

Primaries and Candidate Polarization: Behavioral Theory and Experimental Evidence

Jonathan Woon*
University of Pittsburgh

July 9, 2018[†]

Abstract

Do primary elections cause candidates to take extreme, polarized positions? Standard equilibrium analysis predicts full convergence to the median voter's position, but behavioral game theory predicts divergence when players are policy-motivated and have out-of-equilibrium beliefs. Theoretically, primary elections can cause greater extremism or moderation, depending on the beliefs candidates and voters have about their opponents. In a controlled incentivized experiment, I find that candidates diverge substantially and that primaries have little effect on average positions. Voters employ a strategy that weeds out candidates who are either too moderate or too extreme, which enhances ideological purity without increasing divergence. The analysis highlights the importance of behavioral assumptions in understanding the effects of electoral institutions.

(Word count: 10,461)

*Professor, Department of Political Science, Department of Economics (secondary), and Pittsburgh Experimental Economics Laboratory, woon@pitt.edu, 4437 Wesley W. Posvar Hall, Pittsburgh, PA, 15260

[†]Thanks to Keith Dougherty, Sandy Gordon, Greg Huber, Scott Moser, Charlie Plott, Danielle Thomsen, Alan Wiseman, the editor and anonymous reviewers, participants at the Yale CSAP American Politics Conference, seminar participants at Washington University in St. Louis, University of Oxford (Nuffield College CESS), IC3JM (Juan March Institute), and the Pittsburgh Experimental Economics Lab for helpful comments and feedback. I am also indebted to Kris Kanthak for vigorous discussions during the early stages of this project. Previous versions of the paper were presented at the 2014 Annual Meeting of the American Political Science Association, the 2016 Public Choice Society Meeting, and the 2016 Midwest Political Science Association Conference. This research was supported by the National Science Foundation (SES-1154739) and was approved by the University of Pittsburgh Institutional Review Board under protocol PRO14060001.

The partisan primary system, which favors more ideologically pure candidates, has contributed to the election of more extreme officeholders and increased political polarization. It has become a menace to governing.

— Sen. Charles Schumer (D-NY)¹

The divergence between candidates and legislators from the two major parties is an enduring feature of the American political landscape (Ansolabehere, Snyder and Stewart 2001, Bonica 2013, Poole and Rosenthal 1984, 1997), and the fact that polarization is at historically high levels is a significant concern for scholars and observers of democratic governance, representation, and public policy (Hacker and Pierson 2006, Mann and Ornstein 2013, McCarty, Poole and Rosenthal 2006, 2013). Indeed, politicians and the popular press often lay much of the blame for this phenomenon on partisan primary elections, typically employing a simple, intuitively appealing argument: Candidates take extreme positions because they must appeal to partisan primary voters, whose preferences are more extreme than those of voters in the general election.

Political scientists have tested this argument, finding that while there is some evidence to suggest that primary elections promote extremism, the empirical record is mixed. Extremists are more likely to win congressional primaries than moderates (Brady, Han and Pope 2007), and legislators elected under closed primaries take more extreme positions than legislators elected under open primaries (Gerber and Morton 1998). But other analyses find that polarization is largely unrelated to the introduction of direct primaries (Hirano et al. 2010) and to the variation in the openness of primaries across states (McGhee et al. 2014). At best, primaries may cause polarization under limited circumstances (Bullock and Clinton 2011), and despite the divergence of candidate positions, general elections nevertheless exert nontrivial pressure on candidates to moderate (Hall 2015, Hirano et al. 2010). These findings seem puzzling in light of the basic theory of representation at the heart of this literature that the extremity of primary electorates should directly affect the extremity of party candidates.

¹Charles E. Schumer, “Adopt the Open Primary,” *New York Times*, July 21, 2014.

How, then, should we understand the causal relationship between primary elections and candidate positioning? I examine the connection, both theoretically and experimentally, by comparing elections with and without primaries while holding other features of the electoral environment constant, including preferences and information. The analysis focuses on a particular aspect of primary elections—how the introduction of voters in the candidate selection process affects strategic competition between parties—while abstracting away from many other considerations that might also affect polarization.²

The contribution of this paper is to develop a more nuanced theoretical understanding of the relationship between primaries and polarization than is portrayed in the existing literature. The key theoretical innovation is to move beyond preference-based explanations by treating beliefs as a primitive of the model in a way that is ruled out by standard equilibrium analysis in complete information environments. I show that primaries can cause polarization *or* moderation, depending on candidates’ beliefs about opposing voters’ strategic behavior—even when preferences are held constant. To generate this insight, I rely on ideas from behavioral game theory, which retains much of the theoretical apparatus from standard game theory while allowing for key departures (Camerer 2003). Specifically, I allow players to have “incorrect” or “non-equilibrium” beliefs about others’ actions (Crawford, Costa-Gomes and Iriberri 2013), but assume they are nevertheless strategic in the sense that they best respond to what they *think* other players do (Camerer, Ho and Chong 2004, Nagel 1995, Stahl and Wilson 1995). The analysis demonstrates that changes in preferences alone are not the only cause of polarization. Instead, beliefs and expectations about the strategic behavior of others play important roles in conditioning the effect of institutions.

I then turn to the laboratory and conduct a series of experiments to test the effects of primaries on candidate positions and to distinguish between behavioral assumptions. The chief advantage of the laboratory for theory testing is control (Aldrich and Lupia 2011, Falk

²Such considerations include candidate valence, turnout, activists, or campaign contributions (Adams and Merrill 2008, Callander and Wilson 2007, Hirano, Snyder and Ting 2009, Hummel 2013, Meirowitz 2005, Snyder and Ting 2011).

and Heckman 2009, Morton and Williams 2010), so we can be confident that the observed behavior occurs under the conditions specified by the theoretical model. Importantly, subjects face the same key trade-off in the experiment as the actors do in the theoretical model between increasing the favorability of winning outcomes versus increasing the probability of winning. In the lab, theoretically-relevant quantities of interest that are difficult to measure using observational data with any accuracy or without strong assumptions (in particular, preferences and positions) are also known exactly. Furthermore, experimental manipulations permit tests of mechanisms not possible using observational data. Thus, laboratory experiments are ideal for theory testing given their high internal validity.³

The key finding from the experiment is that primaries appear to cause a kind of ideological purity rather than greater extremism. Regardless of whether there is a primary election or not, I find that subjects take positions that diverge significantly from the median voter’s position. This finding lends support for the behavioral theory. However, the extent to which primaries cause polarization is limited. Greater polarization only occurs without feedback such that candidates cannot learn about the behavior of others, and this polarization happens because voters tend to select extremists over moderates, even though candidate positions do not vary with the election format. More precisely, the analysis reveals that

³The main question of interest for theory testing, as Aldrich and Lupia (2011, 90) put it, is “Will people who are in the situations you describe in your model act as you predict?” For related discussions, see Dickson (2011), Palfrey (2006), and especially, Morton and Williams (2010). While the question of external validity (“to what extent can we generalize from a particular sample?”) is an enduring source of controversy in political science, Falk and Heckman (2009) argue in their insightful defense of the value of lab experiments in social science that “Behavior in the laboratory is reliable and real: Participants in the lab are real human beings who perceive their behavior as relevant, experience real emotions, and take decisions with real economic consequences” (536). Indeed, there are many precedents for testing theories of elite behavior using laboratory experiments (e.g., Aragonés and Palfrey 2007, Frechette, Kagel and Lehrer 2003, Morton 1993). Moreover, Druckman and Kam (2011) note that there is nothing *inherently* problematic with using student samples, and there is little evidence to suggest that using undergraduates as stand-ins for elites biases the results in any particular direction (see Morton and Williams 2010, 343–347). For example, Potters and van Winden (2000) find significant, but small, differences between students and lobbyists, Fatas, Neugebauer and Tamborero (2007) find elites do not fit prospect theory as well as students, while studies by Belot, Duch and Miller (2015), Cooper et al. (1999), and Mintz, Redd and Vedlitz (2006) suggest that student samples provide a *lower bound* to departures from rational decision making. Despite the common assertion that politicians must be better decision makers and more strategic than ordinary people (because they have more experience, access to advice, information, etc.), there is surprisingly little evidence to support such claims. To the contrary, recent comparisons by Sheffer et al. (2018) demonstrate that politicians are just as (and sometimes more) susceptible to choice anomalies than ordinary citizens.

voters do not support party extremists or party moderates unconditionally. Instead, they select candidates with intermediate positions, consistent with their own subjective beliefs about optimal candidate positions, approximately halfway between the median voter and their own party’s ideal point. This behavior generates a greater concentration of candidate positions around an average that diverges from the median voter. Hence, greater ideological homogeneity reinforces, rather than exacerbates, polarization.

Related Literature

My analysis follows a long tradition of using spatial voting models to understand elections. Although existing spatial models predict candidate divergence in elections with primaries, they do so in isolation and do not compare them explicitly to elections without primaries (Aronson and Ordeshook 1972, Coleman 1972, Owen and Grofman 2006).⁴ These models also assume that general election outcomes are probabilistic, which is theoretically consequential because the mechanism they rely on to produce divergence is the combination of policy-motivations and uncertainty about which candidate will win the general election—the same forces that generate incentives for candidate divergence in the absence of primaries (Calvert 1985, Wittman 1983). Thus, it is unclear from the literature whether polarization can be traced to any distinctive features of primaries per se, as electoral institutions. By explicitly comparing institutions, my analysis speaks directly to the connection between primaries and polarization.

Existing theoretical models of two-stage elections also typically maintain the assumption that all political actors, candidates as well as voters, are strategic and forward-looking (e.g., Owen and Grofman 2006). Several models consider the issue of raiding and cross-over voting in open primaries (Cho and Kang 2014, Chen and Yang 2002, Oak 2006), which re-

⁴An exception is Jackson, Mathevet and Mattes (2007), who compare alternative nomination systems in a citizen-candidate framework. In their model, primary elections affect whose preference is decisive in nominating candidates and have no effect if party leaders and the median party voter have the same preferences. Other formal models of primary elections largely focus on considerations of voter uncertainty, incomplete information, and signaling along with issues of candidate valence and distributional concerns.

quires a fairly high degree of strategic sophistication, but this kind of behavior is outside the scope of my analysis. My results also differ from Adams and Merrill (2014), who find that strategic versus expressive voting both generate divergence, but in their model candidates are office-motivated and vary in their campaign skills. In contrast to the preponderance of existing formal models, I take a behavioral (i.e., bounded rationality) approach advocated by Simon (1955), Ostrom (1998), Bendor (2010), and others. I do so by explicitly allowing for sincere (myopic) voting as well as subjective beliefs that are inconsistent with observed behavior.

This paper is also related to two distinct literatures in experimental political science. The experimental literature on candidate positioning in two-party elections finds a strong tendency for candidates and election outcomes to converge to the median voter's position and, more generally, to the Condorcet winner under a variety of conditions, including incomplete information (Collier et al. 1987, McKelvey and Ordeshook 1982, McKelvey and Ordeshook 1985). An exception is Morton (1993), in which candidates are ideological and voting is probabilistic. The other related literature, on strategic voting, generally finds little (at best, mixed) evidence for voter sophistication in the early stages of a multi-stage voting agenda or election contest (Cherry and Kroll 2003, Eckel and Holt 1989, Herzberg and Wilson 1988, McCuen and Morton 2010, Plott and Levine 1978, Van der Straeten et al. 2010).⁵ Taken together, these previous studies raise doubts that voters will be highly strategic (even if candidates are), calling into question theories predicated on voter rationality and strategic sophistication.

