking at its personal timeline $R_i(t)$. Updating of all active users is done synchronously. In particular, a user whose timeline shows 30% of posts promoting opinion $A$ and thus 70% promoting opinion $B$, will adopt opinion $B$ with probability 0.7 and opinion $A$ with probability 0.3. In other words, the updates are done proportionally to the share of posts promoting each opinion in the timelines. Although the most common opinion shown to each user will be most likely adopted, the process is not deterministic. In fact, there is a level of randomness which is an important component of opinion formation and dynamics[7]. In the main text, we consider a scenario in which the activation probabilities are extracted from a power-law distribution ($F(p) \sim p^{-1.5}$ with $p \in [0.01, 1]$). Indeed, analyses of real online social networks show that the propensity of each user to start a social interaction follows an heterogeneous distribution[62–64]. In the Supplemental Information (SI), we consider the case of a constant activation probability ($p = 0.1$, which is the average of the heterogeneous case) across the entire population. The results are qualitatively similar to the homogeneous case.
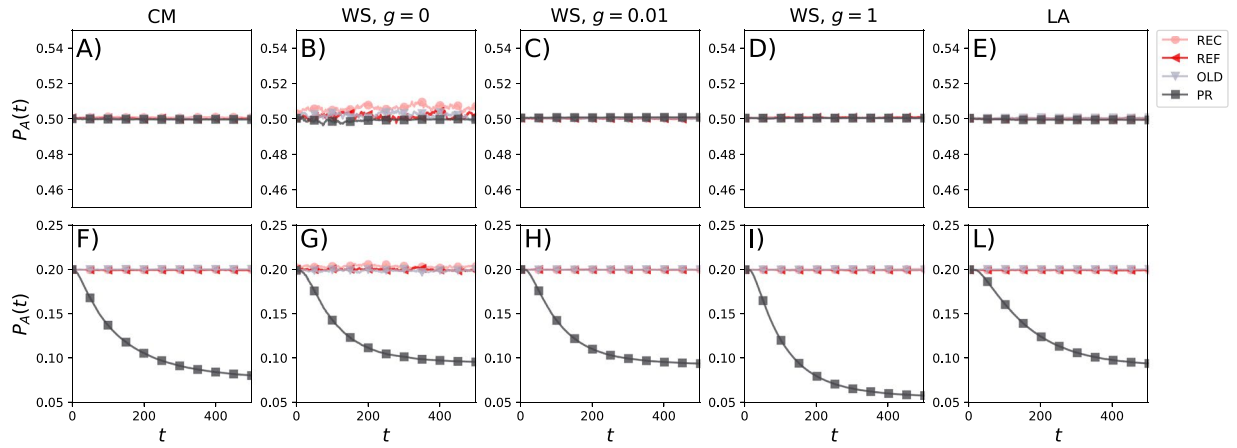
**Algorithmic personalisation and filtering.**    At time $t$, each user $i$ has a hidden list $L_i(t)$ containing all the opinions received from its friends since the last time $t_{last}$ that user $i$ was active, i.e. $L_i(t)$ contains all opinions posted by $i$'s friends during $[t_{last}, t)$, plus the opinions in user $i$ timeline at time $t_{last}$ ($R_i(t_{last})$).

We consider different methods, see Fig. 1D, to select items from the list $L_i(t)$ and thus create the timeline $R_i(t)$: (i) select $Q$ random items from $L_i(t)$ (reference method, named REF), (ii) show older posts first (old method, OLD), (iii) sort the list $L_i(t)$ from the most recent to the oldest received and show most recent posts first (recent method, REC), (iv) sort $L_i(t)$ according to $i$'s current opinion (e.g. if $i$'s opinion is $A$, sort the list as $AAA \ldots BBB$) (preference method, PR), (v) given any of the sorting method just described, a random fraction $z$ of opinions in this list is replaced by a pre-defined opinion centrally chosen (nudge method, NU). In all the scenarios the timeline $R_i(t)$ is a subset of this virtual list $L_i(t)$ and has maximum length $Q$. The remaining of the list is discarded once a timeline is formed at a given time $t$. In the main text we set $Q = 20$, but as we show in the SI its value has little impact on the results.
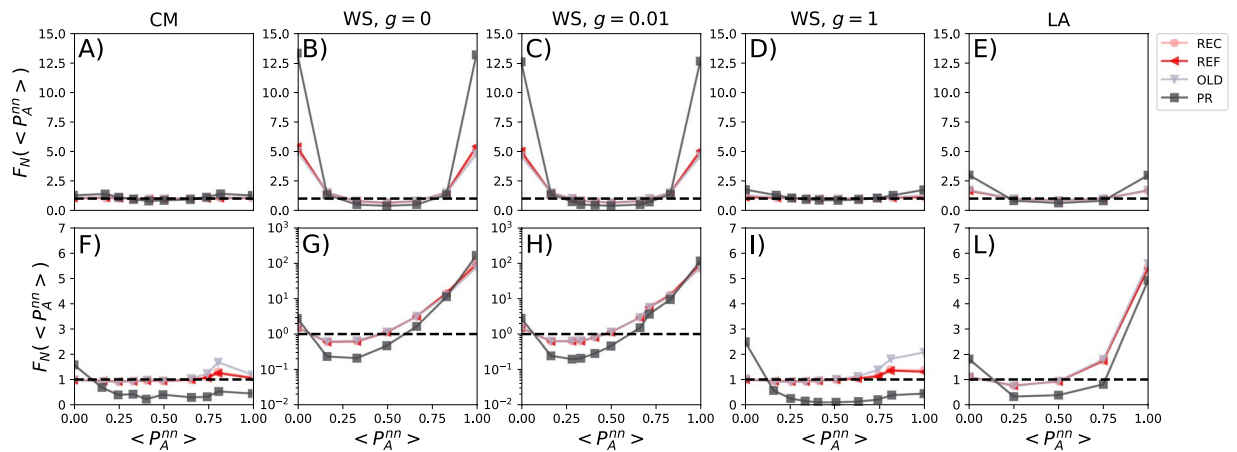
## Results
In this section, we study the effects of sorting algorithms on the group and on the individual opinion dynamics by first looking at methods REF, OLD, REC and PR. Then we investigate the effect of nudging, or centrally controlled bias.

