

Interventions reducing affective polarization do not necessarily improve anti-democratic attitudes

Received: 13 July 2021

Accepted: 16 September 2022

Published online: 31 October 2022

 Check for updates

Jan G. Voelkel¹✉, James Chu², Michael N. Stagnaro³, Joseph S. Mernyk¹, Chrystal Redekopp¹, Sophia L. Pink⁴, James N. Druckman⁵, David G. Rand³ and Robb Willer¹✉

There is widespread concern that rising affective polarization—particularly dislike for outpartisans—exacerbates Americans’ anti-democratic attitudes. Accordingly, scholars and practitioners alike have invested great effort in developing depolarization interventions that reduce affective polarization. Critically, however, it remains unclear whether these interventions reduce anti-democratic attitudes, or only change sentiments towards outpartisans. Here we address this question with experimental tests (total $n = 8,385$) of three previously established depolarization interventions: correcting misperceptions of outpartisans, priming inter-partisan friendships and observing warm cross-partisan interactions between political leaders. While these depolarization interventions reliably reduced affective polarization, we do not find compelling evidence that these interventions reduced support for undemocratic candidates, support for partisan violence or prioritizing partisan ends over democratic means. Thus, future efforts to strengthen pro-democratic attitudes may do better if they target these outcomes directly. More broadly, these findings call into question the previously assumed causal effect of affective polarization on anti-democratic attitudes.

Affective polarization—the tendency of partisans to view opposing partisans negatively and co-partisans positively¹—has been a major focus of research in recent years^{2,3}. In particular, research shows that contemporary US politics is characterized by growing affective polarization^{4,5}. Notably, not only academics but also a large majority of Americans believe the country is extremely divided⁶ and view this division as a serious problem⁷.

There is great concern about rising affective polarization in part because its presumed negative consequences may be uniquely harmful or destabilizing for democratic societies, for example, by stimulating support for undemocratic candidates and practices, or by fomenting

political violence (for example, refs. ^{2,8–13}). In light of the presumed dire consequences of affective polarization, academics and practitioners have invested substantial energy in developing depolarization interventions that reduce affective polarization, typically using outcomes based on sentiment towards opposing partisans (for example, refs. ^{12,14–25}). This body of work has uncovered numerous effective approaches for reducing affective polarization, tools that offer hope for maintaining—or restoring—democratic norms and practices.

Unfortunately, this hope may be premature. This is because much prior work has focused on treating affective polarization itself, and assumed that these interventions would in turn improve downstream

¹Department of Sociology, Stanford University, Stanford, CA, USA. ²Department of Sociology, Columbia University, New York, NY, USA. ³Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA, USA. ⁴Wharton School, University of Pennsylvania, Philadelphia, PA, USA.

⁵Department of Political Science, Northwestern University, Evanston, IL, USA. ✉e-mail: jvoelkel@stanford.edu; willer@stanford.edu

outcomes that pose consequential threats to democracy². Although this assumption may seem reasonable, there is little evidence evaluating its implications for the benefits of depolarization interventions. Here we shed light on the question of whether previously established depolarization interventions have the hoped-for consequence of effectively reducing anti-democratic attitudes.

Researchers who study depolarization interventions frequently propose that reducing affective polarization will indeed have positive effects on democratic outcomes. In line with this, 10 of the 12 papers on depolarization interventions we identified in the literature discuss anticipated effects on democratic outcomes (for quotes, see Supplementary Table 1). Recent review papers^{2,26} also discuss improving democratic outcomes as a goal of depolarization interventions.

The reason depolarization interventions are expected to reduce anti-democratic attitudes is that many researchers assume that affective polarization causes greater anti-democratic attitudes. The first link in this causal chain—the effect of depolarization interventions on affective polarization—is well supported empirically in the published literature. The second link—a causal effect of affective polarization on anti-democratic attitudes—has been suggested by many researchers (for example, refs. ^{2,3,27–31}), and there is little evidence that the relationship between affective polarization and anti-democratic attitudes is contested. In fact, with the exception of ref. ³², we did not find any work suggesting a lack of a causal effect of affective polarization on anti-democratic attitudes.

There are several potential mechanisms through which reducing affective polarization could reduce anti-democratic attitudes. First, reducing affective polarization could reduce the perceived threat of the negative consequences when the outparty wins an election that, in turn, could reduce support for undemocratic inparty candidates. Second, reducing affective polarization could reduce identification with the inparty that, in turn, could decrease the desire to break democratic norms to win at all costs. Third, reducing affective polarization could increase empathy for outpartisans, making it more difficult to justify violence against them.

However, while there are many good reasons to expect that depolarization interventions reduce anti-democratic attitudes, so far there is little empirical evidence in favour of this hypothesis. One paper, using cross-sectional data, found that affective polarization is negatively correlated with support for democratic norms²⁸. Conversely, a recent paper found that manipulating affective polarization had no significant causal effect on accountability, attitudes about democratic norms, or support for partisan violence³². In short, it remains unclear whether commonly used depolarization interventions would have the hoped-for consequence of reducing anti-democratic attitudes.

In this article, we shed light on this question by testing the impact of three previously validated depolarization interventions on a variety of anti-democratic attitudes. In doing so, we advance the literature on affective polarization in three ways. First, we assess the robustness of previous findings by testing the interventions' effects on standard sentiment measure of affective polarization. Second, we extend previous findings by testing the interventions' effects on incentivized behavioural measures of affective polarization to address concerns that standard sentiment measures of affective polarization are inconsequential partisan signalling, with no impact on interpersonal behaviours (for example, refs. ^{1,33,34}).

Third, and most importantly, we test whether the effects of these interventions extend beyond affective polarization to impact three measures of anti-democratic attitudes: support for undemocratic candidates, support for partisan violence, and prioritizing partisan ends over democratic means. The outcomes we study were based on previous research on anti-democratic attitudes^{18,35,36} and chosen because of their relevance for contemporary US politics. If results show that existing depolarization interventions reduce these more consequential outcomes, this would suggest that depolarization researchers and

practitioners should extend their current work towards maximizing the effectiveness of current interventions. However, if we find that existing depolarization interventions do not reduce anti-democratic attitudes, this would suggest that reducing anti-democratic attitudes requires the development of new interventions and direct measurement of anti-democratic attitudes.

Results

We examined whether existing depolarization interventions reduce not only affective polarization but also anti-democratic attitudes in three large-scale experiments. Studies 1 ($n = 2,341$) and 3 ($n = 4,023$) were conducted on non-probability samples that were similar to the national population on several key demographics. Study 2 ($n = 2,021$) was conducted on a convenience sample. Studies 2 and 3 followed pre-registered analysis scripts.

Correlational results

We begin by assessing the correlations between affective polarization and anti-democratic attitudes. In our three studies, we measured affective polarization with the canonical feeling thermometer item that indicates how cold participants felt towards outpartisans. In addition, we measured two behavioural indicators of dislike for outpartisans in Study 1. First, we measured how much money participants would give to an outpartisan in a dictator game (reverse coded). Dictator games were used as a behavioural indicator in a seminal paper on affective polarization¹. Second, we measured how much money participants would spend to take money away from an outpartisan in a 'joy of destruction' game³⁷. This behavioural indicator captures how much people are willing to personally sacrifice to reduce the earnings of outpartisans. Affective polarization was weakly to moderately correlated with withholding money in a dictator game ($r = 0.34, P < 0.001$) and weakly correlated with spending money in a joy-of-destruction game ($r = 0.07, P < 0.001$).

Finally, we measured several anti-democratic attitudes. We measured support for undemocratic candidates—a willingness to sacrifice democratic principles in electoral contexts for inparty victories—with items such as 'How likely would you be to vote for the [Democratic/Republican] candidate if you learned that they support a proposal to reduce the number of polling stations in areas that support the [Republican/Democratic] party?'³⁵. This approach follows work showing partisans are often willing to violate democratic norms (for example, regarding electoral fairness, checks and balances, and civil liberties) to win elections³⁵. We also measured support for partisan violence, another key facet of undemocratic attitudes, with items such as 'How much do you feel it is justified for [Democrats/Republicans] to use violence in advancing their political goals these days?'³⁶. Finally, we measured prioritizing partisan ends over democratic means—a willingness to help the inparty at the expense of the country and/or in contravention of democratic norms—with items such as '[Democrats/Republicans] should redraw districts to maximize their potential to win more seats in federal elections, even if it may be technically illegal'¹⁸.