⁵Smirnov (2009), who studies endogenous agendas and finds behavior consistent with sophisticated expected utility maximization, is an exception. There is stronger experimental evidence for other kinds strategic voting, however, such as coordinating on a less-preferred candidate in multi-candidate contests (Rietz 2008), and in incomplete information pivotal voter settings (e.g., Battaglini, Morton and Palfrey 2010).

Theoretical Framework and Analysis

I consider an environment with two parties, Party L and Party R , competing to win a single office. Candidates choose positions in a one-dimensional policy space, and the winning candidate's position is implemented as the policy outcome. In the electorate, there are an equal number of voters in each party and a set of independent, non-partisan “swing” voters. Candidates and voters alike are entirely *policy-motivated*, caring only about the location of the policy outcome $w \in \mathbb{R}$. The incentive to win office is therefore purely instrumental in this model, which departs from the usual Downsian office motivations. Parties are completely homogeneous, as candidates and voters belonging to the same party both care about policy and have the same ideal point. Thus, there are three ideal points in the model: θ_L for members of Party L , θ_R for members of Party R , and θ_M for the electorate's median voter, where $\theta_L < \theta_M < \theta_R$. I assume that preferences are symmetric and single-peaked. Specifically, in the experimental implementation, all actors have linear loss utility functions, $u_i(w) = K - |w - \theta_i|$, for $i \in \{L, M, R\}$ and some constant K . Preferences are also common knowledge, so the election takes place under conditions of complete information.

There are two types of elections. In *one-stage elections* (1S), each party has one candidate and their positions are c_L and c_R , respectively, and there is one round of majority rule voting to select the winning candidate. In *two-stage elections* (2S), each party has two candidates (denoted c_{L1} and c_{L2} for Party L , c_{R1} and c_{R2} for Party R) who first compete in intra-party elections (the primaries). The candidates who win their respective party primaries then compete in a second round election (the general election) to select the winning policy w . In other words, the parties hold simultaneous “closed” primaries in which the voter with ideal point θ_L chooses $c_L \in \{c_{L1}, c_{L2}\}$ for Party L at the same time that the voter with ideal point θ_R chooses $c_R \in \{c_{R1}, c_{R2}\}$ for Party R . In the general election, the median voter with ideal point θ_M chooses the election outcome from the two candidates selected by the parties' respective median voters, $w \in \{c_L, c_R\}$.

To generate predictions about candidate positioning and to identify the effects of the election format, I consider a variety of alternative behavioral assumptions. I begin with standard game theoretic analysis, applying Nash equilibrium as the solution concept. Since I am interested in making behavioral predictions, the interpretation of Nash equilibrium is worth a brief discussion. One way to interpret Nash equilibrium is to think of it as an idealized set of assumptions such that actors are not only fully rational but also that their rationality is common knowledge (Aumann and Brandenburger 1995). In this interpretation, we can think of political actors as forming beliefs about others' current and future behavior (as well as beliefs about beliefs and rationality, and so on) that are fully consistent with players' actual strategies and behavior. Alternatively, Nash equilibrium can be interpreted as merely representing a stable outcome in which strategies are mutual best responses, without necessarily invoking an epistemic or belief-based justification of how individuals make decisions in games. Such an approach, however, does not make clear cut predictions about how games are played before an equilibrium state is reached. Nevertheless, under a wide variety of learning models, experience can lead play to converge to Nash equilibrium (Fudenberg and Levine 1998), and the role of experience can be investigated experimentally.

Relaxing the Nash assumption of the mutual consistency of beliefs and actions generates an interesting variety of behavioral possibilities. In my analysis, I first explore the implications of voter sophistication for candidate positioning while holding candidate rationality constant. If voting is "sincere," then primary elections produce more polarized candidates than voting that follows an equilibrium strategy. I then consider another departure from standard assumptions: beliefs that some players make mistakes in choosing their positions. They might do so for any number of reasons, such as miscalculating the optimal position, misjudging or underestimating the rationality of others, or having preferences over outcomes of the game that are not fully captured by their material payoffs. Strategically sophisticated players, recognizing that there are other players who make mistakes, will then choose positions that differ from the Nash predictions—in the direction of their

parties' ideal points—but that are optimal given their own beliefs about the distribution of opponents' positions. Introducing noise or the possibility of mistakes generates divergence in both one-stage and two-stage elections, despite complete information about preferences.

With noise, the effect of introducing a primary election is more nuanced. Similar to the case in which candidates do not make mistakes, the optimal positions depend critically on the behavior of primary voters. If voters choose moderate primary candidates, then two-stage elections will generate greater convergence of candidate positions than in one-stage elections. However, if voters choose extreme primary candidates, then candidates in two-stage elections will be more polarized than candidates in one-stage elections. There is also a third possibility: If voters form their own beliefs about the position most likely to maximize their expected utility and vote for candidates closest to this position, then the degree of candidate divergence in two-stage elections is increasing in what we might call voters' *belief-induced ideal points*. Behavioral game theory thus establishes a critical link between candidates' beliefs about opponents' primary voting behavior and the effect of primaries.

Candidate equilibrium with fully strategic voters

Standard equilibrium analysis leads to identical predictions for both one-stage and two-stage elections. This is because, in any equilibrium, the winning candidate's position is the median voter's ideal point. In one-stage elections, the logic is straightforward. The median voter chooses the party candidate closest to his or her ideal point as the winning candidate, so if one candidate adopts θ_M as a campaign position, no other position can defeat it. In the unique equilibrium of the one-stage election game, both parties' candidates must choose $c_L = c_R = \theta_M$. If not, either the winning party's candidate could do better by finding a position closer to her ideal point while still winning the election or the losing candidate can find a position that wins the election, thereby obtaining a better policy outcome for herself. Thus, $w = \theta_M$ is the unique equilibrium policy outcome.

In two-stage elections, the outcome is the same, but the equilibrium strategies of the primary voters must be specified. Given a set of candidate positions and voters' expectations that the general election median voter will choose the more moderate of the parties' candidates, a primary voter's strategy is to choose the candidate closest to her ideal point as long as she believes the candidate will also win the general election (and in equilibrium, the voter's beliefs about which candidate will win are correct). Because candidates and voters have the same preferences, the incentives guiding optimal candidate strategies in the one-stage election are similar to those that guide rational voting behavior in two-stage elections: if offered the same choices, candidates and voters would choose the same position (the only difference is that candidates can choose any position while primary voters' choices are constrained).

In any equilibrium of the two-stage election game, there must be at least one candidate from *each* party located at θ_M , so primary voters will always be observed choosing the moderate candidate along the path of play. If so, both parties' primary voters will select a candidate at the median voter's ideal point and the policy outcome is therefore $w = \theta_M$. Ruling out other possible outcomes then follows from the same logic as in the nonprimary election. Assuming fully strategic behavior from voters therefore predicts full convergence to the median voter's position in both one-stage and two-stage elections.

Prediction 1. *If voters and candidates are rational, forward-looking agents and form correct beliefs about others' behavior, then (a) the moderate candidates from each party will adopt the median voter's position and (b) primaries will have no effect on the polarization of candidates in the general election.*

Candidate equilibrium with sincere voters

I next consider the possibility that primary voters are myopic and vote “sincerely.”⁶ I assume that sincere voters simply vote for the candidate closest to their ideal points, so they are myopic in the sense that they fail to recognize that the candidate’s chances of winning the general election affect the policy outcome (and hence their payoffs). With myopic voters, the two-stage election game has multiple equilibria in which candidates take divergent positions while the equilibrium of the one-stage election game remains the same (full convergence, since there are no primary voters).

In any equilibrium of the two-stage election game with sincere voters, candidates within each party must adopt the same position, and opposing party candidates must be equidistant from the median voter. Specifically, an equilibrium is characterized by the condition that $c_{L1} = c_{L2} = \theta_M - \delta$ and $c_{R1} = c_{R2} = \theta_M + \delta$, where $\delta \geq 0$ denotes some amount of divergence between candidates. The median voter’s strategy is to select the candidate closest to her own ideal point, breaking ties in favor of each party with equal probability.⁷ The result of the general election is therefore a lottery over $w \in \{\theta_M - \delta, \theta_M + \delta\}$, and the expected value of the outcome is the median voter’s position, $E[w] = \theta_M$. Any candidate who adopts a more extreme position would, at best, be able to win their own primary but then would lose the general election with certainty. Moving to a more moderate position would not change the result of the primary and thus would not change the general election result either. Since no candidate can obtain a better policy outcome by unilaterally adopting a different position, campaign promises characterized by intra-party convergence and inter-party symmetric divergence constitute an equilibrium of the primary election game with

⁶While the overall level of voter “rationality” remains an ongoing subject of debate, the assumption that voters are myopic is consistent with recent observational and experimental research on accountability (e.g., Healy and Malhotra 2009, Huber, Hill and Lenz 2012, Woon 2012a). A theory of elections with boundedly rational, behavioral voters is also worked out by Bendor et al. (2011).

⁷Note that it is also possible to construct equilibria in which the median voter has a bias for one of the parties (i.e., breaks ties in favor of one party rather than randomizing), but this would not affect the equilibrium positions of the candidates. Thus, even though the random tie-breaking rule matches the experimental setup, it is not necessary for the results.

sincere voters. The basic intuition underlying this result is that due to the myopic behavior of sincere primary voters, intra-party competition limits any one candidate’s ability to moderate their party’s position in the general election. Thus, in contrast to full convergence in one-stage elections, any amount of divergence can be supported in two-stage elections.