**Opinion dynamics subject to algorithmic personalisation.**    If we start the networks with two opinions $A$ and $B$, in equal proportions ($P_A = 0.5$ and $P_B = 0.5$) and uniformly distributed among the users (i.e. nodes), the prevalence of both opinions remains stable around the starting values for all four sorting methods for various network configurations (first row in Fig. 2, see also SI for longer times). Thus, in these scenarios all the filtering algorithms are not able to break the status quo. Both ideas are able to coexist and the social platform does not drive the system towards a different equilibrium with respect to the initial status. If the opinions are unequal from the start ($P_A = 0.2$ and $P_B = 0.8$), the prevalence also remains stable for the methods REC, REF and OLD during the entire observation period (second row in Fig. 2, see also SI for longer times). Indeed, as the order in which posts appears in the list $L_i$ of each user is random and regulated by the activation process of its neighbours, introducing a temporal bias (in the case of REC and OLD) does not affect the share (20−80) of opinions in each timeline thus preserving the status quo. On the other hand, the preference method (PR) leads to a decrease in the prevalence of the minority opinion ($P_A$) that mainly occurs within the initial 200 time steps across all networks. In an unbalanced situation, the ordering of opinions according to users' preferences, for users holding opinion $B$, might lead to the disappearance of the opposite idea from their timeline. Effectively, such semantic filtering brings to zero the probability of such users to change idea (the *visibility* of the others is zero) thus modifying the initial share and increasing the already dominant position of opinion $B$. It is important to notice that although the less popular opinion decreases its share, it is still able to survive. As shown in the SI, the survival happens even in case of much skewed initial configurations such as 1−99. Results are similar for all network configurations, but for the

**Figure 2.** Evolution of group opinion in the various network models. The prevalence $P_A(t)$ of opinion $A$ over time for starting $P_A(0) = 0.5$ and $P_B(0) = 0.5$ (first row) and $P_A(0) = 0.2$ and $P_B(0) = 0.8$ (second row). Each column describes the results for one of the networks. In each plot, we show results for the four sorting algorithms. Each plot is the average of $10^2$ independent simulations and to improve the visualisation we are showing the data points every 50 time steps. In all scenarios we set $Q = 20$.



**Figure 3.** Opinions of friends. We show the distribution of the fraction of friends (nearest neighbours, nn) of $i$ with the same opinion $A$ at $t = 500$. We normalise the $y$-axis by dividing for the same quantity computed at $t = 0$. In the first row we show the results for starting conditions $P_A(0) = 0.5$ and $P_B(0) = 0.5$. In the second row, for starting conditions $P_A(0) = 0.2$ and $P_B(0) = 0.8$. Each column describes the results for a particular network. Each plot is the average of $10^2$ independent simulations. In each plot, we consider the four ranking algorithms and set $Q = 20$.

WS network with $g = 1$ the decrease in $P_A$ occurs earlier in time and $P_A$ stabilises at lower values ($P_A(500) = 0.06$, see Fig. 2I). The CM network reaches larger values ($P_A(500) = 0.08$) which are however slightly smaller than the other three networks ($P_A(500) = 0.09$). Altogether, these results point to the fact that the PR filtering mechanism, in case of unbalanced initial conditions, causes more biases in networks characterised by null clustering, short path length, and homogeneous degree distributions. The heterogeneity in contact patterns hampers the reduction of the less popular opinion due to the presence of hubs which limit the effectiveness of the algorithm to hide the subordinate opinion. High values of clustering both in combination with low and high average path length behave as in the case of regular lattice which features null clustering, high spatial correlations and high average path length. This result suggests that correlations in the connectivity patterns hamper the effectiveness of semantic filtering with respect to uncorrelated cases, which is qualitatively in line with previous research on recommender systems mentioned above[50].

The interplay between the features of the sorting mechanisms and the network structure might induce the formation of opinion clusters (echo chambers) and polarisation in the population. Figure 3 shows the distribution of the fraction of users $i$ whose friends (1st nearest neighbours, nn for short) share the opinion $A$ at $t = 500$, $\langle P_A^{nn} \rangle$, irrespective of $i$'s opinion (as shown in the SI, plots done conditioning on the opinion of $i$ are very similar but by construction they do not allow to observe in a single plot the polarisation around both opinions). In these plots, the region of the $x$-axis close to one describes neighbourhoods formed by a majority (minority) of users sharing opinion $A$ ($B$). Conversely, the region of the $x$-axis close to zero describes neighbourhoods in which none, or very
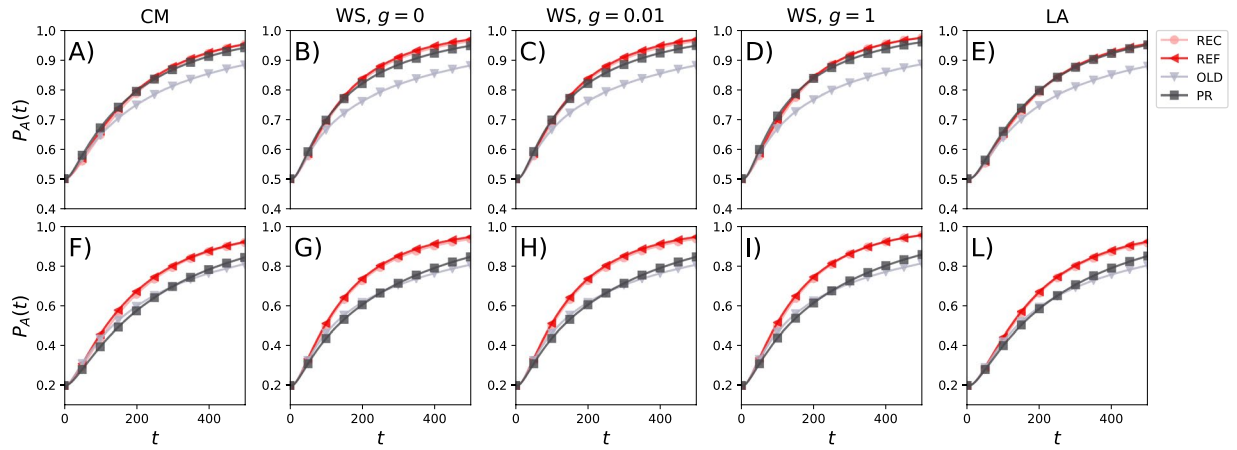
5

