

Game Theory Based Community-Aware Opinion Dynamics

Shanfan Zhang¹, Yongyi Lin², Xiaoting Shen³, Zhan Bu^{4*}, Yuan Rao^{1*}

¹School of software Engineering, Xi'an Jiaotong University, Xi'an, China

²School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an, China

³School of Automation, Nanjing University of Information Science & Technology, Nanjing, China

⁴School of Computer Science, Nanjing Audit University, Nanjing, China

zhangsfxjtu@gmail.com, linyongyi@stu.xjtu.edu.cn, njdnsxt@126.com, zhanbu@nau.edu.cn, raoyuan@mail.xjtu.edu.cn

Abstract

Understanding opinion evolution in complex social networks is crucial for modeling social influence and predicting collective behavior. Yet, most models overlook how community structures shape opinion updates, often assuming homogeneous influence. This abstraction neglects individuals' stronger responsiveness to intra-community peers—an empirically observed driver of localized consensus and inter-group polarization. We propose *GCAOFP*, a co-evolutionary framework that jointly models opinion dynamics and community formation as an integrated process. In *GCAOFP*, agents strategically alternate between two coupled modules: (1) a *Community Dynamics Module*, where agents play a non-cooperative game to optimize community memberships based on opinion alignment and structural cohesion; and (2) an *Opinion Adjustment Module*, where agents revise opinions via a bounded-confidence mechanism modulated by community-induced influence weights. This dual-stage process captures the feedback loop between structure and opinion. We prove that *GCAOFP* converges to stable equilibria, ensuring intra-community consensus and inter-community diversity—dynamics that standard models fail to replicate. Experiments on real-world networks show that *GCAOFP* better reproduces localized opinion clusters, while offering strong scalability and interpretability, illuminating the strategic foundations of polarization.

Datasets & Code & Extended version —

<https://github.com/ZINUX1998/GCAOFP-CODE>

Introduction

Opinion dynamics in social networks lies at the intersection of sociology, computer science, and systems engineering, aiming to understand how individual beliefs evolve through decentralized, interpersonal influence (Shi, Altafini, and Baras 2019; Peralta et al. 2021; Aghbolagh et al. 2023; Xu et al. 2025a). A core challenge is to model how local agent interactions yield emergent global phenomena such as consensus, polarization, and fragmentation (Dong et al. 2018; Zhang et al. 2020; Xu et al. 2025b).

Agent-based models flexibly capture localized opinion updates and their systemic consequences (Li et al. 2020;

Bindel, Kleinberg, and Oren 2015; Zhou et al. 2020; Hunter and Zaman 2022). Agents iteratively update their opinions by aggregating neighbor signals, reflecting the decentralized and recursive nature of social influence (Olfati-Saber, Fax, and Murray 2007; Hassani et al. 2022; Dong, Zhang, and Kong 2024). Classical formulations such as the *DeGroot* model (Degroot 2020) assume that repeated averaging over social links drives global consensus (Olfati-Saber and Murray 2004; Carli et al. 2010). However, real-world opinion landscapes often exhibit persistent disagreement and fragmentation, revealing the limitations of consensus-driven approaches. Bounded-confidence models, e.g., Hegselmann–Krause (*HK*) (Hegselmann and Krause 2002), restrict influence to peers with similar views, producing outcomes from consensus to fragmentation controlled by a confidence bound (Battiston et al. 2010).

Despite progress, most opinion dynamics models assume uniform influence or ignore latent community structures, hindering reproduction of empirical opinion clusters at the meso scale. In practice, opinion formation exhibits locality, as individuals are more influenced by peers within the same community due to stronger ties (Wu and Huberman 2004; Sunstein 2006). Xing et al. (Xing et al. 2022) report strong alignment between opinion clusters and community modularity, highlighting the role of meso-scale structures. Classical models focus on local averaging, overlooking mesoscale topologies mediating cross-community influence (Iñiguez et al. 2009). Recent studies emphasize that modular structures and sparse inter-community bridges shape global opinion patterns (Arruda et al. 2022; Anagnostopoulos et al. 2022; Ding et al. 2024). Few models jointly capture strategic behavior and opinion-community co-evolution, leaving a gap in modeling how individual incentives and structural adaptations shape collective opinion landscapes.

To bridge this gap, we introduce *GCAOFP*—a *Game Theory-Based Community-Aware Opinion Formation Process*—that integrates community-sensitive utility optimization into individual opinion updates. Unlike prior approaches that either assume fixed community partitions or exploit opinions merely for community inference (He et al. 2020; Ren et al. 2022), *GCAOFP* jointly models the co-evolution of opinions and latent community affiliations under a fixed interaction topology. Agents assign higher weights to intra-community opinions and iteratively

*Corresponding author.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

adjust their affiliations according to opinion similarity, forming a bi-level feedback loop formalized as a potential game. The mechanism embeds the bounded-confidence principle (Bernardo et al. 2024), limiting interactions to opinion-similar neighbors and reflecting cognitive constraints such as limited attention and selective exposure. This framework offers a controlled and interpretable setting to analyze micro-level opinion clustering and is particularly suited for environments where relational structures remain stable while ideological landscapes evolve, such as online forums, political subgroups, and long-lived digital communities. The main contributions are:

(1) **Community Susceptibility Metric:** We propose a novel metric measuring structural asymmetries in social influence, reflecting stronger intra-community trust.

(2) **Game-Theoretic Framework with Bounded Confidence:** Co-evolving opinions and communities are modeled as a non-cooperative game, with updates restricted to bounded confidence neighborhoods.

(3) **Efficient and Provably Convergent Algorithm:** A linear-time algorithm is proposed, with theoretical convergence guarantees for opinions and community assignments.

(4) **Comprehensive Empirical Validation:** Experiments on multiple real-world networks demonstrate that *GCAOFP* outperforms state-of-the-art baselines in community detection and opinion clustering.

