

# Generalizability of Conformist Social Influence Beyond Direct Reference

Ryutaro Mori<sup>1,2</sup>, Hidezo Suganuma<sup>1</sup>, and Tatsuya Kameda<sup>3</sup>

ryutau.mori@gmail.com, suganuma.hiz@gmail.com, tkameda@mi.meijigakuin.ac.jp

<sup>1</sup>The University of Tokyo, <sup>2</sup>Japan Society for the Promotion of Science, <sup>3</sup>Meiji Gakuin University

## Abstract

Conformity refers to phenomena where people match their behavior to others. Much research has focused on cases where people observe others in identical situations, saying little about its depth or generalizability. When conforming, do people revise behaviors only in that specific situation, or do they update more deeply to maintain consistent behaviors across situations? Using simulations, we first show that deep and shallow conformity leads to contrasting group dynamics; only with deep conformity can groups accumulate improvements beyond individual lifespans. We further conduct an experiment using an estimation task to examine the depths of conformity in humans. People generally extended conformist social influence to new situations without direct reference to others. However, those who simply averaged their answer with that of the direct reference showed notable failures in this generalization. Collectively, our research highlights the importance of distinguishing different depths of conformity when studying social influence and resulting group outcomes.

**Keywords:** conformity; social learning; cultural evolution; generalization

## Introduction

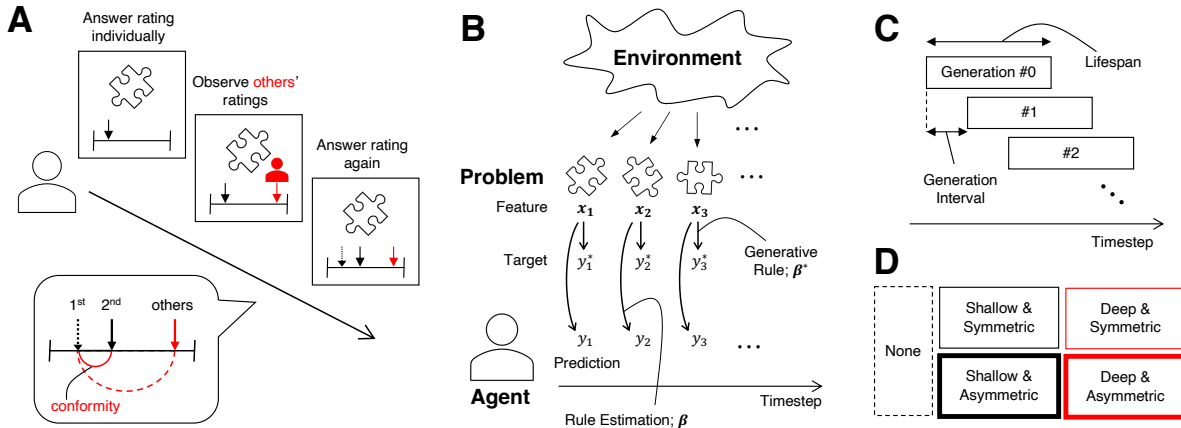
Humans routinely conform their behaviors to those of other individuals in the same situation. For example, in making everyday consumer choices, we often follow sales rankings or friends' reviews rather than independently selecting from available options. While earlier studies typically associated such conformity with irrational acts that make us elude objective criteria (Asch, 1956), it has also been argued that conformity can be understood as informationally rational, with people using social cues to acquire information regarding the adaptive behaviors in given situations (Kameda & Nakanishi, 2002; Toelch & Dolan, 2015). Given that there is a limitation in an individual's knowledge of what is adaptive in the environment, behaving as others do, or at least as similar to others, can provide a rational solution to avoid costly mistakes or save time. Conformity by individuals also catalyzes a variety of collective phenomena such as efficient collective search (List, Elsholtz, & Seeley, 2009) or erroneous informational cascades (Hung & Plott, 2001). Interestingly, small differences in how individuals conform to each other may amplify to shape drastically different dynamics at the collective level.