Prediction 2. *If candidates are rational and forward-looking but primary voters “sincerely” select candidates closest to their own ideal points, then (a) candidates from each party will take positions that diverge from the median voter by the same amount in two-stage elections, and (b) winning candidates will be weakly more polarized in two-stage elections than in one-stage elections, while candidates in the latter will converge to the median voter.*

Candidate best responses to out-of-equilibrium beliefs

The previous sections assumed that candidates correctly anticipate whether voters use either Nash or sincere voting strategies and that their beliefs about other candidates are consistent with those candidates’ actual behavior. That is, if candidate j chooses the platform c_j , then candidate i must believe with certainty that c_j must really be j ’s position. However, this mutual consistency of candidates’ beliefs and actions might break down in a number of ways. Candidates are likely to face cognitive constraints, they may engage in incomplete strategic reasoning, or they may doubt the rationality of other candidates. In this section, I apply the notion of limited strategic sophistication motivated by level- k models in behavioral game theory (Crawford 2003, Nagel 1995, Stahl and Wilson 1995), positing that candidates have some (possibly arbitrary) beliefs and analyze the best response to such beliefs.⁸

To model this, let candidate i ’s beliefs about the positions of candidates from the opposing party $j \neq i$ be given by the cumulative distribution $F(c_j)$. Importantly, these beliefs need not be accurate. For instance, if j ’s true position is $c_j = 0$, candidate i might believe that c_j is uniformly distributed between -1 and 1 . We can think of the distribution $F(c_j)$

⁸While level- k models are a subset of the class of models that assume out-of-equilibrium beliefs, my theory does not rely on different levels of sophistication or reasoning as modeled explicitly in the level- k framework.

as representing *subjective beliefs* that will typically not satisfy the equilibrium consistency requirement.⁹

By relaxing the standard equilibrium assumption of belief consistency, an otherwise expected utility maximizing candidate will choose a position that diverges from the median voter's ideal point. The reasoning is as follows. If a candidate believes there is *some* possibility that the opposing candidate's position diverges from the median voter, then it cannot be optimal for a policy-motivated candidate to choose a platform exactly at the median voter's ideal point. Instead, the candidate will choose a position that trades off some probability of winning against potential policy gains obtained from choosing a position closer to his or her own ideal point. To illustrate this concretely, suppose that $\theta_R = 1$, the left party's ideal point is $\theta_L = -1$, the median is $\theta_M = 0$, and $F(c_L)$ is a uniform random variable, $c_L \sim U[-1, 0]$. With linear loss utility, the optimal position that balances this trade-off is $c_R^* = \frac{1}{3}$. This is illustrated by the solid line showing the expected utility function $EU(c_R)$ in Figure 1.¹⁰ While this logic is similar to the trade-off found in Calvert (1985) and Wittman (1983), the important distinction is that the source of uncertainty in this model is entirely about opponents' behavior rather than about voters or preferences. Moreover, candidate positions are responsive to beliefs such that when a candidate is more likely to expect her opponent to be extreme (i.e., when $F(c_j)$ puts more weight on extreme positions), then she herself will take a position with greater divergence from the median voter in response.

Next, I consider how these beliefs about opposing candidates' positions interact with the type of election. The main result is that the effect of primaries will depend on the

⁹I assume that the density $f(c_j)$ has full support over the interval between median voter θ_M and the opposing party θ_j . The distribution $F(c_j)$ can also be interpreted as an objective probability distribution if candidates' choices are noisy and $F(c_j)$ reflects the true distribution of candidate positions.

¹⁰Formally, given beliefs with density $f(c_L)$, the expected utility function is given by

$$EU(c_R) = \int_{\theta_L}^{2\theta_M - c_R} u(c_R) f(c_L) dc_L + \int_{2\theta_M - c_R}^{\theta_M} u(c_L) f(c_L) dc_L$$

where the integral on the left is the expected utility if c_R is closer to the median voter and wins while the integral on the right is the expected utility if the opposing candidate c_L is closer to the median.

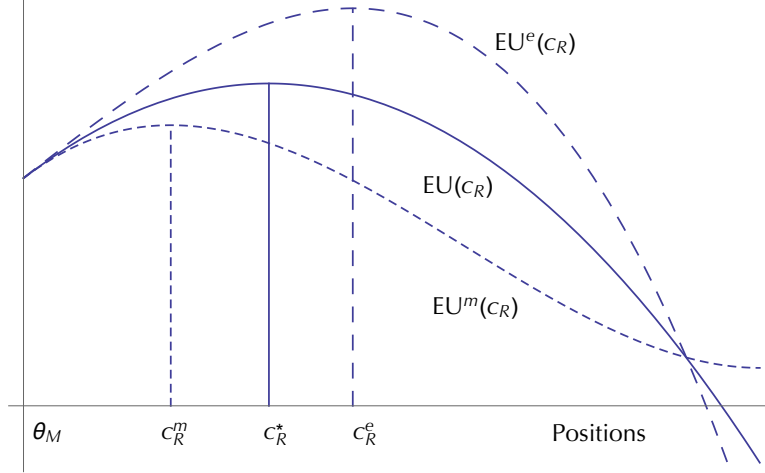


Figure 1: Comparison of candidate expected utility in one-stage and two-stage elections as a function of out-of-equilibrium beliefs and opponents' primary voting behavior: $EU(c_R)$ corresponds to 1S elections, $EU^m(c_R)$ corresponds to 2S elections against moderates, and $EU^e(c_R)$ corresponds to 2S elections against extremists.

candidates' beliefs about the opposing party's primary voters. The baseline for comparison is a one-stage election with opponents drawn from the belief distribution $F(c_j)$. For the purposes of exposition, suppose that $F(c_j)$ is uniform as in the example just given and as shown in the left side of Figure 2, so the candidate's best position is the one that maximizes the same expected utility function $EU(c_R)$ mentioned previously in Figure 1.

In a two-stage election, it is not the original distribution of candidates $F(c_j)$ that matters, but beliefs about which candidate will emerge from the primary election. Let $G(c_j)$ denote this latter set of beliefs about the candidate selected by the opposing party's primary—the candidate that i expects to face in the general election. We can think of the primary election as a selection mechanism or filtering process that affects whether a party's candidate is systematically more or less extreme than the party's initial set of candidates.

More precisely, suppose that both of the opposing party's candidates are independently drawn from $F(c_j)$. Now consider how primary voting behavior affects $G(c_j)$ and, in turn, candidates' positions. If j 's primary voters unconditionally select the more extreme candidate (as they would if they voted sincerely), then party j 's candidate in the general

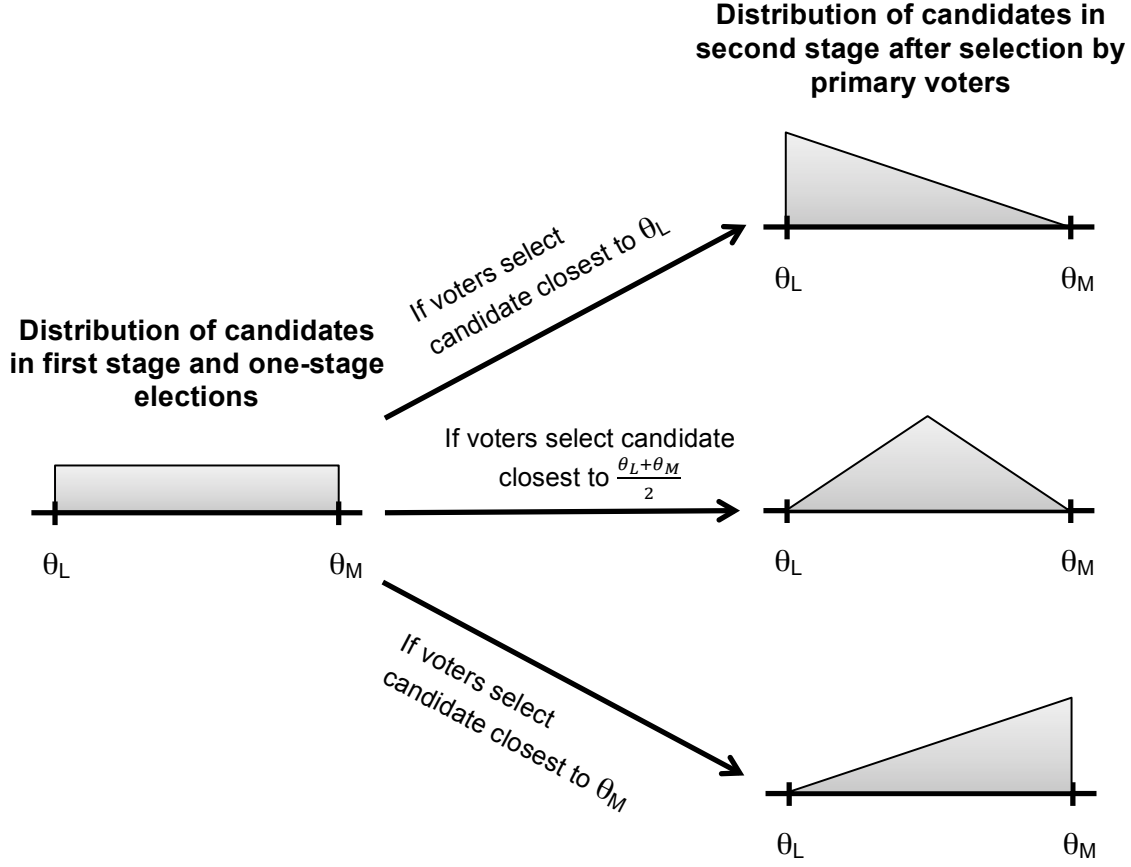


Figure 2: Comparison of beliefs in one-stage and two-stage elections as a function of opponents' primary voting behavior

election will be the more extreme of two independent draws from $F(c_j)$. This results in a distribution $G(c_j)$ that is skewed more towards j 's own ideal point than $F(c_j)$, as shown by the triangular distribution in the upper-right of Figure 2 when $F(c_j)$ is uniform.¹¹ When voters choose extremists, primaries generate incentives for greater *extremism* than in one-stage elections, as illustrated by the upper-dashed expected utility function $EU^e(c_R)$ in Figure 1.

The flip-side of this is that if j 's primary voters select the more moderate candidate (as they would in equilibrium, they generate incentives for greater *moderation* than in one-stage elections. This is because party j 's general election candidate will be the more moderate of

¹¹If L party voters select the extremist, then $c_L = \min\{c_{L1}, c_{L2}\}$ and $G^e(c_L)$ is first order stochastically dominated by $F(c_L)$. If c_{L1} and c_{L2} are independently drawn from $U[-1, 0]$, then $G^e(c_L)$ has density $g^e(c_L) = -2c_L$ for $c_L \in [-1, 0]$.

two independent draws from $F(c_j)$, resulting in a distribution of beliefs $G(c_j)$ that is skewed more towards the median voter than $F(c_j)$. This is shown in Figure 2 by the triangular distribution in the bottom-right.¹² When the probability of facing an extremist opponent is lower, a candidate must moderate their position in response, which is shown by the lower-dashed expected utility function $EU^m(c_R)$ in Figure 1.

These are not the only possibilities, as primary voters might also behave in other ways. For example, a fairly sophisticated voter might reason in the same way as a candidate and form the same beliefs $G(c_i)$ based on expectations about opposing primary voters. To generalize this idea, suppose that a voter has a *belief-induced ideal point* c_j^* and always votes for the candidate in the primary whose position is closest to c_j^* (sometimes this will be the moderate and sometimes the extremist). The result is a distribution of primary candidates $G(c_j)$ that has greater mass closer to c_j^* than $F(c_j)$ does. If c_j^* happens to be the midpoint between the voter's ideal point and the median voter, $G(c_j)$ will be the symmetric triangular distribution in the middle-right of Figure 2. Note that while the mean of this distribution is the same as the original distribution $F(c_j)$, it has lower variance.¹³ Primary elections may therefore also have the effect of reinforcing ideological purity (i.e., increasing homogeneity) within parties even when there is no discernible effect on average candidate positions (i.e., absent changes in polarization).