Methodology

We consider opinion dynamics over a directed social network $\mathcal{N} = \{1, 2, \dots, n\}$, where influence asymmetry is captured by a non-negative matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$. Each entry \mathbf{W}_{ij} represents agent i 's confidence in agent j 's opinion, with no self-loops ($\mathbf{W}_{ii} = 0$), and the induced graph is assumed strongly connected (i.e., \mathbf{W} is irreducible). At any time $T \in \mathbb{R}^+$, the opinion profile is denoted as $\mathbf{x}(T) \in [0, 1]^n$, where $\mathbf{x}_i(T)$ is agent i 's belief. Community structure is encoded by a binary matrix $\mathbf{S}(T) \in \{0, 1\}^{n \times K}$ under single-membership constraint: $\forall i \in N, \sum_{k=1}^K s_{ik}(T) = 1$. The row vector $\mathbf{s}_i(T)$ indicates agent i 's membership, with $\mathbf{s}_i(T) \mathbf{s}_j^T(T) = 1$ if agents i and j belong to the same community. To reflect intra-community affinity, we define a diagonal trust amplification matrix $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, where $\lambda_i \in (1, +\infty)$ quantifies agent i 's increased trust in same-community neighbors (inter-community trust is normalized to 1). Standard notation applies: $\mathbf{1}_n$ is the all-one column vector, $\mathbf{1}^n = \mathbf{1}_n \mathbf{1}_n^T$, and $\text{tr}(\cdot)$ denotes the trace operator.

GCAOFP iteratively updates agent opinions and community affiliations. Agents play a non-cooperative game optimizing community membership and aligning individual incentives with collective cohesion. Concurrently, opinion evolution follows a bounded-confidence rule, limiting interactions to cognitively plausible intra-group neighbors. This coupled process captures the reciprocal dynamics of belief revision and structural adaptation, explaining emergent phenomena such as polarization and fragmentation. See **Appendix A** for game-theoretic analysis and **Appendix B** for algorithmic details and convergence guarantees.

Community Dynamics Process

During community updates, each agent strategically selects a community label in a non-cooperative game, aiming to maximize opinion similarity with neighbors and promote partitions exhibiting high intra-community coherence.

Definition 1. [Community Dynamics Process (CDP)] Given the agents set \mathcal{N} , opinion vector $\mathbf{x}(T)$, influence matrix \mathbf{W} and relative trust matrix Λ , the **Community Dynamics Process** is defined as the tuple, denoted as $\Omega(t, \mathbf{S}(t), \mathcal{U}(\cdot))$, where

- * $t = 0, 1, 2, \dots$ denotes the discrete time step;
- * $\mathbf{S}(t) \in \{0, 1\}^{n \times K}$ is the community membership matrix at time t , with initial membership $\mathbf{S}(0) = \mathbf{s}(T)$.
- * $\mathcal{U}(\cdot)$ is the community update function that determines $\mathbf{S}(t+1)$ based on current opinions and memberships.

At each period t , agent i maximizes a utility function $u_i(t)$ defined as follows, reflecting key behavioral principles: (1) **Opinion alignment:** Preference for neighbors with similar opinions, quantified by $|\mathbf{x}_i(t) - \mathbf{x}_j(t)|$; (2) **Community trust bias:** Stronger weighting of intra-community ties via λ_i ; (3) **Bounded confidence:** Interactions only considered when opinion differences are below cognitive threshold ψ ; (4) **Intra-community coherence:** Penalty for opinion divergence within same community captured by $u_{i1}(t)$; (5) **Cross-community interaction:** Selective inter-community influence modeled by $u_{i2}(t)$, activated when opinions align sufficiently (ReLU threshold).

$$\begin{aligned} u_i(t) &= u_{i1}(t) + u_{i2}(t) \\ u_{i1}(t) &= \lambda_i \sum_{j \neq i} \mathbf{S}_i(t) \mathbf{S}_j^T(t) \cdot [\psi - |\mathbf{x}_i(t) - \mathbf{x}_j(t)|] \cdot \mathbf{W}_{ij} \\ u_{i2}(t) &= \sum_{j \neq i} (1 - \mathbf{S}_i(t) \mathbf{S}_j^T(t)) \cdot \Xi[\psi - |\mathbf{x}_i(t) - \mathbf{x}_j(t)|] \\ &\quad \cdot \mathbf{W}_{ij} \end{aligned} \quad (1)$$

Let $\mathbf{D}(T) \in [0, 1]^{n \times n}$ denotes the pairwise bounded confidence matrix, where $\mathbf{D}_{ij}(T) = \psi - |\mathbf{x}_i(T) - \mathbf{x}_j(T)|$. $\Xi(\cdot)$ denotes the ReLU function. This encodes opinion proximity subject to a cognitive threshold ψ . The utility matrix across all agents at time t is then defined as:

$$\begin{aligned} \mathbf{U}(t) &= \Lambda \mathbf{S}(t) \mathbf{S}^T(t) \odot \mathbf{D}(T) \odot \mathbf{W} \\ &\quad + (\mathbf{1}^n - \mathbf{S}(t) \mathbf{S}^T(t)) \odot \Xi(\mathbf{D}(T)) \odot \mathbf{W} \\ &= \mathbf{B}(T) \odot \mathbf{W} + \mathbf{S}(t) \mathbf{S}^T(t) \odot (\Lambda \mathbf{D}(T) - \mathbf{B}(T)) \odot \mathbf{W} \end{aligned} \quad (2)$$

where $\mathbf{D}(T) = \Psi - \mathbf{E}(T)$, $\Psi = \psi \mathbf{1}^n$, $\mathbf{E}(T) = |\mathbf{x}(T) \mathbf{1}_n^T - \mathbf{1}_n \mathbf{x}^T(T)|$, $\mathbf{B}(T) = \Xi(\mathbf{D}(T))$ and \odot denotes element-wise multiplication. The system-level objective is to maximize social welfare, the aggregate utility:

$$\begin{aligned} sw(t) &= \sum_i u_i(t) = \sum_{ij} \mathbf{U}_{ij}(t) \\ &= \text{tr} \{ \mathbf{S}^T(t) [(\Lambda \mathbf{D}(T) - \mathbf{B}(T)) \odot \mathbf{W}] \mathbf{S}(t) \} \\ &\quad + \text{tr}(\mathbf{W} \mathbf{B}(T)). \end{aligned} \quad (3)$$