few users, adopt opinion $A$ (or neighbourhoods featuring a majority of nodes holding opinion $B$). In order to remove spurious effects, on the $y$-axis we show $F_N(\langle P_A^{nn} \rangle)$ that is obtained dividing the average distribution of opinions in each neighbourhood at time $t = 500$ with its correspondent value at $t = 0$. In doing so, to avoid issues with the binning, we imposed 10 equispaced bins ($\Delta x = 0.1$) in all cases. It is important to stress how we consider evidence of polarisation, or formation of echo-chambers, a higher concentration of neighbourhoods sharing the same opinion respect to an initial random distribution of the two ideas in the system. Other definitions, potentially based on the share of ideas in the various communities, are computationally more expensive, dependent on the community detection algorithm, and will be explored in future work.

For equal starting conditions ($P_A = 0.5$ and $P_B = 0.5$) (first row in Fig. 3), we notice that CM and WS ($g = 1$) behave similarly. The distribution of opinions in the neighbourhood of all nodes is similar to the initial (random) configuration. The PR mechanism, in the case of WS with $g = 1$, produces two little spikes at zero and one which indicate the emergence of low levels of polarisation in the network. Thus, the semantic filtering, in networks characterised by homogeneous degree distributions, null clustering, and short average path length, is able to produce few patches like-minded people which however do not form in case of heterogeneous connectivity patterns (CM network). The polarisation effects of the filtering algorithms, in particular PR, are much stronger in the other networks. Indeed, in highly clustered networks, independently of the average path length (WS graphs with $g = 0$ and $g = 10^{-2}$) the sorting mechanisms induce the emergence of strong opinion polarisation. As shown in the SI, if the size of the timelines is comparable with the degree (i.e. $Q = 5$), REC and REF filtering methods produce stronger polarisation effects in combination with high clustering in the contacts patterns. This is due to the fact that the PR mechanism in these cases might over represent the local minority opinion, thus hampering the convergence to a local majority. Echo chambers appear, although to smaller degree, also in case of strong spatial correlations and null clustering (LA network). These observations provide further evidence that topological as well as spatial correlations are important features for algorithmic filtering and that heterogeneity in contact patterns has a slightly stronger potential to hamper the formation of each chambers in comparison to homogeneous patterns. In fact, in combination with the filtering mechanisms, these features induce a re-organisation of the distribution of opinions creating *compact domains* made of like-minded nodes (echo chambers). This observation aligns with previous results in the literature of complex contagion processes on complex networks[65]. The polarisation emerges even if the total fraction of the population with opinion $A$ or $B$ is the same and stable as function of time (Fig. 2). Thus, although the interplay between the topological features and algorithm filtering does not change the status quo, it might introduce high levels of polarisation and echo chambering.
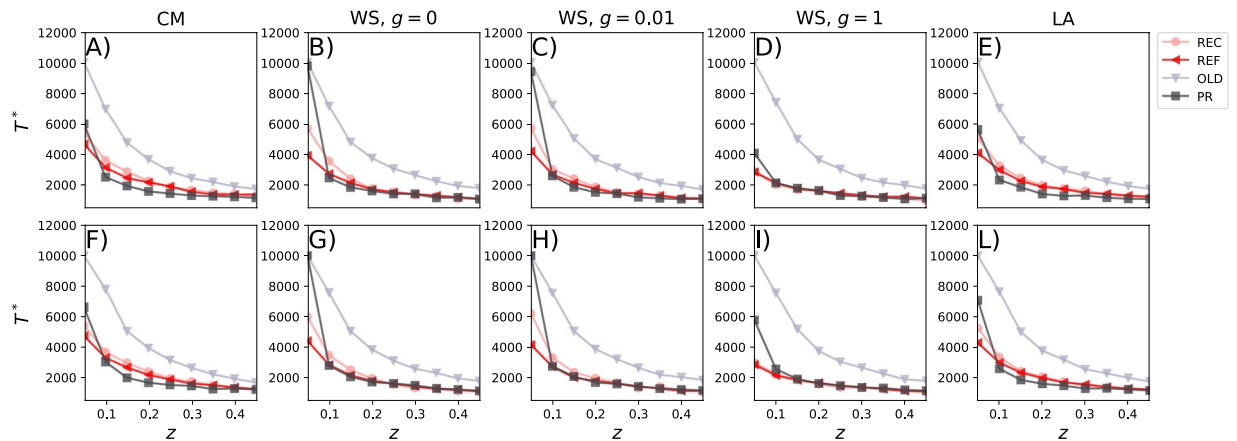
The case of unbalanced starting conditions (i.e. $P_A = 0.2$ and $P_B = 0.8$) confirms this picture (second row in Fig. 3). In fact, the behaviour of CM and WS ($g = 1$) networks is still very similar. The absence of topological clustering and the presence of short cuts between parts of the network (small-world phenomenon) does not allow the emergence of echo chambers around the dominant opinion. This is somehow surprising considering that opinion $B$ reaches around 90% of the population. There is an exception however, the semantic filtering (PR method) fights against this tendency. As mentioned above, this personalisation algorithm reduces the *visibility* of the few nodes holding opinion $A$ that might be present in neighbourhoods of $B$ nodes. The effective reduction of the share of nodes with opinion $A$ in such configuration brings the system to a different, and lower, equilibrium as eventually such nodes change opinion to $B$. This is confirmed by the fact that across the entire $x$-axis (with exception of $x = 0$) the average number of neighbours holding some fraction of opinion $A$ decreases for all network models. The presence of heterogeneous connectivity patterns in CM networks reduces the effectiveness of such effect. Figure 3F–L shows that the three other sorting algorithms, OLD above all, induce a relatively higher number of echo chambers around the subordinate opinion in WS ($g = 1$) networks. As such effects are reduced in CM networks, we can speculate that this is due to the interaction between the algorithms and homogeneous contact patterns.

