

Surname Diversity, Social Ties and Innovation^{*}

Max Posch[†] Jonathan Schulz[‡] Joseph Henrich[§]

July 11, 2023

Abstract

A growing body of evidence suggests that innovation may be fueled by social interactions between individuals with different skills, expertise and ways of thinking. Exploiting the fact that the extent of such diverse interactions varies among communities and across time, this paper investigates how such interactions influenced innovation in U.S. counties from 1850 to 1940. We introduce and validate a new measure of social interactional diversity based on the distribution of surnames: lower surname diversity indicates that social interactions are more concentrated among like-minded people. Leveraging quasi-random variation in counties' surname compositions—stemming from the interplay between historical fluctuations in immigration and local factors that attract immigrants—we find that surname diversity increases both the quantity and quality of innovation. The results support the view that social interactions between diverse minds are key drivers of innovation.

Keywords: Innovation, social interactions, surname diversity, immigration

JEL Classification: O33, R11, N92, J15, Z13

^{*}Max Posch and Jonathan Schulz share co-first authorship. For comments and helpful discussions, we thank Ran Abramitzky, Alberto Alesina, Mike Andrews, Pablo Balan, Tyler Cowen, Klaus Desmet, Dan Fetter, Erik Hornung, Chad Jones, Ross Mattheis, Petra Moser, Nathan Nunn, Tzachi Raz, Slava Savitskiy, and seminar and workshop participants at NYU, UBC, Cologne & Bonn, Lund University and University of Louvain. We thank Enrico Berkes for generous data sharing. Research reported in this publication was supported by the John Templeton Foundation under Award Number 62161.

[†]Department of Economics, University of Exeter. (email: m.posch@exeter.ac.uk)

[‡]Department of Economics, George Mason University. (email: jonathan.schulz77@gmail.com)

[§]Department of Human Evolutionary Biology, Harvard University. (email: henrich@fas.harvard.edu)

It is hardly possible to overrate the value [...] of placing human beings in contact with persons dissimilar to themselves, and with modes of thought and action unlike those with which they are familiar. [...] Such communication has always been, and is peculiarly in the present age, one of the primary sources of progress.

John Stuart Mill
Principles of Political Economy

1 Introduction

At least since John Stuart Mill (1871), economists and other scholars have argued that social interactions among diverse minds encourage innovation and creativity (e.g., [Jacobs, 1969](#); [Glaeser et al., 1992](#); [Weitzman, 1998](#); [Muthukrishna and Henrich, 2016](#); [Galor, 2022](#); [Andrews, 2023](#)). In this view, innovations emerge primarily through the recombination of ideas created as individuals with diverse knowledge, skills and perspectives interact and share their ideas ([Jones, forthcoming](#)). Recent research has revealed remarkable variation in the degree to which interactions among diverse individuals occur across different countries and even among communities within the same country ([Banfield, 1958](#); [Alesina and Giuliano, 2014](#); [Schulz et al., 2019](#); [Enke, 2019](#); [Henrich, 2020](#); [Ghosh et al., 2023](#)). Some communities exhibit concentrated social networks, often structured around strong family ties, which can cultivate an inward-looking psychology marked by mistrust of outsiders and reluctance to engage beyond one’s family or in-group. If this mindset hampers the free exchange of knowledge and ideas beyond one’s family or in-group, then less concentrated, or more diverse, social interactions should foster innovation. However, quantifying the impact of this diversity on innovation has proven elusive due to empirical challenges related to measurement—the need for a finely grained measure of the diversity of social interactions, and causal identification—individuals and their social interactions are not randomly distributed across locations or time.

In this paper, we propose a novel approach to measuring the diversity of social interactions based on the diversity of surnames within U.S. counties, and then use this to assess the existence of a causal relationship with innovation, as measured using data on patents. Theoretically, we develop our diversity measure to test the hypothesis that social interaction among diverse minds spurs innovation, largely because many new ideas arise from the recombination of existing ideas. The potential informational exchanges that occur during these social interactions are assumed to depend on two components: an informational element, in which different individuals possess distinct knowledge, and a

social-psychological aspect, in which individuals must trust each other and be willing to share their ideas. Surnames in the U.S. are well suited for capturing both components. From an informational perspective, surnames, which are typically patrilineally inherited in the U.S., can serve as markers of different kinds of knowledge, skills and perspectives that originate from learning within familial, ethnic, or professional networks and persist across generations because of intergenerational cultural transmission (Boyd and Richerson, 1985; Bisin and Verdier, 1998). From a social-psychological perspective, since individuals with the same surname are more likely to be related, low diversity in surnames may indicate the prevalence of family networks, fostering an inward-looking orientation.¹ Conversely, high surname diversity suggests less concentrated social networks, which makes broader engagement with strangers relatively more beneficial, nurturing the cultivation of impersonal trust.

In our empirical approach, we utilize all surnames reported in the full-count U.S. Census data available from 1850 to 1940. We use this data to compute the diversity of surnames across U.S. counties, which are presumed to be the primary locations of social interactions during this period. While counties might not encapsulate every social interaction, particularly in today's highly interconnected world, they provide a reasonable approximation in this pre-1950 historical context.

We validate the use of surname diversity as a measure of the diversity of social interactions using two approaches. First, to evaluate the informational aspect of social interaction diversity, we analyze the extent to which surnames cluster within specific occupations and immigrants' ancestral regions. Using the Census data, we find that two individuals selected at random, who share the same surname, have a higher than random probability of also sharing an occupation or having origins in the same country or subnational region. This finding aligns with recent studies on social mobility (Clark, 2014; Güell et al., 2015; Bell et al., 2019; Barone and Mocetti, 2021), bolstering the assertion that surnames capture distinct facets of knowledge tied to specific occupations and ethnicities. Second, for the psychological aspect of social interaction diversity, we probe the predictive power of county-level surname diversity for contemporary impersonal trust. In addition, we scrutinize the correlation between a historical measure of the strength of family ties (Raz, 2023) and surname diversity, since previous work has found that strong family ties are tightly linked to lower trust (Alesina and Giuliano, 2014; Schulz et al., 2019). We find a strong association with trust as well as a high correlation with the strength of family ties. We do

¹In population genetics, low surname diversity is often used as a marker of cousin marriage and other forms of inbreeding (e.g., Barraï et al., 1996). This approach has recently been adopted by the economics literature (Ghosh et al., 2023).

not find these patterns for the more conventional measures of diversity based on country of birth and race, which is consistent with much prior evidence (Glaeser et al., 2000; Alesina and La Ferrara, 2000, 2002; Ashraf and Galor, 2013) and underlines the distinctive nature of surname diversity compared to these other types of diversity. Taken together, these results support the view that surname diversity is an effective measure of social interactional diversity, encapsulating both informational and psychological components.

To measure innovation, we rely on two indicators based on U.S. patents. First, we calculate the total number of patents per capita for each U.S. county for 5 or 10-year periods from the 1850s to the 1940s based on the Comprehensive Universe of U.S. Patents (Berkes, 2018). Second, we use the breakthrough patent indicator created by Kelly et al. (2021) to capture highly important patents. Breakthrough patents are identified via textual similarity to previous and subsequent patents: breakthrough patents have low similarity to previous patents but high similarity to subsequent ones.

Using these diversity and innovation measures, our analyses proceed as follows: First, we study the correlation between surname diversity and both patents and breakthrough patents per capita across U.S. counties from the 1850s to the 1940s. We find positive and economically important relationships between surname diversity and both innovation measures: a one standard deviation increase in surname diversity within a county is associated with approximately 78% more patents per capita and a slightly larger increase in breakthrough patents per capita. These relationships are also remarkably stable over time—our sample spans almost a century of U.S. innovation—and hold when controlling for county and period-state fixed effects, population-scale effects, and the composition of immigrants and race within these counties. We find that the more conventional country-of-birth diversity measure is also a robust predictor of patents in most specifications. However, surname diversity offers additional explanatory power even when controlling for country-of-birth diversity and country-of-birth-specific immigrant shares.

Second, we provide evidence that a greater diversity of social interactions results in faster innovation. As noted, such causal evidence has proven elusive because individuals do not allocate randomly across space but tend to move into innovative regions, possibly creating a spurious correlation between surname diversity and patents. To address this concern, we employ an instrumental variable (IV) strategy, building on the approach developed by Burchardi et al. (2019). This strategy leverages historical immigration patterns as a significant determinant of surname diversity in U.S. counties.

Migration, beyond births and deaths, is the key driver of the composition of surnames. However, immigration does not monotonically increase surname diversity. Its impact largely depends on the pre-existing surname distribution in specific counties. Influxes

of individuals carrying rare or new surnames in the county increase surname diversity, while inflows of people bearing locally common surnames results in the opposite effect. We hypothesize that this relationship between immigration and surname diversity substantially affects both the informational and psychological channels (detailed further in section 2). When individuals carrying locally rare surnames arrive, they enhance surname diversity, and in turn, may create opportunities for diverse social interactions, knowledge acquisition, and the cultivation of trust towards individuals with differing cultural and family backgrounds. On the other hand, an inflow of individuals bearing locally common surnames decreases surname diversity. This movement of individuals, who are culturally and genealogically related to the dominant groups within counties, may limit opportunities for novel knowledge acquisition, strengthen family ties or bonds among culturally homogeneous groups, and nurture a low-trust mentality towards outsiders.

The IV approach isolates quasi-random variation in counties' surname composition, which stems solely from the historical interplay of two forces: (i) the staggered arrival of migrants with different surnames and (ii) temporal variation in the relative attractiveness of different destination counties for the average migrant arriving at the time. The interaction of these two historical forces enables us to isolate variation in surname distributions across counties that is essentially inherited from plausibly exogenous shocks to historical migrations dating back far into the 19th century.

Using data across counties from 1900 to 1940, we find that a one-standard deviation increase in surname diversity raises patents and breakthrough patents (per 1,000 people) by 76-93% and 144-149%, respectively.

These results hold across key robustness checks. First, to scrutinize the potential for reverse causality—that is, an increase in a county's innovation leading to increased diversity—we perform a falsification exercise in which we regress past patents on future surname diversity. The coefficients from this exercise are near zero, or even negative, and statistically insignificant, providing strong evidence against the concern that reverse causality might confound our results.

Second, to address the potential influence of scale effects, including through immigration, we control for quasi-random variation in population size isolated by the IV procedure. Once again, the estimates align with our primary findings.

Third, we confront the possibility that our results are region-specific within the U.S., given factors like racial segregation that could simultaneously affect innovation, surname diversity and immigration. Estimating the impact of surname diversity on patents across the four major U.S. census regions (Northeast, Midwest, South, and West), we find consistently positive coefficients, most of which are accurately estimated.

Fourth, recent contributions to research on social mobility (e.g., [Clark, 2014](#); [Barone and Mocetti, 2021](#)) raise the concern that unobserved characteristics embedded in specific (rare) surnames, such as abilities, interests, or knowledge drive the results rather than diversity per se. To explore this, we change the unit of observation from county-period to surname-county-period and include surname-fixed effects in our specifications to absorb any surname-specific traits. We find that our estimates remain highly significant across all specifications and are minimally affected by the inclusion of these fixed effects.

Last, we address the potential concern that a direct effect of immigration, that is not channeled through diversity, confounds our estimates. Replicating our analysis within a subsample of U.S.-born individuals, we find effects of surname diversity on patents and breakthrough patents that are virtually unchanged, reinforcing the interpretation that it is diversity rather than immigration per se driving our results.

Next, we examine the mechanisms underlying our findings. First, we use the more than 140,000 technology codes assigned by the United States Patent and Trademark Office (USPTO) to categorize patents into three distinct types: (1) novel technologies, (2) novel combinations, and (3) reuse/refinement combinations (we follow the methodologies developed by [Strumsky et al. \(2011\)](#) and [Akcigit et al. \(2013\)](#)). We consider a patent a novel technology for a given county if any of its technology codes appear for the first time in that county and the grant year of the patent. If a patent is not classified as a novel technology but includes a unique pairwise combination of technologies appearing for the first time in the county and grant year, we categorize it as a novel combination. Any remaining patents are classified as reuse/refinement combinations. We find robust positive effects of surname diversity across all three patent types, with a particularly noteworthy impact on novel combinations. This finding is consistent with the hypothesis that diversity spurs the recombination of existing ideas.

Second, to shed more light on the social-psychological channel, we investigate the causal relationship between surname diversity and the strength of family ties. We have already identified a negative correlation between these two factors, but our examination now moves beyond this to uncover causal effects. We find that an immigration-induced increase in surname diversity leads to weaker family ties. This finding supports the notion that immigration shrinks the relative sizes of family networks, increasing the relative benefits of interacting with non-family members. This effect has downstream effects on impersonal trust, which promote the exchange of ideas among diverse individuals, increasing the probability of innovation.

Finally, we explore geographic spillover effects. Not only does a county's own surname diversity matter, but we also find effects from the diversity of neighboring counties within

a 100-mile radius. This suggests that both local diversity and the diversity in neighboring communities play a significant role in fostering innovation. People can derive inspiration, knowledge, and novel ideas from individuals they regularly observe and interact within their everyday activities.

Taken together, these results indicate that the diversity of social interactions is causally linked to the rate, quality and type of patenting over much of U.S. history. The findings point to an important role of social interactions between diverse individuals in driving innovation.

1.1 Contributions and Related Literature

Understanding the drivers of innovation is central to many lines of research in economics, from endogenous growth (Romer, 1990; Galor and Weil, 2000) to the origins of the industrial revolution (Mokyr, 2002). Here, we focus narrowly on those recent lines of research that connect most closely with our efforts.

To start, our paper picks up on ideas related to the impact of cities on innovation and the role of agglomeration (Carlino and Kerr, 2015; Glaeser et al., 1992; Glaeser, 2011). Research in this area emphasizes the importance of skill complementaries, localized knowledge spillovers and other information transfers. Consistent with our approach, several studies link innovation the formation of immigrant clusters and to a greater diversity of social interactions (Kerr, 2010, 2008). This paper extends these observations and insight both more broadly—across the entire U.S. and back to the mid-19th century—and offers a viable approach to measuring the diversity of social interactions across many contexts.

Further, our findings directly add to an emerging empirical literature on social interaction and innovation, which explores how social institutions and organizations spur innovation. For example, the closure of saloons during Prohibition reduced patenting rates (Andrews, 2023), demonstrating the role of social establishments in innovation. Similar mechanisms operate today, as illustrated by evidence suggesting that the spread of coffee shops spurred innovation (Andrews and Lensing, 2020). Potential tapping the same mechanism, the historical rise of economic societies in Germany reduced information access costs, thereby fostering innovation (Cinnirella et al., 2022). Further back in time, de la Croix et al. (2018) emphasize the role of pre-industrial apprenticeship institutions in Western Europe, including journeymanship, that facilitated the exchange of knowledge and ultimately contributed to Europe’s growth. These studies, among others, underscore the premise that social interactions stimulate knowledge diffusion, contributing to human

capital and innovation-based growth (Akçigit et al., 2018).

Our work intersects with studies exploring how various forms of diversity shape economic prosperity. Galor (2022) emphasizes the importance of cultural diversity for innovation and stresses the importance of diversity as a key factor underlying cross-societal differences in economic prosperity. Using measures of genetic diversity and appropriate proxies, Ashraf and Galor (2013) provide cross-country evidence that genetic diversity fosters innovation while it decreases trust, resulting in an inverse U-shaped relation between genetic heterogeneity and economic prosperity. Conceptually, our measure of surname diversity is related to genetic heterogeneity since, like genes, surnames are transmitted vertically from parent to offspring, and research in population genetics has shown that under certain conditions genetic heterogeneity can be approximated using surname diversity (Barrai et al., 1996). Our paper supports Ashraf and Galor (2013)'s findings by providing causal evidence on the role of social interaction diversity on innovation within U.S. counties. Furthermore, consistent with Ashraf and Galor (2013), we empirically establish that the results are only driven by diversity per se by using surname-fixed effects to rule out that specific surnames or any genes associated with such surnames are influencing the results. That is, when comparing people with the same surnames, those located in counties with greater social interaction diversity are more innovative.

Additionally, previous studies have highlighted the positive effects of birthplace or country-of-ancestry diversity on local economic growth or wages, both within the U.S. (Ottaviano and Peri, 2006; Ager and Brückner, 2013; Docquier et al., 2020; Fulford et al., 2020) and across countries (Alesina et al., 2016). Our use of surname diversity complements Buonanno and Vanin (2017), who used it as a measure of social closure, although their focus was on crime.

Our paper also enriches the literature connecting migration to innovation and economic prosperity (Abramitzky and Boustan, 2017). Drawing on historical data from 1850 to 1920, Sequeira et al. (2020) show how rising flows of immigrants into U.S. counties resulted in faster rates of patenting. Based on an analysis of foreign patents and consistent with the social interaction diversity hypothesis, the authors argue that much of this effect occurred through making native-born Americans more creative—or at least more likely to patent. Similarly, focusing on the period from the mid-1920s to the mid-1960s in the U.S., Moser and San (2020) show how anti-immigration policies in the form of quotas seeking to preserve the ethnic homogeneity of 1890 America reduced the inflow of migrants from Eastern and Southern Europe, which in turn stifled the production of innovations in the scientific fields favored by such immigrants from these countries prior to the quotas. Revealing the importance of social tie diversity, their work finds a 62% decline

in patenting in these particular fields by native-born American scientists. The authors argue that resident scientists lost the mentorship and fresh approaches that inevitably flow in with researchers trained elsewhere. Similarly, [Abramitzky et al. \(2023\)](#) show that quotas did not benefit US-born workers. On the flip side, exploiting the United States' relative openness to immigrants fleeing Germany and Austria prior to World War II, [Moser et al. \(2014\)](#) also demonstrate the impact of Jewish immigrant scientists on U.S. patents. Their analysis reveals not only how refugee chemists stimulated innovation and interest among native-born individuals, but even how their impact reverberated through the social networks to impact the patenting of collaborators of the immigrants' collaborators. Our work supports these findings by highlighting an important channel through which immigration acts on innovation, via increasing the diversity of social interactions.