$\mathcal{N}_i^T = \{j \mid j \in N_i \wedge |\mathbf{x}_i(T) - \mathbf{x}_j(T)| < \beta\gamma^T\}$ denotes the set of the trusted neighbors of user i at time T . $\beta \in (0, 1)$, $\gamma \in (0, 1)$ are the tuning parameters for the size of \mathcal{N}_i^T . Obviously, as the number of iterations increases, the agents' opinions will gradually be consolidated and their confidence level will continue to decrease. The opinion weight matrix during time T , denoted by $\Phi(T)$, can be written as

$$\Phi(T) = \Theta(T) \oslash [\Theta(T) \mathbf{1}^n + \mathbf{1}^n], \quad (11)$$

$$\Theta(T) = [\mathbf{1}^n + (\Lambda - \mathbf{I}) \mathbf{1}^n \mathbf{s}(T) \mathbf{s}^\top(T)] \odot \mathbf{W} \odot \Xi(\beta\gamma^T \mathbf{1}^n - \mathbf{E}(T)) \quad (12)$$

where $\mathbf{I} = \text{diag}(1, 1, \dots, 1) \in \mathbb{R}^{n \times n}$ denotes the identity matrix, \oslash means element-wise divide, i.e., $\Phi_{ij}(T) = \Theta_{ij}(T) \div (\sum_{j \neq i} \Theta_{ij}(T) + 1)$. During every time T , each agent i is also associated with an intrinsic confidence level, denoted by $\delta_i(T)$, which is defined as

$$\delta_i(T) = \frac{1}{A_i}, \quad \text{where} \quad A_i = \sum_{j \neq i} [1 + (\lambda_i - 1) \mathbf{s}_i(T) \mathbf{s}_j^\top(T)] \cdot \mathbf{W}_{ij} \cdot \Xi[\beta\gamma^T - |\mathbf{x}_i(T) - \mathbf{x}_j(T)|] + 1 \quad (13)$$

Therefore, the opinion update rule Eq. (9) can be written in matrix form as

$$\mathbf{x}(T+1) = \mathbf{1}_n \oslash [\Theta(T) \mathbf{1}_n + \mathbf{1}_n] \odot \mathbf{x}(T) + \Phi(T) \mathbf{x}(T) \quad (14)$$

Lemma 2. $\mathbf{x}(T+1)$ must be in the interval of $[0, 1]$, as $\forall i \in \mathcal{N}$, $\delta_i(T) + \sum_{j \neq i} \Phi_{ij}(T) = 1$.

Remark 1. *GCAOFP integrates community-aware mechanisms into the opinion update process, enabling joint modeling of community formation and opinion evolution. Leveraging the bottom-up architecture of CDP, it uncovers latent community structures in a self-organizing manner—without prior knowledge such as the number of communities—while simultaneously simulating localized opinion dynamics.*

Theorem 4. *If the opinion dynamics follows the function in Eq. (9), for every user i and every time $T \geq 0$, we have*

$$|\mathbf{x}_i(T+1) - \mathbf{x}_i(T)| < \beta\gamma^T \quad (15)$$

Therefore, we claim that every agents i 's opinion vector will converge to a relatively stable state in finite periods T^* , such that $\forall i \in \mathcal{N}$, $|\mathbf{x}_i(T^*) - \mathbf{x}_i(T^* - 1)| < \epsilon$, where ϵ is some very small positive threshold.

Theorem 5. *From any arbitrary initial distribution of opinions $\mathbf{x}(0)$ and any arbitrary initial community membership matrix $\mathbf{s}(0)$, the CDP introduced in Definition 1 will converge to the equilibrium state where the community membership matrix keeps unchanged.*

Theorem 6. *Given that there are M directed connections among agents, from arbitrary initial opinion vector $\mathbf{x}(0)$ and community membership matrix $\mathbf{s}(0)$, the time complexity of GCAOFP introduced is upper-bounded by $\mathcal{O}(t^{Exp} M)$. (Detailed experimental validation is provided in the Appendix C.1).*

Comparative Discussion

LPA-HK (Peng, Zhao, and Hu 2023) updates community labels by considering both network topology and agent opinions. It assigns higher influence weights to neighbors whose opinions fall within a bounded confidence threshold ψ . The corresponding utility function is defined as:

$$u'_i(T) = \sum_{j \neq i} \mathbf{s}_i(T) (\mathbf{s}_j^*(T))^\top \varepsilon(\psi - |\mathbf{x}_i(T) - \mathbf{x}_j^*(T)|) \cdot (1 - |\mathbf{x}_i(T) - \mathbf{x}_j^*(T)|) \mathbf{W}_{ij}. \quad (16)$$

Can be implemented in a simpler way:

$$u'_i(T) = \sum_{j \neq i} \mathbf{s}_i(T) (\mathbf{s}_j^*(T))^\top \cdot \Xi(\psi - |\mathbf{x}_i(T) - \mathbf{x}_j^*(T)|) \mathbf{W}_{ij} \quad (17)$$

where $\varepsilon(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{if } x \leq 0 \end{cases}$, $\mathbf{s}_j^*(T) = \mathbf{s}_j(T)$ and $\mathbf{x}_j^*(T) = \mathbf{x}_j(T)$ if the opinion and label of agent j have already been updated, otherwise $\mathbf{s}_j^*(T) = \mathbf{s}_j(T-1)$ and $\mathbf{x}_j^*(T) = \mathbf{x}_j(T-1)$. However, Eq.17 exclusively accounts for neighbors within the confidence bound, which can result in undue dominance by the closest opinion neighbor. As illustrated in Fig.1 (Left), with $\psi = 0.4$, agent i favors the left community driven by a single nearby neighbor, despite the majority lying outside the confidence threshold. In contrast, the utility function proposed in Eq. 1 introduces penalties for intra-community neighbors whose opinions exceed the confidence bound, mitigating this bias. For the same scenario, agent i prefers the right community with a utility of 0.60 versus -0.4 for the left (with $\lambda = 1.5$), effectively correcting the misleading aggregation under local confidence.