Also for unbalanced initial conditions, topological and spatial correlations induce a much stronger polarisation (notice the log scale on the $y$-axis). In particular, in networks with high values of clustering, the number of $A$ echo chambers is two orders of magnitude larger than in the initial configuration. Spatial correlations introduce less stronger effects which are however of the order of a factor 6. Overall, these results suggest that correlations interact with the sorting algorithms creating compact domains that, in this case, boost the survival of the less common opinion with respect to the other two networks.

### Nudging scenario.
We now consider the same filter algorithms studied in the previous section with the added feature that a random fraction $z$ of opinions, in the list $L_i(t)$ shown to each user $i$, is replaced by a pre-defined opinion centrally chosen and fixed in time (see Fig. 1D). In other words, we are investigating a scenario in which the social platform, where the interactions take place, nudges the users towards a particular opinion. In doing so, a fraction of the opinions shown to a user is substituted with ideas that might or might not be representative of what its nearest neighbours are posting. In the following, we set $z = 0.1$, $Q = 20$ and chose $A$ as the centrally set opinion. Independently of the starting conditions (balanced or unbalanced prevalence of $A$ and $B$), the group opinion moves towards $A$ for all network configurations and sorting algorithms (Fig. 4). In the case of unbalanced starting conditions (i.e. $P_A = 0.2$ and $P_B = 0.8$), opinion $A$ may become dominant as early as $t \sim 100$ time steps and typically no later than $t \sim 150$. At $t = 500$ almost the entire population turned to opinion $A$. Across the board, the REF method is the most efficient at $t = 500$ to nudge the population opinion towards $A$. This is more evident with unbalanced initial conditions. Arguably, this is due to the fact that this filtering method (REF) is the least biased towards the current and past opinions hold by each user and its neighbours. The OLD method instead, in the balanced case, is significantly slower than the rest. By nature, this sorting algorithm slows down the overall drift of the system by showing to each user old neighbours' beliefs. In contrast to the results in the previous section, method PR performs in between the other methods.
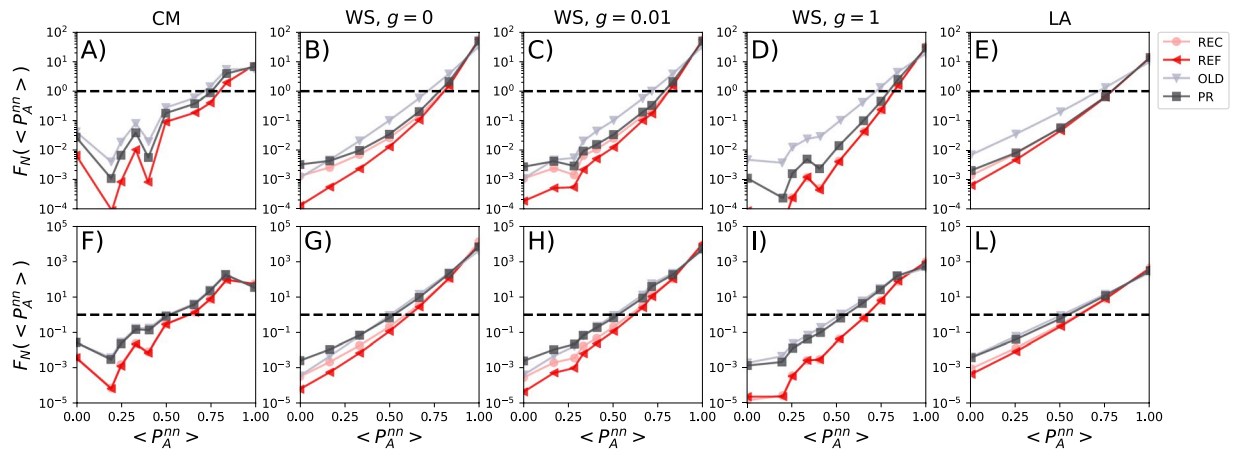
**Figure 4.** Evolution of population opinion with central influencer (nudging). The prevalence $P_A(t)$ of opinion $A$ over time for starting conditions $P_A(0) = 0.5$ and $P_B(0) = 0.5$ in the first row; $P_A(0) = 0.2$ and $P_B(0) = 0.8$ in the second row. Each column describes the results for a particular network. In each plot, we set $Q = 20$ and show results for the four sorting algorithms. Each plot is the average of $10^2$ independent simulations and to improve the visualisation data points are shown every 50 time steps.