After the U.S.'s broad immigration quotas were lifted in 1964, [Burchardi et al. \(2021\)](#) provide causal evidence that by the mid-1970s, American innovation was again powerfully fueled by immigrants, now coming from places like Mexico, China, India, Philippines, and Vietnam. In our paper, we follow their instrumental variable approach to provide causal evidence on the role of social interaction diversity for innovation. Importantly, while we do not directly focus on the role of migration, our approach suggests that immigration fuels innovation through its effect on the diversity of social interactions.

2 Concepts and Measurement

In this section, we first describe the recombinative process that arguably underlies much innovation and then highlight supporting lines of evidence. Next, we detail our measure of surname diversity, explain how and why it proxies for diversity of social interactions, and then empirically demonstrate the key conceptual linkages using census information on occupations and ancestral regions along with measures of psychological openness and family ties. Finally, we discuss how we use U.S. patents as a measure of innovation.

2.1 Diversity of Social Interactions and Innovation

The notion that innovation emerges from the recombination of ideas, propelled by social interaction, has venerable lineages in both economics ([Schumpeter, 1983](#)) and history ([Usher, 2013](#)), and has received persistent attention ever since ([Jacobs, 1969](#); [Glaeser et al., 1992](#); [Henrich, 2009](#); [de la Croix et al., 2018](#); [Ridley, 2020](#); [Johnson, 2011](#); [Mokyr, 2015](#); [Olsson and Frey, 2002](#); [Jones, forthcoming](#); [Nunn, 2021](#); [Lucas Jr and Moll, 2014](#); [Akcigit et al., 2018](#)). At the level of a population, the meeting and merging of people and

ideas involves both an informational component—different people possess distinct ideas, approaches, skills, and perspectives—and a psychological component—individuals have to interact and ideally be willing to share their thoughts. Both elements are required since a population of diverse minds that never interact will not generate any recombinations, and a group of cognitive clones who freely interact but all have the same mentality will also fail to generate recombinations. Thus, conceptually, social interaction diversity – the extent to which the free flow of ideas among diverse individuals occurs – should capture the capacity of local populations to generate novel recombinations, some of which will turn into inventions (and for us, patents).

The plausibility of the process—as a hypothesis—is supported by three separate lines of research that (1) suggest a central role of recombination in innovation, (2) reveal the importance of cultural, genetic, disciplinary, and occupational diversity on innovation, and (3) demonstrate the role of sociality on innovation by focusing on institutions that provide opportunities for social interaction or on the role of trust or other psychological factors that shape social interaction and exchange. We briefly discuss each of these research lines.

Empirically, the idea that all or most innovations are recombinations has been explored both within economics and in related fields. Using 1.8 million U.S. patents from 1975-2004 and their citations to other patents, [Acemoglu et al. \(2016\)](#) model the linkages among patents to show how the production of new patents in particular technological domains depends on developments in other related domains. That is, developments in linked technological domains supply the constituent elements or insights for new patents—supplying the fuel for recombination. Complementing this work and using the full U.S. patent database, [Youn et al. \(2015\)](#) and [Akcigit et al. \(2013\)](#) use the detailed patent class codes to show that most patents are indeed recombinations, drawing from different technological classes. Pushing this further, [Clancy \(2018a,b\)](#) fits a recombinative model of innovation that captures both the ‘fishing out’ of obvious recombinations and the innovation-generating impact of each new recombinative idea (patent). The model’s predictions are consistent with patterns found in U.S. Patents. Using scientific citations to assess recombination, [Uzzi et al. \(2013\)](#) find that the highest-impact scientific papers drew on journals rarely referenced by others in the same journal but were, in the main, otherwise highly conventional in their referencing patterns. Finally, using detailed analyses of 21,745,538 lines of computer code based on entries in programming competitions over 14 years, [Miu et al. \(2018\)](#) shows that entries largely copied prior earlier leading entries, which were publicly available, and then added novelty by recombining code drawn from other prior entries. Recombination was, by far, the key element that led to the gradual improvement of these algorithms over time. Finally, [Thagard \(2012\)](#) coded lists of the top

100 most important inventions and scientific discoveries of all time and found them all to involve conceptual recombinations. Based on work in cognitive science, he argues that all creativity arises from recombination based on neuroscientific models of how brains actually form new ideas.

Alongside such evidence for the centrality of recombination for innovation, many researchers have studied the impact of diversity on innovation. As summarized above, researchers have linked measures of genetic, birthplace, academic discipline, and ethnic diversity to measures of innovation (Ashraf and Galor, 2013; Alesina et al., 2016; Docquier et al., 2020; Suedekum et al., 2014; Ozgen et al., 2014; Page et al., 2019). AlShebli et al. (2018), for example, show how both the ethnic and disciplinary diversity of coauthors are linked to scientific impact. Conceptually, our approach suggests that such diversity fuels innovation because these factors are associated with individuals possessing different skills, techniques, knowledge (explicit beliefs), tacit know-how, intuitions and perspectives.

Finally, both social institutions and psychological traits that facilitate the exchange of ideas have been linked to innovation. As noted above, saloons, cafés and knowledge societies have all been linked to innovation (Mokyr, 1995; Andrews, 2023; Andrews and Lensing, 2020; Cinnirella et al., 2022; Henrich, 2020). Similarly, psychological traits that motivate people to (1) tolerate, trust and cooperate with strangers and (2) reveal non-conforming ideas, views and perspectives have been linked to innovation. For example, focusing on trust at the levels of countries and U.S. states, Algan and Cahuc (2014) reveal positive correlations between impersonal trust, based on the generalized trust question, and three measures of innovation. Similarly, using U.S. firm-level data, Nguyen (2021) shows that more trusting CEOs, based on their national ancestry (marked with surnames), generate an uptick in innovation upon their arrival. Finally, using national-level measures, Gorodnichenko and Roland (2016) link individualism to innovation. Conceptually, these social institutions and aspects of psychology foster the flow of ideas among diverse minds, increasing the likelihood of useful recombinations.

The economics literature often describes the flow of information or exchange of skills among minds as ‘knowledge spillovers’ or ‘skill complementarities.’ While our approach here certainly includes these, we think it is important to consider a broader class of cultural and cognitive diversity (Muthukrishna and Henrich, 2016; Page et al., 2019). Across populations, people rely on different languages, thinking styles, decision heuristics, reading preferences, metaphors, attentional biases and ritual practices (Henrich, 2020; Nisbett, 2003). Work in cognitive science, for example, indicates that speaking and thinking in different languages has consequences for people’s perceptions, attention and reasoning (Blasi et al., 2022). Indeed, we show that our measure of surname diversity accounts for

substantial variation in innovation across U.S. counties even when occupational diversity is held constant (Appendix Table B4).

As such, in this paper we take a broad perspective on the link between diversity and innovation. A more diverse local population may increase the diversity of the workforce which then fuels innovation through skill complementarities in production teams. Yet, casual observations of and interactions with people in non-professional contexts can likewise inspire and fuel recombinative processes and more so in highly diverse local populations. So, while we do not disentangle these different channels, the fact that most patents are attributed to single inventors suggests that at the turn of the 19th century, the latter channel—interactions other than within production teams—was likely important. The average number of inventors per (breakthrough) patent is roughly 1.4 and remarkably flat over most of the period of our analysis, as shown in Appendix Figure C7.²

2.2 Operationalizing the Diversity of Social Interactions with Surnames

To conceptualize the diversity of social interactions, consider a subpopulation consisting of N individuals, K groups, and each individual belongs to exactly one group, k , with size N_k such that $\sum_{k=1}^K N_k = N$. Each group carries unique information (e.g., skills, know-how, metaphors), labeled as $s_k \in \{a, b, c, d, \dots\}$ and $s_k \neq s_h$ for all $k \neq h$. In practice, these groups could be extended families, occupations, castes, clans, ethnicities, birth country, or other geographical units, though below we will explain why surnames are a particularly potent partition. When individuals from different groups meet, the likelihood of recombination and innovation increases. Information theory tells us that the average informational content (or the innovation potential) of such a population (Shannon, 1948):

$$E = - \sum p_k \log_2 p_k \quad (1)$$

where $p_k = \frac{n_k}{N}$ is the probability that a person with group affiliation k is drawn and $\log_2 p_k$ is the informational content embedded in this individual (expressed in bits). This is a version of Shannon entropy.

Shannon entropy is a central concept in information theory and is widely used in many scientific disciplines. The term $-\log_2 p_k$ is the self-information of subgroup k and captures

²At the same time, our data do permit us to test whether patenting teams (those with multiple authors) are more likely to be breakthrough innovations when the surname diversity of the team is greater. Indeed, as Appendix Table C2 shows, if we focus only on patents with multiple inventors (2+), those with only inventors who carry the same surname are less likely to generate breakthroughs. This holds with both year and technology category fixed effects. While this cannot explain our overall effects, because most patents are solos, it does suggest that surname diversity operates at the level of the individual patent.

the level of surprise (or the informational content of a specific outcome). The negative log reflects that rarely-encountered groups carry more surprise (or more information) compared to more frequent encounters. To arrive at Shannon entropy, the self-information is weighted by the probability of its occurrence and summed over all possible outcomes. For example, if the population only consists of one group k , the outcome of a draw is not surprising (the outcome can be predicted perfectly), the entropy is 0, and thus, no recombinations can arise through social interaction. On the other hand, entropy for a population with a fixed number of groups is maximized if all groups are equal in size. An individual who goes out for a random social interaction in such a diverse population is most likely to observe someone different from themselves. A random draw will thus have more informational content (in expectation), which is reflected by higher entropy.³

Conceptually, we expect surname diversity to proxy for both the informational and psychological aspects of social interaction diversity. Informationally, we hypothesize—and later provide evidence—that much of the important informational diversity occurs among groups identified by surnames. The idea here is that people who hold the same surnames share much information among themselves, which can occur through vertical as well as both horizontal and oblique forms of cultural transmission (Cavalli-Sforza and Feldman, 1981; Bell et al., 2019). At the same time, people who hold different surnames differ in their socialization and their social networks as well as their cultural heritage; ultimately, this results in the clustering of many kinds of information—skills, know-how, customs, languages and thinking styles—within groups defined by surnames.

Socially and psychologically, we hypothesize—and later test empirically—that surname diversity is associated with weaker family ties and lower trust in strangers. Consistent with the literature on the impact of kinship and family ties on sociality (Enke, 2019; Alesina and Giuliano, 2014; Schulz et al., 2019), we suspect this holds because rising surname diversity both constrains people’s ability to meet their needs within their own (shrinking) group and implies more opportunities for exchanges outside of one’s group.

Notably, this proposal may run counter to the intuitions of many readers, who suspect that rising diversity will create social miscoordinations or activate in-group biases that inhibit social interactions, thereby thwarting increases in the opportunities for recombinative innovation created by informational diversity. First, if this does indeed play an

³In economics, a Herfindahl measure, which population geneticists call Isonomy, is frequently used to capture diversity. For our purposes, however, Shannon’s entropy has several advantages to conceptualize informational diversity and has favorable mathematical properties (Carcassi et al., 2021). In particular, a Herfindahl approach underweights the importance of rare surnames (p_k vs. $\log_2 p_k$), i.e., under the assumption that surnames carry unique pieces of information, rare surnames are more “valuable”—they carry a higher expected surprise. Empirically, the two measures are highly correlated in our setting ($\rho = 0.80$, see Appendix Table B1).

important role, we should detect it empirically by observing a negative or concave relationship between our surname entropy measure and innovation. As we will show, we do not see this. Second, conceptually, it is important to distinguish between finely-grained partitions like surnames from larger-scale divisions like tribes, ethnicities and nationalities. As the size of extended families and surname groups shrink relative to the entire population, interactions outside of one’s own group—even if riskier—become both more necessary and more profitable. So, unlike larger-scale measures of diversity, greater surname diversity is more likely to also create greater opportunities for beneficial exchange with those outside one’s group. Relatedly, a large body of anthropological and other evidence indicates that the scales of preferential interactions occur at the boundaries of languages, ethnicities, social norms and religions (Handley and Mathew, 2020; White et al., 2021; Desmet et al., 2017). Consistent with this, we will see that diversity measures based on national origins and race both do not create the same convex relationship with innovation, and at least in the case of race diversity it is inversely U-shaped. Indeed, consistent with our view of surname diversity (but not other forms of diversity), we will show that, *ceteris paribus*, surname diversity is associated with possessing more sociality as captured by generalized trust. This suggests that it is important to distinguish between different forms of diversity.

2.3 U.S. Surname Diversity 1850-1940

Based on the U.S. census, we calculate a measure of surname diversity for each U.S. county, i , for each time-period, t :

$$E_{i,t} = - \sum_{k=1}^K p_{k,i,t} \log_2 p_{k,i,t} \quad (2)$$

where $p_{k,i}$ denotes the fraction of people with surname k who live in county i and K denotes the total number of distinct surnames.

Our data source is the full-count Integrated Public Use Microdata Series (IPUMS) compiled by Ruggles et al. (2021) and available on the NBER servers. We use the nine waves from 1850 to 1940 which contain the variable `namelast` of all individuals and county identifiers.⁴ We implement the Philips (1990) phonetic algorithm *metaphone* to deal with misspellings in the name string. Following Burchardi et al. (2021), we also obtain

⁴We harmonize all historical Census data to the 2000 boundaries of U.S. counties using the Ferrara et al. (2021) crosswalks. Specifically, we use the M4 weights that account for urban and rural areas and topographic suitability. We use 2000 as the reference year because the patent dataset is geocoded to 2000 county boundaries. The harmonization procedure sometimes results in counties with very few people, predominantly in the Midwest and West, and for Census years before 1900. As a remedy, we winsorize all harmonized variables from the lower tail at the 1% level.

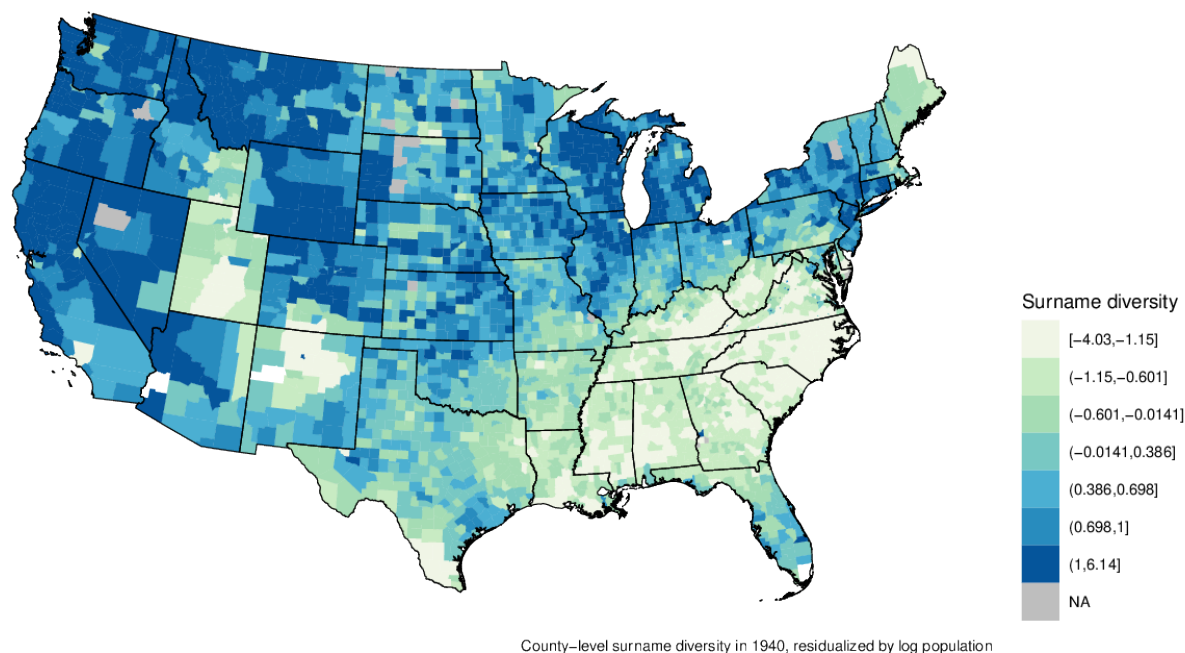


Figure 1: Geographic variation in surname diversity in 1940

Notes: The figure shows standardized county-level surname diversity residualized by log county population in 1940.

the variables `age` and `yr immig` (the year of immigration) to estimate surname diversity for the mid-decades 1895, 1905, 1915, and 1925 by removing all individuals who were born or immigrated after the mid-decade. Ideally, we would also remove all individuals who moved to the county after the midyear, but this information is not available.

Figure 1 maps our measure of surname diversity for U.S. counties in the year 1940, partialing out population size.⁵ Clear geographical patterns emerge. While counties in California and most of the Northeast score high on surname diversity (independent of population size), Utah and the Southern states score substantially lower, i.e., they are more homogeneous with regard to surnames and hence less culturally diverse. Surname diversity correlates moderately with more common diversity measures based on country of birth and occupation (see Appendix Table D1). This moderate correlation provides support for the notion that surname diversity not only encompasses variation within these broader categories but also captures a so far unobserved slice of diversity at the level of families.

⁵Appendix Figure D1 reports the spatial variation in surname diversity unconditional on log county population.

2.4 Surnames Capture Relevant Social Interactional Diversity

Here, we empirically establish that our measure of surname diversity captures both aspects of the informational dimension—i.e., surnames are indicative of occupation and ancestral origins—and the psychological aspect—i.e., surname diversity predict contemporary survey measures of impersonal trust and historical strength of family ties. Our endeavor here is not to establish causal linkages, but merely to reveal the kinds of empirical relationships one would expect if surname diversity indeed captures diversity of social interactions, i.e., the degree that exchange of ideas among diverse people occurs.

Recent work in the social mobility literature demonstrates that surnames capture unique skills, socialization, and know-how (Clark, 2014; Güell et al., 2015; Barone and Mocetti, 2021). This is not surprising given that surnames are indicative of the ancestral heritage of origin regions, professions, and (family) lineages. And even though very frequent surnames will capture family-, or profession-specific traditions to a lesser degree, these surnames nevertheless still capture unique knowledge. For example, “Smith” and “Johnson” are the most frequent surnames in the 1880 census. The Smiths outnumber the Johnsons by a factor of about 1.65 times. However, among individuals who reported blacksmith as their occupation, there are 2.46 times more Smiths than Johnsons. This makes sense given that the surname "Smith" has a long history of association with metalworking and blacksmithing, and the data show that it is still a relatively more common surname among metalworkers in 1880.