Moreover, Eq.1 explicitly models inter-community influence, ensuring agents with divergent views are less likely to be clustered together. As shown in Fig.1 (Right), while intra-community utilities $\lambda_i \sum_{j \neq i} \mathbf{S}_i(t) \mathbf{S}_j^\top(t) \cdot [\psi - |\mathbf{x}_i(T) - \mathbf{x}_j(T)|] \cdot \mathbf{W}_{ij}$ are zero for both communities, the inter-community utility $\sum_{j \neq i} (1 - \mathbf{S}_i(t) \mathbf{S}_j^\top(t)) \cdot \Xi[\psi - |\mathbf{x}_i(T) - \mathbf{x}_j(T)|] \cdot \mathbf{W}_{ij}$ favors the right community ($0.3 > 0.2$), leading to more semantically coherent partitions. Overall, *GCAOFP* achieves finer-grained community separation by jointly modeling agreement and disagreement across and within communities—capabilities lacking in *LPA-HK*—resulting in clusters with higher intra-community opinion cohesion.

Experiments

In this section, we will conduct comprehensive experiments using real-world networks to verify the *effectiveness*, *scalability* and *parameter sensitivity* of *GCAOFP*.

Experiment Setup

Comparison Methods We compare *GCAOFP* against representative opinion dynamics models grouped as follows: classical methods ignoring community structure—*DeGroot*, *Friedkin-Johnsen (FJ)* (Friedkin and Johnsen 2010), *DeGroot-Friedkin (DF)* (Jia et al. 2015), and *HK*; and recent

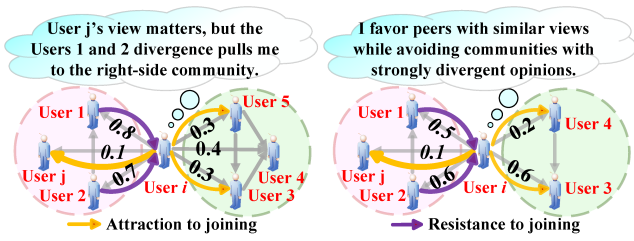


Figure 1: **Left:** Strategy selection contrast between *GCAOFP* and *LPA-HK*. **Right:** Motivation for including inter-community neighbor influence in Eq. (1). Edge weights are uniform, E_{ij} values for agent i are shown, and node colors denote communities.

co-evolutionary approaches—*GK-Means* (Bu et al. 2020) and *LPA-HK*, which jointly update opinions and community labels. In addition, we consider the *Coevolutionary Model (CEM)* (Aghbolagh et al. 2023), a game-theoretic framework where agents simultaneously update actions and opinions across influence and communication layers, with guaranteed convergence to Nash equilibria. Details on repeatability are provided in **Appendix C.2**.

Network Datasets We conducted case studies on eight real-world networks listed in Table 1. Preprocessing involved removing self-loops and extracting the largest strongly connected components. All datasets are publicly accessible, and the processed versions are provided.

Evaluation Metrics We define three metrics to assess the effectiveness of *GCAOFP*. *OSWG* measures the overall gain in social welfare from initial to steady state. *RCSW* quantify the final consensus level. *ACL* captures intra-community agreement at convergence. These metrics jointly capture global welfare gains and local consensus. (Computational details are provided in **Appendix C.3**)

$$OSWG = \frac{sw(T^*) - sw(0)}{\lambda W} \quad RCSW = \frac{sw(T^*)}{\lambda \psi W}$$

$$ACL = \frac{1}{k} \sum_k CL_k = \frac{1}{k} \sum_k \sum_{i \in C_k} \frac{1 - \|\mathbf{x}_i^* - \mathbf{c}_k^*\|}{|C_k|}$$

where $W = \sum_{i \in \mathcal{N}} \sum_{j \neq i} \mathbf{W}_{ij}$, C_k indicates the k -th community in the equilibrium state and $\mathbf{c}_k^* = \frac{1}{|C_k|} \sum_{i \in C_k} \mathbf{x}_i^*$ denotes the average opinion of all members within C_k .

Experimental Results

Effectiveness Analysis We evaluate the effectiveness of *GCAOFP* on multiple real-world and synthetic networks, each initialized with 100 randomized opinion vectors and community assignments. As shown in Fig.2, our method consistently achieves significant improvements in overall social welfare—measured via *OSWG*—across diverse initial conditions. Interestingly, we observe that denser networks exhibit markedly lower variance in outcomes. This stabilization arises from higher average degrees (d), which promote richer opinion exchange and enhance alignment robustness.

Network	D	W	# Nodes	# Edges	$\langle d \rangle$	# GT
Karate	×	×	34	78	4.59	2
Dolphins	×	×	62	159	5.13	2
PolBooks	×	×	105	441	8.40	3
PolBlogs	✓	×	793	15,783	19.90	2
BitAlpha	✓	✓	3,235	23,299	7.20	×
BitOTC	✓	✓	4,709	33,461	7.11	×
GRQC	×	×	4,158	13,425	6.58	×
HEPH	×	×	8,638	24,806	5.74	×

Table 1: Statistics of real-world networks. (D: Directed, W: Weighted, #GT: Ground-Truth community numbers)

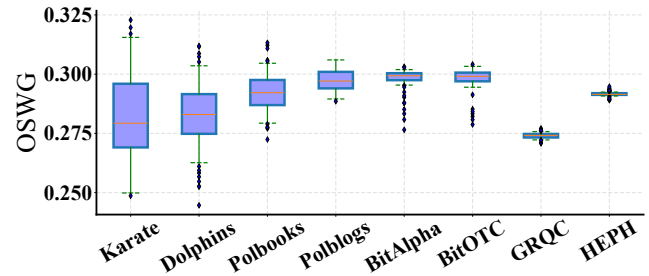


Figure 2: Boxplots of *OSWG* from 100 runs of *GCAOFP* on real-world networks ($\lambda = 1.4$, $\psi = 0.4$).