Previous research considers such anti-democratic attitudes to be direct consequences of affective polarization, leading to assumptions that interventions that reduce affective polarization should also reduce these anti-democratic attitudes (for example, ref. ²). Accordingly, we expected that affective polarization would be correlated with anti-democratic attitudes. However, across our three studies, affective polarization was not reliably correlated with support for undemocratic politicians, support for partisan violence, or prioritizing partisan ends. Affective polarization was positively correlated with support for undemocratic politicians in studies 1 ($r = 0.19, P < 0.001, P_{\text{Holm-Bonferroni adjustment (HBA)}} < 0.001$) and 2 ($r = 0.26, P < 0.001, P_{\text{HBA}} < 0.001$) but negatively correlated in study 3 ($r = -0.05, P = 0.001, P_{\text{HBA}} = 0.001$). Affective polarization was positively correlated with prioritizing partisan ends over democratic means in studies 1 ($r = 0.05, P = 0.011, P_{\text{HBA}} = 0.011$) and 2 ($r = 0.16, P < 0.001, P_{\text{HBA}} < 0.001$) but negatively correlated in study 3

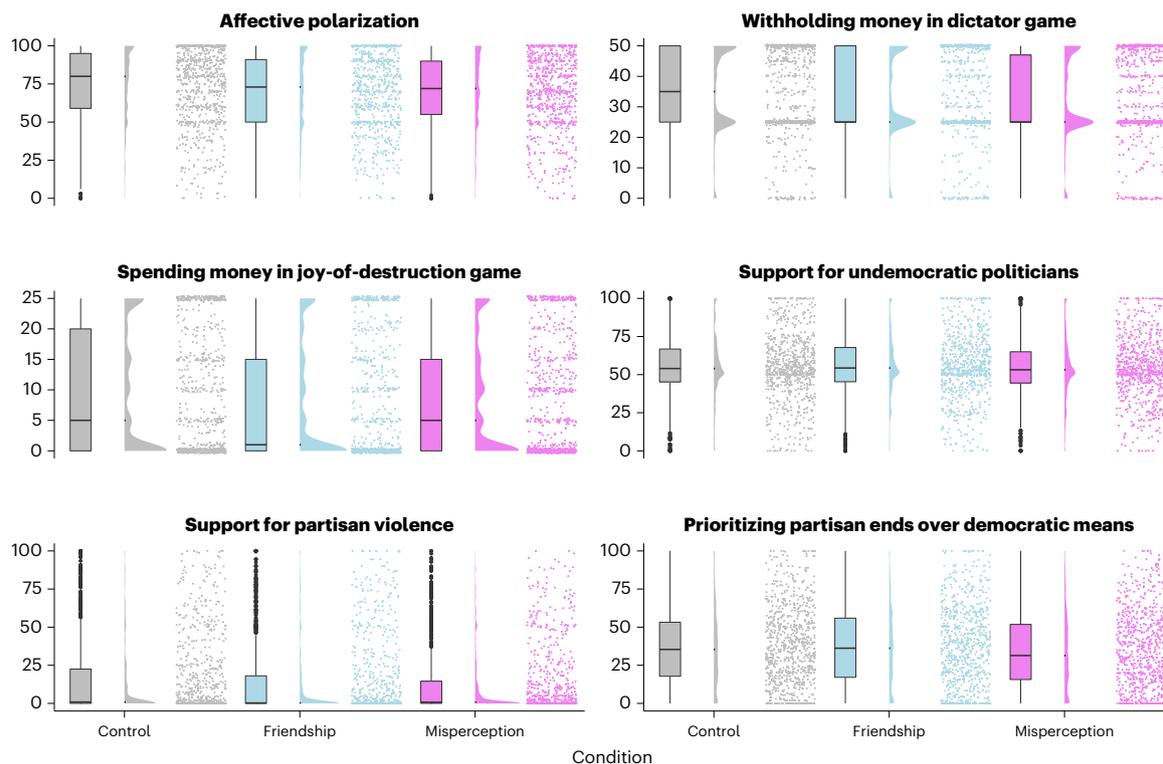


Fig. 1 | Effects of the friendship intervention and the misperception correction intervention on affective polarization and anti-democratic attitudes, study 1. For each condition and outcome, the figure shows a box plot (left), a halved violin plot (middle) and a point cloud (right). The box of the box plot shows the 25th percentile, the median and the 75th percentile. The length of the whiskers is $1.5 \times$ interquartile range (IQR) unless the minimum/maximum falls within $1.5 \times$ IQR of the quartiles. For affective polarization and

withholding money in dictator game: $n_C = 835$, $n_F = 754$, $n_M = 752$; spending money in joy-of-destruction game: $n_C = 829$, $n_F = 751$, $n_M = 751$; support for undemocratic politicians: $n_C = 806$, $n_F = 736$, $n_M = 736$; support for partisan violence: $n_C = 812$, $n_F = 744$, $n_M = 740$; prioritizing partisan ends over democratic means: $n_C = 833$, $n_F = 751$, $n_M = 750$. Please note that the two economic games used different ranges (0–25 and 0–50) from the other dependent variables (0–100).

($r = -0.23$, $P < 0.001$, $P_{\text{HBA}} < 0.001$). Affective polarization was negatively correlated with support for partisan violence in studies 1 ($r = -0.22$, $P < 0.001$, $P_{\text{HBA}} < 0.001$) and 3 ($r = -0.36$, $P < 0.001$, $P_{\text{HBA}} < 0.001$) and uncorrelated in study 2 ($r = 0.02$, $P = 0.463$, $P_{\text{HBA}} = 0.463$). However, because the observed correlations could be suppressed by unobserved variables, we conducted experiments to estimate the causal effect of depolarization interventions on anti-democratic attitudes.

Experimental results

We now turn to our main focus, testing the causal effects of promising depolarization interventions on both (1) measures of affective polarization and (2) measures of anti-democratic attitudes. We chose a varied set of interventions that we perceived to be particularly promising ways to reduce affective polarization.

In study 1, we tested two recently proposed interventions. The first intervention made an outparty friendship salient¹², following the logic that thinking of a friend who supports the opposing party will generate more positive and/or less threatening feelings about the other party. The second intervention corrected exaggerated misperceptions about the extent of outparty opposition to inparty attempts to pass a policy^(16,19); see¹⁴ for another similar intervention). This misperception correction intervention makes clear that the other party is less of a threat to the agenda of the participant's favoured party as if often believed. Both interventions were compared with a control group. Afterwards, we measured the variables described above in the correlational analyses: affective polarization, withholding money from an outpartisan in a dictator game, spending money to take money away from an outpartisan in a joy-of-destruction game, support for undemocratic candidates,

support for partisan violence, and prioritizing partisan ends over democratic means. The main findings are shown in Fig. 1.

First, both interventions reduced affective polarization as measured by cold feelings towards outpartisans. Participants in the friendship intervention condition (mean 69.9, standard deviation (s.d.) 24.7) and in the misperception correction intervention condition (mean 70.3, s.d. 24.0) reported significantly lower levels of affective polarization than participants in the control condition (mean 74.1, s.d. 24.3) (for the friendship intervention condition: $b = -3.86$, standard error (s.e.) = 1.20, $t(2,326) = -3.20$, $P = 0.001$, Cohen's $d = -0.16$, 95% confidence interval (CI) for $b = (-6.22, -1.49)$; for the misperception correction intervention condition: $b = -3.66$, s.e. = 1.20, $t(2,326) = -3.04$, $P = 0.002$, Cohen's $d = -0.15$, 95% CI for $b = (-6.02, -1.30)$).

Second, both interventions also reduced behavioural indicators of affective polarization. Participants in the friendship intervention condition (mean 32.9, s.d. 14.2) and in the misperception correction intervention condition (mean 32.3, s.d. 13.6) withheld significantly less money from an outpartisan in a dictator game than participants in the control condition (mean 35.0, s.d. 14.2) (for the friendship intervention: $b = -2.14$, s.e. = 0.70, $t(2,326) = -3.07$, $P = 0.002$, $P_{\text{HBA}} = 0.004$, Cohen's $d = -0.15$, 95% CI for $b = (-3.51, -0.77)$; for the misperception correction intervention: $b = -2.75$, s.e. = 0.70, $t(2,326) = -3.94$, $P < 0.001$, $P_{\text{HBA}} < 0.001$, Cohen's $d = -0.20$, 95% CI for $b = (-4.12, -1.38)$). Participants in the friendship intervention condition (mean 8.0, s.d. 9.8) also spent significantly less money to take money away from an outpartisan in a joy-of-destruction game than participants in the control condition (mean 9.1, s.d. 10.1) ($b = -1.12$, s.e. = 0.48, $t(2,316) = -2.32$, $P = 0.020$, $P_{\text{HBA}} = 0.020$, Cohen's $d = -0.11$, 95% CI for $b = (-2.07, -0.17)$).