**Figure 5.** Time for convergence. We show the time $T^*$ that the nudged opinion needs to converge (i.e. be shared by all nodes) as function of $z$. We set $10^4$ as maximal time of convergence of the simulations. In he first row we show the results for $P_A(0) = 0.5$ and $P_B(0) = 0.5$; in the second row we show the results for $P_A(0) = 0.2$ and $P_B(0) = 0.8$. In each plot, we set $Q = 20$ and each column describes the results for a particular network model. Each plot is the average of 10 independent simulations.

In order to gather a better understanding of the impact of different sorting algorithms and networks topologies in the nudging process, Fig. 5 shows the time $T^*$ that the opinion $A$ needs to take over the entire population. To reduce the computational cost of the simulations, we set $10^4$ as time limit, thus the actual $T^*$ for the curves that reach such threshold is higher. The first observation is that the OLD sorting mechanism is, across the board, the slowest to reach convergence. The OLD method introduces a delay by promoting the oldest posts of each neighbours, which, over time, might sponsor a different opinion with respect to the opinion the focus nodes currently hold. The second observation is that for $z = 0.05$, thus when only one opinion is centrally replaced in the list of each user (since $Q = 20$), the PR sorting algorithm is generally faster than the OLD method but slower than the others. This is particularly visible in the case of networks with topological (WS networks with $g = 0$ and $g = 10^{-2}$) or spatial correlations (LA network). As mentioned above, correlations interact with the semantic filtering (PR) creating relatively strong echo chambers which shield its members protecting them from the nudged opinion. The third observation is that for $z \geq 0.15$, the time for convergence for the REF, REC, and PR methods becomes a very weak function of $z$ and of the initial share of the two opinions. Fourth, as also noted in Fig. 4, the REF filtering mechanism, which does not introduce any type of bias in the ordering of neighbours' posts, is the fastest to reach convergence in all topologies. In the case of networks with negligible clustering and small average shortest path (WS networks with $g = 1$) the convergence, for $z = 0.05$, is reached faster than in the other topologies. The comparison with the CM networks points one more time to the fact that heterogeneous connectivity patterns hamper both local (as seen before) and global convergence (as seen now) to a single opinion.

**Figure 6.** Opinions of friends. The distribution of the fraction of friends of $i$ with the opinion $A$ at $t = 500$. The $y$-axis is normalised by dividing by the same quantity computed at $t = 0$. In the first row, we show the results for starting conditions $P_A(0) = 0.5$ and $P_B(0) = 0.5$. In the second row, we show the results for starting conditions $P_A(0) = 0.2$ and $P_B(0) = 0.8$. Columns show the results for different networks. Each plot is the average of $10^2$ independent simulations and we set $Q = 20$.

Figure 6 shows the normalised distribution of the fraction of 1$^{st}$ nearest neighbours sharing opinion $A$ at $t = 500$ for both starting conditions irrespective of the opinion of the focal node in each neighbourhood. We set $z = 0.01$ thus 1% of opinions in this list of each node is replaced with opinion $A$. Since opinion $B$ is almost extinguished in both cases, the distributions are all asymmetric and skewed towards values of $\langle P_A^{nn} \rangle$ close to one. Across all networks, the system is converging towards opinion $A$ and a large number of neighbourhoods share all the same nudged opinion. Indeed, the reduction of all the other values of $\langle P_A^{nn} \rangle$ is dramatic (note the logarithm scale on the $y$). The comparison between the CM and WS ($g = 1$) networks shows how the heterogeneity in connectivity patterns alone hampers the formation of echo chambers around the dominant opinion.

## Discussion

Technology has dramatically reduced the cost of communications and vastly increased their reach. As consequence, the amount of information produced at any given moment far surpasses our ability to consume it. Individuals, companies, and institutions alike compete for our attention that has acquired a real economic value. For example, online social platforms such as Twitter and Facebook profit from our engagement with their systems. Thus, optimising our experience is a key part of their business model. To this end, they adopt algorithmic personalisation which broadly describes a range of automated methods that filter the information to which we are exposed to. The filtering is not random, but targeted to stimulate our interests and online activity. Generally speaking the aims of algorithmic filtering are benign. However, in our hyper-connected society they might have unintended and unanticipated consequences. One area of concern is the possibility that such methods, which often are black boxes, affect opinion formation and dynamics. As many of our social interactions are mediated by such platforms and since much of the information gathering and exposure happen online, the perception of our friends as well as of specific issues might be distorted by such algorithms. These concerns are now mainstream and the words *echo chamber*, *confirmation bias*, and *filter bubbles* have become a staple in the news cycle. Despite the heated debate the effect of algorithmic filtering on opinion dynamics is still largely not understood.

In this paper, we studied the effect of different filtering algorithms on the opinion dynamics of a population of interconnected individuals. We modelled their interactions having in mind prototypical online platforms where users post their opinions which appear in friends' timelines. Due to the limited attention of each user, timelines are curated by means of algorithmic personalisation. We considered different scenarios where posts are selected (i) randomly, (ii) as function of the time when they were posted, and (iii) as function of previous posts of users (semantic filtering). Even more, we studied how the features of the connectivity patterns mediating the interactions between users interplay with filtering algorithms and possibly affect opinions. In particular, we considered three different types of networks which allow us to isolate and contrast three main properties of real networks (i) heterogeneity in degree distributions, (ii) high values of clustering, and (iii) small-world phenomenon. Finally, we investigated another scenario in which the social platform where interactions take place tries to nudge the population towards a particular opinion by manipulating a fraction of the posts in users' timelines.