To gain a systematic sense of the degree to which surnames reflect unique knowledge in our data set, we calculated Herfindahl concentration measures that capture how strongly surnames cluster in several different domains including occupations, country or region of origin, and technology categories of patents. The construction of the concentration measure for each domain proceeds in two steps. For example, in the case of occupation, we first calculate a normalized Herfindahl index for each surname across all occupations. This gives us a measure of how strongly a specific surname clusters in occupations. We normalized this measure such that it is zero in case of a uniform surname distribution and one if the surname is only found within one single occupation. Second, we average the surname-specific Herfindahl indices across all surnames, weighted by the number of people with a given surname. This averaged index reveals the overall surname concentration in occupations based on the U.S. population. In the same way, we construct the concentration measures for the other domains.

Table 1 reports the surname concentration in the first row. For occupations, columns 1 to 4 reveal substantial concentrations of specific surnames in 1880 and 1940 based on both the full sample (columns 1 and 2) and an immigrant sub-sample (columns 3 and 4).

Table 1: Surnames cluster in occupations, birthplaces, and patent fields

Sample:	Occupation				Country of origin		Region of origin	USPO tech category	
	All		Immigrants				Germans	Inventors	
	1880	1940	1880	1940	1880	1940	1880	1880-9	1940-9
Year:	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Surname	0.117	0.045	0.097	0.065	0.393	0.227	0.189	0.092	0.068
U.S. county of residence	0.171	0.071	0.153	0.080	0.288	0.130	0.159	0.014	0.017
Country of origin	0.120	0.043	0.096	0.060					
Age	0.154	0.049	0.101	0.048					

Notes: The table reports normalized Herfindahl indices, where larger values indicate greater concentration. The indices are calculated as the average Herfindahl indices of the variable in the header computed for each value of the variables on the left. For example, we calculate the normalized Herfindahl index of occupations for each surname and then average over all surnames using the number of individuals with a given surname as weights. Column 1 (2) includes all individuals in the 1880 (1940) census. Columns 3 and 5 (4 and 6) include all immigrants in the 1880 (1940) census. Column 7 restricts the sample to German immigrants in 1880. Column 8 (9) includes all inventors of patents issued from 1880 to 1889 (1940 to 1949). In column 7, we use the 31 subnational regions of origin of German immigrants recorded by the Census (the variable `bp1d` with codes 45301 to 45361).

Focusing on ancestral origin, Columns 5 to 7 reveal immigrants’ surname concentration across originating countries (columns 5 and 6) and regions within Germany (column 7); the latter is restricted to German immigrants because fine-grained subregional birthplace data are available for this country. Lastly, columns 8 and 9 report surname concentrations across technology categories (USPO).⁶

The concentration measures in Table 1 show a consistent pattern: surnames are not distributed uniformly. The concentration indices are well above zero. For example, in the year 1880, two people with the same surname have a roughly 12% (above chance) probability of holding the same occupation (out of 249 possible occupations), or two same-surname immigrants have about a 39% probability of being from the same country of origin. Further, Column 7 shows that surname even indicates the region of a country where people come from. Same-surname immigrants from Germany have a 19% probability of being from the same inner-German region (out of 31 regions).

Having established that certain surnames concentrate in occupations, originating countries, regions, and patent categories, we can put the concentration indices into context by comparing them to measures of residence-county (row 2), country of origin (row 3), and

⁶The patent data set does not allow us to uniquely identify inventors. Hence, we are unable to detect inventors who file multiple patents in the same technology category, which will bias the concentration upwards. We still report this statistic because this bias is likely small given the low level of regional clustering in this variable (row 2), where we would expect a similar upward bias if regional mobility among inventors is not very high.

age (row 4) concentration. In the year 1880, occupations are relatively more concentrated in counties compared to surnames, though this difference markedly narrows in the year 1940 (row 2). Surnames are substantially more indicative of originating countries and regions compared to immigrants' residence counties. Compared to country of origin and age, surnames are about equally indicative of occupation (rows 3 and 4).

A potential concern is that, although surnames may often be nested within coarser categories like country of origin, regional birthplace and race, there have been historical processes that muddy this hierarchical nesting. For example, many formerly enslaved Africans carry the European surnames of their enslavers (Cook et al., 2022). Surname diversity may thus underestimate the diversity stemming from African cultural heritage.

To address this, we construct a more finely-grained measure and check how it relates to our main indicator. This measure creates additional 'surname categories' based on race-surname combinations. For example, the number of white 'Jacksons' enters the diversity indicator as a separate category from the number of black 'Jacksons'. Similarly, we calculate a surname diversity indicator that further differentiates along country of birth. Appendix table B1 shows that the main surname diversity indicator in 1940 is almost perfectly correlated with those more finely-grained diversity measures.⁷

Moving now to focus on the social-psychological side of surname diversity, we analyze the correlation between surname diversity in 1940 with responses to the generalized trust question in the General Social Survey (waves from 1972 to 2016). Since the observations of surname diversity in 1940 precede the trust data by several decades, we supplement this analysis with an analysis of the correlations between surname diversity and the strength of family ties from 1860 to 1940 (Raz, 2023).

Appendix Figure C1 depicts the relationship between surname diversity in 1940 and generalized trust (1972-2016 in the GSS). A positive relationship emerges both in the raw data and conditional on state fixed effects and log county population, indicating that greater surname diversity is associated with more trust. Examining the strength of this relationship over time, the bottom panel in Appendix Figure C1 reveals that this relationship is quite stable from 1870 to 1940 with a magnitude of roughly 0.1.

Appendix Figure C2 shows a negative relationship between surname diversity and the strength of family ties.⁸ Individuals in counties with higher surname diversity tend

⁷Furthermore, we obtain similarly high correlation coefficients between the main surname diversity indicator and indicators that are based on (i) phonetically uncorrected surnames, (ii) surnames of men only, and (iii) surnames of whites only.

⁸Here, following Raz (2023), the strength of family ties is captured by the first principal component of four underlying variables: (i) the divorce-to-marriage ratio, (ii) the share of elderly people living without a relative, (iii) the share of people living with at least one person who is not their relative, and (iv) the mean size of families.

to have weaker family ties than those in counties with lower surname diversity in 1940. The relationship becomes even stronger when we control for log county population size and state fixed effects. Investigating the strength of this relationship over time, Appendix Figure C2 shows how the magnitude of the relationship increases over time, from around -0.2 between 1860-80 to roughly -0.60 by 1940.

Notably, none of these patterns exist for the more conventional measures of diversity based on country of birth and race (Appendix Figure C3 to C6). These findings underline the distinctive nature of surname diversity compared to these other types of diversity.

Overall, the analyses presented in this section provide *prima facie* evidence that our measure of surname diversity captures both the informational and social components of recombinative innovation.

2.5 Measuring Innovation

To measure innovation, we rely on patent data. Our first measure is the total number of patents per 1,000 individuals. We calculate this measure for each U.S. county for 5 and 10-year periods from 1850 to 1940 based on the Comprehensive Universe of U.S. Patents (CUSP) data set compiled by Berkes (2018). The primary source of this data set is Google Patents supplemented with information from other sources.

To address the concern that this indicator captures a host of patents that do not meaningfully push the knowledge frontier forward (and is thus noisy), we rely on a second indicator, breakthrough patents (per 1,000 individuals). This indicator is created by Kelly et al. (2021) and captures highly important patents. Applying textual analysis, they compare a given patent to previous and subsequent ones. Breakthrough patents fulfill two criteria: they are distinct from previous patents (capturing novelty) and have a high similarity to subsequent ones (capturing impact). For more details on both indicators, see Appendix A.

Of course, as a measure of innovation, patent data suffers from a number of well-known and widely discussed shortcomings (Griliches, 1990; Moser, 2013; Lerner and Seru, 2022), including the fact that many innovations are never patented, industries vary widely in their tendencies to patent, types of inventions vary in how easily they can be patented, and that more patenting in a particular technology category may inhibit innovation rates. Though, as we show in the appendix our results hold in specifications that control for technology category fixed effects (see Appendix Table B14).

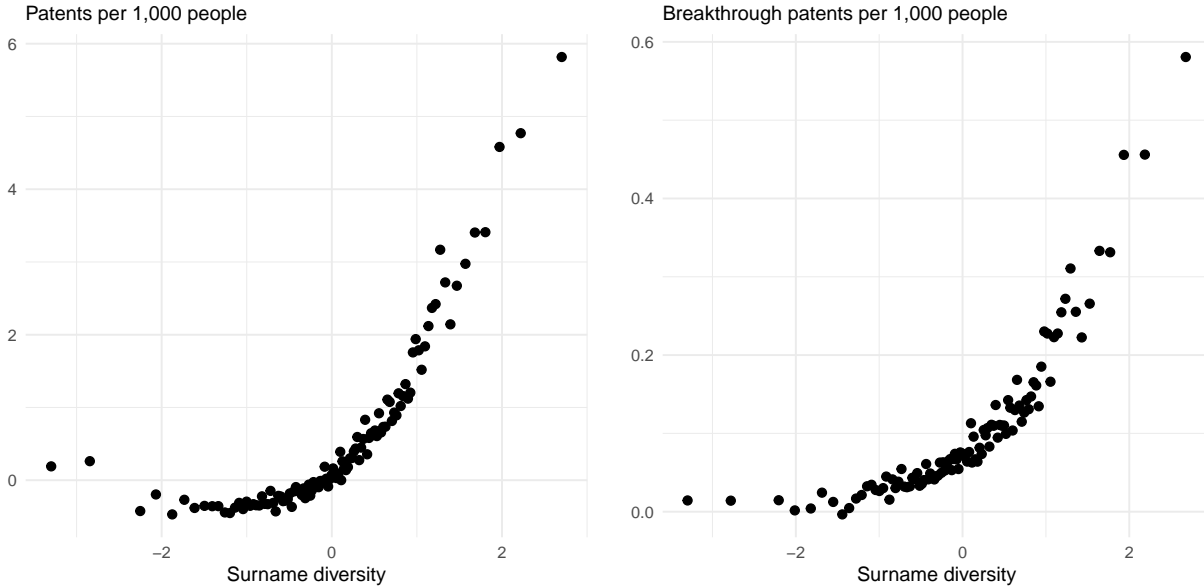


Figure 2: Bivariate relationships between surname diversity and our patent-based innovation measures

Notes: County-level data from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850 and residualized by census year fixed effects. Binscatter plot created using the R package written by Cattaneo et al. (2019).

3 Least-squares regressions

Our analysis is structured into four sections. Here, in section 3, we begin by reporting the estimates from a battery of least-squares regressions to establish a robust positive relationship between sociocultural diversity and innovation using the full decadal data set from 1850 to 1940. In the next section 4, we provide causal evidence on this relationship based on our instrumental variable strategy that exploits the pseudo-random nature of immigration flows into U.S. counties. Section 5 supplies a variety of additional robustness and sensitivity checks including the use of surname-fixed effects. Finally, in Section 6, we delve into the mechanism underlying our main result.

Displaying the relationship between surname diversity and innovation visually, Figure 2 reveals tight, strong relationships. These patterns indicate few outliers and offer no hint of a ‘hump shape’; instead, the plots suggest convexity in the relationship between surname diversity and innovation—that is, greater surname diversity is associated with relatively more (breakthrough) patents per capita. By contrast, the relationship between country of birth diversity and innovation is more linear (Appendix Figure B3), and the relationship between race diversity and innovation is concave and hump-shaped (Appendix Figure B4).

Table 2 reports the least-squares estimates of the relationship between surname diversity and innovation of both patents per 1,000 people and breakthrough patents per 1,000 people in Panels A and B, respectively. Column 1 reports the bivariate relationships between the two innovation outcomes and surname diversity (including period year fixed effects to control for trends affecting surname diversity and innovation across all counties). We find positive and significant relationships: a one standard deviation higher surname diversity in a county is associated with approximately 1.76 more patents per capita (Panel A) and 0.15 breakthrough patents per capita (Panel B). Relating these coefficients to the mean of the dependent variable suggests that a one standard deviation increase in surname diversity is associated with roughly 80% more patents and slightly more breakthrough patents.

For ease of exposition, Table 2 reports linear specifications which do not take into account the non-linearities which are visible in Figure 2. However, the size of the estimated linear approximations (expressed in percent) is very similar in case we apply the inverse hyperbolic sine (IHS) transformation to the outcome variable (see Appendix Table B6). A visual check confirms that such a transformation of the dependent variable linearizes the relationship to some degree, which is displayed in Figure B6. Throughout the paper we report on the untransformed dependent variable because the IHS transformation is not scale invariant (Aihounton and Henningsen, 2020). In any case, all results are qualitatively and quantitatively similar when using IHS transformed dependent variables.

In column 2 of Table 2, we report estimates of the relationship between innovation and (standardized) country of origin diversity, a more conventional measure of cultural diversity (e.g., Ottaviano and Peri, 2006; Ager and Brückner, 2013; Alesina et al., 2016). We find a qualitatively similar relationship with innovation suggesting that country-of-birth diversity likewise fuels the recombinative processes.

In column 3, we include both diversity measures in the regression. The coefficient of surname diversity is significant; surname diversity hence predicts innovation conditional on controlling for country of origin diversity. This suggests that it captures a distinct type of diversity.

In column 4, we include the share of immigrants from each of the 59 originating countries consistently recorded in the census data between 1850 and 1940 to account for the role of immigration from specific countries. We find that the inclusion of these variables hardly changes the coefficients on surname diversity. This suggests that recent immigration does not confound the relationship between innovation and surname diversity.

In column 5, we add period-state fixed effects. The coefficients on surname diversity are virtually unchanged; thus, the relationship between surname diversity and innovation

Table 2: Least-squares estimates: surname diversity and innovation from 1850s to 1940s

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i>						
	Patents per 1,000 people (mean = 2.26, sd = 2.58)					
Surname diversity	1.76*** (0.175)		0.735*** (0.156)	0.705*** (0.160)	0.856*** (0.176)	0.502*** (0.168)
Country of origin diversity		1.64*** (0.093)	1.10*** (0.181)	1.18*** (0.175)	1.07*** (0.232)	0.165 (0.176)
R ²	0.503	0.550	0.574	0.607	0.691	0.864
<i>Panel B:</i>						
	Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.24)					
Surname diversity	0.154*** (0.021)		0.064*** (0.012)	0.059*** (0.012)	0.064*** (0.015)	0.044** (0.021)
Country of origin diversity		0.144*** (0.013)	0.098*** (0.019)	0.106*** (0.018)	0.113*** (0.024)	0.004 (0.015)
R ²	0.416	0.459	0.480	0.524	0.619	0.787
Immigrant shares by country of origin (59 shares)				✓	✓	✓
Period fixed effects	✓	✓	✓	✓		
Period-State fixed effects					✓	✓
County fixed effects						✓
Observations	22,222	22,222	22,222	22,222	22,222	22,222

Notes: The table reports estimates of least-squares regressions of innovation outcomes on surname diversity and immigrant diversities. In Panel A (Panel B), the outcome is number of (breakthrough) patents issued in a given period per 1,000 people. The unit of observation is a county-period from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

is not driven by persistent or time-varying differences across states, including North-South differences (e.g., [Cook, 2014](#)).

In the final specification in column 6, we include county fixed effects to examine the relationship between *changes* in surname diversity and innovation. This specification harnesses only the variation within each county over time that is distinct from other counties within the same state and is not due to immigration. As reported in column 6, the estimates for surname diversity shrink but still remain large: a one standard deviation increase in surname diversity is associated with an increase of roughly 23% in both patents and breakthrough patents per 1,000 people. Interestingly, while the coefficient on surname diversity remains large and well estimated, the coefficients on country-of-origin diversity are no longer significant when surname diversity is held constant.

Although our specification focuses on innovations per capita, there remains a concern that this specification is restrictive and does not adequately capture population-scale

effects. To address this, Appendix Table B2 controls for counties' population size. The coefficients on surname diversity remain large, well-estimated, and robust across specifications for both of our measures of innovation.

To gain further insights into the relationship between surname diversity and other diversity measures, Appendix Table B4 reports on regressions that also include diversity measures based on race and occupations. In all specifications, we find that surname diversity is consistently and at least weakly significantly positively associated with patents—even in the case of occupational diversity which likely overlaps with surname diversity.

In mechanism section 6.4, we discuss education and show that surname diversity is a significant predictor of innovation controlling for the average education of a county's population in the year 1940.

Lastly, to check the stability of the relationship between surname diversity and innovation over time, we estimate a coefficient for surname diversity for each decade separately (in a specification that otherwise parallels column 1 of Table 2). Appendix Figure B1 shows that all the estimated coefficients for surname are positive, significant and sizeable.

4 IV estimates

The previous section established that surname diversity strongly correlates with the quantity and quality of patenting across counties from 1850 to 1940. However, reverse causality or unobserved factors may bias these estimates. For example, it is possible that migrants (with rare surnames) tend to move to innovative counties and that this immigration then increases surname diversity (though we controlled for immigration in some specifications). Similarly, highly-skilled people, who are more likely to innovate, might prefer to live in more diverse counties. In both cases, we would observe a relationship between surname diversity and innovation even if no causal relationship exists.

Given such concerns, we implement an estimation strategy that isolates quasi-random variation in surname diversity. We use this variation as an instrumental variable (IV) for actual surname diversity to provide estimates that allow us to shed light on the causal effect of a change in surname diversity on innovation in a county-level panel from 1900 to 1940. As we will detail below, the benefit of this strategy is that it accounts for reversed causality and local unobservable factors that simultaneously affect innovation and surname diversity.

4.1 Construction of the Instrument and Estimating Equations

The central idea underlying our IV strategy arises from the observation that immigration is a major determinant of the change in counties' surname diversity. This fact allows us to build on recent advances in the immigration literature to isolate variation in surname diversity that is independent of any unobserved determinants of innovation that may bias our estimates. This approach implies that we are estimating the local average treatment effect (LATE) of the change in surname diversity on innovation that is induced by immigration to the US. As such, we are not capturing the average treatment effect that would also take changes in surname diversity into account that stem from births, deaths and domestic migration.