The resulting convergence is less sensitive to initial configurations, a phenomenon further validated in **Appendix C.4** through controlled synthetic experiments that isolate the effect of network density.

We evaluate dynamic opinion models on the *Karate* and *PolBooks* networks, focusing on convergence behavior and the quality of induced community structures. Fig. 3, 4 and 5 show representative opinion trajectories and corresponding *ACL* dynamics. (1) **Fast Convergence.** *GCAOFP* exhibits rapid convergence comparable to baseline methods. Except for *FJ*, most models drive global consensus with *ACL* values approaching 1, indicating high alignment of individual opinions. (2) **Community-Aware Equilibrium.** Unlike models enforcing global uniformity, *GCAOFP* facilitates localized convergence: agents within the same community reach consensus, while opinions diverge across communities. This results in interpretable steady states closely reflecting known social partitions, as observed in the *Karate* network (Fig. 7 and 8). (3) **Instability of *GK-Means*.** The susceptibility parameter λ critically affects *GK-Means* performance. On *Karate*, a high $\lambda = 0.9$ fosters strong intra-group cohesion but weakens inter-group separation, blurring boundaries. Conversely, on *PolBooks*, a lower $\lambda = 0.4$ improves partition accuracy but reduces intra-community agreement, lowering *ACL* scores. (**Appendix C.5**) (4) **Polarization Threshold Effects in *LPA-HK*.** The bounded confidence parameter δ significantly shapes opinion dynamics. Increasing δ leads the system from pluralism to polarization, and eventually to global consensus. At $\delta = 0.95$, both networks converge to uniform opinions but fail to recover meaningful community structures (**Appendix C.5**), demonstrating a trade-off between consensus and community fi-

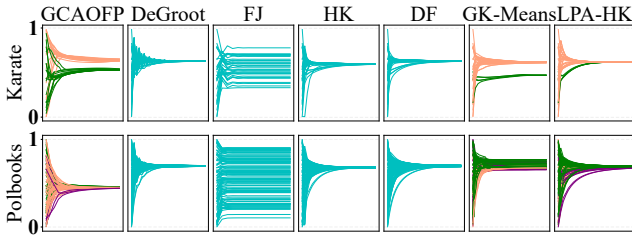


Figure 3: Opinion dynamics on *Karate* and *PolBooks*, with curve colors indicating agent community labels per model.

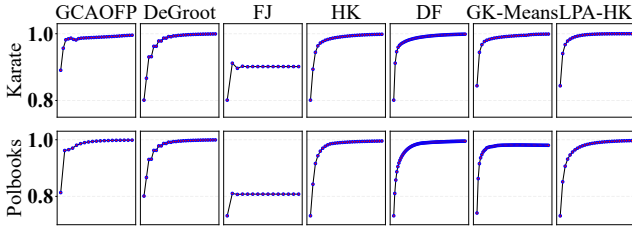


Figure 4: Dynamics of *ACL* for different opinion dynamics models on *Karate* and *PolBooks*.

delity. (5) **Opinion Polarization in CEM.** The coevolutionary framework explicitly partitions the network into two opposing groups, without enforcing consensus. Under this model, opinions progressively polarize toward two extremes (+1 and -1). In contrast to other methods, the *ACL* steadily decreases over time, reflecting the emergence of strong opinion polarization.

Utility-Driven Structural Adaptation. While *GK-Means* and *LPA-HK* promote intra-community agreement, their reliance on static or heuristic rules often leads to suboptimal partitions. In contrast, *GCAOFP* tightly integrates opinion dynamics with community evolution via a welfare-guided update mechanism, where $\mathbf{D}_{ij}^+(T)$ denotes $\mathbf{D}_{ij}(T) > 0$, $\mathbf{D}_{ij}^-(T)$ denotes $\mathbf{D}_{ij}(T) \leq 0$:

$$\mathbf{Y}_{ij}(T) = \begin{cases} [(\lambda_i - 1) \mathbf{W}_{ij} + (\lambda_j - 1) \mathbf{W}_{ji}] \mathbf{D}_{ij}(T), \mathbf{D}_{ij}^+(T) \\ (\lambda_i \mathbf{W}_{ij} + \lambda_j \mathbf{W}_{ji}) \mathbf{D}_{ij}(T), \mathbf{D}_{ij}^-(T) \end{cases} \quad (19)$$

The interaction score $\mathbf{Y}_{ij}(T)$ increases with opinion distance $\mathbf{D}_{ij}(T)$, modulated by susceptibility weights. Agents iteratively reassign community labels to maximize local utility (Eq. 6), reinforcing intra-community cohesion. This coevolution yields stable, utility-consistent partitions robust to initialization and aligned with social semantics.

Parameter Analysis The parameter λ quantifies agents' sensitivity to community identity, with higher values accelerating intra-community opinion convergence and producing finer-grained community partitions. Conversely, the parameter ψ , analogous to the confidence threshold ϱ in the *HK* model, regulates the acceptance range for neighbors' opinions. Larger values increase opinion receptivity, thereby promoting global consensus and resulting in fewer, larger communities. We systematically evaluate the effects of λ

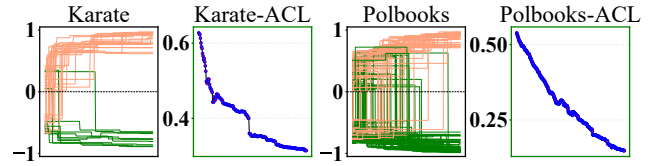


Figure 5: Opinion trajectories and *ACL* evolution of *CEM*.