For simplicity we investigated a bipartisan system in which users can adopt one of two opinions, i.e. $A$ or $B$. We studied two very different initial configurations where the share of opinion $A$ was equal (50–50) or much lower (20–80) than $B$. In these settings, we found that, independently of the social network structure, the sorting mechanisms are not able to change the status quo when the prevalence of opinions is equally distributed in the population. Despite the overall share of opinions does not change, we observed that some features of the networks interact with the filtering algorithms inducing polarisation effects. In particular, we found that topological (high values of clustering) and to a lesser extent spatial (proximity of nodes in a regular lattice) correlations create compact domains (echo chambers) formed only by like-minded individuals. Despite the algorithmic personalisation, absence of correlations, especially in conjunction with heterogeneity in the contacts patterns, hamper

the formation of such echo chambers. In case of unbalanced initial conditions, we found that semantic filtering induces a further increase of the predominant opinion because of positive reinforcement. Also in this scenario correlations play a significant role hampering the reinforcement by creating compact domains around the subordinate opinion that effectively is able to protect such minorities. As mentioned above the absence of topological, or the presence of spatial, correlations hinders the formation of echo chambers, even those centred around the dominant opinion.

If a centralised bias mechanism (i.e. nudging) is introduced, population opinion can move towards the nudged opinion relatively fast and eventually completely switch the status quo. In this case, filtering algorithms that randomly select friends' opinions or that give preference to the most recent ones speedup convergence to the centrally desired opinion. Conversely, methods biased towards older posts are the slowest to reach convergence. Semantic filtering, in case only a single post in users timelines is manipulated, slows down convergence especially in networks characterised by topological or spatial correlations. In these cases, the emergence of echo chambers opposes the change towards the nudged opinion.

Overall, our results support the view that opinion dynamics can be manipulated by algorithmic personalisation methods. Furthermore, they highlight the interplay between filtering mechanisms and the features of the networks where interactions take place. In particular, correlations in contact patterns are key ingredients that might lead to the emergence of polarisation and echo chambers even in case of equal share between opinions. Our results show that filtering algorithms might exacerbate polarisation but the organisation of social ties is a key factor. Conversely, the absence of correlations, especially in conjunction with heterogeneous distributions of contacts hampers the creation of compact domains. These results are in line with the literature of complex contagion processes on complex networks[8,54,55]. The study of the nudging scenario shows how the dominant opinion in the population can completely switch within short time by moderately pushing a desired opinion. This idea is supported by the Nudge Theory in behavioural economics originally developed by Richard Thaler[66] and by recent experiments on the emergence of conventions[67]. The vulnerability of social media to manipulations can be disastrous because nudging for example can be easily introduced by the company controlling the media, by external players through paid advertisement or by the use of bots[68]. There are claims that such strategies have been used to bias political public opinion in recent years[69,70]. Not least, since positive reinforcement mechanisms may increase opinion polarisation according to our study, the spread of fake news, that tend to spread faster than true news[71], may generate local convergence of opinion that later leads to the emergence of bubbles and group isolation[72]. In light of democratic access to information, a common and intuitive idea is that social media developers should make efforts to minimise opinion biases by increasing the diversity of opinions received from friends and counter-balancing the content of advertisements. However, recent empirical evidence points out that the actual dynamics at play are much more complex[47]. Indeed, the exposure to opposite view points might further increase polarisation[47]. Thus, the research on the feedback between opinion dynamics and algorithmic personalisation is just in its infancy.

Our model aims to capture key mechanisms driving online social interactions and opinion dynamics. It has however some limitations since it disregards other potentially relevant mechanisms. In particular, our model does not consider that constantly diverging friends may stop following each other and thus their opinions are not shared. This manual curation, complementary to the automated curation included in our model, might further contribute to the formation of opinion bubbles and likely boosts biases towards the users' own opinion. Future models should include manual curation taking into account the user's own memory of his or her friends opinions. In the language of modern Network Science this will imply the development of adaptive and time-varying networks models that couple the dynamics of the networks (link formation) with the dynamics on the networks (opinion dynamics). Another limitation is that users have some level of resistance to change and thus some people need reinforced exposure to an opinion before switching. This may affect the timing, possibly delaying, and the likelihood to switch the prevalent opinion in the population. Also, friends do not necessarily have the same authority and possibly information coming from specific friends, e.g. close offline friends or relatives, may have higher importance on shaping the user's opinions. Another limitation is the type of filtering mechanisms considered. Future work should consider more complex and realistic models possibly linked to the history and features of each user as well as the popularity of specific memes or information circulating in the system. From the analysis perspective, the numerical results indicate that opinions converge fast and are stable over time. In the future, an analytic treatment of the model, e.g. using master equation stability analysis, could provide stronger evidence on the stability of these solutions. Finally, here we considered only some prototypical networks models as way to isolate three of the main features of real networks. Other properties such as the presence of communities, high-order correlations, and temporal connectivity patterns will be matter of future explorations.

## References

1. Latane, B. The psychology of social impact. *Am. Psychol.* **36**, 343–356 (1981).
2. Isenberg, D. J. Group polarization: A critical review and meta-analysis. *J. Pers. Soc. Psychol.* **50**, 1141–1151 (1986).
3. D'Onofrio, B., Eaves, L., Murrelle, L., Maes, H. & Spilka, B. Understanding biological and social influences on religious affiliation, attitudes, and behaviors: A behavior genetic perspective. *J. Pers.* **67**, 953–984 (1999).
4. Olson, J. M., Vernon, P. A., Harris, J. A. & Jang, K. L. The heritability of attitudes: A study of twins. *J. Pers. Soc. Psychol.* **80**, 845–860 (2001).
5. Watts, D. J. & Dodds, P. S. Influentials, networks, and public opinion formation. *J. Consumer Res.* **34**, 441–458 (2007).
6. Muchnik, L., Aral, S. & Taylor, S. J. Social influence bias: A randomized experiment. *Sci.* **341**, 647–651 (2013).
7. Castellano, C., Fortunato, S. & Loreto, V. Statistical physics of social dynamics. *Rev. modern physics* **81**, 591 (2009).
8. Baronchelli, A. The emergence of consensus: a primer. *Royal Soc. open science* **5**, 172189 (2018).
9. Christakis, N. A. & Fowler, J. H. The spread of obesity in a large social network over 32 years. New Engl. *J. Medicine* **357**, 370–379 (2007).
10. Paluck, E. L., Shepherd, H. & Aronow, P. M. Changing climates of conflict: A social network experiment in 56 schools. *Proc. Natl. Acad. Sci. USA* **113**, 566–571 (2016).