We note that it is crucial to distinguish between estimating the immigration-induced LATE of diversity on innovation and investigating the impact of immigration (measured by the inflow of immigrants) on innovation. While we utilize the changes in surname composition resulting from immigration, it does not imply a straightforward monotonically increasing relationship between immigration inflow and surname diversity. In some instances, an influx of immigrants might decrease surname diversity, while in others, it may increase it. To reinforce this point, we demonstrate that even when controlling for the migration-induced changes in population, our IV results on surname diversity remain consistent (see section 4.2 for a detailed discussion on identification). The key aspect of our IV strategy is that migration affects surname composition (in complex ways), and we can capture quasi-random variation in immigration to create an instrument of surname diversity.

Specifically, we adapt the IV strategy developed by [Burchardi et al. \(2019\)](#) to our context. The construction of this instrument requires two steps. First, we isolate quasi-random variation in the stock $N_{k,i}^t$ of each surname k residing in county i in census year t based on historical migration patterns. Second, we compute the instrument for surname diversity by calculating diversity based on these (predicted) quasi-random stocks of surnames $\hat{N}_{k,i}^t$. We now explain the details of these two steps.

Step 1: Isolating Quasi-random Variation in Counties' Surname Stocks We adopt [Burchardi et al. \(2019\)](#)'s historical push-pull approach to isolate quasi-random variation in the composition of surnames in U.S. counties. The approach assumes that a combination of push and pull factors jointly determines the allocation of immigrants with given surnames to counties and that the historical interactions of these two factors generate quasi-random variation in surname stocks that persists over time.

Empirically, the push factor is summarized by the total number of immigrants with a

given surname entering the U.S. during a given period; the pull factor is the attractiveness of a county in this period, which is operationalized by the share of immigrants (out of all immigrants entering the US) who settle in this county in the same period. These two factors vary over time, and their interactions, which we can trace back to 1880, generate quasi-random variation in a county's distribution of surnames.

Formally, we predict the stock of people $N_{k,i}^t$ (in thousands) with surname k residing in county i in year t by estimating the following zero-stage equation:

$$N_{k,i}^t = \delta_i + \delta_{k,r(i)} + \sum_{\tau=1880}^{t-1} b^\tau \underbrace{I_{k,-r(i)}^\tau}_{\text{Push}} \underbrace{\frac{I_{i,-k}^\tau}{I_{-k}^\tau}}_{\text{Pull}} + \sum_{\tau=1880}^{t-1} d^\tau \frac{I_{i,-k}^\tau}{I_{-k}^\tau} + u_{i,k}, \quad (3)$$

where i indexes counties, k denotes surnames, t indexes census years from 1900 to 1940, including the midyears, and $r(i)$ denotes the census region containing county i . The variable $I_{k,-r(i)}^\tau$ is the push factor in the period ending in year τ (1880, 1895, 1900, 1905, 1910, 1915, 1920, 1925, 1930). It is given by the total number of migrants (in thousands) with surname k who arrive in the U.S. during this period and settle *outside* the region containing county i . The pull factor captures the relative attractiveness of a specific county i during the period ending in τ . It is given by the share of migrants a county attracts $\frac{I_{i,-k}^\tau}{I_{-k}^\tau}$, where $I_{i,-k}^\tau$ is the total number of migrants who settle in county i during this period and who do not have surname k , and $I_{-k}^\tau = \sum_i I_{i,-k}^\tau$ is the total number of migrants who settled in the U.S. during the same period and who do not have surname k .⁹

Core to the identification strategy are the historical interactions between the push and pull factors in each period τ (up to period $t - 1$). We estimate a coefficient for this interaction, b^τ , for each period stretching back to the year 1880 (the earliest period for which we have data on immigrants or their parents). That is, equation (3) attributes the stock of each name in a county (in a given year t) to the past inflow of migrants who are allocated according to the push-pull factors over the course of several decades.

In addition to the push-pull factors, equation (3) also includes the term $\sum_{\tau=1880}^{t-1} d^\tau \frac{I_{i,-k}^\tau}{I_{-k}^\tau}$, i.e., the relative share of migrants which settle in a county in each period τ . This term captures the time-varying relative attractiveness of a county in the past. It isolates the push-pull instruments from county-level conditions that drew migrants in each period

⁹We follow [Burchardi et al. \(2019\)](#) and estimate equation (3) using a leave-out approach. That is, we exclude migrants with surname n from the pull factor (denoted by $-k$), and we exclude the census regions r that county i is located in from the push factor (denoted by $-r(i)$). This leave-out approach ensures that our estimates are not driven by the settlement outcomes of migrants with surname k who settled in region $r(i)$. We note, though, that at the level of surnames, this is likely less of a concern because the fractions of surnames relative to all migrants are small.

τ up to $t - 1$, which may affect innovation still in period t . Moreover, δ_i , denotes county fixed effects. They remove any time-invariant factors that make specific counties more attractive to all migrants. $\delta_{k,r(i)}$ are name-region fixed effects. They remove time-invariant unobserved factors that may make specific census regions more attractive to migrants with certain surnames.

Based on equation (3) we estimate the coefficients \hat{b}^τ for each period τ and then calculate the predicted stocks of name k in county i at time t as

$$\hat{N}_{k,i}^t = \sum_{\tau=1880}^{t-1} \hat{b}^\tau \left(I_{k,-r(i)}^\tau \frac{I_{i,-k}^\tau}{I_{-k}^\tau} \right)^\perp$$

where \hat{b}^τ is the estimate of b^τ from equation (3), and $\left(I_{k,-r(i)}^\tau \frac{I_{i,-k}^\tau}{I_{-k}^\tau} \right)^\perp$ are residuals of a regression of the push-pull interaction, $I_{k,-r(i)}^\tau \frac{I_{i,-k}^\tau}{I_{-k}^\tau}$, on δ_i , $\delta_{k,r(i)}$ and $\frac{I_{i,-k}^\tau}{I_{-k}^\tau}$. This residualization ensures that the predicted stock of each name $\hat{N}_{k,i}^t$ relies on the component of the push-pull factors that is orthogonal to the control variables included in equation (3). This orthogonalization is particularly useful with regard to $\frac{I_{i,-k}^\tau}{I_{-k}^\tau}$, because it ensures that the instrument is orthogonal to the past attractiveness of a county, which could be driven by an underlying factor that also determines innovation decades later.

Step 2: Calculating the Instrument for Surname Diversity In step 2, we compute the instrument for surname diversity by applying the entropy formula on the predicted stock of each surname $\hat{N}_{k,i}^t$:

$$\widehat{\text{Surname diversity}}_i^t = - \sum_k \left(\frac{\hat{N}_{k,i}^t}{\sum_k \hat{N}_{k,i}^t} \log \left(\frac{\hat{N}_{k,i}^t}{\sum_k \hat{N}_{k,i}^t} \right) \right)$$

We repeat steps 1 and 2 eight times to obtain an instrument for diversity in each of the eight periods (ranging from $t = 1900$ to $t = 1940$) that form part of our panel analysis.

Step 3: IV Estimating Equations We implement our IV procedure using 2SLS. The equations are given by equations (4) and (5), where equation (4) is the first stage and equation (5) is the second stage.

$$\widehat{\text{Surname diversity}}_i^t = \gamma \widehat{\text{Surname diversity}}_i^t + \mu_{t,s(i)} + \mu_i + v_i^t \quad (4)$$

$$Y_i^t = \beta \widehat{\text{Surname diversity}}_i^t + \alpha_{t,s(i)} + \alpha_i + \varepsilon_i^t \quad (5)$$

where i indexes counties, s states, and t census years (including the midyears). Y_i^t is county i 's number of (breakthrough) patents per 1,000 people in the five-year period starting in t . Surname diversity $_i^t$ is county i 's surname diversity in t ; and $\widehat{\text{Surname diversity}}_i^t$ county i 's predicted surname diversity in t described above. The coefficient β is our main interest.

Equations (4) and (5) also include state-period fixed effects, $\mu_{t,s(i)}$ and $\alpha_{t,s(i)}$, and county fixed effects, μ_i and α_i . By including these fixed effects, β is estimated from changes in surname diversity within the same county over time while controlling for persistent and time-varying differences across states.

In addition, several specifications include county-specific linear time trends such that β captures the relationship between deviations in the changes in diversity and innovation within counties over time relative to their overall trend. Comparing the estimates of these specifications to the baseline estimates of equation (5) provides another exogeneity check of the instrument. If the estimates remain similar, this suggests that the instrument is orthogonal to persistent or gradually growing county-level confounding factors.

4.2 Identification

Our identification strategy is valid if $\widehat{\text{Surname diversity}}_i^t$ is truly exogenous in the specification of equation (5). A sufficient condition for this to hold is

$$\left(\begin{array}{c} I_{k,-r(i)}^\tau \\ I_{i,-k}^\tau \end{array} \right)^\perp \perp \varepsilon_i^t.$$

It requires that any factor that affects counties' innovation in t is independent of the interaction between the orthogonalized historical push-pull factors. If this condition holds, the predicted stocks of surnames are exogenous to innovation (Step 1), and so is the instrument for diversity (Step 2).

An important question regarding the validity of this empirical strategy is whether past push-pull factors are independent of a county's future innovative capacity. For example, it is possible that migrants preferred to settle in counties that were more innovative in the past, likely increasing their diversity, and those same counties are subsequently still more innovative. This possibility would give rise to reverse causality. More generally, (persistent) unobserved factors may determine both the past pull factors and future innovation, which may create a correlation between the push-pull instrument and the error term.

In their paper, [Burchardi et al. \(2019\)](#) detail why this possibility is unlikely given the substantial variation in the push-pull factors over time and space. Empirically, we address this concern in three ways: First, we orthogonalize our push-pull instrument with regard

to the historical attractiveness of a county as captured by the fraction of immigrants who settled there over time (see our zero-stage equation (3)). Consequently, our IV estimates do not reflect unobserved persistent factors that had already manifested themselves in immigrants' past settlement decisions. Second, in several specifications, we control for county-specific linear time trends. To the degree that these linear time-trends capture persistent unobserved factors, they will mitigate concerns of estimation bias. Lastly, and most importantly, we conduct a falsification exercise and regress previous-period innovation on subsequent surname diversity. We do not find any evidence for reverse causality, i.e., a shock to surname diversity is statistically unrelated to previous-period innovation. Instead, we find that a shock to diversity has a lasting impact on innovation (see Section 5.1). Therefore, our estimates are unlikely to be driven by reverse causality or persistent unobserved confounders.

Another concern is that unobserved individual characteristics co-determine settlement patterns and innovation. For example, people with a high (unobserved) propensity to innovate may prefer to settle in relatively more diverse counties. In this case, the observed relationship between diversity and innovation would be due to settlement preferences of individuals with high innovative capacity and not due to diversity per se. The IV approach addresses this concern because the predicted surname stocks in a county are solely determined by the interaction of the historical push and pull factors, i.e., the allocation of migrants to counties does not rest on individual preferences.¹⁰ This push-pull instrument is orthogonal to county fixed effects and surname-region fixed effects. Thus, unobserved stable settlement preferences of people with a certain surname cannot bias the estimates. In addition, in section 5.4, we further address this concern by devising specifications with surname-fixed-effects. This specification will absorb any genetic, environmental, or acquired characteristic embodied in surnames and, thus, captures the 'pure' diversity effect that is independent of the type of information embedded in surnames. Taken together, our results are unlikely to be biased due to individual characteristics that co-determine settlement patterns and innovation.

A final concern is the possibility that there is a direct effect of immigration (other than through diversity) that confounds the estimates. Though, while we establish that immigration impacts counties' surname composition, conceptually there is not a simple

¹⁰The exclusion restrictions could be violated if the push-pull factors primarily reflect the migration decisions (= preferences) of people with a specific surname. Yet, this is unlikely because any specific surname makes up only a tiny fraction of all people entering the U.S. in a given period (the push factor) and a small fraction of migrants settling in a county (the pull factor). Nevertheless, we follow [Burchardi et al. \(2019\)](#) and report leave-out estimators such that the push factor does not contain individuals with surname n and the pull factor does not contain regions in which a county is located in $r(i)$.

monotonously increasing relationship between immigration and surname diversity. The extent to which immigrants impact surname diversity depends on the surname composition of the immigrants in comparison to the local population. For example, if migrants predominantly hold the same names as the dominant local groups then immigration will decrease surname diversity. The fact that the least-squares estimates in Table 2 hardly change when we added controls for the share of immigrants from different countries likely attest to these considerations and provides evidence that the relationship between surname diversity and innovation is not confounded by a direct effect of immigration that goes through other channels than diversity.

In the IV-specifications, we avoid controlling for endogenous variables such as the number of immigrants and our push-pull IV-strategy (which is based on surname stocks and not flows) does not allow for a straightforward way to instrument it. Rather, we rely on alternative strategies. First, as explained above, our instrument is orthogonal to each county's immigration history, which mitigates endogeneity concerns. Second, we conduct a robustness check and control for instrumented population. The instrument for population is likewise based on the push-pull approach (calculated by taking the sum over all predicted stocks of surnames). As such, the IV-specification controls for the *migration-induced* change in counties' population which is a proxy variable for migration. We find that our IV-results hold controlling for instrumented population. Finally, we also conduct a robustness check and estimate IV specifications based on a sub-sample that does not contain any immigrant innovators. Even though this is a conservative approach (see Section 5.5 for details), this restricted sample of only US-born individuals likewise provides IV-evidence that diversity increases innovation.

4.3 Zero-stage Estimates

We report the zero-stage estimates of equation (3) in Table 3. These estimates allow us to obtain predicted stocks for each surname in each county, which we will use to compute the diversity instrument. In total, we estimate equation (3) eight times, once for each period from 1900 to 1940. The reported standard errors are clustered at the surname level.

The results indicate that we identify variation in the stock of surnames based on the push-pull factors stretching across the full range of periods in our sample. For example, the estimates reported in column 8 suggest that push-pull factors as far back as 1880 and all the way up to 1930 are significant predictors of the stock of surnames in 1940. Qualitatively, our results are broadly similar to [Burchardi et al. \(2019\)](#), who estimate the push-pull factor at the level of originating country (and not surnames). For example, they

Table 3: Zero-stage panel estimates

	No. of people in county i with surname n in year:							
	1900 (1)	1905 (2)	1910 (3)	1915 (4)	1920 (5)	1925 (6)	1930 (7)	1940 (8)
$I_{n,-r(i)}^{1880} \times \frac{I_{-n,i}^{1880}}{I_{-n}^{1880}}$	3.181*** (0.113)	3.192*** (0.121)	3.683*** (0.130)	3.982*** (0.095)	4.386*** (0.113)	4.516*** (0.136)	4.945*** (0.108)	5.030*** (0.101)
$I_{n,-r(i)}^{1895} \times \frac{I_{-n,i}^{1895}}{I_{-n}^{1895}}$	1.962*** (0.274)	2.336*** (0.336)	2.902*** (0.311)	2.276*** (0.125)	2.588*** (0.120)	4.071*** (0.199)	4.471*** (0.243)	4.551*** (0.290)
$I_{n,-r(i)}^{1900} \times \frac{I_{-n,i}^{1900}}{I_{-n}^{1900}}$		-2.405 (3.101)	-0.262 (3.497)	-6.766** (2.632)	-5.988** (2.742)	-9.868*** (3.227)	-10.297*** (3.513)	-8.711*** (3.312)
$I_{n,-r(i)}^{1905} \times \frac{I_{-n,i}^{1905}}{I_{-n}^{1905}}$			12.702*** (0.885)	17.984*** (0.911)	20.613*** (1.110)	27.031*** (1.401)	30.116*** (1.219)	32.922*** (1.250)
$I_{n,-r(i)}^{1910} \times \frac{I_{-n,i}^{1910}}{I_{-n}^{1910}}$				14.937*** (2.540)	16.972*** (2.918)	24.672*** (3.346)	26.990*** (3.213)	28.736*** (2.902)
$I_{n,-r(i)}^{1915} \times \frac{I_{-n,i}^{1915}}{I_{-n}^{1915}}$					8.268*** (0.803)	8.023*** (0.830)	9.419*** (0.546)	10.294*** (0.602)
$I_{n,-r(i)}^{1920} \times \frac{I_{-n,i}^{1920}}{I_{-n}^{1920}}$						3.659* (2.198)	5.087*** (1.181)	9.319*** (1.584)
$I_{n,-r(i)}^{1925} \times \frac{I_{-n,i}^{1925}}{I_{-n}^{1925}}$							25.957*** (1.293)	31.662*** (1.490)
$I_{n,-r(i)}^{1930} \times \frac{I_{-n,i}^{1930}}{I_{-n}^{1930}}$								-29.396*** (2.731)
Observations	5,933,320	7,336,530	7,806,098	7,365,155	7,448,013	7,977,906	8,053,917	8,811,918
R ²	0.710	0.687	0.698	0.701	0.711	0.660	0.706	0.690
County fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
Surname-Region fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
$I_{-n,i}^t/I_{-n}^t$ controls	✓	✓	✓	✓	✓	✓	✓	✓

Notes: This table reports OLS estimates for the specification described in equation (3), corresponding to step 1 of the instrument construction. An observation is a surname-county in a period from 1900 to 1940. Standard errors clustered at the surname level. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

likewise obtain a negative coefficient for the interaction for the period ending in 1930, a period with a high degree of out-migration.

Based on the results of Table 3 we calculate the predicted (and orthogonalized) stock of each surname in each county for each of the eight periods. Finally, we compute the instrument for surname diversity by applying the entropy formula to the predicted surname stocks.

4.4 IV Estimates

We now turn to the IV estimation. Table 4 reports first-stage, reduced-form, and second-stage estimates. Starting with the first-stage estimates reported in Panel D, we find that the instrument for surname diversity is strongly correlated with actual surname diversity, with a Kleibergen-Paap F -statistic of around 51 in our baseline specification in columns

2 and 5. The F -statistic shrinks to roughly 28 if we additionally control for county-specific linear time trends in columns 3 and 6. The point estimates imply that a one standard deviation increase in the instrument is associated with 0.43 standard deviation greater surname diversity in the baseline (columns 2 and 5) and with a 0.39 standard deviation greater surname diversity when conditioning on county-specific linear time trends (columns 3 and 6). Taken together, the first-stage relationship of the instrument is highly significantly related to surname diversity, and the F -statistics of the excluded instrument in all specifications are above the conventional thresholds commonly used in the literature.