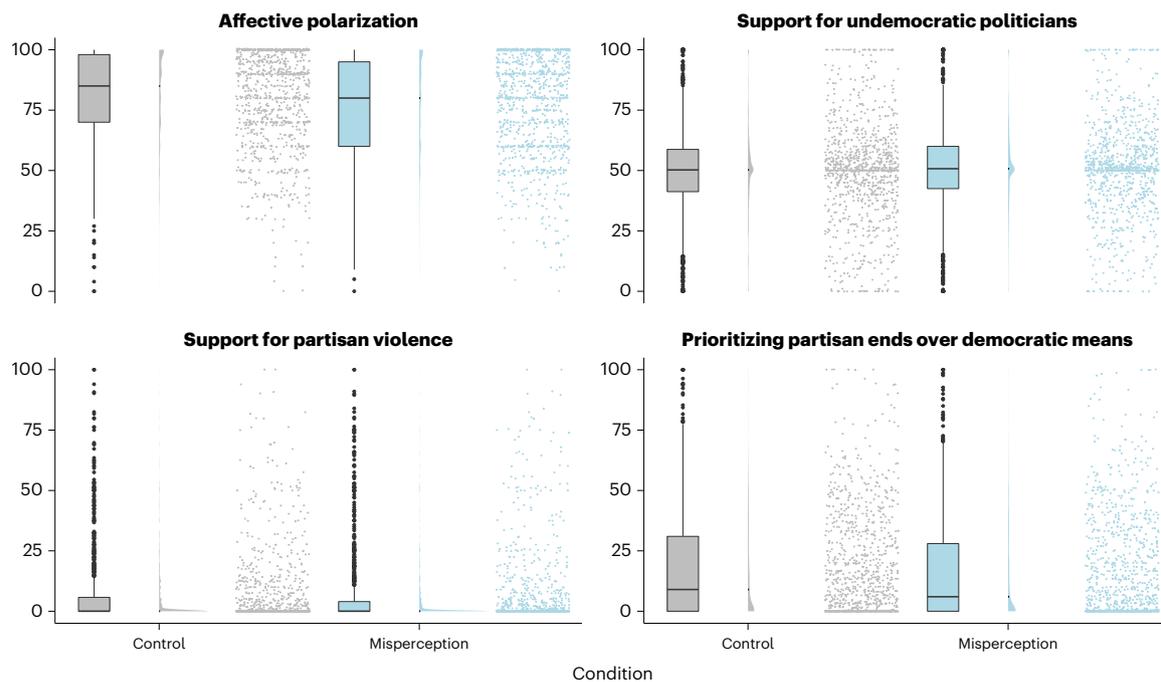


Fig. 2 | Effects of the misperception correction intervention on affective polarization and anti-democratic attitudes, study 2. For each condition and outcome, the figure shows a box plot (left), a halved violin plot (middle) and a point cloud (right). The box of the box plot shows the 25th percentile, the median and the 75th percentile. The length of the whiskers is $1.5 \times \text{IQR}$ unless

the minimum/maximum fall within $1.5 \times \text{IQR}$ of the quartiles. For affective polarization, $n_C = 1,016$, $n_M = 1,005$; support for undemocratic politicians: $n_C = 1,010$, $n_M = 1,002$; support for partisan violence: $n_C = 1,012$, $n_M = 1,004$; prioritizing partisan ends over democratic means: $n_C = 1,013$, $n_M = 1,005$.

The misperception correction intervention (mean 8.3, s.d. 9.4) did not significantly reduce spending in the joy-of-destruction game ($b = -0.74$, s.e. = 0.48, $t(2,316) = -1.53$, $P = 0.127$, $P_{\text{HBA}} = 0.127$, Cohen's $d = -0.08$, 95% CI for $b = (-1.68, 0.21)$).

Critically, however, neither intervention significantly reduced any of the downstream measures. Participants in the friendship intervention condition (mean 55.9, s.d. 21.1) and in the misperception correction intervention condition (mean 54.2, s.d. 20.6) did not report significantly less support for undemocratic candidates than participants in the control condition (mean 55.2, s.d. 21.9) (for the friendship intervention: $b = 0.47$, s.e. = 0.99, $t(2,263) = 0.48$, $P = 0.631$, $P_{\text{HBA}} = 1$, Cohen's $d = 0.02$, 95% CI for $b = (-1.46, 2.41)$); for the misperception correction intervention: $b = -0.75$, s.e. = 0.98, $t(2,263) = -0.76$, $P = 0.447$, $P_{\text{HBA}} = 0.447$, Cohen's $d = -0.04$, 95% CI for $b = (-2.68, 1.18)$). Participants in the friendship intervention condition (mean 15.1, s.d. 25.7) and in the misperception correction intervention condition (mean 13.2, s.d. 22.8) did not report significantly less support for partisan violence than participants in the control condition (mean 15.6, s.d. 25.1) (for the friendship intervention: $b = -0.84$, s.e. = 1.15, $t(2,281) = -0.73$, $P = 0.468$, $P_{\text{HBA}} = 1$, Cohen's $d = -0.03$, 95% CI for $b = (-3.09, 1.42)$); for the misperception correction intervention: $b = -1.52$, s.e. = 1.15, $t(2,281) = -1.32$, $P = 0.185$, $P_{\text{HBA}} = 0.436$, Cohen's $d = -0.06$, 95% CI for $b = (-3.78, 0.73)$). Participants in the friendship intervention condition (mean 37.2, s.d. 25.8) and in the misperception correction intervention condition (mean 35.0, s.d. 24.9) did not report significantly lower levels of prioritizing partisan ends over democratic means than participants in the control condition (mean 37.3, s.d. 25.3) (for the friendship intervention: $b = -0.22$, s.e. = 1.18, $t(2,319) = -0.19$, $P = 0.851$, $P_{\text{HBA}} = 1$, Cohen's $d = -0.01$, 95% CI for $b = (-2.53, 2.09)$); for the misperception correction intervention: $b = -1.71$, s.e. = 1.18, $t(2,319) = -1.46$, $P = 0.145$, $P_{\text{HBA}} = 0.436$, Cohen's $d = -0.07$, 95% CI for $b = (-4.02, 0.59)$).

Bayesian analyses provided further evidence that neither intervention affected any of the downstream measures. Whereas non-significant

P values in null hypothesis significance testing cannot be interpreted as evidence in favour of null effects, the Bayes factor quantifies the relative predictive performance of the null hypothesis of no difference relative to the alternative hypothesis³⁸. Bayesian analyses require a choice of the prior of the relative plausibility of the two hypotheses before looking at the data. Here we used the prior that the null hypothesis and alternative hypothesis are equally likely to be true. Such a relatively high prior for the alternative hypothesis is justified because we and other researchers believed that depolarization interventions would reduce anti-democratic attitudes. The Bayes factor describes to what extent the data warrant a change in the relative plausibility of the two hypotheses from this prior³⁸.

We found strong evidence for null effects of the friendship intervention on support for undemocratic candidates ($\text{BF}_{01} = 14.06$), support for partisan violence ($\text{BF}_{01} = 16.43$), and prioritizing partisan ends over democratic means ($\text{BF}_{01} = 17.66$). We found strong evidence for a null effect of the misperception correction intervention on support for undemocratic candidates ($\text{BF}_{01} = 11.26$). The data were also more consistent with a null effect of the misperception correction intervention on support for partisan violence ($\text{BF}_{01} = 2.88$) and prioritizing partisan ends over democratic means ($\text{BF}_{01} = 3.53$), but this evidence was weak and moderate, respectively. Taken together, the results from study 1 show that depolarization interventions can reduce both attitudinal and behavioural indicators of affective polarization without necessarily reducing anti-democratic attitudes. This calls into question the commonly held assumption that anti-democratic attitudes are downstream consequences of affective polarization.