11. Aral, S. & Walker, D. Identifying influential and susceptible members of social networks. *Sci.* **337**, 337–341 (2012).
12. Bond, R. M. *et al.* A 61-million-person experiment in social influence and political mobilization. *Nat.* **489**, 295–298 (2012).
13. Kim, D. *et al.* Social network targeting to maximise population behaviour change: a cluster randomised controlled trial. *The Lancet* **386**, 145–153 (2015).
14. Watts, D. J. A simple model of global cascades on random networks. *Proc. Natl. Acad. Sci.* **99**, 5766–5771 (2002).
15. Falkinger, J. Attention economies. *J. Econ. Theory* **133**, 266–294 (2007).
16. Weng, L., Flammini, A., Vespignani, A. & Menczer, F. Competition among memes in a world with limited attention. *Sci. reports* **2**, 335 (2012).
17. Wu, F. & Huberman, B. A. Novelty and collective attention. *Proc. Natl. Acad. Sci.* **104**, 17599–17601 (2007).
18. Leskovec, J., Backstrom, L. & Kleinberg, J. Meme-tracking and the dynamics of the news cycle. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, 497–506 (ACM, 2009).
19. Lerman, K. & Ghosh, R. Information contagion: An empirical study of the spread of news on digg and twitter social networks. *Icwsm* **10**, 90–97 (2010).
20. Ratkiewicz, J., Fortunato, S., Flammini, A., Menczer, F. & Vespignani, A. Characterizing and modeling the dynamics of online popularity. *Phys. review letters* **105**, 158701 (2010).
21. Onnela, J.-P. & Reed-Tsochas, F. Spontaneous emergence of social influence in online systems. *Proc. Natl. Acad. Sci.* **107**, 18375–18380 (2010).
22. Gonçalves, B., Perra, N. & Vespignani, A. Modeling users' activity on twitter networks: Validation of dunbar's number. *PloS one* **6**, e22656 (2011).
23. Dunbar, R. I. The social brain hypothesis. *Evol. Anthropol. Issues, News, Rev. Issues, News, Rev.* **6**, 178–190 (1998).
24. Dunbar, R. I., Arnaboldi, V., Conti, M. & Passarella, A. The structure of online social networks mirrors those in the offline world. *Soc. Networks* **43**, 39–47 (2015).
25. Weng, L., Karsai, M., Perra, N., Menczer, F. & Flammini, A. Attention on weak ties in social and communication networks. In *Complex Spreading Phenomena in Social Systems*, 213–228 (Springer, 2018).
26. Simon, H. A. Designing organizations for an information-rich world. *Comput. communications, public interest* (1971).
27. Davenport, T. H. & Beck, J. C. *The attention economy: Understanding the new currency of business* (Harvard Business Press, 2001).
28. Bozdag, E. Bias in algorithmic filtering and personalization. *Ethics information technology* **15**, 209–227 (2013).
29. Backstrom, l. a window into news feed, https://www.facebook.com/business/news/News-Feed-FYI-A-Window-Into-News-Feed. Accessed: 2018-07-20.
30. Rader, E. & Gray, R. Understanding user beliefs about algorithmic curation in the facebook news feed. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*, 173–182 (ACM, 2015).
31. DeVito, M. A. From editors to algorithms: A values-based approach to understanding story selection in the facebook news feed. *Digit. Journalism* **5**, 753–773 (2017).
32. Möller, J., Trilling, D., Helberger, N. & van Es, B. Do not blame it on the algorithm: An empirical assessment of multiple recommender systems and their impact on content diversity. *Information, Commun. & Soc.* **21**, 959–977 (2018).
33. Bozdag, E. & van den Hoven, J. Breaking the filter bubble: democracy and design. *Ethics Inf. Technol.* **17**, 249–265 (2015).
34. Pariser, E. *The filter bubble: What the Internet is hiding from you* (Penguin UK, 2011).
35. Del Vicario, M. *et al.* Echo chambers: Emotional contagion and group polarization on facebook. *Sci. reports* **6**, (37825 (2016).
36. Bucher, T. Want to be on the top? algorithmic power and the threat of invisibility on facebook. *new media & society* **14**, 1164–1180 (2012).
37. Rader, E. Examining user surprise as a symptom of algorithmic filtering. *Int. J. Human-Computer Stud.* **98**, 72–88 (2017).
38. Ciampaglia, G. L., Nematzadeh, A., Menczer, F. & Flammini, A. How algorithmic popularity bias hinders or promotes quality. *Sci. reports* **8**, 15951 (2018).
39. Friedman, B. & Nissenbaum, H. Bias in computer systems. ACM Transactions on. *Inf. Syst. (TOIS)* **14**, 330–347 (1996).
40. Gillespie, T. The relevance of algorithms. *Media technologies: Essays on communication, materiality, society* **167** (2014).
41. Stevenson, D. & Pasek, J. Privacy concern, trust, and desire for content personalization. In *TPRC 43: The 43rd Research Conference on Communication, Information and Internet Policy Paper* (2015).
42. Mittelstadt, B. D., Allo, P., Taddeo, M., Wachter, S. & Floridi, L. The ethics of algorithms: Mapping the debate. *Big Data & Soc.* **3**, 2053951716679679 (2016).
43. Sandvig, C., Hamilton, K., Karahalios, K. & Langbort, C. Can an algorithm be unethical? Ann. *Arbor* **1001**, 1285 (2015).
44. Liu, C., Belkin, N. J. & Cole, M. J. Personalization of search results using interaction behaviors in search sessions. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, 205–214 (ACM, 2012).
45. Goldman, E. Search engine bias and the demise of search engine utopianism. In *Web Search*, 121–133 (Springer, 2008).
46. Nikolov, D., Lalmas, M., Flammini, A. & Menczer, F. Quantifying biases in online information exposure. *J. Am. Soc. for Inf. Sci. Technol.* (2018).
47. Bail, C. A. *et al.* Exposure to opposing views on social media can increase political polarization. *Proc. Natl. Acad. Sci.* **115**, 9216–9221 (2018).
48. Rossi, W. S., Polderman, J. W. & Frasca, P. The closed loop between opinion formation and personalised recommendations. *arXiv preprint arXiv:1809.04644* (2018).
49. Geschke, D., Lorenz, J. & Holtz, P. The triple-filter bubble: Using agent-based modelling to test a meta-theoretical framework for the emergence of filter bubbles and echo chambers. *Br. J. Soc. Psychol* (2018).
50. Bressan, M., Leucci, S., Panconesi, A., Raghavan, P. & Terolli, E. The limits of popularity-based recommendations, and the role of social ties. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 745–754 (ACM, 2016).
51. Fleder, D. & Hosanagar, K. Blockbuster culture's next rise or fall: The impact of recommender systems on sales diversity. *Manag. science* **55**, 697–712 (2009).
52. Dandekar, P., Goel, A. & Lee, D. T. Biased assimilation, homophily, and the dynamics of polarization. *Proc. Natl. Acad. Sci.* **110**, 5791–5796 (2013).
53. Spinelli, L. & Crovella, M. Closed-loop opinion formation. In *Proceedings of the 2017 ACM on Web Science Conference*, 73–82 (ACM, 2017).
54. Guilbeault, D., Becker, J. & Centola, D. Complex contagions: A decade in review. In *Complex Spreading Phenomena in Social Systems*, 3–25 (Springer, 2018).
55. Barrat, A., Barthelemy, M. & Vespignani, A. *Dynamical processes on complex networks* (Cambridge university press, 2008).
56. Baronchelli, A., Felici, M., Loreto, V., Caglioti, E. & Steels, L. Sharp transition towards shared vocabularies in multi-agent systems. *J. Stat. Mech. Theory Exp.* **2006**, P06014 (2006).
57. Barabási, A. -L. *et al. Network science* (Cambridge University press, 2016).
58. Centola, D. & Baronchelli, A. The spontaneous emergence of conventions: An experimental study of cultural evolution. *Proc. Natl. Acad. Sci.* **112**, 1989–1994 (2015).
59. Newman, M. E. The structure and function of complex networks. *SIAM review* **45**, 167–256 (2003).
60. Watts, D. J. & Strogatz, S. H. Collective dynamics of 'small-world' networks. *nature* **393**, 440 (1998).