Appendix Figure B7 reports binscatter plots that show the first-stage relationship between the instrument and actual surname diversity, both with and without controls for county-specific time trends. They demonstrate that the relationship is strong, linear and not driven by a small set of observations.

We next turn to the estimates relating surname diversity to innovation. Table 4 reports the estimates for both main outcome variables—patents per 1,000 people (columns 1 to 3) and breakthrough patents per 1,000 people (columns 4 to 6). Panel A reports least-squares estimates for comparison, Panel B reports reduced-form estimates and Panel C reports the IV estimates. All specifications control for county fixed effects and period fixed effects (column 1 and 4) or state-period fixed effects (columns 2 to 3 and 5 to 6). In addition, the specification reported in columns 3 and 6 controls for county-specific linear time trends. We report estimates for weighted regressions, with the weights determined by a county’s population in 1900. The reported standard errors are clustered at the state level. The least-squares estimates reveal a significantly positive relationship between surname diversity and both patents and breakthrough patents. In column 2 and 4, a one standard deviation increase in a county’s diversity is associated with 1.5 more patents per 1,000 people, roughly a 74% increase relative to the sample mean, and 0.15 more breakthrough patents per 1,000 people, a 104% increase in breakthrough patents.¹¹ These estimates are similar to those reported in Table 2 in the previous section even though they cover a different time span.

Turning to the IV specifications, Panel B shows statistically significant reduced-form relationships between the dependent variables and our historical push-pull instrument (predicted surname diversity based on the zero stage). Appendix Figure B8 reports partial correlation plots to visualize these relationships. Finally, Panel C reports on the IV estimates. The coefficients for surname diversity are positive and highly significant in all specifications and for both innovation outcomes. The estimates in the baseline

¹¹To obtain these percentages, we divide the coefficients by the sample means reported in the table.

Table 4: Panel estimates of the effect of surname diversity on innovation

	Patents per 1,000 people (mean = 2.04, sd = 2.60)			Breakthrough patents per 1,000 people (mean = 0.14, sd = 0.24)		
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Least-squares estimates</i>						
Surname diversity	1.528*** (0.332)	1.511*** (0.361)	1.319*** (0.320)	0.183*** (0.042)	0.146*** (0.043)	0.131*** (0.046)
<i>Panel B: Reduced-form estimates</i>						
Surname diversity (push-pull IV)	0.686*** (0.204)	0.773*** (0.165)	0.734*** (0.173)	0.090*** (0.018)	0.085*** (0.018)	0.080*** (0.024)
<i>Panel C: Instrumental-variable estimates</i>						
Surname diversity	1.542*** (0.382)	1.794*** (0.378)	1.902*** (0.543)	0.202*** (0.043)	0.197*** (0.055)	0.208** (0.090)
Kleibergen-Paap <i>F</i> -statistic	63.280	51.050	28.341	63.280	51.050	28.341
<i>Panel D: First-stage estimates</i>						
Surname diversity (push-pull IV)	0.445*** (0.056)	0.431*** (0.060)	0.386*** (0.073)	0.445*** (0.056)	0.431*** (0.060)	0.386*** (0.073)
County fixed effects	✓	✓	✓	✓	✓	✓
Period fixed effects	✓			✓		
State-Period fixed effects		✓	✓		✓	✓
County-specific linear time trends			✓			✓
Observations	23,660	23,660	23,660	23,660	23,660	23,660

Notes: The table reports least-squares, reduced-form, and instrumental-variable (IV) estimates for the specifications described in equation (5) and first-stage estimates for equation (4). An observation is a county in a period from 1900 to 1940. Observations are weighted by county population in 1900. The endogenous variable is county-level surname diversity in t . In columns 1 to 3, the dependent variable is number of patents filed in the county in the five-year period starting in t divided by county population size in 1900. In columns 4 to 6, the dependent variable is number of breakthrough patents filed in the county in the five-year period starting in t divided by county population size in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

specifications (columns 2 and 5) suggest that a one standard deviation increase in a county’s surname diversity increases patents (per 1,000 inhabitants) by about 88% relative to the sample mean and breakthrough patents by about 141%. When we control for county-specific linear trends, the effect of surname diversity on patents remains remarkably stable at around 93% for patents (column 3) and around 149% for breakthrough patents (column 6), bolstering our confidence that the instrument for diversity is orthogonal to persistent or gradually growing county-level confounding factors. Overall, the estimates suggest that surname diversity has large positive effects on both the quantity and quality of innovation.

Comparing the least-square (Panel A) with the IV estimates (Panel C) reveals that the latter are slightly larger in magnitude, but simple t -tests suggests that the differences are not significant. We believe that a likely explanation for this observation is that the IV estimates capture the local average treatment effect (LATE), i.e., the average effect of a change in diversity that is due to immigration. This contrasts with the average treatment effect that would take into account the overall change in diversity (e.g., including change brought about by births, deaths, and internal migration). It is plausible that the component of surname diversity driven by recent immigration has a larger impact on innovation than other sources of surname diversity—even considering two people who share the same surname—because recent immigrants might have more distinct ideas, knowledge, skills and perspectives that are more valuable in the recombination process.¹²

5 Robustness and Sensitivity Checks

We now check for the robustness of our estimates in five ways: we examine (1) a placebo test that regresses past innovation on surname diversity, (2) controlling or instrumenting for population size in our IV estimation, (3) heterogeneity across the four major census regions, (4) an approach that uses surname fixed effects to assess whether surname-specific traits affect the results, and (5) the removal of immigrant innovators from our analyses, so that we only consider how greater surname diversity impacts innovation among native-born Americans. We also examine the robustness of our results to the use of alternative procedures to construct an instrument for surname diversity. These results are reported in Appendix Section B.4.

¹²Another explanation for the smaller least-square estimates is that they are biased downward due to negative selection, i.e., migrants settled in regions that were doing economically worse. [Sequeira et al. \(2020\)](#) discuss the relevant literature and conclude that the negative selection of immigrants into poor regions is indeed a possibility.

5.1 Placebo Tests and Reverse Causality

A potential concern with our results is a form of reverse causality, i.e., that innovative counties attract relatively more immigrants which then increases diversity. This possibility is unlikely because our instrument is orthogonal to a county’s past attractiveness as captured by the (time-varying) shares of immigrants who settled in a county over the course of several decades (see section 4.1 for details). Moreover, we examined specifications that include county-specific linear time trends, which absorb the effects of trending unobserved factors associated with innovation and migration.

Nevertheless, to further challenge the validity of our instrument, we conduct a placebo exercise of whether contemporaneous diversity affects past innovation activity. A significant estimate in such a regression would be evidence of reverse causality, i.e., that innovative counties attract immigrants, and this increases surname diversity. In columns 1 and 2 of Table 5, we regress our measures of innovation for one (in $t - 1$) and two (in $t - 2$) periods prior to our measure of surname diversity in the current period (in t). Column 3 replicates our IV specification in which we regress innovation on same-period diversity (reported in Table 4, column 5). The estimates demonstrate that there is no significant positive relationship between previous periods’ innovation and subsequent diversity. This is the case for both patents (Panel A) and breakthrough patents (Panel B). When we regress patenting in t (i.e., that occurs between t and $t + 1$) on diversity in period t , the coefficients increase in size and become significantly different from zero (column 3). In short, we find no support for the importance of reverse causality in our identification strategy.

In addition, we investigate the persistence of the impact of a diversity shock on innovation by regressing the one-period lead (innovation in $t + 1$, column 4), the two-period lead ($t + 2$, column 5), and the three-period lead ($t + 3$, column 6) on surname diversity (in t). The estimates in columns 4 to 6 of Table 5 suggest that the impact of diversity on patenting in the two following periods is significantly positive and decreasing in magnitude over time (panel A). The effect on breakthrough patents shows similar persistence (panel B). We fail to detect an effect in period $t + 3$, suggesting that the effect of surname diversity on innovation persists for about 15 years.

5.2 Sensitivity of Estimates to Population Size

In our IV specifications, population enters through the dependent variables, which is per capita (per the population in 1900), and each county is weighted based on the population size in 1900. Here, we follow the literature (e.g., [Burchardi et al., 2021](#)) and choose a time-invariant base year because population growth is likely endogenous—innovative

Table 5: Placebo test and persistence: IV estimates of the effect of surname diversity on past and future innovation

	$t - 2$	$t - 1$	t	$t + 1$	$t + 2$	$t + 3$
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i>						
	Patents per 1,000 people					
Surname diversity	0.132 (0.323)	-0.047 (0.218)	1.902*** (0.543)	1.348*** (0.451)	1.000** (0.429)	0.097 (0.173)
<i>Panel B:</i>						
	Breakthrough patents per 1,000 people					
Surname diversity	-0.019 (0.067)	-0.041 (0.083)	0.208** (0.090)	0.252** (0.116)	0.179** (0.073)	-0.098 (0.073)
Kleibergen-Paap F -statistic	24.092	23.922	28.341	20.842	19.700	16.600
County fixed effects	✓	✓	✓	✓	✓	✓
State-Period fixed effects	✓	✓	✓	✓	✓	✓
County-specific linear time trends	✓	✓	✓	✓	✓	✓
Observations	17,743	17,746	23,660	17,746	17,743	14,785

Notes: The table reports IV estimates of the leads and lags of innovation outcomes on surname diversity for the specifications described in equation (5). Columns 1 and 2 use the two-period and one-period lag of the dependent variables, respectively. Column 3 repeats the baseline specification (contemporaneous values of the dependent variables). Columns 4 to 6 use the one-period, two-period and three-period lead of the dependent variables, respectively. Observations are county-periods and weighted by county population in 1900. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

regions may attract more people. A concern with this specification is that it may not be able to adequately capture scale effects due to an increasing population. The least-squares estimates reported in Appendix Table B2 suggest that this is unlikely to be the case—the coefficients on surname diversity are not sensitive to controlling for counties’ population size. Here, we assess the robustness of our IV estimates to the inclusion of population size in the specification.

The results of this robustness check are reported in Appendix Table B8. We use actual population in Panel A, predicted population in Panel B, and we instrument population with predicted population in Panel C. The construction of predicted population parallels the construction of the surname-diversity instrument: based on the historical push-pull interaction of the zero stage (see Section 4), we obtain the predicted stock of each surname in a county in a given period. By aggregating these stocks, we obtain quasi-random

estimates of county population at a given point in time.

Appendix Table B8 shows that our IV results are robust to controlling for population. All estimates hardly change compared to our baseline estimates.

5.3 Estimates for Sub-regions

In Appendix Table B10, we explore whether the relationship between diversity and innovation holds in each of the four major U.S. regions (the Midwest, Northeast, South, and West). The estimates suggest that this is indeed the case; all region-specific estimates are positive and in almost all cases the coefficients are statistically significant. The coefficients for the Midwest tend to be larger than for the other regions, however the estimates are not precise enough to draw strong conclusions about regional differences. We also note that some IV estimates may suffer from a weak first stage.

5.4 Surname Fixed Effects

Another potential concern with the interpretation of our findings is that surname-specific traits, such as abilities, interests, or knowledge, drive innovation rather than the diversity of these traits. For example, Clark (2014) and Barone and Mocetti (2021) find that rare surnames are proxies for the vertical transmission of traits, and these traits might affect innovation. We assess this concern by estimating specifications that include surname fixed effects which absorb any surname-specific traits. This requires us to change the unit of observation from county-period to surname-county-period. The estimating equations are given by equations (6) and (7), where equation (6) is the first stage and equation (7) is the second stage.

$$\text{Surname diversity}_i^t = \gamma \widehat{\text{Surname diversity}_i^t} + \mu_{t,s(i)} + \mu_i + \mu_{t,k} + v_{i,k}^t \quad (6)$$

$$Y_{i,k}^t = \beta \text{Surname diversity}_i^t + \alpha_{t,s(i)} + \alpha_i + \alpha_{t,k} + \varepsilon_{i,k}^t \quad (7)$$

where i indexes counties, s states, t census years (including the midyears), and k surnames. As before, $\text{Surname diversity}_i^t$ is county i 's surname diversity in t , and $\widehat{\text{Surname diversity}_i^t}$ is county i 's predicted surname diversity in t . $Y_{i,k}^t$ now is the number of (breakthrough) patents filed by people with surname k residing in county i per 1,000 of these residents in the five-year period starting in t . For example, 36,984 individuals with the surname Johnson resided in Cook County (IL) in 1940 and filed about 105 patents and 11 breakthrough patents between 1940 and 1944. Therefore, the innovation outcomes vary at the surname-

Table 6: Surname fixed effects

	Patents per 1,000 people (mean = 2.06, sd = 44.86)				Breakthrough patents per 1,000 people (mean = 0.16, sd = 9.29)			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
<i>Panel A: Least-squares estimates</i>								
Surname diversity	1.694*** (0.454)	1.736*** (0.466)	1.667*** (0.504)	1.111*** (0.297)	0.262*** (0.079)	0.267*** (0.081)	0.195** (0.074)	0.132*** (0.039)
<i>Panel B: Reduced-form estimates</i>								
Surname diversity (push-pull IV)	0.639*** (0.236)	0.657*** (0.230)	0.655*** (0.223)	0.815** (0.337)	0.137*** (0.039)	0.140*** (0.040)	0.126*** (0.028)	0.117** (0.053)
<i>Panel C: Instrumental-variable estimates</i>								
Surname diversity	1.477*** (0.458)	1.524*** (0.451)	1.571*** (0.483)	2.205** (0.941)	0.316*** (0.091)	0.325*** (0.094)	0.301*** (0.091)	0.317** (0.156)
Kleibergen-Paap F-statistic	61.575	61.793	49.541	31.909	61.575	61.793	49.541	31.909
<i>Panel D: First-stage estimates</i>								
Surname diversity (push-pull IV)	0.433*** (0.055)	0.431*** (0.055)	0.417*** (0.059)	0.370*** (0.065)	0.433*** (0.055)	0.431*** (0.055)	0.417*** (0.059)	0.370*** (0.065)
County fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
Period fixed effects	✓				✓			
Surname-Period fixed effects		✓	✓	✓		✓	✓	✓
State-Period fixed effects			✓	✓			✓	✓
County-specific linear time trends				✓				✓
Observations	30,416,997	30,416,997	30,416,997	30,416,997	30,416,997	30,416,997	30,416,997	30,416,997

Notes: The table reports least-squares, reduced-form, and instrumental-variable (IV) estimates for the specifications described in equation 7 and first-stage estimates for equation 6. An observation is a surname in a given county in a period from 1900 to 1940. Observations are weighted by the surname population in a given county in the year 1900. In columns 1 to 3, the dependent variable is number of patents filed by individuals with surname n residing in county i in the five-year period starting in t divided by surname population size in county i in 1900 (multiplied by 1,000). The dependent variable in columns 4 to 6 is the corresponding number of breakthrough patents. Standard errors are two-way clustered on states and surnames and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

county-period level, while surname diversity remains defined at the county-period level.¹³ Crucially, we can now include surname-period fixed effects, denoted by the parameter $\alpha_{t,k}$, which implies we non-parametrically control for surname-specific traits across periods (i.e., traits specific to all individuals named Johnson in 1940). The remaining parameters and variables are as in equations (6) and (7). As before, the coefficient of interest is β . Observations are weighted by the number of people in a county carrying the surname in the year 1900. Standard errors are clustered in two ways, on states and surnames.

The results reported in Table 6 show that the estimates are highly significant in all specifications and the inclusion of surname fixed effects (in columns 2 to 4 and 6 to 8) hardly changes the estimates. The IV specification suggests that a one standard deviation increase in a county's surname diversity increases patent filings per 1,000 people with the same surname by around 1.6 patents (column 3) and 0.3 breakthrough patents (column 7).

¹³Consequently, the number of observations increases because they are now determined by the total number of unique surnames in a given county. Consistent with our baseline specification, we normalize the number of patents and breakthrough patents by the surname population in the year 1900. If a surname does not exist in a given county in 1900, we drop it from the sample. See Appendix A for all the details on how we construct the sample.

Expressed as percent (by relating it to mean values of the dependent variable) this is an increase of 76% in patents (column 3) and 188% in breakthrough patents (column 7). Even though we changed the unit of observation, the estimates are comparable to the county-level estimates for patents reported in Table 4, though they are larger for breakthrough patents. Crucially, in this specification, the surname fixed effects additionally ensure that the estimates are independent of any unobserved surname-specific characteristics.¹⁴

5.5 Immigrant Innovators

A final concern is that immigration per se, rather than diversity, biases our estimates. To tackle this in the county-level least-squares specifications reported in Table 2 (Section 3), we regress innovation on diversity while controlling for immigrant shares (separately by each country of origin). The coefficients on surname diversity are remarkably robust to the inclusion of these control variables. This finding is evidence against the possibility that immigration per se biases our estimates because, for example, immigrants may be more highly skilled, entrepreneurial, or possess novel patentable knowledge. Furthermore, the concern that our IV estimates are confounded by immigration is mitigated by the fact that our instrument is orthogonal to counties' past immigration history (see Section 4.1 for details on the construction of the instrument) and are robust to controlling for migration-induced population changes (see Section 5.2).

Nevertheless, to drive home the point that immigration per se does not fully explain our findings, we conduct a further robustness check. In the surname-fixed-effects specification, the unit of observation is surname-county-period. This allows us to drop all names for which we know that at least one immigrant also holds this name.¹⁵ We then estimate equations (6) and (7) with this non-immigrant sample. We acknowledge that this approach is coarse because we drop many native innovators who happen to share their surname with an immigrant living in the same county at the same time. Consequently, the sample size decreases substantially by roughly 44% from 30.4 to 17.1 million observations. This will skew the estimates towards zero because the sample average of the number of patents per 1,000 surnames in the full sample is 2.06, while it is only 0.64, or a third, in the non-immigrant sub-sample. This difference is even more pronounced for breakthrough

¹⁴Using a similar approach (for details see Appendix Section B.5), we estimate specifications that include patent technology class fixed effects which absorb any patent class-specific factors. These specifications address the concern that systematic variation in patenting practices across technologies may bias our results. The estimates, reported in Appendix Table B14, show that the results hold with this patent class fixed effects specification. This suggests that patent class-specific factors do not bias our estimates.