In study 2 we sought to address a limitation of study 1: the lack of full randomization of the order of the dependent variables. In study 1, support for partisan violence and support for undemocratic candidates were always presented after the feeling thermometer, the two behavioural measures of affective polarization and prioritizing partisan ends over democratic means. Thus, the lack of effects on

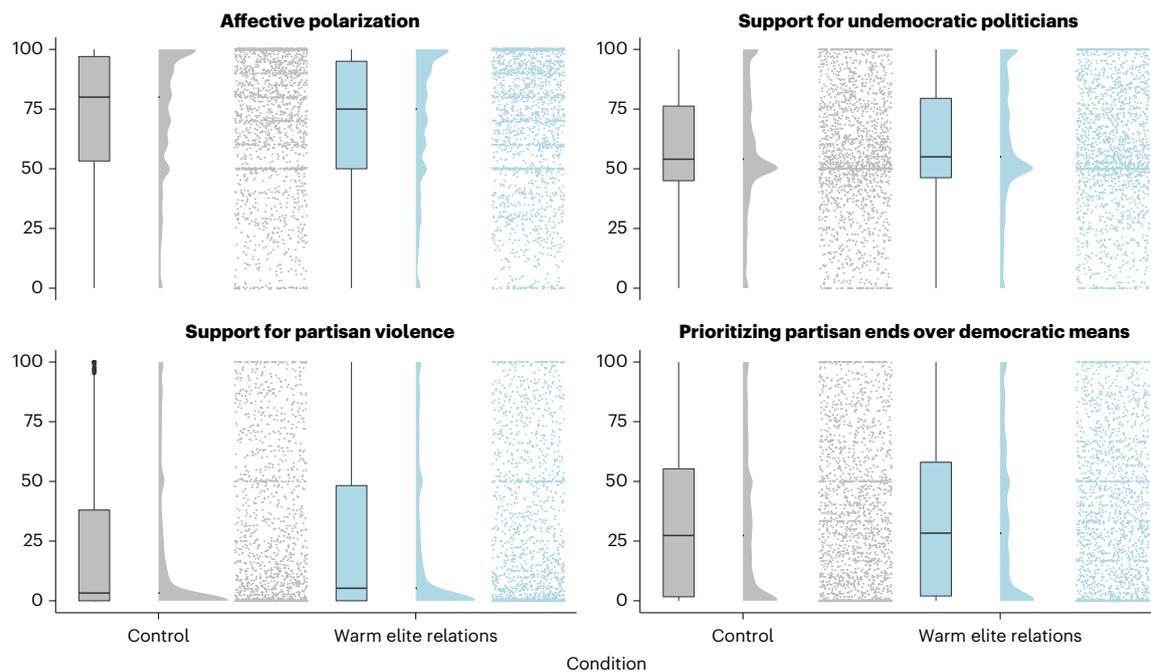


Fig. 3 | Effects of the warm elite relations intervention on affective polarization and anti-democratic attitudes, study 3. For each condition and outcome, the figure shows a box plot (left), a halved violin plot (middle) and a point cloud (right). The box of the box plot shows the 25th percentile, the median, and the 75th percentile. The length of the whiskers is $1.5 \times$ IQR unless

the minimum/maximum fall within $1.5 \times$ IQR of the quartiles. For affective polarization, $n_c = 1,998$, $n_w = 2,025$; support for undemocratic politicians: $n_c = 1,962$, $n_w = 2,006$; support for partisan violence: $n_c = 1,993$, $n_w = 2,017$; prioritizing partisan ends over democratic means: $n_c = 1,990$, $n_w = 2,016$.

support for partisan violence and support for undemocratic candidates could have been because the effect of the treatment decreased over time (owing to participant fatigue or something else). In study 2, we fully randomized the order of four dependent variables: the feeling thermometer, support for undemocratic candidates, support for partisan violence, and prioritizing partisan ends over democratic means. We focused on the comparison of the control condition and the misperception correction intervention, as this intervention showed larger (yet statistically non-significant) effects on anti-democratic attitudes than the friendship intervention in study 1. We did not include the behavioural measures of affective polarization in study 2, given their secondary interest.

The results of study 2 were similar to the results of study 1 (Fig. 2). Once again, participants in the misperception correction intervention condition (mean 76.7, s.d. 21.0) reported significantly lower levels of affective polarization than participants in the control condition (mean 80.2, s.d. 19.5) ($b = -2.90$, s.e. = 0.88, $t(2,007) = -3.29$, $P = 0.001$, Cohen's $d = -0.14$, 95% CI for $b = (-4.63, -1.17)$). However, the misperception correction intervention did not significantly reduce any of the negative downstream outcomes. Participants in the misperception correction intervention condition (mean 51.0, s.d. 19.4) did not report significantly less support for undemocratic candidates than participants in the control condition (mean 50.2, s.d. 20.5) ($b = 0.90$, s.e. = 0.86, $t(1,998) = 1.05$, $P = 0.295$, $P_{\text{HBA}} = 0.591$, Cohen's $d = 0.04$, 95% CI for $b = (-0.78, 2.57)$). Participants in the misperception correction intervention condition (mean 6.7, s.d. 16.0) did not report significantly less support for partisan violence than participants in the control condition (mean 7.1, s.d. 15.4) ($b = -0.37$, s.e. = 0.69, $t(2,002) = -0.54$, $P = 0.590$, $P_{\text{HBA}} = 0.591$, Cohen's $d = -0.02$, 95% CI for $b = (-1.71, 0.98)$). Participants in the misperception correction intervention condition (mean 16.1, s.d. 21.1) did not report significantly lower levels of prioritizing partisan ends over democratic means than participants in the control condition (mean 17.7, s.d. 21.4) ($b = -1.42$, s.e. = 0.93, $t(2,004) = -1.53$, $P = 0.127$, $P_{\text{HBA}} = 0.381$, Cohen's $d = -0.07$, 95% CI for $b = (-3.24, 0.40)$).

Bayesian analyses provided further evidence that the misperception correction intervention did not affect the downstream measures. We found strong evidence for a null effect of the misperception correction intervention on support for undemocratic candidates ($\text{BF}_{01} = 13.03$) and on support for partisan violence ($\text{BF}_{01} = 18.13$). The data were also consistent with a null effect of the misperception correction intervention on prioritizing partisan ends over democratic means, but this evidence is only moderately strong: $\text{BF}_{01} = 4.68$). Taken together, the results from study 2 replicate the finding that interventions reducing affective polarization do not necessarily reduce anti-democratic attitudes.

In study 3, we tested an elite-focused intervention. While the friendship and misperception correction interventions tested in studies 1 and 2 used content about outpartisan voters, study 3 tested a warm elite relations intervention highlighting the friendship between a Republican politician (John McCain) and a Democratic politician (Joe Biden). This intervention is similar to the warm elite relations treatments used in previous research¹⁵ but avoids deception.

The results of study 3 again support the idea that interventions reduce affective polarization but do not reduce anti-democratic attitudes (Fig. 3). Participants in the warm elite relations intervention condition (mean 69.6, s.d. 28.5) reported significantly lower levels of affective polarization than participants in the control condition (mean 72.0, s.d. 27.6) ($b = -2.04$, s.e. = 0.83, $t(4,008) = -2.45$, $P = 0.014$, Cohen's $d = -0.07$, 95% CI for $b = (-3.67, -0.40)$).

However, the warm elite relations intervention did not significantly reduce any of the measures of anti-democratic attitudes. On the contrary, participants in the warm elite relations intervention condition (mean 58.8, s.d. 26.0) reported significantly more support for undemocratic candidates than participants in the control condition (mean 57.4, s.d. 25.6) ($b = 1.49$, s.e. = 0.73, $t(3,954) = 2.04$, $P = 0.042$, $P_{\text{HBA}} = 0.083$, Cohen's $d = 0.06$, 95% CI for $b = (0.06, 2.93)$); although this result was not robust to a Holm–Bonferroni adjustment for multiple hypothesis testing. Participants in the warm elite relations intervention

condition (mean 24.4, s.d. 32.1) also reported significantly more support for partisan violence than participants in the control condition (mean 22.0, s.d. 31.0) ($b = 2.20$, s.e. = 0.88, $t(3,996) = 2.51$, $P = 0.012$, $P_{\text{HBA}} = 0.036$, Cohen's $d = 0.07$, 95% CI for $b = (0.48, 3.92)$). Participants in the warm elite relations intervention condition (mean 34.4, s.d. 31.8) did not report significantly lower levels of prioritizing partisan ends over democratic means than participants in the control condition (mean 33.4, s.d. 31.5) ($b = 0.76$, s.e. = 0.90, $t(3,991) = 0.84$, $P = 0.401$, $P_{\text{HBA}} = 0.401$, Cohen's $d = 0.02$, 95% CI for $b = (-1.01, 2.53)$).

Bayesian analyses provided further evidence that the warm elite relations intervention did not decrease the downstream measures. We found weak to moderate evidence for a null effect of the warm elite relations intervention on support for undemocratic candidates ($\text{BF}_{01} = 6.18$) and on support for partisan violence ($\text{BF}_{01} = 1.52$). Note that, even if these effects were not null, the direction of these effects was positive, not negative. We found strong evidence for a null effect of the warm elite relations intervention on prioritizing partisan ends over democratic means ($\text{BF}_{01} = 17.76$). Taken together, the results of study 3 provide further support that interventions that reduce affective polarization do not reduce anti-democratic attitudes. We even found some evidence that the warm elite relations intervention can increase anti-democratic attitudes.