One prominent experimental paradigm to investigate conformity is a "rating paradigm (Figure 1A)," in which participants initially take independent behaviors (rating or estimation), then observe other individuals' behaviors as social information, and finally arrive at a revision (Figure 1A;

Klucharev, Hytönen, Rijpkema, Smidts, & Fernández, 2009; Schultze, Rakotoarisoa, & Stefan, 2015; Jayles, Kim, Escobedo, Cezera, Blanchet, Kameda, Sire, & Theraulaz, 2017). By comparing the initial and the revised behaviors, researchers can determine the extent to which the individuals updated their behaviors under conformist social influence. Human individuals are known to conform to social information across domains, from facial attractiveness ratings (Klucharev et al., 2009; Izuma & Adolphs, 2013) to gambling choices (Suzuki, Jensen, Bossaerts, & O'Doherty, 2016) and perceptual estimations (Molleman, Kurvers, & van den Bos, 2019; Kuroda, Ogura, Ogawa, Tamei, Ikeda, & Kameda, 2022), and more so, for example, if others to whom they refer is more reliable (Toelch, Bruce, Newson, Richerson, & Reader, 2014; Kameda, Toyokawa, & Tindale, 2022) or have more things in common with themselves (Baron, Kerr, & Miller, 1992).

Most of the present research about conformity focuses on how people conform to others in direct reference, a situation where they can observe others in the same specific situation (but see Nook, Ong, Morelli, Mitchell, & Zaki, 2016 for a notable exception). Here, we argue that this leaves a crucial question unanswered: How much does conformist social influence generalize across similar but distinct situations even when others are no longer observable (i.e., without direct reference)? Suppose a child conforms to their parent refraining from talking at a funeral. Does this only drive the imitation of a particular behavior to stay silent at the specific funeral, or does it generalize more broadly across similar yet distinct situations, such as subdued behavior in other formal settings? The generalization of conformity can be restated as the difference in the depths at which it operates. If conformity is *shallow*, affecting only observable behaviors (e.g., staying calm at that funeral with the parent), changes will be confined to specific situations where direct reference is possible. Conversely, if it is *deep* and influences the underlying generative function of behaviors (e.g., behave reservedly if a situation is formal), then agents will also update their behavioral patterns in line with others in other situations without direct reference.

Our study explores these different depths of conformity from two perspectives. First, we ask how the distinction is relevant in terms of collective phenomena: When and how do agents with deep and shallow conformity shape different dynamics in groups and yield different macro patterns? Drawing from the literature on cultural evolution, we suggest that the depth of social influence determines what a group can accumulate over time. Theories of cultural evolution



**Figure 1:** **A:** Illustration of the “re-rating” paradigm. **B, C, D:** Schematic diagram of the simulation setup. Agents exist within an environment that continuously poses them a new problem, which consists of a feature vector and a target value (B). Agents share overlapping generations (C). We explored five different types of conformity (D).

maintain that conformist social learning plays a key role in enabling the spread of rare effective inventions across people with high fidelity while eliminating the inefficient “re-invention of wheels” within a group (e.g., Lala, 2017). Here, we show that for conformity to work in that way, it must be *deep*. We introduce a simple model, in which agents of overlapping generations learn adaptive behaviors in a vast environment. We manipulate how individual agents conform to each other and compare the resulting collective dynamics.

The next question we ask is empirical: Do people actually exhibit such deep conformity? As we have argued, most existing research addressing people’s conformity behavior, including the rerating paradigm, failed to distinguish the different depths of influence. This is because participants could observe the behavior of others in the same situation every time they were asked to revise their behavior. Critically, regardless of whether they are merely aligning behaviors in the focal situation (i.e., shallow conformity) or updating a deeper function that produces behaviors across situations (i.e., deep conformity), their behavioral outputs are the same unless they face a new situation in which they have never directly referred to others’ behaviors. Therefore, we extend the rerating paradigm so that participants can observe others in only half of the situations. By comparing their revision patterns with and without direct social reference, we examine whether, and if so how, conformist social influence is deep and generalizes to new situations. We find some evidence that people generalize conformity beyond direct reference, but only partially. We further examine the reasons explaining this partiality.

### Simulation

In this simulation, we explore and demonstrate how conformity in varying depths at the individual level shapes different dynamics at the collective level. We set up a simple

model that captures an environment in which agents strive to learn its true generative rule to behave adaptively. Although the lifespan of an agent is too short to reach optimal performance individually, as their lives partially overlap, they may be able to accumulate improvements across generations through social influence. Critically, we systematically vary how an individual conforms to other agents and compare the resultant collective outcomes.