In contrast to standard equilibrium analysis, which predicts full convergence, a simple model with out-of-equilibrium beliefs generates divergence in candidate positions, even in the absence of primaries and with complete information about preferences. Moreover, the effect of primaries varies with candidates' expectations about the opposing party's voting behavior. Primaries can indeed cause greater polarization (but only if primary voters select sufficiently

¹²In selecting the moderate, $c_L = \max\{c_{L1}, c_{L2}\}$ and so $G^m(c_L)$ first order stochastically dominates $F(c_L)$. In the example where the uniform distribution is $U[-1, 0]$, $G^m(c_L)$ has density $g^m(c_L) = 2c_L + 2$ for $c_L \in [-1, 0]$.

¹³The best response, however, is not exactly the same in the 2S election for the symmetric triangular distribution as it is for the uniform distribution in the 1S election. Nevertheless, there does exist a belief-induced ideal point (via an intermediate value theorem argument, which generates an asymmetric distribution) such that candidates' optimal positions diverge from the median voter by the same amount (i.e., the best responses are identical) in both 1S and 2S elections.

extreme candidates), cause greater moderation (if primary voters select moderates), or cause increased intra-party homogeneity (if voters select on the basis of intermediate belief-induced ideal points). Thus, the behavioral theory identifies how the effect of primaries depends on the connection between beliefs and behavior rather than on preferences alone.

Prediction 3. *If candidates have out-of-equilibrium beliefs about the distribution of opposing candidates, then (a) candidate positions will diverge from the median voter’s ideal point in both one-stage and two-stage elections, (b) the direction of the effect of primary elections on candidate polarization depends on expectations about voting behavior, and (c) polarization in two-stage elections is increasing in the expected extremity of candidates selected by the opposing party’s primary voters.*

Experimental Analysis

The theoretical analysis generated a set of competing predictions about the effect of primaries as a function of alternative behavioral assumptions. If all players are fully strategic, then we should observe full convergence to the median voter’s position and primaries should have no effect. If candidates are strategic but voters are not, then we should observe candidate divergence only in two-stage elections but not in one-stage elections. If the behavioral theory has merit and candidates have subjective beliefs about their opponents’ positions, then polarization, moderation, or increased homogeneity are possible depending on voter behavior. Which set of assumptions is a better reflection of how humans behave is ultimately an empirical question, and thus, I turn to the lab to distinguish between the competing theories.¹⁴

¹⁴See Woon (2012*b*) for an extended discussion of why laboratory experiments are well-suited for behavioral inference in the context of formal models.

Procedures

The experiment was conducted at the Pittsburgh Experimental Economics Laboratory and involved a total of 182 participants drawn primarily from the university’s undergraduate population. Each session involved 14 participants, and each subject participated in one session of either the one-stage (1S) election treatment (6 sessions) or the two-stage (2S) election treatment (7 sessions). At the beginning of each session, following standard laboratory procedures, subjects gave informed consent, the instructions were read out loud to induce public knowledge, and subjects answered a set of questions about the rules on their computers to ensure comprehension.¹⁵ The interface was computerized and programmed using the software z-tree (Fischbacher 2007). Each session took about an hour and a half to complete, and subjects earned an average of \$21.05 (including a \$7 show-up fee).

Subjects participated in a total of 40 elections, and the instructions emphasized that each election was to be treated as a “separate decision task.” For each election, subjects were divided into two groups of seven participants, and every member of a group had the same payoff function and ideal point.¹⁶ Throughout the experiment, the policy space was the set of integers from 1 to 200, and payoffs were given by the linear loss function $200 - |w - \theta_i|$.¹⁷ The parties’ ideal points were located symmetrically from the median voter’s ideal point θ_M such that $\theta_L = \theta_M - d$, $\theta_R = \theta_M + d$, and $d \in \{50, 75\}$. The numerical value of θ_M varied

¹⁵See the Supplementary Materials for the full text of the experimental instructions. Comprehension of the instructions was high. The percentage of correct responses for individual questions ranged from 81% to 94%, and 69% answered all 4 questions correctly while only 8% missed more than one question. These figures likely underestimate the overall degree of comprehension since subjects read explanations of the correct answers before playing the game.

¹⁶We can think of each group as a party, although I was careful to avoid using the term “party” when describing the game to subjects. Groups were randomly reassigned between rounds in two sessions of each treatment, while the remaining sessions involved fixed groups. The method of group assignment did not affect the results, so I ignore the distinction and pool the data in the analysis.

¹⁷Note that with a linear loss function (in contrast to quadratic loss), every possible policy outcome between the parties’ ideal points generates an equal amount of total social welfare, making it unlikely that risk neutral, altruistic subjects will want to choose the midpoint between parties to maximize the total social monetary payoffs of both groups. However, to the extent that subjects’ preferences for money exhibit risk aversion (and they expect this of other subjects), total social welfare will be maximized at the midpoint between parties, which would bias the results *toward* median convergence. Similarly, inequity aversion would also bias choices towards convergence to the median.

from election to election, while the exact sequence of values was identical across sessions and treatments.¹⁸ Payoffs were denominated in “points” and converted to cash by dividing by 10 and rounding to the nearest quarter. Given this conversion rate, the range of possible monetary payoffs for each election was between \$0 and \$20 dollars. The final payment was determined by randomly selecting one election to count from the entire session and then adding the show-up fee.

At the beginning of each election period, subjects first learned the position of every player’s ideal point. Every subject then chose a policy position (referred to as their “campaign promise”), and they were informed that if their campaign promise was selected as the winning position, it would affect every other subject’s payoff. After subjects chose their campaign promise, the computer then randomly selected candidates from each group: one candidate from each group in the 1S election and two candidates from each group in the 2S election, with each group member equally like to be selected and the selection of candidates independent across election periods. The rest of the subjects were assigned to the role of a voter in that election. Thus, at the beginning of each election, every subject was a potential candidate and did not know whether he or she was a candidate until after submitting a campaign promise.¹⁹

Once the candidates were selected, the game proceeded to the voting stages. In the 2S election, voters first chose between one of their group’s two candidates by majority rule. Each primary (first stage) vote is held simultaneously, and neither party knew the positions of the other group’s candidates while voting. Abstentions were not allowed. After each group selected its nominee, a second round of voting took place to choose the winning policy from the two groups’ nominees. All voters participated in this second round, which was effectively

¹⁸To determine the sequence of values, I randomly selected the median’s position, θ_m , from the integers between 51 and 150 for $\delta = 50$ and between 76 and 125 when $\delta = 75$. I varied the numerical values in order to encourage subjects to pay attention and think about their relative, rather than absolute, positions.

¹⁹This method of role assignment is similar in spirit to the strategy method and maximized the number of observed positions in the experiment given that one of the primary goals of the experiment is to measure and test candidate positioning behavior.

the “general election.”²⁰ In contrast to the 2S election treatment, the 1S election treatment featured only one round of voting in which every voter participated.

The median voter in the general election in both the 1S and 2S election treatments was a computer voter who had a distinct ideal point and, as the instructions explained to subjects (similar to Morton 1993), was “like a robot programmed to always vote for the candidate whose campaign promise gives it the higher payoff value.” In the case of ties, the computer voted for each candidate with equal probability. The subjects were informed of the computer voter’s ideal point before every election.

The 40 elections within each session were divided into two parts, where each part varied the type of feedback subjects received. Part 1 consisted of 10 elections without any feedback.²¹ Part 2 consisted of 30 elections with feedback provided to subjects after each election. The information subjects received included the positions of the subjects who were selected as candidates, the number of votes for each candidate, the winning position, and the payoff from the final outcome. Within each part, I varied the distance between the groups’ ideal points by dividing each set of elections into two halves. In the first half, the left and right groups’ ideal points were 100 units apart ($d = 50$), while they were 150 units apart ($d = 75$) in the second half. Note that both of these within-subjects manipulations varied ancillary assumptions (feedback and distance between ideal points) and therefore serve as robustness checks. The experimental manipulation of theoretical interest is the between-subjects manipulation of the electoral institution.

²⁰To avoid priming subjects’ political attitudes regarding primaries, I avoid referring to the two rounds of voting as a “primary” and “general” election but instead refer to them as the “first voting stage” and the “second voting stage.”

²¹The fact that the game is sequential means that it would be impossible to prevent learning across elections if subjects completed each election game before proceeding to the next. I solved this problem by implementing a procedure similar in spirit to the strategy method whereby the game was divided into stages and subjects made their decisions for all elections in one stage before moving to the next stage. That is, subjects first chose their positions for all 10 elections in Part 1, subjects in the 2S treatment then voted for their party’s candidates in all 10 primary elections, and then subjects in both treatments voted in all 10 general elections.

Electoral Dynamics

To get a sense of the kinds of promises candidates make and whether moderates or extremists win elections, Figure 3 presents the sequence of candidate positions and outcomes for selected sessions (2 one-stage sessions and 2 two-stage sessions). The horizontal axis indicates the election, and the vertical axis indicates the promises of the subjects selected as the candidates. These positions are adjusted (centered) so that the general election median voter’s position is 0. The vertical lines indicate when the electoral conditions change in terms of feedback and the distance between the parties’ ideal points. General election candidates are depicted using solid shapes (candidates in one-stage elections and the primary winners in two-stage elections) while primary candidates who lost the first stage election are depicted with hollow shapes. The winning position of the general election is shown by the solid line. Although the dynamics of each session differ, these plots reveal several noteworthy patterns.

First, the positions of candidates from the two parties clearly diverge from the median voter’s position. This is true for both one-stage and two-stage elections, and it appears to persist over the course of the experiment even after subjects gain considerable experience. In session 10 (1S), for example, the candidates from each party choose positions close to their own ideal points, and polarization between the candidates’ positions increases when the underlying preference polarization increases. Along with divergence, there also appears to be substantial heterogeneity and fluctuation in candidate positions.²²

Second, while the general election candidate closer to the median voter’s position generally wins, it is rare for the winning candidate to be located exactly at the predicted equilibrium position. Even in session 4 (one stage), in which the electoral outcome appears most frequently near the median voter’s position, the winning candidate is located at the median’s position in only 3 elections (in another 8 elections, the winning candidate is ± 1 from the median voter’s position). In session 10, the winning candidate usually appears to be just barely closer to the median voter than the losing candidate.

²²The figures also reveal that candidates and voters sometimes make mistakes. For example, in election 1 in session 4, *both* parties’ candidates are located to the left of the median voter, with the party R candidate located at leftmost position in the policy space.