Meta-analyses

To provide tests of the effect of the depolarization interventions on anti-democratic attitudes with the largest possible sample size, we conducted random-effects meta-analyses³⁹ using the effect sizes from the three studies and from two additional pilot tests we conducted (described in more detail in the supplementary information on pilot studies 1 and 2). According to power analyses conducted with metapower⁴⁰, these tests had approximately 80% power to detect $|\text{Cohen's } d| \geq 0.06$ for support for undemocratic candidates, $|\text{Cohen's } d| \geq 0.08$ for support for partisan violence, and $|\text{Cohen's } d| \geq 0.07$ for prioritizing partisan ends over democratic means. The differences in the power analyses result from using the observed heterogeneity which was different for the different dependent variables.

The meta-analytic results support the idea that the tested interventions reliably reduced affective polarization. Participants in the intervention conditions reported significantly lower levels of affective polarization than participants in the control conditions (Cohen's $d = -0.13$, s.e. = 0.02, $z = -5.77$, $P < 0.001$, 95% CI for Cohen's $d = (-0.17, -0.09)$). However, we did not find evidence that the tested interventions reduced anti-democratic attitudes. Specifically, the depolarization interventions did not significantly reduce support for undemocratic candidates (Cohen's $d = 0.03$, s.e. = 0.02, $z = 1.28$, $P = 0.200$, 95% CI for Cohen's $d = (-0.01, 0.06)$), nor support for partisan violence (Cohen's $d = -0.01$, s.e. = 0.03, $z = -0.30$, $P = 0.767$, 95% CI for Cohen's $d = (-0.06, 0.05)$), nor prioritizing partisan ends over democratic means (Cohen's $d = -0.03$, s.e. = 0.02, $z = -1.13$, $P = 0.259$, 95% CI for Cohen's $d = (-0.07, 0.02)$).

These results do not rule out that depolarization interventions may have small effects on anti-democratic attitudes. As the effect of the interventions on affective polarization is Cohen's $d = -0.13$ and our power analysis suggested that we have 80% power to detect $|\text{Cohen's } d| \geq 0.06$ –0.08, we were powered to detect downstream effects of an approximately 2:1 ratio. This means that we cannot rule out that there are downstream effects via affective polarization with a larger ratio. However, our results suggest that depolarization interventions usually only modestly reduce affective polarization itself. Thus, even if there were downstream effects via affective polarization of a 3:1 or 4:1 ratio, the downstream effects of depolarization interventions would be very small.

Nonetheless, our best estimate is that the effect of depolarization interventions on anti-democratic is essentially null. This is partly based on the estimated effect sizes. In addition, Bayesian random-effects meta-analyses found strong evidence for a null effect of the depolarization interventions on support for undemocratic candidates

($\text{BF}_{01} = 24.35$), support for partisan violence ($\text{BF}_{01} = 24.59$) and prioritizing partisan ends over democratic means ($\text{BF}_{01} = 16.89$). Finally, the results of instrumental variable analyses were also consistent with the conclusion that reducing affective polarization does not reduce anti-democratic attitudes (Supplementary Table 21).

Possible mediators and moderators

These results beg the question, where does the causal chain from depolarization interventions to anti-democratic attitudes break? In additional exploratory analyses, we did not find evidence that any of the tested interventions reduced potential mediators for the link between affective polarization and anti-democratic attitudes (for example, strength of inparty identification, and empathy for outpartisans; Supplementary Table 6). We also failed to find evidence for reliable moderation effects by party identity (Republican versus Democrat; Supplementary Tables 7–13) or strength of inparty identification (Supplementary Tables 14–16).

Discussion

Our findings call into question whether depolarization interventions developed to reduce affective polarization also reduce anti-democratic attitudes. Across three experiments, we successfully replicate previous research, finding that three depolarization interventions reliably reduced self-reported affective polarization. We also extend past work by showing that these interventions also impact behavioural indicators of affective polarization, thereby demonstrating that depolarization interventions can influence behaviours with real, monetary stakes for outpartisans. Critically, however, the depolarization interventions did not reliably reduce any of three measures of anti-democratic attitudes: support for undemocratic candidates, support for partisan violence, and prioritizing partisan ends over democratic means. Even the correlational associations between affective polarization and these anti-democratic attitudes were not reliable. Thus, we conclude that many researchers (including ourselves, for example, refs.^{2,22}) may have substantially overestimated the effects of depolarization interventions on anti-democratic tendencies.

Our paper has several important limitations. First, we focused on estimating the effects of depolarization interventions on anti-democratic attitudes. The observed null effects do not imply that depolarization interventions cannot have effects on other important measures, such as economic, social or romantic discrimination against outpartisans. For example, we find that depolarization interventions can increase giving in a dictator game and reduce spending in a joy-of-destruction game (see also refs.^{32,41,42}). Depolarization interventions may also have other political effects. For example, previous research has found that reducing affective polarization increases willingness to compromise¹². Another limitation of our paper is that our studies do not use probability samples and, thus, are not truly representative of the US population. However, two of our three studies were conducted on samples that were similar to the national population on several key demographics. Another limitation is that we tested only a subset of existing depolarization interventions (although we tested interventions that involved both voters and elites). Future researchers should test additional depolarization interventions to identify what kind of interventions—if any—simultaneously reduce affective polarization and anti-democratic attitudes. Future researchers should also think carefully about what kind of affective polarization they want to measure (see also ref.⁴³). It could be that the joy-of-destruction game is better suited to measure hate while feeling thermometers and the dictator game capture dislike. A final limitation of our paper is that we tested only a limited set of potential mediators and moderators. Future research should examine if other constructs (for example, anti-establishment orientations⁴⁴) moderate the effects of depolarization interventions and test which mediators interventions need to move to reduce anti-democratic attitudes.

Our findings are important because they replace an old assumption, that depolarization interventions will reduce anti-democratic attitudes (for example, refs. ^{2,18,21,22,25}), with an empirical-based default, that depolarization interventions do not reduce anti-democratic attitudes. That is, researchers and practitioners who are concerned about anti-democratic attitudes should not presume that treating affective polarization will impact those outcomes. Instead, they should see affective polarization and anti-democratic attitudes as two separate classes of outcomes that require distinct interventions.

More generally, our results suggest that research on depolarization interventions should avoid assuming downstream consequences on these interventions and instead measure these potential downstream consequences directly. That is, researchers and practitioners who are interested in interventions targeting anti-democratic attitudes such as support for undemocratic politicians, support for partisan violence, and prioritizing partisan ends over democratic means should not focus on treating affective polarization and begin developing more direct interventions—trends that run counter to most current work.

From a broader theoretical perspective, our results raise serious questions about whether a causal link from affective polarization to anti-democratic attitudes actually exists. Future research is needed to examine whether variables that have been identified as leading to affective polarization^{45–48} also affect anti-democratic attitudes. Future research should also examine whether anti-democratic attitudes may affect affective polarization. For example, testing whether interventions designed to reduce anti-democratic attitudes also reduce affective polarization would provide insights into whether there is no causal relationship between the two constructs or whether anti-democratic attitudes actually cause affective polarization. For these reasons, identifying factors that can reduce anti-democratic attitudes should be a priority for future research (for example, ref. ⁴⁹).

Methods

Ethics statement and reproducibility

All studies were approved by the institutional review board at Stanford University. All participants provided informed consent and were compensated for their participation. Materials, anonymized data (including descriptions of how the original files were anonymized) and analysis code for both studies are available via <https://osf.io/n5u9d/>. The pre-registrations for studies 2 and 3 are available via <https://osf.io/a2ukg/> and <https://osf.io/rmhct/>.

Samples

All studies were conducted in Qualtrics. In study 1, we collected data from 2,341 participants who were recruited from an Internet panel provided by Bovitz between 28 October 2020 to 3 November 2020. Bovitz maintains an online panel of approximately one million respondents recruited through random digit dialling and empanelment of Americans with Internet access. In study 2, we collected data from 2,021 participants who were recruited from a large panel of previously recruited Amazon Mechanical Turk workers between 23 March 2021 and 5 April 2021. In study 3, we collected data from 4,023 participants from Lucid between 18 December 2021 and 2 January 2022. The samples in studies 1 and 3 were quota-matched so that they were similar to the national population on key demographics. All samples were recruited with soft quotas for participants' self-identified partisanship (including learners): 50% Democrats and 50% Republicans.