### Method

**Environment, Problem, and Agent** Our model consists of three main classes of objects (Fig. 1B): the environment, the problem, and the agent. *Agents* exist within an *environment* that continuously poses them *problems*. A new problem  $j$  consists of  $n_{\text{dim}}$ -dimensional features;  $\mathbf{x}_j \in \mathbb{R}^{n_{\text{dim}}}$ , and a corresponding target value;  $y_j^* \in \mathbb{R}$ . The environment has true generative rule parameters;  $\boldsymbol{\beta}^* \in \mathbb{R}^{n_{\text{dim}}}$ , that specify the relationship between the feature and the target value of each problem such that  $y_j^* = \mathbf{x}_j^T \boldsymbol{\beta}^*$ , for any  $j^1$ . At each time step, the environment generates a new problem  $j$ . An agent  $i$  can only observe the problem’s feature vector and must predict its target value,  $y_{i,j}$ . We assume agents make rule-based predictions. Specifically, each agent possesses estimations of the rule parameters (beta estimation),  $\boldsymbol{\beta}_i$ , and uses it to generate a prediction,  $y_{i,j} = \mathbf{x}_j^T \boldsymbol{\beta}_i$ . Individual learning occurs when they update their beta estimations after getting feedback from the prediction error,  $y_j^* - y_{i,j}$ . Here we simply assume the gradient descent:  $\boldsymbol{\beta}_i \leftarrow \boldsymbol{\beta}_i - \alpha_i (-2) \mathbf{x}_j^T (y_j^* - y_{i,j})$ , where  $\alpha_i$  is a learning parameter. Together, our setting provides a minimal model of an environment producing a constant flow of similar but distinct situations, and agents learning to behave adaptively in those situations.

<sup>1</sup> For simplicity, here we assume that the true generative rule of the environment is a linear combination of features, and agents’

primary task is to determine the specific parameters of this combination.

**Cohort Structure of Agents** The agents are not permanent fixtures within the environment but instead exit after a certain number of timesteps,  $T_{\text{lifespan}}$ , which we posit is too short for individual agents to independently develop optimal beta estimations (Fig. 1C). Agents experience shared lifespans with others born at the same timestep, forming a generation that consists of  $N_{\text{generation}}$  individual agents. A new generation of agents is introduced at every  $T_{\text{interval}}$ , timesteps. Crucially, when  $T_{\text{interval}} < T_{\text{lifespan}}$ , two or more generations partially overlap, meaning that agents at different ages of their lifespans coexist in the environment at certain timestep.

**Five Conformity Types** The lifetimes of agents are neither entirely synchronous nor completely separate, allowing each agent to observe and potentially be influenced by others. We model this situation in a similar way to the rerating paradigm. At each timestep, agents can observe other agents and possibly revise their predictions and beta estimations in light of those of the agents they referred to (Figure 1A). A crucial aspect of our simulation is examining various nuances of conformist social influence operating in this revision process. First, we categorize the depths of conformity in three levels:

- **No Conformity (“None”)**: Agents do not adjust their predictions or beta estimations in reference to others.
- **Shallow Conformity**: Agents conform only to the observable behaviors of other agents. This is implemented by adjusting the agent’s ( $i$ ) prediction to align with the prediction of the referred agent ( $k$ ):  $y_i \leftarrow y_i + s_i \times (y_k - y_i)$ , where  $s_i$  is a sensitivity parameter that determines how much the agent’s prediction shifts towards that of the referred agent.
- **Deep Conformity**: Beyond conforming in terms of predictions for a specific problem ( $y_{i,j}$ ), agents also adjust their internal models (i.e., beta estimations) to match those of the referred agents:  $\beta_i \leftarrow \beta_i + s_i \times (\beta_k - \beta_i)^2$ . This deeper influence generalizes across problems because the agent uses the updated beta estimations to generate predictions for new problems in subsequent timesteps.

Additionally, in line with the transmission bias in cultural evolution literature (e.g., Mesoudi, 2011), we varied whether agents distinguish relative expertise between the target agents and themselves, in both shallow and deep scenarios. Specifically, we assume agents select the agents they refer to either from the entire population (i.e., all the living agents other than themselves) or exclusively from those who are older than themselves. This results in conformity working either *symmetrically* or *asymmetrically* among agents.

In sum, we explored five different types of conformity, each in a separate simulation (Figure 1D).

**Other Simulation Details** The parameters of the true generative rule, agents’ beta estimations, and feature values of each problem were independently initialized using an  $n_{\text{dim}}$ -dimensional vector, with each feature value sampled from a uniform distribution between 0 and 1. For the main simulation analysis (Figure 4), parameters were set to  $n_{\text{dim}} = 40$ ,  $T_{\text{lifespan}} = 20$ ,  $T_{\text{interval}} = 10$ ,  $N_{\text{generation}} = 10$ , and  $\alpha_i = 0.05$  and  $s_i = 0.3$  for all agents. We confirmed that the results remained qualitatively robust when varying each parameter within broader reasonable ranges<sup>3</sup>.