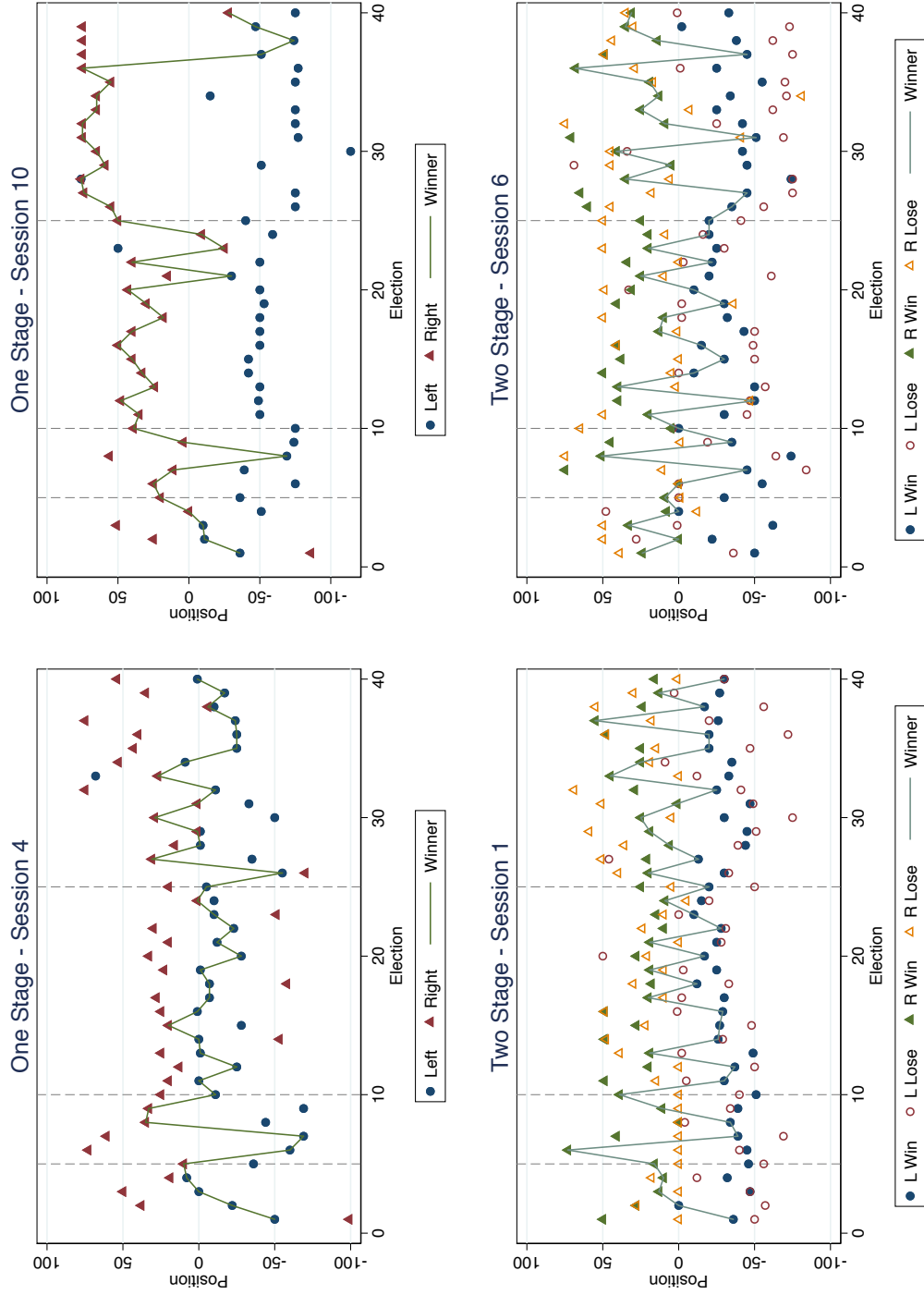


Figure 3: Sample session dynamics

Third, primary voters are inconsistent in selecting either extreme or moderate candidates. Notably, there are several candidates in two-stage elections who locate at exactly the median voter’s position yet lose the primary. In session 1, there were 12 out of 14 such candidates, and in session 6, there were 6 out of 10. While this could suggest that primary voters prefer extremists, there are also many elections in which the more moderate candidate wins. For example, in election 11 of session 6, the left party candidate at -30 defeated the candidate at -45 , and the right party candidate at 20 defeated the candidate located at 50 , with the right party candidate (who is closer to the median voter) winning the general election. Indeed, Figure 3 depicts losing candidates in primary elections on either side of the parties’ winning candidates (indicated by the fact that the hollow candidate markers appear both above and below the solid ones).

These sample dynamics suggest that standard game theoretic analysis poorly predicts candidate positions and voting behavior in the experiment. Whereas equilibrium predicts complete candidate convergence in both one-stage and two-stage elections, I find that candidates’ positions instead diverge. The considerable heterogeneity in candidate positions and the selection of extreme candidates by primary voters indicate that behavioral game theory and non-equilibrium analysis may be useful tools for understanding the consequences of electoral institutions. Of course, Figure 3 only provides a snapshot of experimental behavior. The remainder of the analysis demonstrates that many of the patterns described above generalize across subjects and sessions.

Candidate Positions

Figure 4 shows the average positions over time and by election format for all candidates (top panel) and for winning candidates (bottom panel). In the remainder of the analysis, I measure the extremity of a candidate’s position (vertical axis) by normalizing positions so that a subject’s own ideal point is 1 and the median voter’s ideal point is 0 (so the opposing party’s ideal point is -1 on this transformed scale). The top panel of Figure 4 shows that

candidate positions clearly diverge from the median voter’s position throughout the experiment regardless of the election format. This divergence also appears to persist over time and with no apparent effect of primary elections on polarization. The average normalized position across all rounds is 0.452 in the 1S condition and 0.456 in the 2S conditions. Subjects choose positions only slightly closer to the median voter than the midpoint between their group’s ideal point and the median voter’s ideal point. While the bottom panel shows less stability in the positions of winning candidates due to the fact that there are a small number of sessions per treatment, there are some differences across the feedback conditions. Without feedback, there is slight convergence of winning candidates to the median voter’s position in 1S elections and an increase in divergence once feedback is introduced. In 2S elections, however, the positions of winning candidates remain polarized throughout the experiment.

Table 1 presents a series of ordinary least squares regressions to measure the effect of primaries on candidate divergence while controlling for feedback and experience.²³ The estimates generally reinforce the visual interpretation of the data displayed in Figure 4. Positions diverge (as measured by the intercept) and do not change over time (as the coefficients on *Experience* are small and insignificant across the models). Although primary elections have no effect on the positions chosen by all candidates (column 1), they do have a statistically significant effect on the divergence between party candidates (those standing for election in the second voting stage, column 2) in the absence of feedback. In 1S elections, the divergence of party candidates from the median voter is 0.4 on the normalized scale (i.e., 40% of the distance between the median and the party ideal point) and increases by a fairly substantial 0.175 in 2S elections (to 57.5% of the distance between the median and party ideal point). The natural consequence of this divergence in party candidates is that election outcomes are more extreme in 2S elections than in 1S elections (column 3).

The effect of primary elections disappears, however, when feedback is introduced, as none of the treatment effects in columns (4), (5), or (6) are statistically significant.

²³*Increased Polarization* is an indicator for elections where the distance between parties is 150, and *Experience* counts the number of previous elections for a given condition.

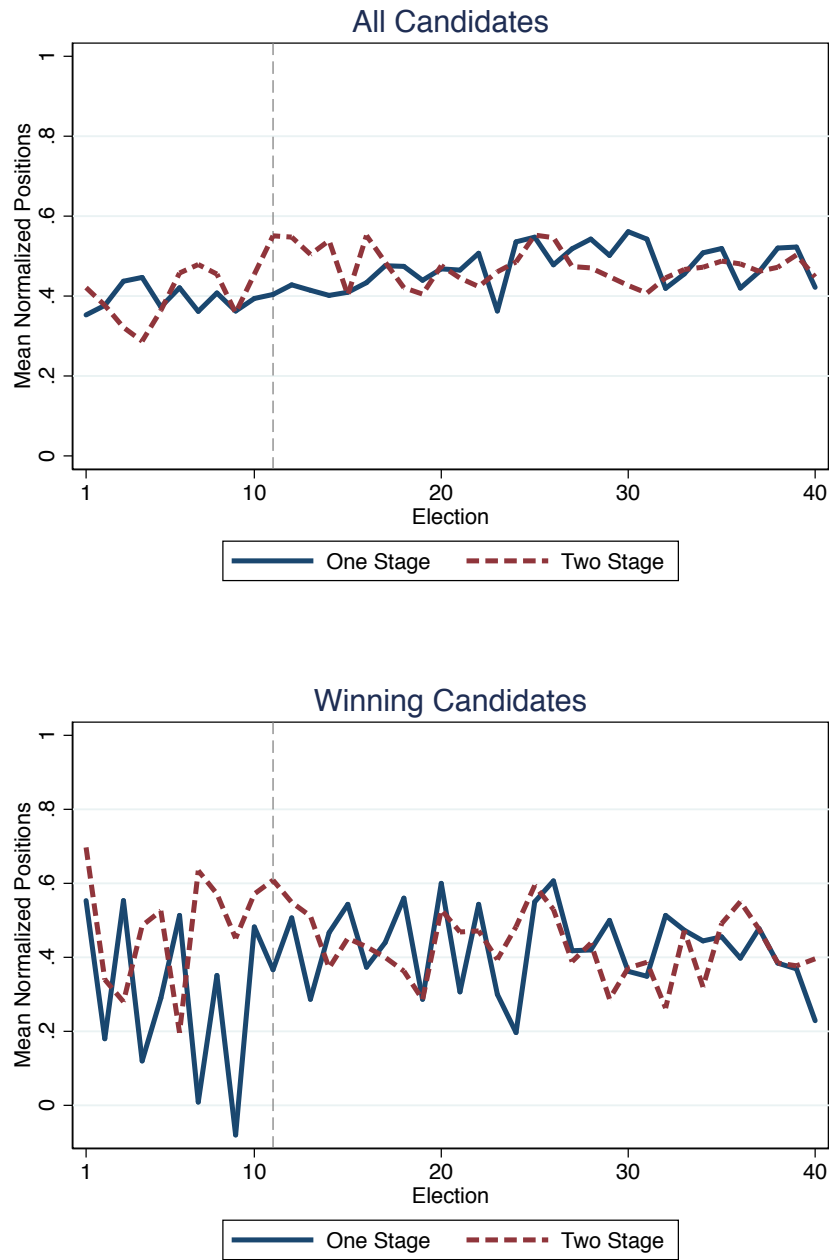


Figure 4: Average positions and outcomes

Table 1: Regression analysis of candidate positions

	No feedback (elections 1-10)			Feedback (elections 11-40)		
	(1) All	(2) Party	(3) Winner	(4) All	(5) Party	(6) Winner
Primary (2S) Elections	0.004 (0.060)	0.175* (0.067)	0.178** (0.058)	0.003 (0.042)	0.046 (0.090)	0.011 (0.102)
Increased Polarization	0.081 (0.041)	-0.098 (0.102)	0.007 (0.143)	0.016 (0.022)	0.044 (0.050)	0.022 (0.080)
Experience	-0.008 (0.009)	0.013 (0.017)	-0.007 (0.023)	-0.000 (0.001)	-0.002 (0.004)	-0.003 (0.004)
Constant	0.387** (0.058)	0.400** (0.051)	0.326** (0.062)	0.469** (0.042)	0.499** (0.103)	0.459** (0.111)
N	1820	260	130	5460	780	390
R^2	0.001	0.028	0.047	0.0001	0.004	0.004

* $p < .05$ ** $p < .01$, OLS regressions with robust standard errors in parentheses, clustered by subjects in (1), (4) and sessions in (2), (3), (5), and (6).