Participants were excluded if they had duplicate IDs (keeping only the first case), did not consent to participate, were underage or did not provide their age, failed attention checks (completed pre-treatment assignment), identified as neither Democrat nor Republican, or left the study before they were randomly assigned to a condition. We used pairwise deletion for missing data. According to power analyses conducted with G*Power⁵⁰, all studies had at least 80% power to detect even small effect sizes ($|Cohen's d| \geq 0.14$ in study 1, $|Cohen's d| \geq 0.12$

in study 2, and $|Cohen's d| \geq 0.09$ in study 3) and 95% power to detect $|Cohen's d| \geq 0.18$ in study 1, $|Cohen's d| \geq 0.16$ in study 2 and $|Cohen's d| \geq 0.11$ in study 3. More details on sample characteristics, exclusion rules and demographics are provided in Supplementary Tables 2 and 3.

Interventions

We chose interventions from the literature that we perceived as particularly promising ways of reducing affective polarization. Our focus was more on testing whether these interventions would reduce anti-democratic attitudes than on replicating the original effects as closely as possible. Therefore, our design differed in several ways from the original studies, including the control conditions, the dependent variables and the analyses conducted.

Participants were randomly assigned to one of three (study 1) or two (studies 2 and 3) conditions. The first intervention we tested was a friendship intervention¹². In this condition, which we included only in study 1, participants received the following instructions: 'Although you are [a Democrat/an Independent who is closer to the Democratic Party/an Independent who is closer to the Republican Party/a Republican], you probably know people who are [Republicans/Democrats]. Think about one such [Republican/Democrat] that you like and respect a great deal. This person could be a friend, relative, neighbour, co-worker, or just someone that you know. Please explain why you feel this way about this person.'

The second intervention we tested was a misperception correction intervention. The original paper¹⁶ tested two versions of this intervention, and we selected the hypocrisy prevention intervention because it had a descriptively bigger effect size (although not significantly different). In this condition, which was included in studies 1 and 2, participants were presented with a scenario, randomly chosen from five scenarios, for example, 'A state [Democratic/Republican] party in control of the state legislature has drafted a proposal to streamline the appointment of judges where judges would be nominated and voted on in groups, not individually. This would reduce the workload of state legislators and make the process more efficient, however it may make it more difficult for the party in the minority, the [Republicans/Democrats], to object to the appointment of individual judges' (all scenarios are available at <https://osf.io/n5u9d/>). Then participants rated how much they believed an outpartisan would (a) dislike and (b) oppose this action, and (c) find this action politically unacceptable. After partisans provided their own beliefs, their potential misperceptions were corrected by presenting the real responses from outpartisans and the real responses from inpartisans who had a read a similar scenario where outpartisans were taking the action. The real responses were based on a nationally representative survey¹⁶.

The third intervention we tested was a warm elites relation intervention. This intervention is similar to the warm elite relations treatments used by previous research¹⁵ but avoids deception. In this condition, which was included in study 3, participants were asked to watch a video about the friendship between the Democratic politician Joe Biden and the Republican politician John McCain. The video is available at shorturl.at/jrOP6.

We used two different control conditions. In the first two studies, we used a null control. That is, in the control condition, participants moved immediately towards the section with the measures of the dependent variable. In the third study, participants watched a video about the history of neckties.

Measures

We measured affective polarization using a feeling thermometer rating for outpartisans: 'We would like to get your feelings toward both Democrats and Republicans. We would like you to rate them using something we call the feeling thermometer. Ratings between 50 degrees and 100 degrees mean that you feel favourable and warm toward them. Ratings between 0 degrees and 50 degrees mean that you don't feel favourable

toward them and that you don't care too much for them. You would rate them at the 50 degree mark if you don't feel particularly warm or cold toward them.' We used the reversed-coded feeling thermometer towards outpartisans so that higher scores indicate colder feelings.

An alternative operationalization of affective polarization is the difference score of feelings towards inpartisans and outpartisans. We found that the interventions tested in studies 1 and 2 significantly reduced the difference score (all $P \leq 0.001$), and the intervention tested in study 3 reduced the difference score directionally but not significantly ($P = 0.136$). Detailed results are reported in Supplementary Table 4. We report the outpartisan feeling thermometer as our main measure in the manuscript for several reasons. First, the difference score could be decreased by (a) decreasing cold feelings towards outpartisans or (b) by increasing cold feelings towards inpartisans. As increasing cold feelings towards inpartisans is not clearly normatively desirable, we focused on decreasing cold feelings towards outpartisans. Second, previous work suggests that feelings towards outpartisans animosity is the major source of change in affective polarization over time⁵¹. Third, we pre-registered the outpartisan feeling thermometer as our main measure in studies 2 and 3. In study 3, we also examined whether the intervention would impact feelings towards outpartisan voters and politicians differently (Supplementary Table 5).

Withholding money in a dictator game was measured with the following item: 'You have been anonymously and randomly matched with another participant who identifies as a [Republican/Democrat]. You have just been given 50 cents. You will now decide how to split these 50 cents between yourself and the [Republican/Democratic] participant. You can give any amount between 0 cents and 50 cents to the other participant. The other participant cannot affect the outcome you choose. How many cents (if any) will you give to the [Republican/Democratic] participant?'. Withholding money in a dictator game is a behavioural measure of affective polarization used in a seminal article on the topic¹.

Spending money in a joy-of-destruction game was measured with the following item (based on ref. ³⁷): 'You have been anonymously and randomly matched with another participant who identifies as a [Republican/Democrat]. Both you and the other participant have just each been given 50 cents. You will now decide whether to leave the [Republican/Democratic] participant's payment unchanged or take away part or all of their 50 cents. For every 1 cent you pay, you remove 2 cents from the [Republican/Democratic] participant. You can pay any amount between 0 cents and 25 cents. The other participant cannot affect the outcome you choose. How many cents (if any) do you want to pay to remove the [Republican/Democratic] participant's earnings? (remember, for every 1 cent you pay, the [Republican/Democratic] participant will lose 2 cents)'. Although spending money in a joy-of-destruction game is not a traditional measure of affective polarization, we included it because it measures hate—a stronger form of animosity—rather than either the outpartisans feeling thermometer or withholding money in the dictator game, which capture dislike (for more information, see Supplementary Information on correlational statistics for the joy-of-destruction game).

Support for undemocratic candidates was measured on a 101-point scale ranging from 'extremely likely to vote for the [Republican/Democratic] candidate' to 'extremely likely to vote for the [Democratic/Republican] candidate' with the following items (based on 35): (1) How likely would you be to vote for the [Democratic/Republican] candidate if you learned that they said they would ignore unfavourable court rulings by [Republican/Democratic]-appointed judges?, (2) How likely would you be to vote for the [Democratic/Republican] candidate if you learned that they support a proposal to reduce the number of polling stations in areas that support the [Republican/Democratic] party?, (3) How likely would you be to vote for the [Democratic/Republican] candidate if you learned that they support a redistricting plan that gives [Democrats/Republicans] ten extra seats despite a decline in the polls?, (4) How likely would you be to vote for the [Democratic/

Republican] candidate if you learned that they said that [Democrats/Republicans] should not accept election results if they do not win?, (5) How likely would you be to vote for the [Democratic/Republican] candidate if you learned that they said they would prosecute journalists who accuse them of misconduct if the journalists won't reveal their sources?, and (6) How likely would you be to vote for the [Democratic/Republican] candidate if you learned that they said they would ban far-[right/left] group rallies on the state capital grounds?. Items 5 and 6 were not included in studies 2 and 3. The items formed a reliable scale in all studies (study 1: Cronbach's $\alpha = 0.90$; study 2: Cronbach's $\alpha = 0.88$; study 3: Cronbach's $\alpha = 0.91$).

Support for partisan violence was measured with the following four items:³⁶ (1) When, if ever, is it OK for [Democrats/Republicans] to send threatening and intimidating messages to [Republican/Democratic] party leaders?, (2) When, if ever, is it OK for an ordinary [Democrat/Republican] in the public to harass an ordinary [Republican/Democrat] on the Internet, in a way that makes the target feel frightened?, (3) How much do you feel it is justified for [Democrats/Republicans] to use violence in advancing their political goals these days?, and (4) How much do you feel it is justified for [Democrats/Republicans] to use violence if the [Republican/Democratic] party wins more races in the next election?. Items 1 and 2 used a 101-point scale ranging from 'never' to 'always', and items 3 and 4 used a 101-point scale ranging from 'not at all justified' to 'extremely justified'. The items formed a reliable scale in all studies (study 1: Cronbach's $\alpha = 0.95$; study 2: Cronbach's $\alpha = 0.93$; study 3: Cronbach's $\alpha = 0.96$).