**Outcome Measures** We focus on two measures to assess the collective performance of agents in each generation. First is the mean lifetime performance of individual agents. At each timestep (i.e., for each new specific problem), each agent’s performance is scored by the negative squared error between their prediction and the true target value. Lifetime performance for an agent is the accumulated sum of this score:  $\sum_{j=t}^{t+T_{\text{lifespan}}-1} -(y_j^* - y_{i,j})^2$ , where  $t$  is the timestep when the agent was born. Higher (closer to zero) lifetime performance indicates the agent behaving more adaptively in the environment on average across their lifespan. A generation’s mean performance is calculated by averaging the lifetime performance of all the agents in the generation.

Another measure is based on the errors in agents’ beta estimations ( $\beta_i$ ) compared to the true generative rule of the environment ( $\beta^*$ ), defined by the Euclidian distance in the  $n_{\text{dim}}$ -dimensional space. At each time step, agents individually update their estimations by locally adjusting parameters after observing prediction errors using gradient descent. Furthermore, when an agent ( $i$ ) conforms “deeply” to another agent ( $k$ ),  $i$  observes  $k$ ’s beta estimations and adjusts their own estimations accordingly. Lower beta gaps indicate that the agent possesses better estimates of the environment. We calculate the average learning curve of the agents (i.e., the decrease in the beta gap over their lifetime) for each generation and investigate whether within-individual improvements differ across generations in a way that allows later generations to learn more quickly by leveraging the learning of earlier generations.

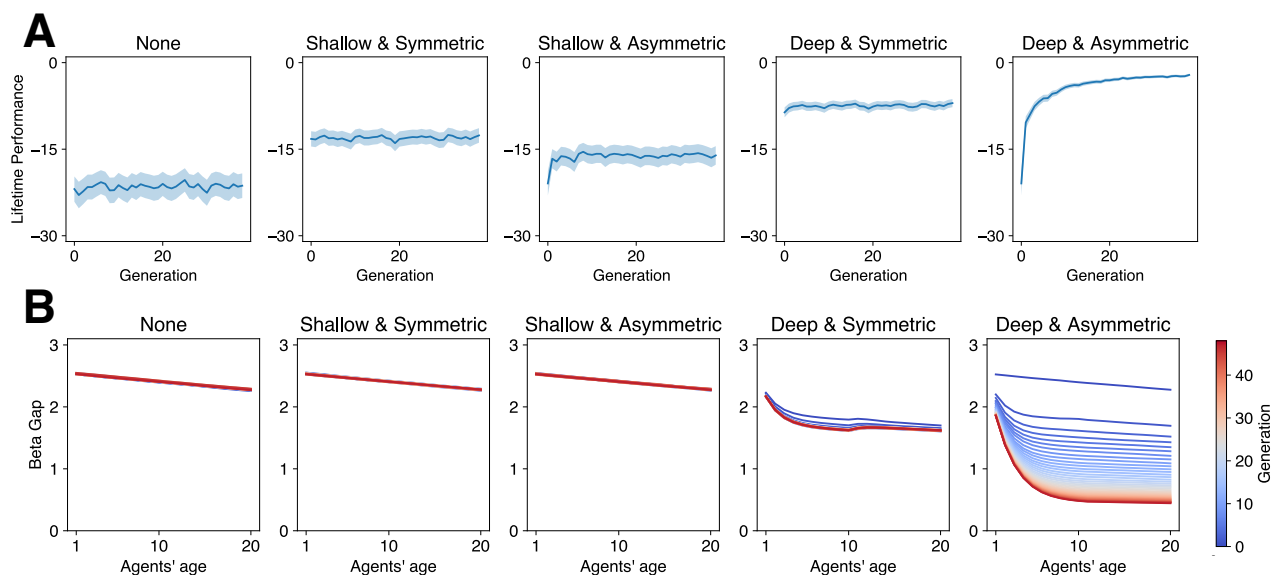
## Results

Figure 2A shows the agents’ mean lifetime performance as a function of their generations, under different types of conformity. The solid lines indicate the mean lifetime performance of each generation while the shaded areas indicate the standard deviations. Comparing “None” to other types of conformity, the presence of any type of conformity leads to better performance. This could be attributed to noise reduction, where adapting predictions to those of other agents reduces random noise and improves the accuracy of predictions. Under “Shallow” conformity, we find that

<sup>2</sup> Here, we assume that agents can observe the internal model of the target agents in each timestep, which may not be realistic. Although we believe it is necessary to assume that learning about others’ beta estimation is at least easier than learning about the

environment’s true beta, how much we can relax this assumption is an interesting future direction.