Comparing the intercepts with and without feedback suggests that this is because candidates in 1S elections take more extreme positions once feedback is introduced.²⁴ Indeed, in 1S elections the average party candidate's position is 0.400 without feedback and increases to 0.469 in elections with feedback. In 2S elections, feedback appears to have the opposite effect with average positions starting at 0.575 without feedback and decreasing to 0.545 with feedback. The effect of primaries on party candidate divergence thus disappears as the result of countervailing effects of feedback across institutions.²⁵

²⁴Additional analysis (see the Supplementary Materials) also suggests that candidate positions are responsive to the extremity of positions in the immediately preceding election.

²⁵The persistence of candidate divergence in a one-dimensional spatial setting is surprising given that previous experiments find a consistent tendency for candidates to converge to the median voter's position (Collier et al. 1987, McKelvey and Ordeshook 1985, Morton 1993) or for outcomes to converge to the Condorcet winner (Fiorina and Plott 1978, McKelvey and Ordeshook 1982, Palfrey 2006). The difference may have to do with the fact that candidates are policy-motivated in my experiment rather than office-motivated in most previous experiments, but there are also a number of other differences between my design and previous experiments, including the use of linear instead of quadratic utility, random role assignment, the strategy method, and the varying of the numerical value of players' ideal points. Isolating the exact cause of the difference would be interesting, but doing so is beyond the scope of this paper. Nevertheless, I conducted a modified version of the experiment (detailed in the Supplementary Materials) in which I increase the salience of players' decisions by assigning fixed roles. The main result that candidates diverge in both 1S

Looking only at average positions obscures the effects of primary elections on other characteristics of the distribution of candidate positions. Although the effect of primaries on average positions is limited to elections without feedback, I find that primaries cause candidate positions to become more tightly centered around the mean—that is, less dispersed. Figure 5 plots the standard deviation of candidate and winning positions over the course of the experiment. The graphs reveal two interesting patterns in candidate dispersion. First, in the upper plot, we see dispersion decreasing steadily over time for all candidates. This implies that because average positions remain unchanged, positions converge, not to the median voter’s position but to the mean position in both 1S and 2S elections. Second, we observe a clear effect of primary elections on dispersion. Variation in positions is consistently lower in 2S elections than in 1S elections for all candidates (upper plot) as well as for winning candidates (lower plot).²⁶ Primary elections therefore appear to reinforce candidate polarization by reducing intra-party heterogeneity.

Voting Behavior

The sample dynamics and analysis of candidate positions suggest that, rather than causing or increasing candidate divergence, primaries instead help to maintain polarization by playing a role in the selection of candidates, weeding out party candidates who are either too extreme or too moderate. In this section, I examine voting behavior in primaries by assessing the extent to which primary voters prefer moderates or extremists and by determining the behavioral rule that best fits the experimental data.

Voters tend to select the more extreme candidate, but it is not an overwhelming preference. Overall, voters sincerely choose the extremist in 57% of the elections in the data.

and 2S elections holds up in this modified Fixed Roles Experiment, albeit with some differences, including movement towards the median over the course of the experiment and greater polarization in 2S than in 1S, but the magnitude of the effect is modest.

²⁶The standard deviations of all candidates’ positions aggregated across elections are 0.550 in 1S elections and 0.467 in 2S elections. For winning candidates, the standard deviations are 0.555 in 1S elections and 0.292 in 2S elections. Variance ratio tests show there are statistically significant differences in these variances ($p < 0.01$ in both cases).

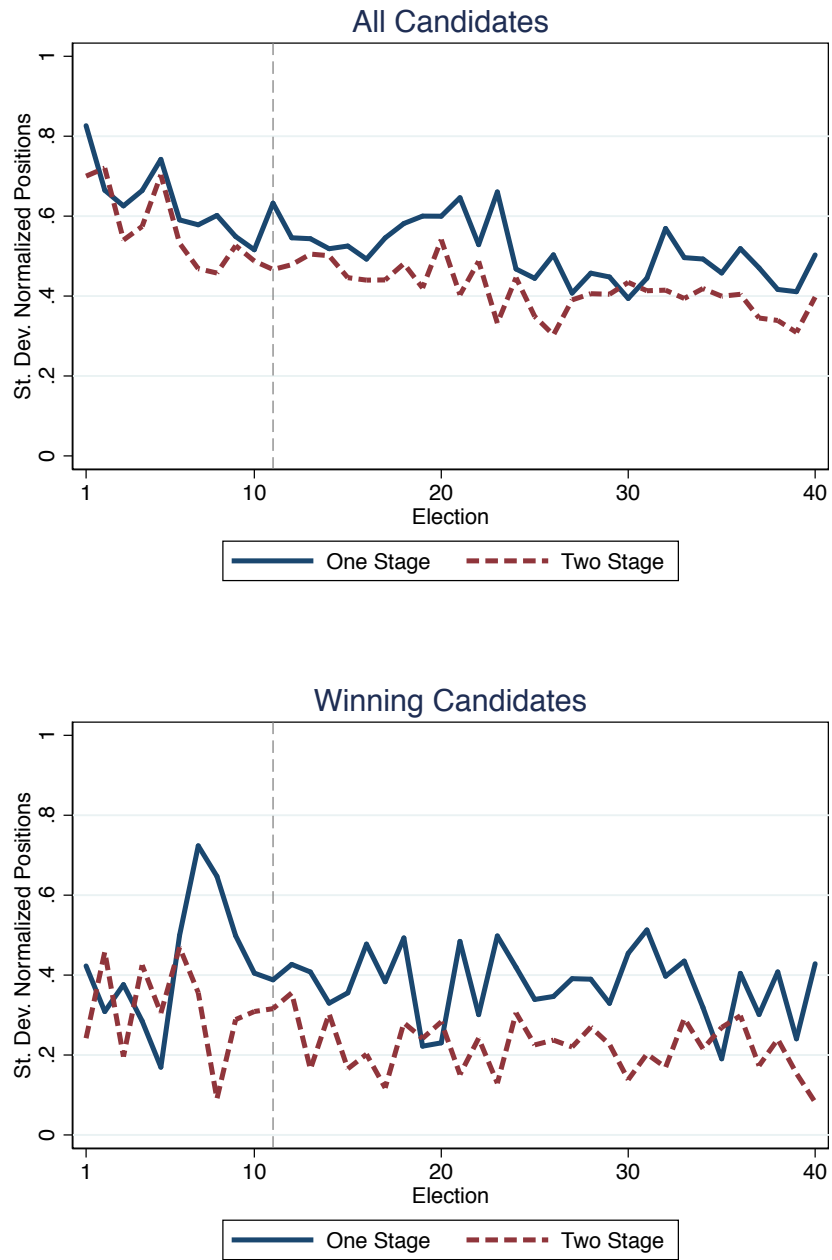


Figure 5: Dispersion of positions and outcomes

When the moderate candidate’s position is closest to the median voter (between 0 and 0.2 on the normalized scale), voters overwhelmingly choose the extremist (69%), especially when both candidates are close to the median voter (90%). Yet, there is an asymmetry because when the extremist candidate’s position is extreme (when the extremist is closest to the voter’s own ideal point between 0.8 and 1 on the normalized scale), the choice between the moderate and extremist is essentially a coin flip (the extremist wins 51% of the time). The more moderate candidates do best when they locate near the midpoint between their party’s ideal point and the median voter and when the extremist is more extreme, but even then, the moderate does not do much better than a coin flip, winning elections at most 57% of the time (when the moderate is between 0.4 and 0.6 and the extremist is between 0.6 and 0.8).²⁷

In Table 2, I characterize voting in each set of 10 elections according to three possible behavioral rules. The first row shows the percentage of strategic voting for the moderate candidate.²⁸ Notice that fewer than half of such votes favor the moderate candidate, 36.6% in the first 10 elections without feedback, increasing slightly to 43.8–46.1% in elections with feedback. The slight increase in voting for moderates appears to lend some support for the theoretical framework, as the change in positioning behavior when feedback is introduced is consistent with the change in voting behavior. Without feedback, more votes are cast for extremists than moderates (63.4% versus 36.6%). If candidates expected this, then their best responses would have been to take more extreme positions, which is consistent with the effect of primaries in elections 1-10. When feedback is introduced, there is an up-tick in voting for moderate candidates, which would lead candidates to expect less extreme opponents and hence to moderate their own behavior.

It is clear, however, that primary voters do not express unconditional preferences for either moderate or extremist candidates in their parties. Table 2 therefore characterizes

²⁷See the Supplementary Materials for additional details about voting behavior as a function of the candidates’ positions.

²⁸This excludes cases where the moderate was the sincere preference due to one or both candidates being outside the region between the median voter and the party’s own ideal point.

Table 2: Classification of behavioral rules for primary voting

Voting rule	Elections				Total
	1-10	11-20	21-30	31-40	
Strategic choice of moderate	36.6% (525)	43.8% (575)	45.7% (630)	46.1% (635)	43.3% (2,365)
Closer to midpoint	67.0% (645)	65.0% (640)	63.7% (680)	64.6% (655)	65.0% (2,620)
Closer to own promise	66.8% (648)	78.9% (644)	79.9% (671)	81.4% (652)	76.8% (2,615)

Total number of votes cast in parentheses, excluding elections where the rule cannot distinguish between candidates.

whether voting is consistent with two additional behavioral rules. The second row shows that a simple strategy where voters select the candidate closest to the midpoint between the median voter and their party’s position is a better description of behavior than either strategic or sincere voting. Roughly two-thirds of votes (overall 65.0%) are consistent with this rule, whereas 43.3% of votes are consistent with strategic voting for moderates and 57.7% of votes are consistent with myopically voting for extremists.

The third row of Table 2 suggests that voters behave as if they have heterogeneous “belief-induced ideal points.” This rule appears to be the most consistent with the data.²⁹ It assumes that each voter has an individual belief that a candidate located at v_i^* maximizes their expected utility and therefore votes for the candidate closest to v_i^* . In the experiment, subjects effectively express such belief-induced ideal points when they choose campaign promises at the beginning of each election, so I use a subject’s campaign promise as a measure of their belief-induced ideal point. Overall, this voting rule attains the highest rate of classification success (76.8%), and outperforms the simple midpoint rule in elections

²⁹Analyzing the results at the subject level (see the Supplementary Material) reinforces this conclusion.

with feedback. By the last 10 elections, 81.4% of votes are consistent with voting for the candidate closest to the belief-induced ideal point (one’s own promise earlier in the election), compared to 46.1% strategic voting for moderates, 53.9% sincere voting for extremists, and 64.6% for the midpoint. Because campaign promises and belief-induced ideal points diverge from the median voter’s position, this voting rule has the effect of reducing variance in candidate positions and reinforcing candidate polarization (as discussed in the theoretical analysis and consistent with the patterns shown in the middle triangular distribution in Figure 5).

A Direct Test of Beliefs and Behavior

The experimental findings that candidate positions diverge from the median voter’s position in both 1S and 2S elections supports the behavioral theory predicated on out-of-equilibrium beliefs (Prediction 3) over the competing predictions based on fully strategic candidate behavior with either forward-looking voters (Prediction 1) or myopic voters (Prediction 2). This inference is indirect, however, because beliefs are neither measured nor manipulated in the experiment. To generate a more direct test of the connection between beliefs and behavior, I conducted another version of the experiment in which beliefs are more carefully controlled and manipulated.³⁰ In this version of the experiment, subjects play as candidates in the 2S election game. Greater control over beliefs is achieved by having subjects play against computer opponents rather than other subjects. This ensures that the distribution of positions is known and exogenous. Variation in beliefs is induced by providing truthful information about whether the opposing party’s candidate is moderate or extreme.