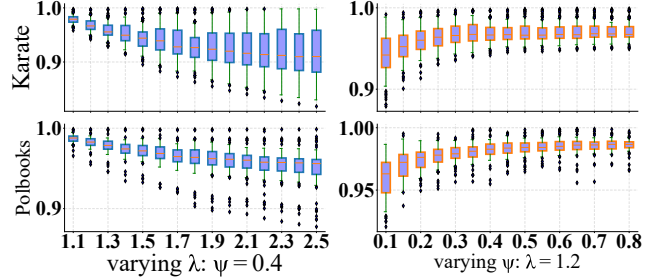


Figure 6: Boxplots of *RCSW* (100 runs) of *GCAOFP* on *Karate* under varying λ and ψ .

and ψ on the *Karate* and *PolBooks* networks, conducting 100 random initializations for each setting. As illustrated in Fig. 6, at fixed ψ , lower λ leads to higher overall social welfare (*RCSW*) and reduced variance across runs, indicating stronger global consensus and coarser community structures. Increasing λ reverses these trends, fostering opinion fragmentation and yielding finer community partitions. At fixed λ , increasing ψ similarly enhances *RCSW* and decreases variance, reflecting improved opinion coherence and reduced fragmentation. Notably, when λ is high, appropriate tuning of ψ enables *GCAOFP* to strike a balance between localized convergence and effective community detection. Additional sensitivity analyses are presented in **Appendix C.6**.

Community detection We evaluate *GCAOFP* with opinion- or community-based initialization (*GCAOFP-O* and *GCAOFP-C*) on four benchmark networks, compared with *LPA* (Raghavan, Albert, and Kumara 2007), *k-clique* (Palla et al. 2005), *GLEAM* (Bu et al. 2018), *GK-Means*, *LPA-HK*. Performance is measured by Adjusted Rand Index (ARI) (Hubert and Arabie 1985) and Adjusted Mutual Information (AMI) (Vinh, Epps, and Bailey 2010), averaged over 100 randomized runs. Fig. 7 highlights five key findings: 1) Opinion-aware models match or outperform structure-only baselines; 2) With well-informed opinion vectors, *GCAOFP* achieves the most accurate recovery of ground-truth communities; 3) Performance depends on initialization, highlighting the need for effective strategies; 4) Sensitivity to opinion initialization decreases with network size, indicating robustness; 5) *GCAOFP* is more robust to community than opinion initialization, showing opinion initialization dominates convergence; **Appendix C.7** examines *GCAOFP*'s sensitivity to λ , ψ , and initialization.

We visualize community detection results on the *Karate* and *PolBooks* networks (Fig. 8), where node size encodes final opinions, color denotes detected communi-

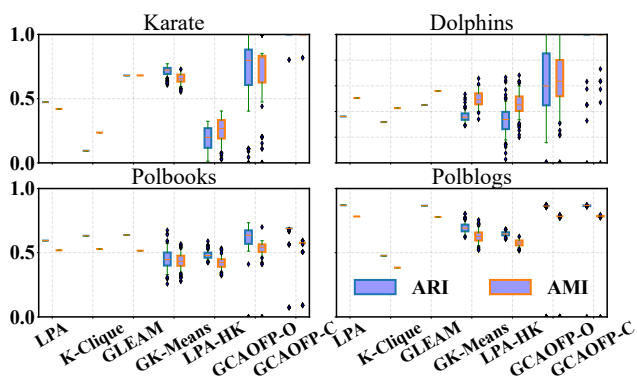


Figure 7: Performance comparison on real-world networks with ground-truth communities, averaged over 100 runs.

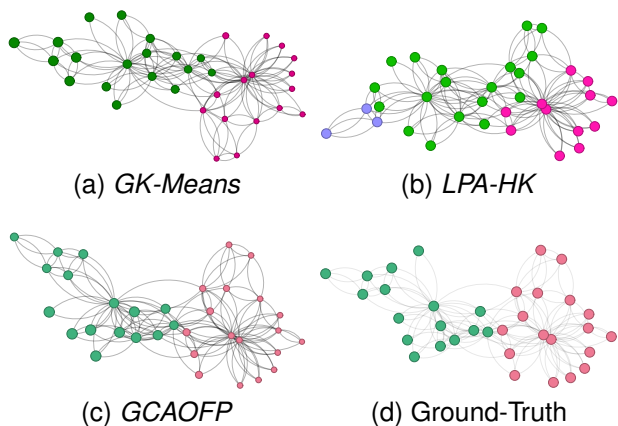


Figure 8: Community detection results on *Karate* network. Complete visualizations with full details, as well as results for *PolBooks*, are provided in **Appendix C.8**.

ties, and edge thickness represents influence. Ground-truth plots color nodes by their true community labels for comparison. On *Karate*, *GCAOFP* perfectly recovers ground-truth communities (ARI/AMI = 1.0), outperforming *GK-Means* (0.772/0.727) and *LPA-HK* (0.324/0.403). On *PolBooks*, *GCAOFP* misclassifies 12 nodes (ARI/AMI = 0.727/0.637), fewer than *GK-Means* (16) and *LPA-HK* (21), with errors concentrated in a sparse cluster that weakens intra-community influence. These results demonstrate that *GCAOFP* achieves more accurate and interpretable community separation, particularly in densely connected regions.