61. Barrat, A. & Weigt, M. On the properties of small-world network models. *The Eur. Phys. J. B-Condensed Matter Complex Syst.* **13**, 547–560 (2000).
62. Rocha, L. E. C., Liljeros, F. & Holme, P. Simulated Epidemics in an Empirical Spatiotemporal Network of 50,185 Sexual Contacts. *PLoS Comput. Biolology* **7**(3), e1001109 (2011).
63. Perra, N., Gonçalves, B., Pastor-Satorras, R. & Vespignani, A. Activity driven modeling of time varying networks. *Sci. reports* **2**, 469 (2012).
64. Ribeiro, B., Perra, N. & Baronchelli, A. Quantifying the effect of temporal resolution on time-varying networks. *Sci. reports* **3**, 3006 (2013).
65. Dall'Asta, L., Baronchelli, A., Barrat, A. & Loreto, V. Nonequilibrium dynamics of language games on complex networks. *Phys. Rev. E* **74**, 036105 (2006).
66. Thaler, R. H. & Sunstein, C. R. *Nudge: Improving Decisions About Health, Wealth and Happiness* (Yale University Press, 2008).
67. Centola, D., Becker, J., Brackbill, D. & Baronchelli, A. Experimental evidence for tipping points in social convention. *Sci.* **360**, 1116–1119 (2018).
68. Ferrara, E., Varol, O., Davis, C., Menczer, F. & Flammini, A. The rise of social bots. *Commun. ACM* **59**, 96–104 (2016).
69. Bessi, A. & Ferrara, E. Social bots distort the 2016 us presidential election online discussion. *First Monday* **21** (2016).
70. Stukal, D., Sanovich, S., Bonneau, R. & Tucker, J. A. Detecting bots on russian political twitter. *Big Data* **5** (2017).
71. Vosoughi, S., Roy, D. & Aral, S. The spread of true and false news online. *Sci.* **359**, 1146–1151 (2018).
72. Nikolov, D., Oliveira, D. F. M., Flammini, A. & Menczer, F. Measuring online social bubbles. *PeerJ Comput. Sci.* **1**, e38 (2015).

## Acknowledgements

## Author Contributions

N.P. and L.R. conceived the model, N.P. conducted the experiments, N.P. and L.R. analysed the results, wrote the manuscript and approved the manuscript.

## Additional Information

**Supplementary information** accompanies this paper at https://doi.org/10.1038/s41598-019-43830-2.

**Competing Interests:** The authors declare no competing interests.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.