Prioritizing partisan ends over democratic means was measured on a 101-point scale ranging from 'strongly disagree' to 'strongly agree' with the following items:¹⁸ (1) I think the [Democrats/Republicans] should do everything they can to hurt the [Republican/Democratic] party, even if it is at the short-term expense of the country, (2) It's OK to sacrifice US economic prosperity in the short term in order to hurt [Republicans'/Democrats'] chances in future elections, (3) [Democrats/Republicans] should redraw districts to maximize their potential to win more seats in federal elections, even if it may be technically illegal, (4) If [Democrats/Republicans] gain control of all branches of government, they should use the Federal Communications Commission to heavily restrict or shut down [Fox News/MSNBC] to stop the spread of fake news, and (5) I think the [Democrats/Republicans] should do everything in their power within the law to make it as difficult as possible for [Trump to run the government effectively/Democrats to take part in governing the country]. Items 4 and 5 were not included in studies 2 and 3. The items formed a reliable scale in all studies (study 1: Cronbach's $\alpha = 0.84$; study 2: Cronbach's $\alpha = 0.85$; study 3: Cronbach's $\alpha = 0.91$). Please note that, while we refer to this variable in our pre-registration as support for undemocratic practices, we think that labelling the scale prioritizing partisan ends over democratic means is more accurate. Descriptive statistics for these measures are discussed in the supplementary information on descriptive statistics. We also included additional dependent variables that are not relevant for answering the research questions of this paper. The questionnaires for all studies are available via <https://osf.io/n5u9d/>.

Analysis strategy

Our main analysis strategy (pre-registered in studies 2 and 3) was null hypothesis significance testing. For the (non-pre-registered) correlational analyses, we used the Pearson correlation coefficient. For the experimental analyses, we used linear regression analyses, controlling for participants' gender, age, race, education, partisan identity and strength of partisan identity. We used P values from two-tailed tests as our inference criteria. In addition, we conducted a robustness check (pre-registered in study 3) using the Holm–Bonferroni adjustment for multiple hypothesis testing when we used multiple dependent variables for a construct (such as the three measures of anti-democratic attitudes).

For the non-significant effects, we conducted exploratory Bayesian analyses to estimate the strength of the evidence in favour of the null hypothesis. We conducted Bayesian analyses of co-variance in JASP including the same control variables as described above. We used JASP's default settings for priors that the null hypothesis and alternative hypothesis are equally likely to be true (for robustness checks with different priors, see Supplementary Figs. 1–12).

For the meta-analysis, we conducted frequentist random-effects meta-analyses using the R package *metafor*³⁹. Bayesian random-effects meta-analyses were conducted via JASP. We used seven effect sizes from the three studies and from two additional pilot tests we conducted (described in more detail in the supplementary information on pilot studies 1 and 2). As some of these effect sizes rely on comparisons with the same control condition, we conducted robustness checks accounting for this dependency. These robustness checks provided converging results (Supplementary Table 17).

Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data Availability

The data for our studies are openly available via <https://osf.io/n5u9d/>.

Code Availability

The analysis scripts for our studies are openly available via <https://osf.io/n5u9d/>.

References

1. Iyengar, S. & Westwood, S. J. Fear and loathing across party lines: new evidence on group polarization. *Am. J. Pol. Sci.* **59**, 690–707 (2015).
2. Finkel, E. J. et al. Political sectarianism in America: a poisonous cocktail of othering, aversion, and moralization. *Science* **370**, 533–536 (2020).
3. Iyengar, S., Lelkes, Y., Levendusky, M., Malhotra, N. & Westwood, S. J. The origins and consequences of affective polarization in the United States. *Annu. Rev. Polit. Sci.* **22**, 129–146 (2019).
4. Boxell, L., Gentzkow, M. & Shapiro, J. Cross-country trends in affective polarization. *Rev. Econ. Stat.* https://doi.org/10.1162/rest_a_01160 (2022).
5. Iyengar, S. & Krupenkin, M. The strengthening of partisan affect. *Polit. Psychol.* **39**, 201–218 (2018).
6. Partisan antipathy: more intense, more personal. *Pew Research Center* <https://www.pewresearch.org/politics/2019/10/10/the-partisan-landscape-and-views-of-the-parties/> (2019).
7. NBC News/Wall Street Journal Survey, Study #181259. *Hart Research Associates/Public Opinion Strategies* <http://wsj.com/public/resources/documents/181259NBCWSJOctober2018PollFinal.pdf> (2018).
8. Abramowitz, A. I. & Webster, S. The rise of negative partisanship and the nationalization of US elections in the 21st century. *Elect. Stud.* **41**, 12–22 (2016).
9. Diermeier, D. & Li, C. Partisan affect and elite polarization. *Am. Polit. Sci. Rev.* **113**, 277–281 (2019).
10. Hetherington, M. J. & Rudolph, T. J. *Why Washington Won't Work: Polarization, Political Trust, and the Governing Crisis* (Univ. Chicago Press, 2015).
11. Klein, E. *Why We're Polarized* (Simon and Schuster, 2020).
12. Levendusky, M. S. *Our Common Bonds: Using What Americans Share to Help Bridge the Partisan Divide*. Unpublished manuscript, Univ. Pennsylvania (2020).
13. Mason, L. *Uncivil Agreement: How Politics Became Our Identity* (Univ. Chicago Press, 2018).
14. Ahler, D. J. & Sood, G. The parties in our heads: misperceptions about party composition and their consequences. *J. Polit.* **80**, 964–981 (2018).
15. Huddy, L. & Yair, O. Reducing affective polarization: warm group relations or policy compromise? *Polit. Psychol.* **42**, 291–309 (2021).
16. Lees, J. & Cikara, M. Inaccurate group meta-perceptions drive negative out-group attributions in competitive contexts. *Nat. Hum. Behav.* **4**, 279–286 (2020).
17. Levendusky, M. S. Americans, not partisans: can priming American national identity reduce affective polarization? *J. Polit.* **80**, 59–70 (2018).
18. Moore-Berg, S. L., Ankori-Karlinsky, L. O., Hameiri, B. & Bruneau, E. Exaggerated meta-perceptions predict intergroup hostility between American political partisans. *Proc. Natl Acad. Sci. USA* **117**, 14864–14872 (2020).
19. Ruggeri, K. et al. The general fault in our fault lines. *Nat. Hum. Behav.* **5**, 1369–1380 (2021).
20. Simonsson, O. & Marks, J. Love thy (partisan) neighbor: brief befriending meditation reduces affective polarization. *Group Process Intergroup Relat.* **25**, 1577–1593 (2022).
21. Swanson, S. By the people: the role of local deliberative forums in combating affective political polarization. *The Project on International Peace and Security* https://www.wm.edu/offices/global-research/_documents/pips/selene-swanson-whitepaper (2021).
22. Voelkel, J. G., Ren, D. & Brandt, M. J. Inclusion reduces political prejudice. *J. Exp. Soc. Psychol.* **95**, 104149 (2021).
23. Warner, B. R., Horstman, H. K. & Kearney, C. C. Reducing political polarization through narrative writing. *J. Appl. Commun. Res.* **48**, 459–477 (2020).
24. Wojcieszak, M. & Warner, B. R. Can interparty contact reduce affective polarization? A systematic test of different forms of intergroup contact. *Polit. Commun.* **37**, 789–811 (2020).
25. Zoizner, A., Shenhav, S. R., Fogel-Dror, Y. & Sheaffer, T. Strategy news is good news: how journalistic coverage of politics reduces affective polarization. *Polit. Commun.* **38**, 604–623 (2021).
26. Hartman, R. et al. Interventions to reduce partisan animosity. *Nat. Hum. Behav.* **6**, 1194–1205 (2022).
27. Gidron, N., Adams, J. & Horne, W. *American Affective Polarization in Comparative Perspective* (Cambridge Univ. Press, 2020).
28. Kingzette, J. et al. How affective polarization undermines support for Democratic norms. *Public Opin. Q.* **85**, 663–677 (2021).
29. Lees, J. & Cikara, M. Understanding and combating misperceived polarization. *Philos. Trans. R. Soc. B* **376**, 20200143 (2021).
30. McCoy, J. & Sommer, M. Toward a theory of pernicious polarization and how it harms democracies: comparative evidence and possible remedies. *Ann. Am. Acad. Pol. Soc. Sci.* **681**, 234–271 (2019).
31. Orhan, Y. E. The relationship between affective polarization and democratic backsliding: comparative evidence. *Democratization* **29**, 714–735 (2022).
32. Broockman, D. E., Kalla, J. L. & Westwood, S. J. Does affective polarization undermine democratic norms or accountability? Maybe not. *Am. J. Pol. Sci.* <https://doi.org/10.1111/ajps.12719> (2022).
33. Carlin, R. E. & Love, G. J. Political competition, partisanship and interpersonal trust in electoral democracies. *Br. J. Polit. Sci.* **48**, 115–139 (2016).
34. Whitt, S. et al. Tribalism in America: behavioral experiments on affective polarization in the Trump era. *J. Exp. Political Sci.* **8**, 247–259 (2021).
35. Graham, M. H. & Svobik, M. W. Democracy in America? Partisanship, polarization, and the robustness of support for democracy in the United States. *Am. Polit. Sci. Rev.* **114**, 392–409 (2020).