I conducted three sessions of the modified experiment (54 participants, 18 subjects per session).³¹ Each subject played 20 rounds of the 2S game (with feedback) against computer opponents.³² The game was modified so that subjects were informed that the opposing

³⁰I thank an anonymous reviewer for suggesting the need for a direct test that experimentally manipulates beliefs.

³¹These sessions were conducted in the same location (the Pittsburgh Experimental Economics Lab) and with the same subject population (mostly University of Pittsburgh undergraduates) as the original experiment. None of the participants in the additional experiment had participated in the original experiment.

³²The sessions were divided into three parts, with the candidate choices of interest being made in Part

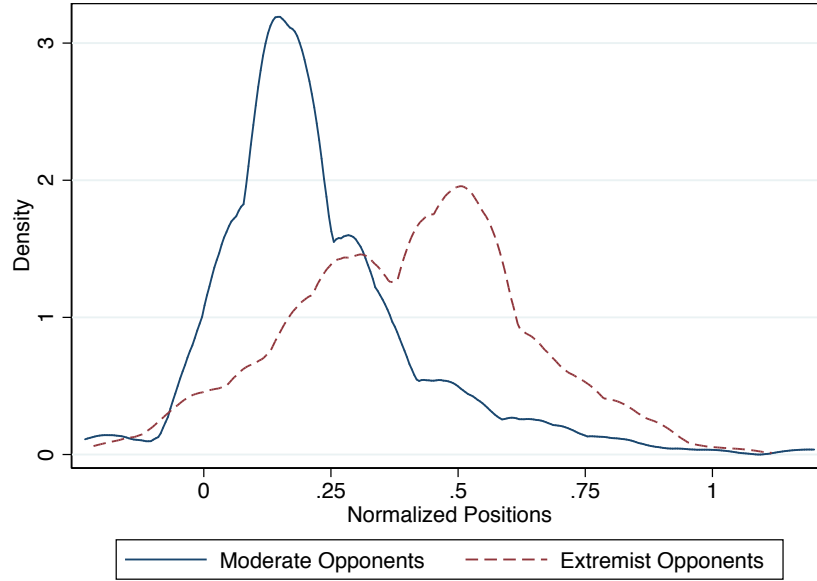


Figure 6: Candidate positions against (computer-selected) moderates versus extremists

candidate's position in the general (second stage) election was stochastically determined by a two-part process. First, two opposing primary candidates' positions were randomly drawn from a uniform distribution over the positions between the median voter's ideal point and the opposing party's ideal point. Second, the opposing party's primary voter (also the computer) randomly selected one of the two positions with equal probability. Information about whether the opposing computer voter chose the moderate or extremist was then provided to the subject. The consequence of this two-part procedure is that candidates faced one of the right-triangular distributions shown in the top-right and bottom-right of Figure 2. When the opponent was known to be more moderate, the belief distribution should have skewed towards the median voter, and when the opponent was known to be more extreme, the distribution should have skewed toward the opposing party's ideal point.

Figure 6 provides a comparison of the distributions that candidates took in response against moderates versus extremists, plotted as kernel densities. The results demonstrate

3. The procedures for Part 1 were the same as in the main experiment for the 2S election game without feedback. This ensures that subjects have the same experience with the game as the subjects in the original experiment. The only minor differences were that groups have 9 players instead of 7 and the distance between parties was held constant at 120. Part 2 involved voting against random opponents, but those data are not analyzed here.

that subjects clearly respond to information about the extremity of their opponents. When computers selected moderate opponents, human candidates took correspondingly moderate positions. Specifically, the majority of such positions (69.6%) fall between 0 and 0.25 on the normalized scale (that is, one-fourth of the distance from the median to their own party’s ideal point). When opponents were extreme, the distribution of positions shifted considerably toward their own party’s ideal point (67.3% of the distribution shifts to *above* 0.25). In addition, the mode increased sharply, from within the interval between 0 and 0.25 to 0.5 on the normalized scale. On average, positions shifted from 0.20 against moderates to 0.35 against extremists, and this difference is statistically significant ($p < 0.01$). These results provide direct evidence that candidates’ positions respond to exogenously induced changes in their beliefs, supporting the belief-based behavioral theory.

Conclusion

The theoretical and experimental models I investigate strip away many of the complexities of real-world elections to focus analytically on how introducing primary voters might affect candidate polarization. By drawing on ideas from behavioral game theory and comparing alternative behavioral assumptions, this approach generates new insights about mechanisms that link voter behavior to strategic expectations and, in turn, candidate positioning. Specifically, when candidates and voters have imperfect, out-of-equilibrium beliefs about the behavior of their electoral opponents, the extremity of the positions and candidates they choose depends on their expectations about the extremity of their opponents. Thus, my analysis suggests that the effect of primary elections should be understood to be conditioned by beliefs and strategic expectations and not simply the product of voter preferences.

Experimentally, I find that the need to win a partisan primary does not affect candidates’ positions—candidates’ positions diverge in elections with and without primaries. This is despite a stark setting in which actors are policy-motivated and have complete information.

Thus, the experimental evidence lends support to the behavioral theory over more standard analytical approaches to studying strategic interaction that predict candidate divergence. Moreover, candidates are responsive to exogenous manipulations of opposing candidates' positions, which provides more direct evidence for the mechanism posited by the theory. To the extent there is any polarization, it tends to occur through the behavior of voters rather than the strategic responses of candidates. The effect is relatively small and limited to settings in which participants cannot learn about the behavior of others from past experience.

Interestingly, the experimental data suggest that primaries have an effect, not by increasing candidate divergence, but by increasing the homogeneity of a party's candidates. Primaries therefore contribute to "ideological purity" in a way that differs from what conventional wisdom would suggest. In the experiment, it is not because primary voters care about ideological purity per se, but because voters use the primary process to weed out candidates both too close to and too far away from the general election median voter's position. Thus, voters in the lab seem to recognize the tension between centrist policies that yield few policy benefits and extreme positions that are unlikely to win the general election, resolving the trade-off by generally splitting the difference. Candidates are responsive to this selective weeding out by voters and, as a consequence, take positions that are more homogeneous in elections with primaries than elections without them. This kind of responsiveness likely reinforces polarization in the long-run. Even though moderates are still more likely to win general elections (Hall 2015), fewer moderates means that parties and candidates are less likely to learn this from experience, which then inhibits convergence through long-run adaptive electoral processes (e.g., Bendor, Mookherjee and Ray 2006, Kollman, Miller and Page 1992).

The behavioral theory helps to make sense of the fact that many people blame partisan primary elections for much of the polarization and dysfunction that afflicts the contemporary American political system but empirical research has not been able to provide compelling

evidence to support the claim. For example, the theory is consistent with the findings that neither the introduction of direct primaries (Hirano et al. 2010) nor the format of primary elections (McGhee et al. 2014) has much to do with increasing polarization. It is also consistent with the fact that polarization has been increasing over time despite the absence of significant changes in electoral institutions.

Future observational research should explore the notion that strategic expectations about increasing polarization may be self-fulfilling. For example, the theory implies that partisans who increasingly perceive the opposing party’s candidates to be more extreme will be emboldened to support more extreme candidates of their own. This may explain Bernie Sanders’ popularity in the 2016 Democratic primary and the popularity of other self-proclaimed “Democratic Socialists” in the aftermath of the 2016 elections. With innovative and appropriate measurement techniques, this hypothesis could be tested both cross-sectionally and over time. It is entirely plausible to the extent that citizens infer extreme ideological positions from their dislike of the opposing party (Brady and Sniderman 1985), given the steady rise of negative partisanship and affective polarization (Abramowitz and Webster 2016, Iyengar, Sood and Lelkes 2012).

References

- Abramowitz, Alan I and Steven Webster. 2016. “The rise of negative partisanship and the nationalization of US elections in the 21st century.” *Electoral Studies* 41:12–22.
- Adams, James and Samuel Merrill. 2008. “Candidate and party strategies in two-stage elections beginning with a primary.” *American Journal of Political Science* 52(2):344–359.
- Adams, James and Samuel Merrill. 2014. “Candidates’ policy strategies in primary elections: does strategic voting by the primary electorate matter?” *Public choice* .
- Aldrich, John H. and Skip Lupia. 2011. Experiments and Game Theory’s Value to Political Science. In *Cambridge Handbook of Political Science*, ed. James N. Druckman, Donald P. Green, James H. Kuklinski and Arthur Lupia. Cambridge University Press.

- Ansolahehere, Stephen, James M Snyder, Jr. and Charles Stewart, III. 2001. "Candidate positioning in US House elections." *American Journal of Political Science* 45(1):136–159.
- Aragones, Enriqueta and Thomas R. Palfrey. 2007. "The Effect of Candidate Quality on Electoral Equilibrium: An Experimental Study." *American Political Science Review* 98(February):77–90.
- Aronson, Peter H. and Peter C. Ordeshook. 1972. Spatial Strategies for Sequential Elections. In *Probability Models of Collective Decision Making*, ed. Richard G. Niemi and Herbert F. Weisberg. Charles E. Merrill pp. 298–331.
- Aumann, Robert and Adam Brandenburger. 1995. "Epistemic conditions for Nash equilibrium." *Econometrica* 63(5):1161–1180.
- Battaglini, Marco, Rebecca B Morton and Thomas R Palfrey. 2010. "The Swing Voter's Curse in the Laboratory." *The Review of Economic Studies* 77(1):61–89.
- Belot, Michele, Raymond Duch and Luis Miller. 2015. "A comprehensive comparison of students and non-students in classic experimental games." *Journal of Economic Behavior & Organization* 113:26–33.
- Bendor, Jonathan B. 2010. *Bounded rationality and politics*. Vol. 6 Univ of California Press.
- Bendor, Jonathan, Daniel Diermeier, David A. Siegel and Michael M. Ting. 2011. *A Behavioral Theory of Elections*. Princeton University Press.
- Bendor, Jonathan, Dilip Mookherjee and Debraj Ray. 2006. "Satisficing and selection in electoral competition." *Quarterly Journal of Political Science* 1(2):171–200.
- Bonica, Adam. 2013. "Ideology and interests in the political marketplace." *American Journal of Political Science* 57(2):294–311.
- Brady, David W., Hahrie Han and Jeremy C. Pope. 2007. "Primary Elections and Candidate Ideology: Out of Step with the Primary Electorate?" *Legislative Studies Quarterly* 32(1):79–105.
- Brady, Henry E and Paul M Sniderman. 1985. "Attitude attribution: A group basis for political reasoning." *American Political Science Review* 79(4):1061–1078.
- Bullock, Will and Josh Clinton. 2011. "More of a Molehill than a Mountain: The Effects of the Blanket Primary on Elected Officials' Behavior from California." *Journal of Politics* 73(3):915–30.
- Callander, Steven and Catherine H. Wilson. 2007. "Turnout, Polarization, and Duverger's Law." *The Journal of Politics* 69(November):1047–1056.
- Calvert, Randall. 1985. "Robustness of the Multidimensional Voting Model: Candidates' Motivations, Uncertainty, and Convergence." *American Political Science Review* 29(1):69–95.