¹⁵We are not able to drop immigrants directly—this individual-level data (in contrast to the surname-level data) would require us to match the patent data to the Census at the individual level. This is currently not feasible at a reasonable level of accuracy.

Table 7: Disentangling diversity from immigration mechanism: Natives subsample

	Patents per 1,000 people (mean = 0.64, sd = 43.76)				Breakthrough patents per 1,000 people (mean = 0.04, sd = 9.61)			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
<i>Panel A: Least-squares estimates</i>								
Surname diversity	0.317*** (0.072)	0.380*** (0.077)	0.331*** (0.081)	0.301*** (0.086)	0.062*** (0.021)	0.063*** (0.020)	0.047*** (0.015)	0.019* (0.010)
<i>Panel B: Reduced-form estimates</i>								
Surname diversity (push-pull IV)	0.213*** (0.064)	0.260*** (0.065)	0.237*** (0.065)	0.246*** (0.066)	0.044*** (0.015)	0.045*** (0.015)	0.039*** (0.013)	0.022* (0.012)
<i>Panel C: Instrumental-variable estimates</i>								
Surname diversity	0.333*** (0.095)	0.406*** (0.095)	0.375*** (0.100)	0.390*** (0.114)	0.069*** (0.023)	0.070*** (0.023)	0.062*** (0.020)	0.035* (0.019)
Kleibergen-Paap F-statistic	460.250	457.336	343.932	217.417	460.250	457.336	343.932	217.417
<i>Panel D: First-stage estimates</i>								
Surname diversity (push-pull IV)	0.638*** (0.030)	0.639*** (0.030)	0.631*** (0.034)	0.631*** (0.043)	0.638*** (0.030)	0.639*** (0.030)	0.631*** (0.034)	0.631*** (0.043)
County fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
Period fixed effects	✓				✓			
Surname-Period fixed effects		✓	✓	✓		✓	✓	✓
State-Period fixed effects			✓	✓			✓	✓
County-specific linear time trends				✓				✓
Observations	17,061,661	17,061,661	17,061,661	17,061,661	17,061,661	17,061,661	17,061,661	17,061,661

Notes: The table reports least-squares, reduced-form, and instrumental-variable (IV) estimates for the specifications described in equation 7 and first-stage estimates for equation 6. An observation is a surname in a given county in a period from 1900 to 1940. The sample is restricted to observations with no immigrants in a given surname-county-period cell. Observations are weighted by the surname population in a given county in year 1900. Standard errors are two-way clustered on states and surnames and reported in parentheses. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

patents (mean of 0.16 in the full vs. 0.04, or a fourth, in the sub-sample). As such, in this robustness check we are less concerned about the effects size—the approach is too coarse—but are interested in whether we can still detect significant effects in this sub-sample.

Table 7 shows that in this sub-sample, the point estimates are all positive and significant; the only exception is the weakly significant estimate in column 8. As expected, the effect sizes are smaller compared to the ones reported in Table 6. Yet, expressed in relation to the mean of the dependent variable they are sizeable. We estimate a 59% increase in patents per 1,000 natives in column 3 and a 155% increase in breakthrough patents per 1,000 natives in column 7. These overall significant and sizable estimates—in a sample that rules out that a patent was filed by an immigrant—are further evidence of diversity’s causal impact on innovation.

6 Mechanisms

Conceptually, following much prior work, we see innovation as arising from the recombination of ideas due to the social interaction occurring among diverse minds. To explore this more deeply, we consider (1) if surname diversity spurs recombination of existing technologies, (2) the impact of surname diversity on the strength of family ties, (3) the extent of spatial spillovers across counties and (4) the interplay between surname diversity and education.

6.1 Novel Combinations Patents

One key question about our results concerns what types of patents are generated by greater surname diversity. Does surname diversity tend to encourage the patenting of novel technology types, potentially imported from other counties? Or, does such diversity encourage either the creation of novel recombinations of existing technologies or the refinement of such combinations? Building on the approach taken by [Strumsky et al. \(2011\)](#) and [Akcigit et al. \(2013\)](#), we use the more than 140,000 technology codes assigned by the United States Patent and Trademark Office (USPTO) to categorize each patent into three distinct types: (1) novel technologies, (2) novel combinations, and (3) reuse/refinement combinations. Patents are considered a novel technology if, for a given county, any of its technology codes appear for the first time in that county in the grant year of the patent. If the patent does not include such novel codes, it is considered as a novel combination if it includes a unique pairwise combination of technologies that appear for the first time in the county and grant year. Any remaining patents are classified as reuse/refinement combinations.

Table 8 show our results. Overall, the impact of surname diversity on patents per capita is similar in magnitude across our three categories, with the influence on patenting novel combinations showing the largest positive effects. These results suggest that diversity increases innovation in different forms, including recombinations of existing technologies.

6.2 Strength of Family Ties

To more directly test for the social-psychological component of our hypothesized mechanism, we now examine the impact of surname diversity on the strength of family ties. In Section 2, we established that the two variables are highly correlated. Here, we investigate causality by estimating regressions that replicate our baseline equation (5) but use the strength of family ties as the dependent variable.

Table 8: Patent types: novel technologies, novel combinations, and reuse/refinements

	Novel technology patents per 1,000 people (mean = 1.46, sd = 1.72)		Novel combination patents per 1,000 people (mean = 0.09, sd = 0.16)		Reuse/refinement patents per 1,000 people (mean = 0.31, sd = 0.47)	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Least-squares estimates</i>						
Surname diversity	1.124*** (0.257)	0.913*** (0.197)	0.059*** (0.019)	0.034*** (0.012)	0.165*** (0.045)	0.135*** (0.045)
<i>Panel B: Reduced-form estimates</i>						
Surname diversity (push-pull IV)	0.526*** (0.113)	0.475*** (0.108)	0.040*** (0.009)	0.028*** (0.008)	0.094*** (0.025)	0.095*** (0.033)
<i>Panel C: Instrumental-variable estimates</i>						
Surname diversity	1.221*** (0.248)	1.229*** (0.348)	0.094*** (0.029)	0.072** (0.032)	0.219*** (0.079)	0.246* (0.124)
Kleibergen-Paap <i>F</i> -statistic	51.050	28.341	51.050	28.341	51.050	28.341
County fixed effects	✓	✓	✓	✓	✓	✓
State-Period fixed effects	✓	✓	✓	✓	✓	✓
County-specific linear time trends		✓		✓		✓
Observations	23,660	23,660	23,660	23,660	23,660	23,660

Notes: The table reports least-squares, reduced-form, and IV estimates for the specifications described in equation (5) of the effect surname diversity on three different types of patents per 1,000 people. We classify a patent as a novel technology if any of the technology codes listed on the patents appear in the grant year of the patent for the first time in the county where the inventor resides. We classify a patent as novel combination of technologies that have previously been used in the county if a pairwise combination of technologies is observed in the county for the first time. The third classification is reuse/refinement patents, where all pairwise combinations of technologies have been observed previously in the county. An observation is a county in a period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

The results are reported in Table 9. The estimates suggest a negative effect that becomes statistically significant when we control for state-period fixed effects (column 2). Hence, greater surname diversity leads to weaker family times, likely because it limits people’s ability to meet their needs within their family and comes with more opportunities for exchanges with unrelated individuals.

6.3 Spatial Spillovers

People acquire inspiration, knowledge and ideas from others they frequently observe and interact with in their daily activities. Mostly, these will be people living in proximity and, consequently, we expect that local diversity drives innovation. (Carlino and Kerr, 2015).

Here, we investigate how local the relationship between diversity and innovation is and whether there are spillovers from nearby counties. To do so, we compute surname diversity among individuals residing in surrounding regions successively further away from the county. Specifically, for each county i at time t , we pool the individuals and

Table 9: Panel estimates of the effect of surname diversity on strength of family ties

	Strength of family ties (mean = -0.14, sd = 0.80)		
	(1)	(2)	(3)
<i>Panel A: Least-squares estimates</i>			
Surname diversity	-0.038 (0.075)	-0.183*** (0.037)	-0.470*** (0.087)
<i>Panel B: Reduced-form estimates</i>			
Surname diversity (push-pull IV)	-0.162* (0.091)	-0.264*** (0.093)	-0.328*** (0.063)
<i>Panel C: Instrumental-variable estimates</i>			
Surname diversity	-0.365 (0.229)	-0.613** (0.277)	-0.851*** (0.266)
Kleibergen-Paap <i>F</i> -statistic	63.349	51.127	28.360
County fixed effects	✓	✓	✓
Period fixed effects	✓		
State-Period fixed effects		✓	✓
County-specific linear time trends			✓
Observations	23,639	23,639	23,639

Notes: The table reports least-squares, reduced-form, and IV estimates for the specifications described in equation (5) with the strength of family ties as the outcome variable. An observation is a county in a period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table 10: Spillover analysis

	Patents per 1,000 people (mean = 2.04, sd = 2.59)				Breakthrough patents per 1,000 people (mean = 0.14, sd = 2.21)			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
<i>Panel A: Least-squares estimates</i>								
Surname diversity	1.474*** (0.382)	1.471*** (0.382)	1.469*** (0.382)	1.180*** (0.297)	0.145*** (0.046)	0.144*** (0.046)	0.144*** (0.046)	0.114** (0.046)
Surname diversity (< 100 miles)	0.872 (0.573)	0.879 (0.579)	0.855 (0.604)	1.505*** (0.487)	0.066 (0.064)	0.067 (0.064)	0.065 (0.067)	0.187** (0.076)
Surname diversity (100 < 200 miles)		-0.522 (0.446)	-0.524 (0.452)	-1.278* (0.661)		-0.053 (0.042)	-0.053 (0.043)	-0.221** (0.093)
Surname diversity (200 < 300 miles)			-0.205 (0.383)	0.121 (0.458)			-0.017 (0.042)	-0.024 (0.084)
<i>Panel B: Reduced-form estimates</i>								
Surname diversity (push-pull IV)	0.780*** (0.161)	0.780*** (0.161)	0.783*** (0.162)	0.751*** (0.176)	0.086*** (0.019)	0.086*** (0.019)	0.087*** (0.019)	0.084*** (0.026)
Surname diversity (push-pull IV, < 100 miles)	0.815*** (0.280)	0.814*** (0.284)	0.837*** (0.304)	0.653*** (0.179)	0.076** (0.031)	0.075** (0.031)	0.081** (0.034)	0.059 (0.036)
Surname diversity (push-pull IV, 100 < 200 miles)		-0.016 (0.174)	0.020 (0.188)	0.303 (0.301)		-0.007 (0.017)	0.002 (0.021)	0.023 (0.037)
Surname diversity (push-pull IV, 200 < 300 miles)			0.106 (0.166)	0.235 (0.194)			0.026 (0.026)	0.046 (0.040)
<i>Panel C: Instrumental-variable estimates</i>								
Surname diversity	1.674*** (0.357)	1.668*** (0.358)	1.668*** (0.357)	1.855*** (0.573)	0.188*** (0.058)	0.187*** (0.058)	0.188*** (0.057)	0.209** (0.099)
Surname diversity (< 100 miles)	2.533** (1.202)	2.558** (1.204)	2.552** (1.190)	1.458** (0.660)	0.217 (0.143)	0.222 (0.146)	0.227 (0.151)	0.115 (0.157)
Surname diversity (100 < 200 miles)		-0.279 (0.552)	-0.286 (0.552)	0.698 (1.029)		-0.054 (0.068)	-0.047 (0.072)	0.050 (0.132)
Surname diversity (200 < 300 miles)			-0.055 (0.679)	0.628 (0.872)			0.053 (0.115)	0.162 (0.168)
F-statistic: Surname diversity	88.191	59.925	44.633	26.361	88.191	59.925	44.633	26.361
F-statistic: Surname diversity (< 100 miles)	17.373	14.392	12.706	37.700	17.373	14.392	12.706	37.700
F-statistic: Surname diversity (100 < 200 miles)		13.568	13.652	25.706		13.568	13.652	25.706
F-statistic: Surname diversity (200 < 300 miles)			6.150	10.814			6.150	10.814
County fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
State-Period fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
County-specific linear time trends				✓				✓
Observations	23,093	23,093	23,093	23,093	23,093	23,093	23,093	23,093

Notes: The table reports least-squares, reduced-form, and instrumental-variable (IV) estimates of regressions of innovation outcomes on surname diversity. The unit of observation is a county-period from 1900 to 1940 (including the midyears). The table sequentially adds surname diversity in areas within 100 miles (excluding i), 100 miles to 200 miles, and 200 miles to 300 miles of county i . Observations are weighted by county population in 1900. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

compute surname diversity and construct a separate instrument for individuals living within 100 miles, excluding i itself, individuals living between 100 miles and 200 miles, and between 200 miles and 300 miles.¹⁶

Table 10 reports the results for patents (columns 1 to 4) and breakthrough patents (columns 5 to 8). According to columns 1 to 4 (panel C), an increase in surname diversity just outside and within 100 miles of county i increases its number of patents. A similar relation is observed for breakthrough patents, though these estimates are less accurate. Moving to regions still further away from county i (i.e., between 100 and 200 miles or 200 and 300 miles), we do not find evidence for spillover effects. Our findings suggest that the causal link between diversity and innovation tends to be local, including the neighboring areas that fall within a 100-mile radius of the county.

6.4 Education

The role of human capital features prominently in the literature on innovation. Here, we explore the interplay between education and diversity. We report on a least-square specification that controls for the average educational attainment of a county's population. In addition, we interact surname diversity with educational attainment. This specification allows us to gain insights on the role of education; in particular, the interaction term reveals whether the relationship between diversity and innovation is more pronounced if average education of the population is higher. Clearly, such a least-square specification has to be interpreted cautiously because education is endogenous. Furthermore, the availability of education data restricts us to a cross-sectional analysis of the year 1940.

With these caveats in mind, table B5 reports the regression results for patents (columns 1–3) and breakthrough patents (columns 4–6). Column 1 and 4 only contain average education as a regressor. In both specification, the coefficient for education is highly significant. Once we add surname diversity and its interaction with education in the remaining columns, the picture changes. Educational attainment is either no longer significant (columns 2 and 5) or decreases in size in the specification that controls for state fixed effects (columns 3 and 6). At the same time, in all specification surname diversity is at least weakly significant and its interaction term with educational diversity is highly significant. A possible interpretation of this finding is that educational attainment plays out its innovative strength mostly in a diverse environment. Without the opportunity to learn from a diverse population educational attainment is less important for innovation.

¹⁶We use the [NBER's County Distance Database](#) to compute these areas for each county.

7 Conclusion

Focusing on the United States, during the period when it rose to dominate global innovation (1850 to 1940), we study the impact of diverse social interactions on innovation. The core idea is that many, if not most, innovations arise from the recombinations of existing ideas, approaches and techniques that come together through the connections among diverse minds. To measure such diversity, we introduce and benchmark an entropic diversity measure that exploits a widely available data source, surnames, obtained from the complete U.S. Census. To measure innovation, we use patents per capita at the county level and a text-based measure of breakthrough patents per capita. In our analysis, we first use least-squares regressions across U.S. counties. These analyses show that surname diversity is robustly correlated with both patent-based innovation measures across the entire period and that it holds controlling for more common diversity measures such as those based on birth country. Next, we employ an instrumental variable approach that uses immigrant flows to extract an exogenous component of surname diversity to examine the effect of surname diversity on our innovation outcomes. These analyses suggest that greater surname diversity causes greater innovation. Third, we subject these results to a battery of robustness and sensitivity checks including a placebo test for reverse causality, explorations of the role of population size, and surname fixed effects, which shows that people with the same surname get more innovative when they live in a more diverse county. Our analysis closes by showing that surname diversity increases novel combinations of existing technologies, weakens family ties and that the impact of surname diversity degrades rapidly with spatial distance and that it is partially driven by weakening family ties.

References

- Abramitzky, Ran, Philipp Ager, Leah Boustan, Elior Cohen & Casper W. Hansen** (2023) “The effect of immigration restrictions on local labor markets: Lessons from the 1920s border closure”, *American Economic Journal: Applied Economics*, 15 (1), pp. 164–191.
- Abramitzky, Ran & Leah Boustan** (2017) “Immigration in American economic history”, *Journal of Economic Literature*, 55 (4), pp. 1311–45.
- Acemoglu, Daron, Ufuk Akcigit & William R. Kerr** (2016) “Innovation network”, *Proceedings of the National Academy of Sciences*, 113 (41), p. 11483–11488.
- Ager, Philipp & Markus Brückner** (2013) “Cultural diversity and economic growth: Evidence from the US during the age of mass migration”, *European Economic Review*, 64, pp. 76–97.
- Aihounton, Ghislain B D & Arne Henningsen** (2020) “Units of measurement and the inverse hyperbolic sine transformation”, *The Econometrics Journal*, 24 (2), pp. 334–351.
- Akcigit, Ufuk, Santiago Caicedo, Ernest Miguelez, Stefanie Stantcheva & Valerio Sterzi** (2018) “Dancing with the Stars: Innovation Through Interactions”, Technical Report w24466, National Bureau of Economic Research, Cambridge, MA.
- Akcigit, Ufuk, William R Kerr & Tom Nicholas** (2013) “The mechanics of endogenous innovation and growth: Evidence from historical U.S. patents”, *Working paper*.
- Alesina, Alberto & Paola Giuliano** (2014) “Chapter 4 - Family Ties”, Philippe Aghion & Steven N. Durlauf eds. *Handbook of Economic Growth*, 2, Elsevier, pp. 177–215.
- Alesina, Alberto, Johann Harnoss & Hillel Rapoport** (2016) “Birthplace diversity and economic prosperity”, *Journal of Economic Growth*, 21 (2), pp. 101–138.
- Alesina, Alberto & Eliana La Ferrara** (2000) “Participation in Heterogeneous Communities”, *The Quarterly journal of economics*, 115 (3), pp. 847–904.
- Alesina, Alberto & Eliana La Ferrara** (2002) “Who trusts others?”, *Journal of Public Economics*, 85 (2), pp. 207–234.
- Algan, Yann & Pierre Cahuc** (2014) “Chapter 2 - trust, growth, and well-being: New evidence and policy implications”, Philippe Aghion & Steven N. Durlauf eds. *Handbook of Economic Growth*, 2 of Handbook of Economic Growth, Elsevier, pp. 49–120.