<sup>3</sup> The data and code to implement all analysis can be accessed at <https://github.com/hiz-repo/depth-of-conformity-exp>.



**Figure 2:** **A:** Lifetime performance of agents as a function of their generations. **B:** Improvements in beta estimations across lifetime of agents in each generation. Blue/red curves indicates agents in earlier/later generations. In both A and B, panels correspond to different types of conformity and results are averaged over 50 independent runs.

“Symmetric” influence outperforms “Asymmetric” influence. This could be because a greater number of agents find their reference agents in the “Symmetric” condition, which simply leads to a better collective performance when the primary benefit of conformity is noise reduction. However, under shallow conformity, the positive effects do not carry over to better beta estimations, by definition. Fig. 2B shows how each generation of agents learns to reduce errors in beta estimations throughout their lifespans, under different types of conformity. Without any conformist social influence (“None”), agents across all generations follow identical learning curves. The “Shallow” conformity results in the identical pattern: As the impact of social learning does not extend to beta estimations, each generation of agents has to re-learn them solely through their individual experiences.

“Deep” conformity differs from “Shallow” conformity in terms of both lifetime performance and beta estimations. First, “Deep” conformity generally yielded better lifetime performance than “Shallow” influence irrespective of whether the direction of conformity was symmetric or asymmetric, suggesting that the deeper social influence had an additional advantage beyond noise reduction. This improvement is explained by the differential learning curves in beta estimations (Fig. 2B): Under “Deep” conformity, agents began to cumulate better beta estimations across generations. With “Symmetric” influence, however, the cumulative advantage is limited to the earliest generations, possibly because adverse effects from the less to the more experienced agents hinder further improvements in later generations. The learning curves of beta estimations consistently improved from earlier to later generations only when conformity was both “Deep” and “Asymmetric” (Fig. 2B, rightmost panel). Consequently, an accumulative

improvement in lifetime performance across generations was observed exclusively under the “Deep” and “Asymmetric” conformity conditions. The lifetime performance of each generation successively improved until it converged to a near-optimal performance of zero (Fig. 2A, rightmost panel). These results underscore the importance of deep (and asymmetric) influence that can be generalized to future new problems in accumulating performance improvements beyond mere noise reduction.

Overall, we have shown that, using a minimal model where overlapping generations of agents learn to behave adaptively in a vast environment, the conformity of different depths shapes distinct collective dynamics. Specifically, cumulative improvements across generations can be achieved only through the combination of deep conformity that reaches the generative model of an agent and asymmetric influence from the more to the less experienced ones.

## Behavioral Experiment

In the following behavioral experiment, we aim to determine whether, and if so how, human participants generalize conformist social influence beyond direct reference. We have adapted the rerating paradigm, presenting participants with a series of problems that are similar yet distinct within a well-defined feature space. Critically, participants are exposed to others’ behaviors in only a subset of the problems. This selective exposure allows us to discern how participants extend the influence of observed behaviors to other similar yet distinct problems without direct reference. Since we are interested in the role of expertise difference in modulating this difference, we also manipulate cover stories associated with social information.

## Method

**Protocol** We administered an in-lab experiment with the University of Tokyo undergraduates who received course credit for their participation. Forty-three students (17 female; age:  $M = 20.3$ ,  $SD = 0.88$ ) participated in the experiment.

The specific domain we employed in an experiment was to estimate the monthly rent of apartments. We confronted participants with twenty-four different apartments, each characterized by three numerical values: age, size, and distance to the nearest station. The task consists of three stages (Figure 3A). In the first stage, participants individually answered their predictions for every twenty-four apartments. In the second stage, for half (twelve) of the apartments randomly chosen from the twenty-four used in the first stage, participants were told to guess “others’ predictions”, followed by feedback. We hereafter call these twelve apartments presented in the second stage as “with reference” apartments. The other twelve apartments were not addressed in the second stage (“without-reference” apartments). Critically, this means that participants could directly observe others’ behaviors in only half of all the situations (apartments). In the third stage, participants could update their predictions of all the twenty-four apartments, including both “with reference” and “without reference” ones. The order of apartments presented in each stage and the division of with- and without-reference apartments were randomized across participants.