36. Kalmoe, N. P. & Mason, L. Lethal mass partisanship: prevalence, correlates, and electoral contingencies. Preprint at National Capital Area Political Science Association American Politics Meeting https://www.dannyhayes.org/uploads/6/9/8/5/69858539/kalmoe__mason_ncapsa_2019_-_lethal_partisanship_-_final_lmedit.pdf (2019).
37. Abbink, K. & Sadrieh, A. The pleasure of being nasty. *Econ. Lett.* **105**, 306–308 (2009).
38. van Doorn, J. et al. The JASP guidelines for conducting and reporting a Bayesian analysis. *Psychonomic Bull. Rev.* **28**, 813–826 (2021).
39. Viechtbauer, W. Conducting meta-analyses in R with the metafor package. *J. Stat. Softw.* **36**, 1–48 (2010).
40. Griffin, J. W. Calculating statistical power for meta-analysis using metapower. *Quant. Method Psychol.* **17**, 24–39 (2021).
41. Gift, K. & Gift, T. Does politics influence hiring? Evidence from a randomized experiment. *Polit. Behav.* **37**, 653–675 (2015).
42. McConnell, C., Margalit, Y., Malhotra, N. & Levendusky, M. The economic consequences of partisanship in a polarized era. *Am. J. Pol. Sci.* **62**, 5–18 (2018).
43. Cassese, E. C. Partisan dehumanization in American politics. *Polit. Behav.* **43**, 29–50 (2021).
44. Uscinski, J. E. et al. American politics in two dimensions: partisan and ideological identities versus anti-establishment orientations. *Am. J. Pol. Sci.* **65**, 877–895 (2021).
45. Bougher, L. D. The correlates of discord: identity, issue alignment, and political hostility in polarized America. *Polit. Behav.* **39**, 731–762 (2017).
46. Mason, L. A cross-cutting calm: how social sorting drives affective polarization. *Public Opin. Q.* **80**, 351–377 (2016).
47. Santos, L. A., Voelkel, J. G., Willer, R. & Zaki, J. Belief in the utility of cross-partisan empathy reduces partisan animosity and facilitates political persuasion. *Psychol. Sci.* **33**, 1557–1573 (2022).
48. Simas, E. N., Clifford, S. & Kirkland, J. H. How empathic concern fuels political polarization. *Am. Polit. Sci. Rev.* **114**, 258–269 (2020).
49. Bartels, L. M. Ethnic antagonism erodes Republicans' commitment to democracy. *Proc. Natl Acad. Sci. USA* **117**, 22752–22759 (2020).
50. Faul, F., Erdfelder, E., Buchner, A. & Lang, A. G. Statistical power analyses using G*Power 3.1: tests for correlation and regression analyses. *Behav. Res. Methods* **41**, 1149–1160 (2009).
51. Groenendyk, E. Competing motives in a polarized electorate: political responsiveness, identity defensiveness, and the rise of partisan antipathy. *Polit. Psychol.* **39**, 159–171 (2018).

Acknowledgements

The authors received funding for this project from the Civic Health Project (R.W.), the Stanford Center on Philanthropy and Civil Society (R.W.) and the Institute for Policy Research, Northwestern University (J.N.D.). This work is supported under a Stanford Interdisciplinary Graduate Fellowship (J.G.V.). The funders had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript.

Author contributions

J.G.V., J.C., M.N.S., J.S.M., C.R., S.L.P., J.N.D., D.G.R. and R.W. designed the studies. J.G.V., J.C., M.N.S., J.S.M., C.R., S.L.P. and R.W. collected the data. J.G.V. analysed the data. J.G.V. and D.G.R. wrote the manuscript. J.C., M.N.S., J.N.D. and R.W. provided comments on the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41562-022-01466-9>.

Correspondence and requests for materials should be addressed to Jan G. Voelkel or Robb Willer.

Peer review information *Nature Human Behaviour* thanks Eric Groenendyk, Omer Yair and Magdalena Wojcieszak for their contribution to the peer review of this work. Peer reviewer reports are available.

Reprints and permissions information is available at www.nature.com/reprints.

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

Springer Nature or its licensor holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

© The Author(s), under exclusive licence to Springer Nature Limited 2022

Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Behavioural & social sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description	All studies were quantitative, survey-based between-subjects experiments that were conducted online.
Research sample	<p>Sample characteristics and demographics are available in Supplementary Tables 2 and 3. In all three studies, we aimed to collect approximately 50% Democrats and 50% Republicans. Our studies do not use probability samples and, thus, are not truly representative of the US population.</p> <p>Study 1: Participants were recruited via Bovitz. Bovitz maintains an online panel of approximately one million respondents recruited through random digit dialing and empanelment of Americans with Internet access. We restricted the sample to self-identified Democrats and Republicans (including leaners). The sample was quota-matched so that they were similar to the national population on key demographics.</p> <p>Study 2: Participants were recruited via a large panel of previously recruited Amazon Mechanical Turk workers. We restricted the sample to self-identified Democrats and Republicans (including leaners).</p> <p>Study 3: Participants were recruited via Lucid between December 18, 2021 and January 2, 2022. We restricted the sample to self-identified Democrats and Republicans (including leaners). The sample was quota-matched so that they were similar to the national population on key demographics.</p>
Sampling strategy	We used quota sampling. We determined our targeted sample sizes via a priori power analyses with G*Power. We aimed to be able to detect even small effect sizes with high power. All studies had at least 80% power to detect even small effect sizes ($ \text{Cohen's } d \geq 0.14$ in Study 1, $ \text{Cohen's } d \geq 0.12$ in Study 2, and $ \text{Cohen's } d \geq 0.09$ in Study 3) and 95% power to detect $ \text{Cohen's } d \geq 0.18$ in Study 1, $ \text{Cohen's } d \geq 0.16$ in Study 2, and $ \text{Cohen's } d \geq 0.11$ in Study 3.
Data collection	Participants completed a survey online via Qualtrics. The questionnaires we used are available via https://osf.io/n5u9d/ . Participants did not interact with the researchers and were blinded to experimental condition and the study hypotheses.
Timing	<p>Study 1: Data collection started on October 28, 2020 and ended on November 3, 2020.</p> <p>Study 2: Data collection started on March 23, 2021 and ended on April 5, 2021.</p> <p>Study 3: Data collection started on December 18, 2021 and ended on January 2, 2022.</p>
Data exclusions	Regarding pre-treatment exclusions: Supplementary Table 2 lists all exclusion criteria and the number of participants who were excluded for each of these criteria. The exclusion criteria were preregistered in Studies 2 and 3.
Non-participation	Regarding post-treatment attrition: As preregistered in Studies 2 and 3, we used pairwise deletion for all dependent variables. Missing values for control variables were replaced by adding an "unanswered" category for categorical variables and the median for continuous variables. The number of people who dropped out for the different dependent variables are provided in Supplementary Table 2.
Randomization	Participants were randomly assigned to experimental conditions.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Human research participants

Policy information about [studies involving human research participants](#)

Population characteristics

Sample characteristics and demographics for the three studies are described in Tables S2 and S3.

Recruitment

Study 1: Participants were recruited by Bovitz. Bovitz maintains an online panel of approximately one million respondents recruited through random digit dialing and empanelment of Americans with Internet access. Samples are drawn such that the demographics of the sample match those of the U.S. population. There might be a self-selection bias and a non-response bias, but it is not clear how this would affect the results.

Study 2: Participants were recruited via MTurk. There might be a self-selection bias and a non-response bias, but it is not clear how this would affect the results.

Study 3: Participants were recruited via Lucid. There might be a self-selection bias and a non-response bias, but it is not clear how this would affect the results.

Ethics oversight

The Institutional Review Board at Stanford University

Note that full information on the approval of the study protocol must also be provided in the manuscript.