- AlShebli, Bedoor K., Talal Rahwan & Wei Lee Woon** (2018) “The preeminence of ethnic diversity in scientific collaboration”, *Nature Communications*, 9 (1), p. 5163.
- Altonji, Joseph G & David Card** (1991) “The effects of immigration on the labor market outcomes of less skilled natives”, *Immigration, Trade and the Labor Market*, Chicago, University of Chicago Press, pp. 201–234.
- Andrews, Michael** (2023) “Bar Talk: Informal Social Interactions, Alcohol Prohibition, and Invention”, *Working paper*.
- Andrews, Michael J. & Chelsea Lensing** (2020) “Cup of joe and knowledge flow: Coffee shops and invention”, *Working paper*.
- Ashraf, Quamrul & Oded Galor** (2013) “The “Out of Africa” Hypothesis, Human Genetic Diversity, and Comparative Economic Development”, *American Economic Review*, 103 (1), pp. 1–46.
- Banfield, Edward C** (1958) *The Moral Basis of a Backward Society*, Glencoe, Ill. : [Chicago, Free Press.
- Barone, Guglielmo & Sauro Mocetti** (2021) “Intergenerational mobility in the very long run: Florence 1427–2011”, *The Review of Economic Studies*, 88 (4), pp. 1863–1891.
- Barrai, I., C. Scapoli, M. Beretta, C. Nesti, E. Mamolini & A. Rodriguez-Larralde** (1996) “Isonymy and the genetic structure of Switzerland I. The distributions of surnames”, *Annals of Human Biology*, 23 (6), pp. 431–455.
- Bell, Alexander, Raj Chetty, Xavier Jaravel, Neviana Petkova & John Van Reenen** (2019) “Who Becomes an Inventor in America? The Importance of Exposure to Innovation”, *Quarterly Journal of Economics*, 134 (2), pp. 647–713.
- Berkes, Enrico** (2018) “Comprehensive universe of US patents (CUSP): Data and facts”, *Working paper*.
- Bisin, A & T Verdier** (1998) “On the cultural transmission of preferences for social status”, *Journal of Public Economics*, 70 (1), p. 75–97.
- Blasi, Damián E., Joseph Henrich, Evangelia Adamou, David Kemmerer & Asifa Majid** (2022) “Over-reliance on english hinders cognitive science”, *Trends in Cognitive Sciences*, 26 (12), p. 1153–1170.

- Boyd, Robert & Peter J Richerson** (1985) *Culture and the Evolutionary Process*, Chicago, IL, University of Chicago Press.
- Buonanno, Paolo & Paolo Vanin** (2017) “Social Closure, Surnames and Crime”, *Journal of Economic Behavior & Organization*, 137, pp. 160–175.
- Burchardi, Konrad B, Thomas Chaney & Tarek A Hassan** (2019) “Migrants, Ancestors, and Foreign Investments”, *The Review of Economic Studies*, 86 (4), pp. 1448–1486.
- Burchardi, Konrad B, Thomas Chaney, Tarek A Hassan, Lisa Tarquinio & Stephen J Terry** (2021) “Immigration, Innovation, and Growth”, *Working paper*.
- Carcassi, Gabriele, Christine A Aidala & Julian Barbour** (2021) “Variability as a better characterization of Shannon entropy”, *European Journal of Physics*, 42 (4), p. 045102.
- Card, David** (2001) “Immigrant inflows, native outflows, and the local labor market impacts of higher immigration”, *Journal of Labor Economics*, 19, pp. 22–64.
- Carlino, Gerald & William R. Kerr** (2015) “Agglomeration and Innovation”, *Handbook of Regional and Urban Economics*, 5, Elsevier, pp. 349–404.
- Cattaneo, Matias D, Richard K Crump, Max H Farrell & Yingjie Feng** (2019) “On bin-scatter”, *arXiv preprint arXiv:1902.09608*.
- Cavalli-Sforza, Luigi Luca & Marc W Feldman** (1981) *Cultural Transmission and Evolution: A Quantitative Approach*, Princeton University Press.
- Cinnirella, Francesco, Erik Hornung & Julius Koschnick** (2022) “Flow of ideas: Economic societies and the rise of useful knowledge”, *CESifo Working Paper*, 9836.
- Clancy, Matthew S** (2018a) “Combinations of Technology in US Patents, 1926-2009: A Weakening Base for Future Innovation?”, *Economics of Innovation and New Technology*, 27 (8), pp. 770–785.
- Clancy, Matthew S.** (2018b) “Inventing by combining pre-existing technologies: Patent evidence on learning and fishing out”, *Research Policy*, 47 (1), pp. 252–265.
- Clark, Gregory** (2014) *The Son also Rises: a Surprising Look how Ancestry still Determines Social Outcomes*, Princeton, Princeton University Press.
- Cook, Lisa D.** (2014) “Violence and economic activity: Evidence from African American patents, 1870–1940”, *Journal of Economic Growth*, 19 (2), pp. 221–257.

- Cook, Lisa D., John Parman & Trevon Logan** (2022) “The antebellum roots of distinctively black names”, *Historical Methods: A Journal of Quantitative and Interdisciplinary History*, 55 (1), pp. 1–11.
- de la Croix, David, Matthias Doepke & Joel Mokyr** (2018) “Clans, guilds, and markets: Apprenticeship institutions and growth in the pre-industrial economy”, *The Quarterly Journal of Economics*, 133 (1), pp. 1–70.
- Desmet, Klaus, Avner Greif & Stephen L. Parente** (2017) “Spatial Competition, Innovation and Institutions: The Industrial Revolution and the Great Divergence”, *SSRN Electronic Journal*.
- Docquier, Frédéric, Riccardo Turati, Jérôme Valette & Chrysovalantis Vasilakis** (2020) “Birthplace diversity and economic growth: Evidence from the US states in the Post-World War II period”, *Journal of Economic Geography*, 20 (2), pp. 321–354.
- Enke, Benjamin** (2019) “Kinship, Cooperation, and the Evolution of Moral Systems*”, *The Quarterly Journal of Economics*, 134 (2), pp. 953–1019.
- Ferrara, Andreas, Patrick Testa & Liyang Zhou** (2021) “New Area-and Population-based Geographic Crosswalks for US Counties and Congressional Districts, 1790-2020”, *SSRN 4019521*.
- Fulford, Scott L., Ivan Petkov & Fabio Schiantarelli** (2020) “Does it matter where you came from? Ancestry composition and economic performance of US counties, 1850–2010”, *Journal of Economic Growth*, 25 (3), pp. 341–380.
- Galor, O & D N Weil** (2000) “Population, Technology, and Growth: From Malthusian Stagnation to the Demographic Transition and beyond”, *The American Economic Review*, 90 (4), pp. 806–828.
- Galor, Oded** (2022) *The Journey of Humanity: The Origins of Wealth and Inequality*, New York, Dutton.
- Ghosh, Arkadev, Sam Il Myoung Hwang & Munir Squires** (2023) “Economic Consequences of Kinship: Evidence from US Bans on Cousin Marriage”, *Working paper, University of British Columbia*.
- Glaeser, Edward** (2011) “Engines of Innovation”, *Scientific American*, p. 7.
- Glaeser, Edward L., Hedi D. Kallal, José A. Scheinkman & Andrei Shleifer** (1992) “Growth in cities”, *Journal of Political Economy*, 100 (6), p. 1126–1152.

- Glaeser, Edward L., David I. Laibson, José A. Scheinkman & Christine L. Soutter** (2000) “Measuring Trust”, *The Quarterly Journal of Economics*, 115 (3), pp. 811–846.
- Gorodnichenko, Yuriy & Gerard Roland** (2016) “Culture, institutions, and the wealth of nations”, *The Review of Economics and Statistics*, 99 (3), pp. 402–416.
- Griliches, Zvi** (1990) “Patent Statistics as Economic Indicators: A Survey”, *Journal of Economic Literature*, 28 (4), pp. 1661–1707.
- Güell, Maia, José V. Rodríguez Mora & Christopher I. Telmer** (2015) “The Informational Content of Surnames, the Evolution of Intergenerational Mobility, and Assortative Mating”, *The Review of Economic Studies*, 82 (2), pp. 693–735.
- Handley, Carla & Sarah Mathew** (2020) “Human large-scale cooperation as a product of competition between cultural groups”, *Nature Communications*, 11 (1), p. 702.
- Henrich, J** (2009) “The evolution of innovation-enhancing institutions”, Stephen J Shennan & Michael J O’Brien eds. *Innovation in Cultural Systems: Contributions in Evolutionary Anthropology*, Cambridge, MIT, pp. 99–120.
- Henrich, Joseph** (2020) *The WEIRDest People in the World: How the West Became Psychologically Peculiar and Particularly Prosperous*, New York, Farrar, Straus and Giroux.
- Jacobs, Jane** (1969) *The Economy of Cities*, New York, Random House.
- Johnson, S.** (2011) *Where Good Ideas Come From: The Natural History of Innovation*, Penguin Publishing Group.
- Jones, Charles I** (forthcoming) “Recipes and economic growth: A combinatorial march down an exponential tail”, *Journal of Political Economy*.
- Kelly, Bryan, Dimitris Papanikolaou, Amit Seru & Matt Taddy** (2021) “Measuring Technological Innovation over the Long Run”, *American Economic Review: Insights*, 3 (3), pp. 303–320.
- Kerr, William** (2008) “The Ethnic Composition of US Inventors”, *HBS Finance Working Paper No. 08-006*.
- Kerr, William R.** (2010) “Breakthrough inventions and migrating clusters of innovation”, *Journal of Urban Economics*, 67 (1), pp. 46–60.
- Lerner, Josh & Amit Seru** (2022) “The use and misuse of patent data: Issues for finance and beyond”, *The Review of Financial Studies*, 35 (6), p. 2667–2704.

- Lucas Jr, Robert E & Benjamin Moll** (2014) “Knowledge growth and the allocation of time”, *Journal of Political Economy*, 122 (1), pp. 1–51.
- Mill, J.S.** (1871) *Principles of Political Economy: With Some of Their Applications to Social Philosophy*, Principles of Political Economy: With Some of Their Applications to Social Philosophy, Longmans, Green, Reader, and Dyer.
- Miu, Elena, Ned Gulley, Kevin N. Laland & Luke Rendell** (2018) “Innovation and cumulative culture through tweaks and leaps in online programming contests”, *Nature Communications*, 9 (2321).
- Mokyr, Joel** (1995) “Urbanization, technological progress, and economic history”, H. Giersch ed. *Urban Agglomeration and Economic Growth*, Berlin; Heidelberg, Springer, pp. 51–54.
- Mokyr, Joel** (2002) *The Gifts of Athena: Historical Origins of the Knowledge Economy*, Princeton, NJ, Princeton University Press.
- Mokyr, Joel** (2015) “ECONOMICS. Intellectuals and the rise of the modern economy.”, *Science (New York, N.Y.)*, 349 (6244), pp. 141–2.
- Moser, Petra** (2013) “Patents and Innovation: Evidence from Economic History”, *The Journal of Economic Perspectives*, 27 (1), pp. 23–44.
- Moser, Petra & Shmuel San** (2020) “Immigration, Science, and Invention. Lessons from the Quota Acts”, *SSRN Electronic Journal*.
- Moser, Petra, Alessandra Voena & Fabian Waldinger** (2014) “German Jewish Émigrés and US Invention”, *American Economic Review*, 104 (10), pp. 3222–3255.
- Muthukrishna, Michael & Joseph Henrich** (2016) “Innovation in the collective brain”, *Philosophical Transactions of the Royal Society B: Biological Sciences*, 371 (1690), pp. 1–14.
- Nguyen, Kieu-Trang** (2021) “Trust and innovation within the firm: Evidence from matched CEO-Firm data”, *Working paper*.
- Nisbett, Richard E** (2003) *The Geography of Thought: How Asians and Westerners Think Differently– and Why*, New York, Free Press.
- Nunn, Nathan** (2021) “Chapter 3 - History as evolution”, Alberto Bisin & Giovanni Federico eds. *The Handbook of Historical Economics*, Academic Press, pp. 41–91.

- Olsson, Ola & Bruno Frey** (2002) “Entrepreneurship as recombinant growth”, *Small Business Economics*, 19 (2), p. 69–80.
- Ottaviano, Gianmarco IP & Giovanni Peri** (2006) “The economic value of cultural diversity: Evidence from US cities”, *Journal of Economic Geography*, 6 (1), pp. 9–44.
- Ozgen, Ceren, Cornelius Peters, Annekatrien Niebuhr, Peter Nijkamp & Jacques Poot** (2014) “Does cultural diversity of migrant employees affect innovation?”, *International Migration Review*, 48 (1), pp. 377–416.
- Page, S.E., N. Cantor & K. Phillips** (2019) *The Diversity Bonus: How Great Teams Pay Off in the Knowledge Economy*, Our Compelling Interests, Princeton University Press.
- Philips, Lawrence** (1990) “Hanging on the metaphor”, *Computer Language*, 7 (12), pp. 39–43.
- Raz, Itzhak T.** (2023) “Soil Heterogeneity and the Formation of Close-knit Communities”, *Working paper*.
- Ridley, Matt** (2020) *How Innovation Works: And Why It Flourishes in Freedom*, London, Harper.
- Romer, Paul M** (1990) “Endogenous Technological Change”, *Journal of Political Economy*, 98 (5, Part 2), pp. S71–S102.
- Ruggles, Steven, Catherine A. Fitch, Ronald Goeken, J. David Hacker, Matt A. Nelson, Evan Roberts, Megan Schouweiler & Matthew Sobek** (2021) *IPUMS Ancestry Full Count Data: Version 3.0 [Dataset]*, Minneapolis, MN, IPUMS.
- Schulz, Jonathan F., Duman Bahrami-Rad, Jonathan P. Beauchamp & Joseph Henrich** (2019) “The Church, Intensive Kinship, and Global Psychological Variation”, *Science*, 366 (6466).
- Schumpeter, J.A.** (1983) *The Theory of Economic Development: An Inquiry Into Profits, Capital, Credit, Interest, and the Business Cycle*, Economics Third World studies, Transaction Books.
- Sequeira, Sandra, Nathan Nunn & Nancy Qian** (2020) “Immigrants and the Making of America”, *Review of Economic Studies*, 87, pp. 382–419.
- Shannon, Claude Elwood** (1948) “A mathematical theory of communication”, *The Bell system technical journal*, 27 (3), pp. 379–423.

- Strumsky, Deborah, Jose Lobo & Sander Van der Leeuw** (2011) “Measuring the relative importance of reusing, recombining and creating technologies in the process of invention”, *SFI Working Paper 2011-02-003*: 23.
- Suedekum, Jens, Katja Wolf & Uwe Blien** (2014) “Cultural Diversity and Local Labour Markets”, *Regional Studies*, 48 (1), pp. 173–191.
- Thagard, Paul** (2012) “Creative combination of representations: Scientific discovery and technological invention.”, *Psychology of Science: Implicit and Explicit Processes.*, New York, NY, US, Oxford University Press, pp. 389–404.
- Usher, A.P.** (2013) *A History of Mechanical Inventions: Revised Edition*, Dover Publications.
- Uzzi, B., S. Mukherjee, M. Stringer & B. Jones** (2013) “Atypical Combinations and Scientific Impact”, *Science*, 342 (6157), pp. 468–472.
- Weitzman, Martin L.** (1998) “Recombinant Growth”, *The Quarterly Journal of Economics*, 63 (2).
- White, Cindel J. M., Michael Muthukrishna & Ara Norenzayan** (2021) “Cultural similarity among coreligionists within and between countries”, *Proceedings of the National Academy of Sciences*, 118 (37), p. e2109650118.
- Youn, H J, D Strumsky, L M A Bettencourt & J Lobo** (2015) “Invention as a combinatorial process: Evidence from US patents”, *Journal of the Royal Society Interface*, 12 (106).

Online Appendix

Surname Diversity, Social Ties and Innovation

Max Posch, Jonathan Schulz, and Joseph Henrich

A Data Sources and Construction

Surname-based surname diversity

To construct county-level surname diversity up until the year 1940, we use the 1850, 1860, 1870, 1880, 1900, 1910, 1920, 1930, and 1940 waves of the full-count Integrated Public Use Microdata Series (IPUMS) compiled by [Ruggles et al. \(2021\)](#) and available on the NBER servers. For each wave, we obtain county identifiers and the variable `namelast` of all individuals. We perform the following steps to clean the surname variable. First, we transform non-ASCII characters into ASCII characters—e.g., we convert characters with accents or umlauts to the closest letter in English. Second, we convert all characters to upper case. Third, we remove all non-alphabetic characters, including all spaces (e.g., ‘MAC ARTHUR’ becomes ‘MACARTHUR’). Fourth, we drop entries with one or fewer letters. Last, we apply the [Philips \(1990\)](#) phonetic algorithm *metaphone* to deal with misspellings.

We harmonize all historical Census data to the 2000 boundaries of U.S. counties using the [Ferrara et al. \(2021\)](#) crosswalks. Specifically, we use the M4 weights, which account for urban and rural areas and topographic suitability. We use 2000 as the reference year because the patent dataset is geocoded to 2000 county boundaries. The harmonization procedure sometimes results in counties with very few people, predominantly in the Midwest and West, and for Census years before 1900. As a remedy, we winsorize all harmonized variables from the lower tail at the 1% level.

Following [Burchardi et al. \(2021\)](#), we also obtain individuals’ age and year of immigration, the variables `age` and `yrimmig`, to estimate surname diversity for the midyears 1895, 1905, 1915, and 1925 by removing all individuals who were born or immigrated after the midyear. Ideally, we would also remove all individuals who moved to the county after the midyear, but this information is unavailable. We also compute alternative measures of surname diversity by interacting surnames with the main categories of race (`race`), birthplace (`bp1`), and education (`higrade`). We recode U.S. states and territories (`bp1` codes <10000) to a single code and use an indicator equaling one if the individual completed elementary school.