Additionally, we manipulated the credibility of social information in a between-participants design by altering the cover stories about “others” in the second stage. Twenty participants were assigned to the high-credibility condition, where social information was presented as the average answers from “senior students actually living in the area.” The other 23 participants were placed in the low-credibility condition, where the social information was attributed to mere “other university students.” To check the manipulation’s effectiveness, we asked participants to assess the accuracy of others’ estimations relative to themselves on a seven-point scale in the postexperimental questionnaire.

**Materials** We retrieved data from existing one-room apartments whose nearest station is Hongo-Sanchome Station in Tokyo. The dataset of twenty-four apartments varied in age (4 - 60 years), size (12 - 38.5 square meters), walking distance from the station (1 - 9 minutes), and monthly rent (59 - 145 thousand JPY). In each stage, participants provided rent estimations on a slider from 0 to 300 thousand JPY.

For the second stage predictions, we used data from a pilot study where four senior University of Tokyo undergraduates, who had house-hunting experiences in the area, estimated the rents for these apartments. To capture the average patterns of their predictions, we used predictive values from linear regression with age, size, and distance as predictors and rent as the outcome. In the analysis, we standardized feature values by dividing them by their respective maximums. The resulting standardized coefficients of the linear model for social information were -0.14 for age, 0.48 for size, -0.057 for distance, and 8.6 for intercept. Importantly, this model can

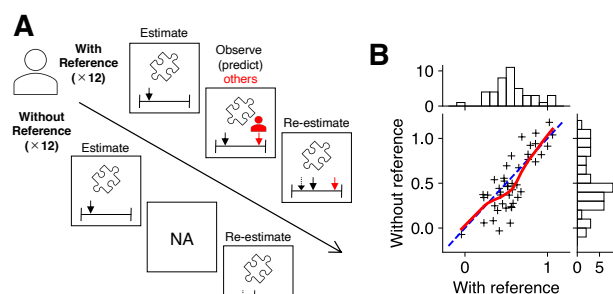
pinpoint social information for all twenty-four apartments, including the twelve apartments where participants did not directly observe the information.

**Analysis** The primary variable of interest is the sensitivity to social influence, denoted as  $s_{i,j}$  for participant  $i$  estimating apartment  $j$ . The sensitivity is defined as  $s_{i,j} = \frac{E_{i,j}^2 - E_{i,j}^1}{M_j - E_{i,j}^1}$ , where  $E_{i,j}^1$  and  $E_{i,j}^2$  are participant  $i$ ’s first and revised estimations of apartment  $j$ , and  $M_j$  is the social information for that apartment. A sensitivity score of 0 means the participant repeated their initial estimation without any conformity, whereas a score of 1 indicates complete adaptation to the social information. Note that  $s$  can become arbitrarily large in its magnitude if the initial estimate and the social information are very close. Therefore, in the analysis, we limited its range to  $-1 \leq s \leq 2$ , which corresponds to a symmetrical range surrounding  $0 \leq s \leq 1$  and contained 91% of the experimental data. Within this range,  $s$  values greater than zero can be seen as expressions of conformity, where larger values indicate stronger influence.

In our experiment, the accuracy of others’ predictions is arbitrarily determined by the researcher and is not self-evident for participants. Therefore, our analysis will focus on *how* participants incorporate social information into their second estimates, that is the degree of sensitivity, rather than looking at the accuracy of their results.

## Results

Do participants show any evidence of conformity for problems with which they did not directly observe social information (i.e., “without-reference” apartments)? If they do so at all, how is the conformity pattern different from conformity for problems with direct reference?



**Figure 3:** **A:** Schematic diagram of the experimental setup. **B:** Histograms and scatter plot comparing participants’ mean sensitivity for problems with (x-axis) and without (y-axis) direct reference to social information (black crosses: individual participants; blue dotted line: the diagonal [i.e., equal sensitivity for both types of problems]; red line: the locally weighted scatterplot smoothing [LOWESS] curve).