Sensitivity to Deviations from Rationality. To assess robustness under bounded rationality, we introduce a stochasticity parameter $\varphi \in [0, 0.09]$. At each iteration, agent i follows the best response (Eq. 7) with probability $1 - \varphi$ and randomly samples from its candidate set $S_i(t)$ with probability φ , injecting controlled noise into the dynamics. As shown in Fig. 9, the impact of stochasticity varies across networks. On *Karate* and *PolBooks*, clustering performance degrades monotonically as φ increases, indicating that even mild deviations from rational updates can disrupt coherent commu-

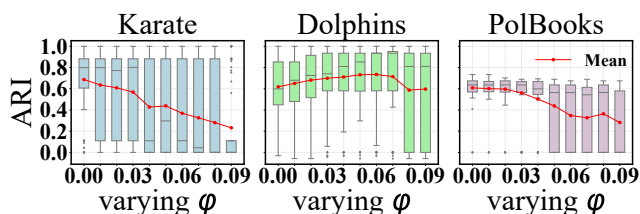


Figure 9: Impact of stochastic community reassignment under varying bounded rationality; 100 runs per parameter.

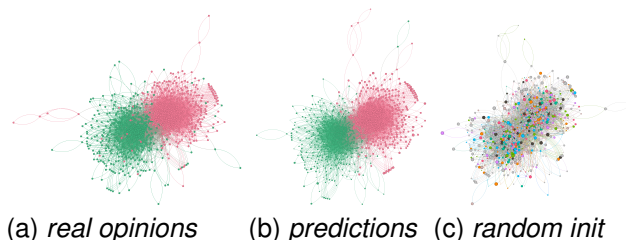


Figure 10: Opinion dynamics on the Blogs dataset.

nity formation. In contrast, *Dolphins* exhibits a slight improvement under small stochasticity before deteriorating at larger φ . This contrast reflects fundamental differences between human and animal social networks. Stochastic deviations disrupt fragile, identity-driven community boundaries in human networks, whereas mild randomness can help refine behaviorally grounded modularity in animal networks.

Evaluation on Real-World Opinion Data. We evaluate *GCAOFP* on the 2005 *Blogs* (Matakos, Terzi, and Tsaparas 2017) network of U.S. political weblogs labeled as Liberal or Conservative. After converting to an undirected form and retaining the largest connected component, the network comprises two communities (636 and 586 nodes) connected by 33,431 edges. Fig. 10(a) visualizes the real opinion distribution using color-coded labels (red: Liberal, green: Conservative). We initialize opinions and community labels randomly based on the network topology (Fig. 10(c)), producing a highly disordered initial state. The predicted opinion states (Fig. 10(b)) exhibit clear opinion polarization, with strong alignment to ground-truth labels (ARI/AMI = 0.8071/0.7132). Moreover, final opinions exhibit strong intra-community coherence and distinct stratification, with red-community nodes showing stronger Liberal tendencies—highlighting *GCAOFP*’s capacity to capture realistic opinion polarization.

Conclusion

We propose *GCAOFP*, a game-theoretic opinion dynamics model integrating bounded confidence with strategic community selection. We introduce *Community Susceptibility* to quantify asymmetric influence and design a scalable algorithm with provable convergence. Extensive experiments validate that *GCAOFP* effectively captures localized consensus and uncovers latent community structures.

Acknowledgments

This research was partially supported by the National Natural Science Foundation of China under Grant No.U22B2036, and in part by the National Natural Science Foundation of China under Grant 72371133, and in part by Shaanxi Provincial Key Research and Development Program - Key Project - "Four-Chains" Integration Project of the Qinchuangyuan General Platform under Grant 2024PT-ZCK-93, and in part by the QingLan Project of Jiangsu Province.