- Camerer, Colin. 2003. *Behavioral game theory: Experiments in strategic interaction*. Princeton University Press.
- Camerer, Colin F., Teck-Hua Ho and Juin-Kuan Chong. 2004. "A Cognitive Hierarchy Model of Games." *Quarterly Journal of Economics* 119(3):861–898.
- Chen, Kong-Pin and Sheng-Zhang Yang. 2002. "Strategic voting in open primaries." *Public Choice* 112(1-2):1–30.
- Cherry, Todd L. and Stephan Kroll. 2003. "Crashing the Party: An experimental investigation of strategic voting in primary elections." *Public Choice* 114:387–420.
- Cho, Seok-Ju and Insun Kang. 2014. "Open primaries and crossover voting." *Journal of Theoretical Politics* .
- Coleman, James S. 1972. The Positions of Political Parties in Elections. In *Probability Models of Collective Decision Making*, ed. Richard G. Niemi and Herbert F. Weisberg. Charles E. Merrill pp. 332–357.
- Collier, Kenneth E, Richard D McKelvey, Peter C Ordeshook and Kenneth C Williams. 1987. "Retrospective voting: An experimental study." *Public Choice* 53(2):101–130.
- Cooper, David J, John H Kagel, Wei Lo and Qing Liang Gu. 1999. "Gaming against managers in incentive systems: Experimental results with Chinese students and Chinese managers." *American Economic Review* 89(4):781–804.
- Crawford, Vincent P. 2003. "Lying for Strategic Advantage: Rational and Boundedly Rational Misrepresentation of Intention." *American Economic Review* 93(1):133–149.
- Crawford, Vincent P, Miguel A Costa-Gomes and Nagore Iriberri. 2013. "Structural models of nonequilibrium strategic thinking: Theory, evidence, and applications." *Journal of Economic Literature* 51(1):5–62.
- Dickson, Eric. 2011. Economics versus Psychology Experiments: Stylization, Incentive, and Deception. In *Cambridge Handbook of Experimental Political Science*, ed. James N. Druckman, Donald P. Green, James H. Kuklinski and Arthur Lupia. Cambridge University Press.
- Druckman, Jamie N. and Cindy D. Kam. 2011. Students as Experimental Participants: A Defense of the "Narrow Data Base". In *Cambridge Handbook of Experimental Political Science*, ed. James N. Druckman, Donald P. Green, James H. Kuklinski and Arthur Lupia. Cambridge University Press pp. 41–57.
- Eckel, Catherine and Charles A Holt. 1989. "Strategic voting in agenda-controlled committee experiments." *American Economic Review* 79(4):763–773.
- Falk, Armin and James J Heckman. 2009. "Lab experiments are a major source of knowledge in the social sciences." *science* 326(5952):535–538.

- Fatas, Enrique, Tibor Neugebauer and Pilar Tamborero. 2007. "How politicians make decisions: A political choice experiment." *Journal of Economics* 92(2):167–196.
- Fiorina, M P and C R Plott. 1978. "Committee decisions under majority rule: An experimental study." *The American Political Science Review* .
- Fischbacher, Urs. 2007. "z-Tree: Zurich Toolbox for Ready-Made Economics Experiments." *Experimental Economics* 10(2):171–8.
- Frechette, Guillaume R, John H Kagel and Steven F Lehrer. 2003. "Bargaining in legislatures: An experimental investigation of open versus closed amendment rules." *American Political science review* 97(2):221–232.
- Fudenberg, Drew and David K Levine. 1998. *The Theory of Learning in Games*. Vol. 2 MIT press.
- Gerber, Elisabeth R. and Rebecca B. Morton. 1998. "Primary Election Systems and Representation." *Journal of Law, Economics, and Organization* 14(2):302–324.
- Hacker, Jacob S and Paul Pierson. 2006. *Off center: The Republican revolution and the erosion of American democracy*. Yale University Press.
- Hall, Andrew B. 2015. "What Happens When Extremists Win Primaries?" *American Political Science Review* 109(1):18–42.
- Healy, Andrew and Neil Malhotra. 2009. "Myopic Voters and Natural Disaster Policy." *American Political Science Review* 103(August):387–406.
- Herzberg, Roberta Q and Rick K Wilson. 1988. "Results on sophisticated voting in an experimental setting." *Journal of Politics* 50(2):471–486.
- Hirano, Shigeo, James M. Snyder, Jr. and Michael M. Ting. 2009. "Distributive politics with primaries." *Journal of Politics* 71(4):1467–1480.
- Hirano, Shigeo, James M. Snyder, Jr., Stephen Ansolabehere and John Mark Hansen. 2010. "Primary Elections and Partisan Polarization in the U.S. Congress." *Quarterly Journal of Political Science* 5(2):169–191.
- Huber, Gregory A., Seth J. Hill and Gabriel S. Lenz. 2012. "Sources of Bias in Retrospective Decision Making: Experimental Evidence on Voters' Limitations in Controlling Incumbents." *American Political Science Review* 106(4):720–741.
- Hummel, Patrick. 2013. "Candidate strategies in primaries and general elections with candidates of heterogeneous quality." *Games and Economic Behavior* 78:85–102.
- Iyengar, Shanto, Gaurav Sood and Yphtach Lelkes. 2012. "Affect, Not Ideology A Social Identity Perspective on Polarization." *Public opinion quarterly* 76(3):405–431.
- Jackson, Matthew O., Laurent Mathevet and Kyle Mattes. 2007. "Nomination processes and policy outcomes." *Quarterly Journal of Political Science* 2(1):67–92.

- Kollman, Ken, John H Miller and Scott E Page. 1992. "Adaptive parties in spatial elections." *American Political Science Review* 86(4):929–937.
- Mann, Thomas E and Norman J Ornstein. 2013. *It's even worse than it looks: How the American constitutional system collided with the new politics of extremism*. Basic Books.
- McCarty, Nolan, Keith T Poole and Howard Rosenthal. 2006. *Polarized America: The dance of ideology and unequal riches*. MIT Press.
- McCarty, Nolan, Keith T Poole and Howard Rosenthal. 2013. *Political Bubbles: Financial Crises and the Failure of American Democracy*. Princeton University Press.
- McCuen, Brian and Rebecca B Morton. 2010. "Tactical coalition voting and information in the laboratory." *Electoral Studies* 29(3):316–328.
- McGhee, Eric, Seth Masket, Boris Shor, Steven Rogers and Nolan McCarty. 2014. "A Primary Cause of Partisanship? Nomination Systems and Legislator Ideology." *American Journal of Political Science* 558(2):337–351.
- McKelvey, Richard D. and Peter C. Ordeshook. 1982. "Two-Candidate Elections without Majority Rule Equilibria An Experimental Study."
- McKelvey, Richard D and Peter C Ordeshook. 1985. "Sequential Elections with Limited Information." *American Journal of Political Science* 29(3):480–512.
- Meirowitz, Adam. 2005. "Informational party primaries and strategic ambiguity." *Journal of Theoretical Politics* 17(1):107–136.
- Mintz, Alex, Steven B Redd and Arnold Vedlitz. 2006. "Can we generalize from student experiments to the real world in political science, military affairs, and international relations?" *Journal of Conflict Resolution* 50(5):757–776.
- Morton, Rebecca B. 1993. "Incomplete Information and Ideological Explanations of Platform Divergence." *American Political Science Review* 87(2):382–392.
- Morton, Rebecca B. and Kenneth C. Williams. 2010. *Experimental Political Science and the Study of Causality: From Nature to the Lab*. Cambridge University Press.
- Nagel, Rosemarie. 1995. "Unraveling in Guessing Games: An Experimental Study." *The American Economic Review* 85(5):1313–1326.
- Oak, Mandar P. 2006. "On the role of the Primary System in Candidate Selection." *Economics & Politics* 18(2):169–190.
- Ostrom, Elinor. 1998. "A Behavioral Approach to the Rational Choice Theory of Coallctive Action: Presidential Address, American Political Science Association, 1997." *American Political Science Review* 92(1):1–22.

- Owen, Guillermo and Bernard Grofman. 2006. "Two-stage electoral competition in two-party contests: persistent divergence of party positions." *Social Choice and Welfare* 26(3):547–569.
- Palfrey, Thomas R. 2006. Laboratory Experiments. In *Oxford Handbook of Political Economy*, ed. Barry R. Weingast and Donald A. Wittman. Oxford University Press pp. 915–936.
- Plott, Charles R and Michael E Levine. 1978. "A model of agenda influence on committee decisions." *American Economic Review* 68(1):146–160.
- Poole, Keith T and Howard Rosenthal. 1984. "The polarization of American politics." *The Journal of Politics* 46(4):1061–1079.
- Poole, Keith T and Howard Rosenthal. 1997. *Congress: A political-economic history of roll call voting*. Oxford University Press.
- Potters, Jan and Frans van Winden. 2000. "Professionals and Students in a Lobbying Experiment: Professional Rules of Conduct and Subject Surrogacy." *Journal of Economic Behavior & Organization* 43:499–522.
- Rietz, Thomas. 2008. "Three-way experimental election results: strategic voting, coordinated outcomes and Duverger's law." *Handbook of Experimental Economics Results* 1:889–897.
- Sheffer, Lior, Peter John Loewen, Stuart Soroka, Stefaan Walgrave and Tamir Sheafer. 2018. "Nonrepresentative Representatives: An Experimental Study of the Decision Making of Elected Politicians." *American Political Science Review* 112(2):302–321.
- Simon, Herbert A. 1955. "A Behavioral Model of Rational Choice." *Quarterly Journal of Economics* 69(1):99–118.
- Smirnov, Oleg. 2009. "Endogenous choice of amendment agendas: types of voters and experimental evidence." *Public Choice* 141(3-4):277–290.
- Snyder, Jr., James M. and Michael M. Ting. 2011. "Electoral selection with parties and primaries." *American Journal of Political Science* 55(4):782–796.
- Stahl, Dale O. and Paul W. Wilson. 1995. "On Players' Models of Other Players: Theory and Experimental Evidence." *Games and Economic Behavior* 10(1):218–254.
- Van der Straeten, Karine, Jean-Francois Laslier, Nicolas Sauger and Andre Blais. 2010. "Strategic, Sincere, and Heuristic Voting under four Election Rules: An Experimental Study." *Social Choice and Welfare* 35(3):435–472.
- Wittman, Donald A. 1983. "Candidate Motivations: A Synthesis of Alternative Theories." *American Political Science Review* 77(1):142–57.
- Woon, Jonathan. 2012a. "Democratic Accountability and Retrospective Voting: A Laboratory Experiment." *American Journal of Political Science* 56(4):913–930.

Woon, Jonathan. 2012*b*. Laboratory Tests of Formal Theory and Behavioral Inference. In *Experimental Political Science: Principles and Practices*, ed. Bernhard Kittel, Wolfgang J. Luhan and Rebecca B. Morton. Palgrave Macmillan.