To analyze the determinants of sensitivity to social influence, we employed a multi-level regression with random intercepts for participants and problems. The dependent

variable was sensitivity, and independent variables included the presence or absence of direct social information, the type of cover story for social information, and their interaction. We found a significantly positive main effect of direct reference to social information ( $\beta = 0.08, p = .040$ ) and a significantly positive intercept ( $\beta = 0.47, p < .001$ ), suggesting that, while direct reference to social information led to stronger conformity, its influence extended to situations where such information was not observable. We did not observe a significant main effect of the cover stories ( $\beta = 0.07, p = .415$ ) or their interaction with direct reference ( $\beta = -0.04, p = .449$ ), indicating that the type of cover story did not significantly influence the degree of conformity. According to the responses to the postexperimental questionnaire, there was no significant difference between the high- and low-credibility conditions in terms of participants' subjective ratings of the accuracy of the social source ( $t(39.9) = 0.96, p = .345$ ). This means that our manipulation of cover stories did not sufficiently shift participants' subjective credibility toward social information. This result is plausible in hindsight: as subjective credibility should be determined by the relative expertise between oneself and others, and given substantial variations in participants' own expertise, manipulation of the cover story alone may not suffice.

Therefore, we have adjusted our approach to take the participants' subjective credibility ratings as an independent measure, replacing the cover story manipulation. In this revised analysis, we again observed a significantly positive main effect of direct reference to social information ( $\beta = 0.06, p = .035$ ) and a significantly positive intercept ( $\beta = 0.50, p < .001$ ), and this time, a significantly positive main effect of credibility ( $\beta = 0.14, p < .001$ ). We, again, did not find the interaction statistically significant ( $\beta = -0.01, p = .618$ ). Collectively, these results suggest that participants not only conformed to the observed behaviors of others but also generalized the influence to situations where they did not directly observe others, albeit to a lesser extent. Moreover, the extent of this conformist influence was larger when they assigned higher subjective credibility to social sources.

Next, we examine the underlying mechanisms behind this "partial" generalization of conformist social influence. Figure 3B shows participants' mean sensitivities for problems with (x-axis) and without (y-axis) reference to others. We highlight several observations. First, most scatters are located along the diagonal line. There was a strong positive rank correlation in mean sensitivities for with- and without-reference problems (Spearman's  $\rho = .72, p < .001$ ), suggesting consistent individual sensitivity levels across problem types, with inter-participant variability. Second, sensitivity for with-reference problems clustered around 0.5, displaying a unimodal distribution, while sensitivity for without-reference problems has more variance, forming a mode below 0.5. Third, the noticeable within-participant drop in sensitivity for without-reference problems (scatters located below the diagonal) primarily came from participants whose with-reference sensitivities are close to

0.5. The red LOWESS curve, capturing the average trend, dips below the diagonal for only those whose initial sensitivities were around 0.5.

These patterns may align well with the idea that two different depths of conformity coexist among participants: One is shallow behavioral alignment in which people utilize others' estimates as additional samples for noise reduction. This mode may be signified by the sensitivity value of 0.5 (i.e., fair averaging) with direct-reference problems and leads to reduced consistency in situations without reference. The other is a deeper level of influence in which people regard social information as a source for a better understanding of the feature-to-prediction mappings, leading to the generalization of conformity using these mapping rules across situations beyond direct reference.

## Discussion

We have examined the generalization property of conformity from two perspectives. On one hand, we formally demonstrated the idea that for agents with overlapping lifetimes in a vast environment to accumulate performance improvements longer than individual lifespans, they need deep and asymmetric conformity—adopting influences at the level of generative function (i.e., mapping between situational features and behaviors) only in the direction from the more to the less experienced agents.

Motivated by this theoretical observation, we conducted a behavioral experiment and found that people do display deep conformity, extending the social influence to situations where direct reference is unavailable. This finding is connected to the literature on function learning, in which people learn a continuous relationship between input and output variables from the training data presented as "truth." However, in our study, social information was not necessarily accurate; in fact, participants did not uniformly adopt social information even when they could directly refer to it. The results suggest that, even in such situations, the influence can be deep enough to generalize to new (similar but distinct) situations.

Interestingly, those who were simply averaging the answers while direct reference was available showed a notable failure in this generalization (i.e., reduced sensitivity for without-reference problems). We conjecture that this suggests that two different depths of conformity coexist among individuals. What determines which mode operates or how we might flexibly arbitrate between the two (Wu, Vélez, & Cushman, 2022) would be an interesting future question.

One important limitation of our study is the limited empirical strength due to a small sample size ( $N = 43$ ) and the exploratory nature of our analysis. As a result, our conclusions are not definitive, and further research is necessary to draw any empirical conclusions. Despite this, our study presents an interesting hypothesis that humans may extend conformist social influence beyond direct reference to others, and such deep (and asymmetric) conformity could underpin the accumulation of improvements longer than individual lifespans, a necessary condition for the cumulation of complex culture or technologies.