References

- Aghbolagh, H. D.; Ye, M.; Zino, L.; Chen, Z.; and Cao, M. 2023. Coevolutionary Dynamics of Actions and Opinions in Social Networks. *IEEE Transactions on Automatic Control*, 68(12): 7708–7723.
- Anagnostopoulos, A.; Becchetti, L.; Cruciani, E.; Pasquale, F.; and Rizzo, S. 2022. Biased opinion dynamics: when the devil is in the details. *Information Sciences*, 593: 49–63.
- Arruda, H. F. d.; Cardoso, F. M.; Arruda, G. F. d.; Hernández, A. R.; Costa, L. d. F.; and Moreno, Y. 2022. Modelling how social network algorithms can influence opinion polarization. *Information Sciences*, 588: 265–278.
- Battiston, F.; Cencetti, G.; Iacopini, I.; Latora, V.; Lucas, M.; Patania, A.; Young, J.-G.; and Petri, G. 2010. Networks beyond pairwise interactions: Structure and dynamics. *Physics Reports*, 874: 1–92.
- Bernardo, C.; Altafini, C.; Proskurnikov, A.; and Vasca, F. 2024. Bounded confidence opinion dynamics: A survey. *Automatica*, 159: 111302.
- Bindel, D.; Kleinberg, J.; and Oren, S. 2015. How bad is forming your own opinion? *Games and Economic Behavior*, 92: 248–265.
- Bu, Z.; Cao, J.; Li, H.-J.; Gao, G.; and Tao, H. 2018. GLEAM: a graph clustering framework based on potential game optimization for large-scale social networks. *Knowledge and Information Systems*, 55: 741–770.
- Bu, Z.; Li, H.; Zhang, C.; Cao, J.; Li, A.; and Shi, Y. 2020. Graph k-means based on leader identification, dynamic game, and opinion dynamics. *IEEE Transactions on Knowledge and Data Engineering*, 32(7): 1348–1361.
- Carli, R.; Fagnani, F.; Frasca, P.; and Zampieri, S. 2010. Gossip consensus algorithms via quantized communication. *Automatica*, 46(1): 70–80.
- Degroot, M. H. 2020. Reaching a consensus. *Journal of the American Statistical Association*, 69(345): 118–121.
- Ding, R.-X.; Cheng, R.-X.; Li, M.-N.; Yang, G.-R.; and Herrera-Viedma, E. 2024. Conflict management-based consensus reaching process considering conflict relationship clustering in large-scale group decision-making problems. *Expert Systems with Applications*, 238: 122095.
- Dong, J.; Zhang, Y.-C.; and Kong, Y. 2024. The evolution dynamics of collective and individual opinions in social networks. *Expert Systems with Applications*, 255: 124813.
- Dong, Y.; Zhan, M.; Kou, G.; Ding, Z.; and Liang, H. 2018. A survey on the fusion process in opinion dynamics. *Information Fusion*, 43: 57–65.
- Friedkin, N. E.; and Johnsen, E. C. 2010. Social influence and opinions. *Journal of Mathematical Sociology*, 15: 193–206.
- Hassani, H.; Razavi-Far, R.; Saif, M.; Chiclana, F.; Krejcar, O.; and Herrera-Viedma, E. 2022. Classical dynamic consensus and opinion dynamics models: A survey of recent trends and methodologies. *Information Fusion*, 88: 22–40.
- He, Q.; Wang, X.; Mao, F.; Lv, J.; Cai, Y.; Huang, M.; and Xu, Q. 2020. CAOM: A community-based approach to tackle opinion maximization for social networks. *Information Sciences*, 513: 252–269.
- Hegselmann, R.; and Krause, U. 2002. Opinion Dynamics and Bounded Confidence Models, Analysis and Simulation. *Journal of Artificial Societies and Social Simulation*, 5(3).
- Hubert, L.; and Arabie, P. 1985. Comparing partitions. *Journal of Classification*, 2: 193–218.
- Hunter, D. S.; and Zaman, T. 2022. Optimizing Opinions with Stubborn Agents. *Operations Research*, 70(4): 2119–2137.
- Iñiguez, G.; Kertész, J.; Kaski, K. K.; and Barrio, R. A. 2009. Opinion and community formation in coevolving networks. *Phys. Rev. E*, 80: 066119.
- Jia, P.; MirTabatabaei, A.; Friedkin, N. E.; and Bullo, F. 2015. Opinion Dynamics and the Evolution of Social Power in Influence Networks. *SIAM Review*, 57(3): 367–397.
- Li, K.; Liang, H.; Kou, G.; and Dong, Y. 2020. Opinion dynamics model based on the cognitive dissonance: An agent-based simulation. *Information Fusion*, 56: 1–14.
- Matakos, A.; Terzi, E.; and Tsaparas, P. 2017. Measuring and moderating opinion polarization in social networks. *Data Min. Knowl. Discov.*, 31: 1480–1505.
- Nash, J. 1951. Non-cooperative games. *Annals of Mathematics*, 54(2): 286–295.
- Olfati-Saber, R.; Fax, J. A.; and Murray, R. M. 2007. Consensus and Cooperation in Networked Multi-Agent Systems. *Proceedings of the IEEE*, 95(1): 215–233.
- Olfati-Saber, R.; and Murray, R. M. 2004. Consensus problems in networks of agents with switching topology and time-delays. *IEEE Transactions on Automatic Control*, 49(9): 1520–1533.
- Palla, G.; Derényi, I.; Farkas, I.; and Vicsek, T. 2005. Uncovering the overlapping community structure of complex networks in nature and society. *Nature*, 435: 814–818.
- Peng, Y.; Zhao, Y.; and Hu, J. 2023. On the role of community structure in evolution of opinion formation: A new bounded confidence opinion dynamics. *Information Sciences*, 621: 672–690.
- Peralta, A. F.; Neri, M.; Kertész, J.; and Iñiguez, G. 2021. Effect of algorithmic bias and network structure on coexistence, consensus, and polarization of opinions. *Phys. Rev. E*, 104: 044312.
- Raghavan, U. N.; Albert, R.; and Kumara, S. 2007. Near linear time algorithm to detect community structures in large-scale networks. *Phys. Rev. E*, 76(3): 036106.

- Ren, R.; Shao, J.; Cheng, Y.; and Wang, X. 2022. Detecting Hierarchical and Overlapping Network Communities Based on Opinion Dynamics. *IEEE Transactions on Knowledge and Data Engineering*, 34(6): 2696–2710.
- Shi, G.; Altafini, C.; and Baras, J. S. 2019. Dynamics over Signed Networks. *SIAM Review*, 61(2): 229–257.
- Sunstein, C. R., ed. 2006. *Infotopia: How many minds produce knowledge*. New York, NY: Oxford University Press.
- Vinh, N. X.; Epps, J.; and Bailey, J. 2010. Information Theoretic Measures for Clusterings Comparison: Variants, Properties, Normalization and Correction for Chance. *Journal of Machine Learning Research*, 11(95): 2837–2854.
- Wu, F.; and Huberman, B. A. 2004. Social structure and opinion formation. arXiv:0407252.
- Xing, Y.; Wang, X.; Qiu, C.; Li, Y.; and He, W. 2022. Research on opinion polarization by big data analytics capabilities in online social networks. *Technology in Society*, 68: 101902.
- Xu, Y.; Liu, S.; Cheng, T.; Feng, X.; Jun, W.; and Shang, X. 2025a. Opinion convergence and management: Opinion dynamics in interactive group decision-making. *European Journal of Operational Research*, 323(3): 938–951.
- Xu, Y.; Liu, X.; Yuan, J.; Luo, J.; Zhou, W.; Yu, M.; and He, Y. 2025b. POMM: A public opinion management model integrating network game and opinion dynamics for social networks. *Knowledge-Based Systems*, 310: 112964.
- Zhang, H.; Zhao, S.; Kou, G.; Li, C.-C.; Dong, Y.; and Herrera, F. 2020. An overview on feedback mechanisms with minimum adjustment or cost in consensus reaching in group decision making: Research paradigms and challenges. *Information Fusion*, 60: 65–79.
- Zhou, Q.; Wu, Z.; Altalhi, A. H.; and Herrera, F. 2020. A two-step communication opinion dynamics model with self-persistence and influence index for social networks based on the DeGroot model. *Information Sciences*, 519: 363–381.