## Acknowledgments

We would like to thank Hye-rin Kim and Mayu Takahashi for helpful discussions. This work was supported by the Japan Society for the Promotion of Science (grant no. JP16H06324 to T.K. and no. JP23KJ0781 to R.M).

## References

- Asch, S. E. (1956). Studies of independence and conformity: I. A minority of one against a unanimous majority. *Psychological Monographs: General and Applied*, 70(9), 1–70.
- Baron, R. S., Kerr, N. L., & Miller, N. (1992). *Group process, group decision, group action*. Thomson Brooks/Cole Publishing Co.
- Hung, A. A., & Plott, C. R. (2001). Information Cascades: Replication and an Extension to Majority Rule and Conformity-Rewarding Institutions. *The American Economic Review*, 91(5), 1508–1520.
- Izuma, K., & Adolphs, R. (2013). Social manipulation of preference in the human brain. *Neuron*, 78(3), 563–573.
- Jayles, B., Kim, H.-R., Escobedo, R., Cezera, S., Blanchet, A., Kameda, T., Sire, C., & Theraulaz, G. (2017). How social information can improve estimation accuracy in human groups. *Proceedings of the National Academy of Sciences of the United States of America*, 114(47), 12620–12625.
- Kameda, T., & Nakanishi, D. (2002). Cost–benefit analysis of social/cultural learning in a nonstationary uncertain environment. An evolutionary simulation and an experiment with human subjects. *Evolution and Human Behavior*, 23(5), 373–393.
- Kameda, T., Toyokawa, W., & Tindale, R. S. (2022). Information aggregation and collective intelligence beyond the wisdom of crowds. *Nature Reviews Psychology*, 1(6), 345–357.
- Klucharev, V., Hytönen, K., Rijpkema, M., Smidts, A., & Fernández, G. (2009). Reinforcement learning signal predicts social conformity. *Neuron*, 61(1), 140–151.
- Kuroda, K., Ogura, Y., Ogawa, A., Tamei, T., Ikeda, K., & Kameda, T. (2022). Behavioral and neuro-cognitive bases for emergence of norms and socially shared realities via dynamic interaction. *Communications Biology*, 5(1), 1379.
- Lala, K. (2017). *Darwin's Unfinished Symphony*. Princeton University Press.
- List, C., Elsholtz, C., & Seeley, T. D. (2009). Independence and interdependence in collective decision making: an agent-based model of nest-site choice by honeybee swarms. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 364(1518), 755–762.
- Mesoudi, A. (2011). An experimental comparison of human social learning strategies: payoff-biased social learning is adaptive but underused. *Evolution and Human Behavior*, 32(5), 334–342.
- Molleman, L., Kurvers, R. H. J. M., & van den Bos, W. (2019). Unleashing the BEAST: a brief measure of human social information use. *Evolution and Human Behavior*, 40(5), 492–499.
- Nook, E. C., Ong, D. C., Morelli, S. A., Mitchell, J. P., & Zaki, J. (2016). Prosocial conformity: Prosocial norms generalize across behavior and empathy. *Personality and Social Psychology Bulletin*, 42(8), 1045–1062.
- Schultze, T., Rakotoarisoa, A.-F., & Stefan, S.-H. (2015). Effects of distance between initial estimates and advice on advice utilization. *Judgment and Decision Making*, 10(2), 144–171.
- Suzuki, S., Jensen, E. L. S., Bossaerts, P., & O’Doherty, J. P. (2016). Behavioral contagion during learning about another agent’s risk-preferences acts on the neural representation of decision-risk. *Proceedings of the National Academy of Sciences of the United States of America*, 113(14), 3755–3760.
- Toelch, U., Bruce, M. J., Newson, L., Richerson, P. J., & Reader, S. M. (2014). Individual consistency and flexibility in human social information use. *Proceedings of the Royal Society B: Biological Sciences*, 281(1776), 20132864.
- Toelch, U., & Dolan, R. J. (2015). Informational and Normative Influences in Conformity from a Neurocomputational Perspective. *Trends in Cognitive Sciences*, 19(10), 579–589.
- Wu, C. M., Vélez, N., Cushman, F. A., Dezza, I. C., & Schulz, E. (2022). Representational Exchange in Human Social Learning. In *The Drive for Knowledge: The Science of Human Information Seeking*. Cambridge University Press.