Construction of the instrumental variable

We build on the [Burchardi et al. \(2019\)](#) approach to construct an instrumental variable for surname-based surname diversity. We identify the number of individuals in a given U.S. county i at the time of each census who immigrated to the U.S. since the prior census and have the surname k . For the 1900 to 1930 census waves, we separate this immigration into five-year periods based on the year each migrant arrived in the U.S. We obtain immigration flows for the following bins: 1881-1895, 1896-1900, 1901-1905, 1906-1910, 1911-1915, 1916-1920, 1921-1925, and 1926-1930. From the 1880 census wave, we count all first- and second-generation immigrants, regardless of the date of arrival in the U.S.

When we predict the stock of people $N_{i,k}^t$ in equation (3), we obtain negative values for some observations. The logarithmic transformation of a negative value is undefined. To obtain Shannon entropy for counties containing $N_{i,k}^t$ with negative values, we truncate those negative values at the smallest positive value we observe in the data in a given year. The resulting variable is highly correlated with the original variable ($\rho = 0.965$).

Construction of other demographic measures

We collect county-level data on population size, ethnic and birthplace diversity, and immigrant shares for each census year from 1850 to 1940. All data are taken from the full-count IPUMS available on the NBER servers. To compute ethnic diversity, we obtain the variable race and use the nine main categories. To compute birthplace diversity and immigrant shares, we draw on the variable bp1 with 188 main categories. As before, we recode U.S. states and territories (codes <100) to a single category, transform the data from each period to 2000 U.S. counties using the M4 weights from the [Ferrara et al. \(2021\)](#) cross-walks, and winsorize all demographic variables from the lower tail at the 1% level.

We follow the instructions in [Raz \(2023\)](#) to construct the strength of family ties measure from the full-count census data for all census waves from 1860 to 1940. The strength of family ties is captured by the first principal component of four underlying variables: (i) the divorce-to-marriage ratio, (ii) the share of elderly people living without a relative, (iii) the share of people living with at least one person who is not their relative, and (iv) the mean size of families. We use the variables age and yr immig to estimate the strength of family ties for the midyears 1895, 1905, 1915, and 1925 by removing all individuals who were born or immigrated after the midyear.

Construction of the innovation measures

We use the *Comprehensive Universe of U.S. Patents* (CUSP) compiled by [Berkes \(2018\)](#). The data set contains U.S. patents from 1836-2015 and is primarily constructed from Google Patents with supplementary information from other sources. For each patent, the data set provides inventor names and location of residence (geocoded to 2000 county boundaries), filing and issuing years of patents, and the U.S. Patent and Trademark Office technology classification.

We also draw on the breakthrough patent indicator created by [Kelly et al. \(2021\)](#). The authors use the text in patent documents to estimate patent quality. They assign a higher quality to patents that are novel in terms of cosine similarities. Patents are considered novel if they have low similarity with the existing stock of patents and are impactful in that they have high similarity with subsequent patents. We use this measure of patent quality rather than the number of citations an individual patent has received because the U.S. Patent and Trademark Office did not consistently begin to record patent citations until after 1947.

We construct the innovation outcome variables at the county-period and surname-county-period levels. The county-period-level outcomes are per capita number of (breakthrough) patents filed by inventors residing in county i during the period starting t . If a patent is filed by more than one inventor, possibly residing in different counties, we divide the patent count by the number of inventors. We use county population sizes at the beginning of t , computed using the full-count IPUMS and the [Ferrara et al. \(2021\)](#) border harmonization procedure (winsorized from the lower tail at the 1% level), to get per capita rates of (breakthrough) patents filed. We use patent issuing years rather than filing years in the least-squares analysis in Section 3, because filing years are not consistently recorded in the CUSP data set before 1870. When the unit of observation is county-period, we winsorize the innovation outcome variables from the upper tail at the 99% level to reduce the influence of outlier counties with a very large number of patents. We do not winsorize the surname-county-period-level outcomes because the number of breakthrough patents filed by inventors with a given surname in a given county during a given period is typically small. We also report results using non-winsorized, inverse hyperbolic sine transformed patent counts.

The surname-county-period-level outcomes are per capita number of (breakthrough) patents filed by inventors with surname k residing in county i during the period starting t . The construction of these outcome variables requires inventor surnames. The CUSP data includes inventor names. This string variable contains the surname, first name, and sometimes middle names or initials. Identifying surnames from this string variable is not

straightforward because the order of first names and surnames is inconsistent: surnames follow first names in some entries but not in others. When a semicolon, colon, or comma delineates surnames from first and middle names, we use these characters to discern the surnames. When the string variable starts with initials followed by a token of two or more characters, or when it ends with a whitespace followed by “DE”, “DU”, “DE LA”, “DI”, “DEL”, “DELLA”, “VAN”, “VON”, “LE”, “LA”, or “ST”, we distinguish the surnames accordingly. For the remaining entries, we tokenize the string variable based on whitespace and keep the first token and the last token, which are the first name and surname in most cases. To determine which of the two tokens is the surname, we compute the frequencies of all names (first name and surname) from the pooled census years 1900, 1910, 1920, 1930, and 1940 and compare which constellation is more common. For example, for the tokens “JOHN” and “PETER”, we identify the surname based on whether there were more individuals named “JOHN PETER” or “PETER JOHN”. Finally, we clean the surname variable following the steps described above.

B Robustness of Main Results

B.1 Alternative Definitions of Surname Diversity

Table B1: Correlations between baseline surname diversity and alternative surname diversity measures

Herfindahl surname	Surname, uncorrected	Surname, men	Surname, whites	Surname-race	Surname-country of birth
0.80	0.96	1.00	0.99	0.98	0.98

Notes: This table reports the correlations between county-level surname diversity (based on Shannon entropy) and (i) a surname-based Herfindahl index, (ii) diversity of surnames that are not phonetically corrected, (iii) surname diversity among men, (iv) surname diversity among white individuals, and (v) alternative diversity measures that interact surnames with race or birthplace. An observation is a county from 1850 to 1940 (excluding the midyears). The sources and construction of all variables are explained in Appendix Section A.

B.2 Least-Squares Results

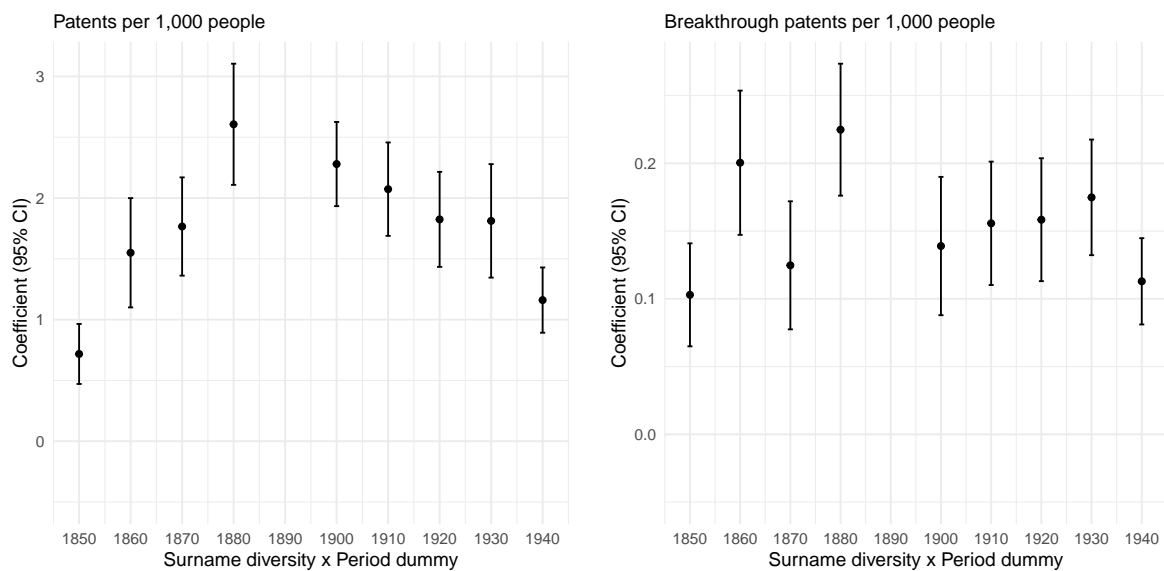


Figure B1: Correlations between surname diversity in years 1850-1940 and innovation
Notes: Each figure shows coefficients of a regression of an innovation outcome on surname diversity interacted with year dummies conditional on year fixed effects.

Table B2: Least-squares estimates: population size

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i>						
	Patents per 1,000 people (mean = 2.26, sd = 2.58)					
Surname diversity		1.51** (0.134)	0.677*** (0.149)	0.665*** (0.154)	0.802*** (0.168)	0.602*** (0.187)
Population	0.345*** (0.050)	0.115*** (0.033)	0.058*** (0.021)	0.043** (0.021)	0.070*** (0.021)	-0.070*** (0.015)
Country of origin diversity			1.03*** (0.177)	1.13*** (0.177)	0.930*** (0.239)	0.147 (0.196)
R ²	0.305	0.521	0.579	0.609	0.695	0.865
<i>Panel B:</i>						
	Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.24)					
Surname diversity		0.123*** (0.014)	0.054*** (0.011)	0.052*** (0.011)	0.057*** (0.014)	0.057** (0.024)
Population	0.033*** (0.005)	0.015*** (0.004)	0.010*** (0.002)	0.007*** (0.002)	0.009*** (0.002)	-0.009*** (0.002)
Country of origin diversity			0.085*** (0.016)	0.098*** (0.016)	0.095*** (0.024)	0.002 (0.018)
R ²	0.283	0.448	0.494	0.530	0.626	0.790
Immigrant shares by country of origin (59 shares)				✓	✓	✓
Period fixed effects	✓	✓	✓	✓		
Period-State fixed effects					✓	✓
County fixed effects						✓
Observations	22,299	22,299	22,299	22,299	22,299	22,299

Notes: The table reports least-squares estimates of regressions of innovation outcomes on surname diversity, immigrant diversities and population size. In Panel A (Panel B), the outcome is number of (breakthrough) patents per 1,000 people. The unit of observation is a county-period from 1850 to 1940 (excluding the midyears). Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table B3: Least-squares estimates: log population size

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i>						
	Patents per 1,000 people (mean = 2.26, sd = 2.58)					
Surname diversity		0.944*** (0.160)	0.143 (0.206)	0.148 (0.206)	0.212 (0.168)	0.747*** (0.257)
Log Population	2.17*** (0.221)	1.16*** (0.317)	0.946*** (0.196)	0.847*** (0.170)	0.947*** (0.126)	-0.432* (0.238)
Country of origin diversity			1.03*** (0.146)	1.14*** (0.152)	0.975*** (0.200)	0.124 (0.172)
R ²	0.503	0.534	0.594	0.619	0.700	0.864
<i>Panel B:</i>						
	Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.24)					
Surname diversity		0.054*** (0.014)	-0.014 (0.016)	-0.017 (0.016)	-0.009 (0.014)	0.071** (0.032)
Log Population	0.200*** (0.024)	0.142*** (0.030)	0.124*** (0.019)	0.114*** (0.016)	0.108*** (0.013)	-0.047* (0.023)
Country of origin diversity			0.087*** (0.015)	0.102*** (0.014)	0.102*** (0.020)	-0.0003 (0.015)
R ²	0.457	0.469	0.520	0.549	0.632	0.788
Immigrant shares by country of origin (59 shares)				✓	✓	✓
Period fixed effects	✓	✓	✓	✓		
Period-State fixed effects					✓	✓
County fixed effects						✓
Observations	22,299	22,299	22,299	22,299	22,299	22,299

Notes: The table reports least-squares estimates of regressions of innovation outcomes on surname diversity, immigrant diversities and log population size. In Panel A (Panel B), the outcome is number of (breakthrough) patents per 1,000 people. The unit of observation is a county-period from 1850 to 1940 (excluding the midyears). Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

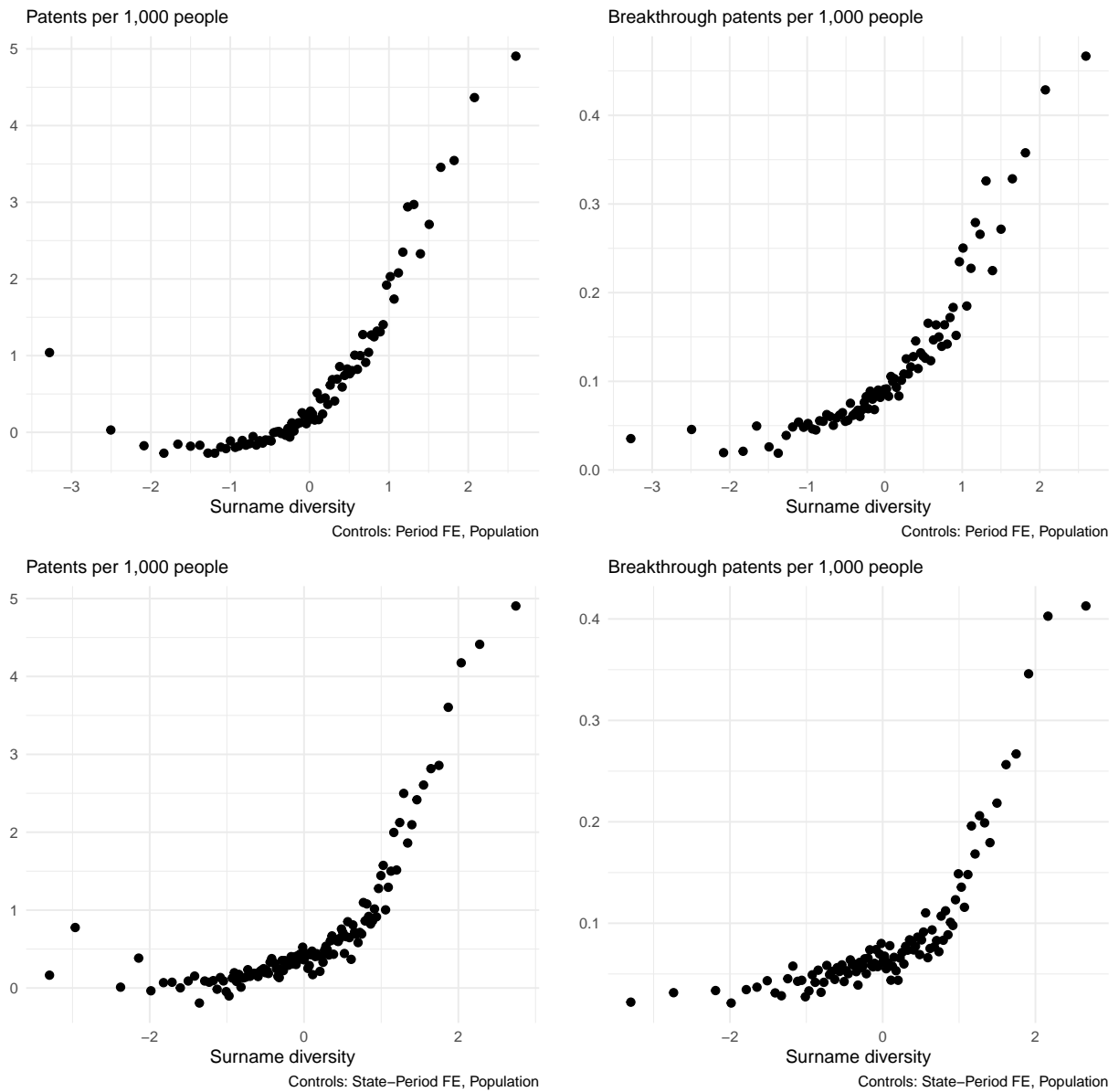


Figure B2: Conditional relationships between surname diversity and (breakthrough) patents

Notes: County-level data from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850 and residualized by census year fixed effects and county population. Bottom graphs: observations are additionally residualized by state-period fixed effects. Binscatter plot created using the R package written by Cattaneo et al. (2019).

Table B4: Least-squares estimates: race and occupational diversity

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
<i>Panel A:</i>							
	Patents per 1,000 people (mean = 1.20, sd = 0.90)						
Surname diversity	1.76** (0.175)		0.734** (0.157)	0.692** (0.161)	0.342** (0.128)	0.577** (0.130)	0.488** (0.203)
Country of origin diversity		1.62** (0.103)	1.10** (0.184)	1.16** (0.178)	1.08** (0.179)	1.03** (0.242)	0.137 (0.194)
Race diversity		-0.086 (0.107)	-0.010 (0.113)	-0.135 (0.094)	-0.146 (0.092)	-0.086 (0.082)	-0.478** (0.148)
Occupational diversity					0.535** (0.095)	0.422** (0.133)	0.180 (0.134)
R ²	0.503	0.550	0.574	0.608	0.620	0.696	0.866
<i>Panel B:</i>							
	Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.27)						
Surname diversity	0.154** (0.021)		0.067** (0.012)	0.060** (0.012)	0.036** (0.009)	0.046** (0.011)	0.044* (0.023)
Country of origin diversity		0.148** (0.015)	0.100** (0.019)	0.108** (0.018)	0.103** (0.019)	0.110** (0.025)	0.002 (0.016)
Race diversity		0.018 (0.011)	0.025** (0.012)	0.015 (0.010)	0.015 (0.011)	0.003 (0.009)	-0.047** (0.017)
Occupational diversity					0.037** (0.009)	0.027** (0.011)	0.013 (0.008)
R ²	0.416	0.461	0.485	0.525	0.531	0.622	0.789
Immigrant shares by country of origin (59 shares)				✓	✓	✓	✓
Period fixed effects	✓	✓	✓	✓	✓		
Period-State fixed effects						✓	✓
County fixed effects							✓
Observations	22,206	22,206	22,206	22,206	22,206	22,206	22,206

Notes: The table reports least-squares estimates of regressions of innovation outcomes on surname diversity and other dimensions of sociocultural diversity, including race and occupational diversity. In Panel A (Panel B), the outcome is number of (breakthrough) patents per 1,000 people. The unit of observation is a county-period from 1850 to 1940 (excluding the midyears). Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

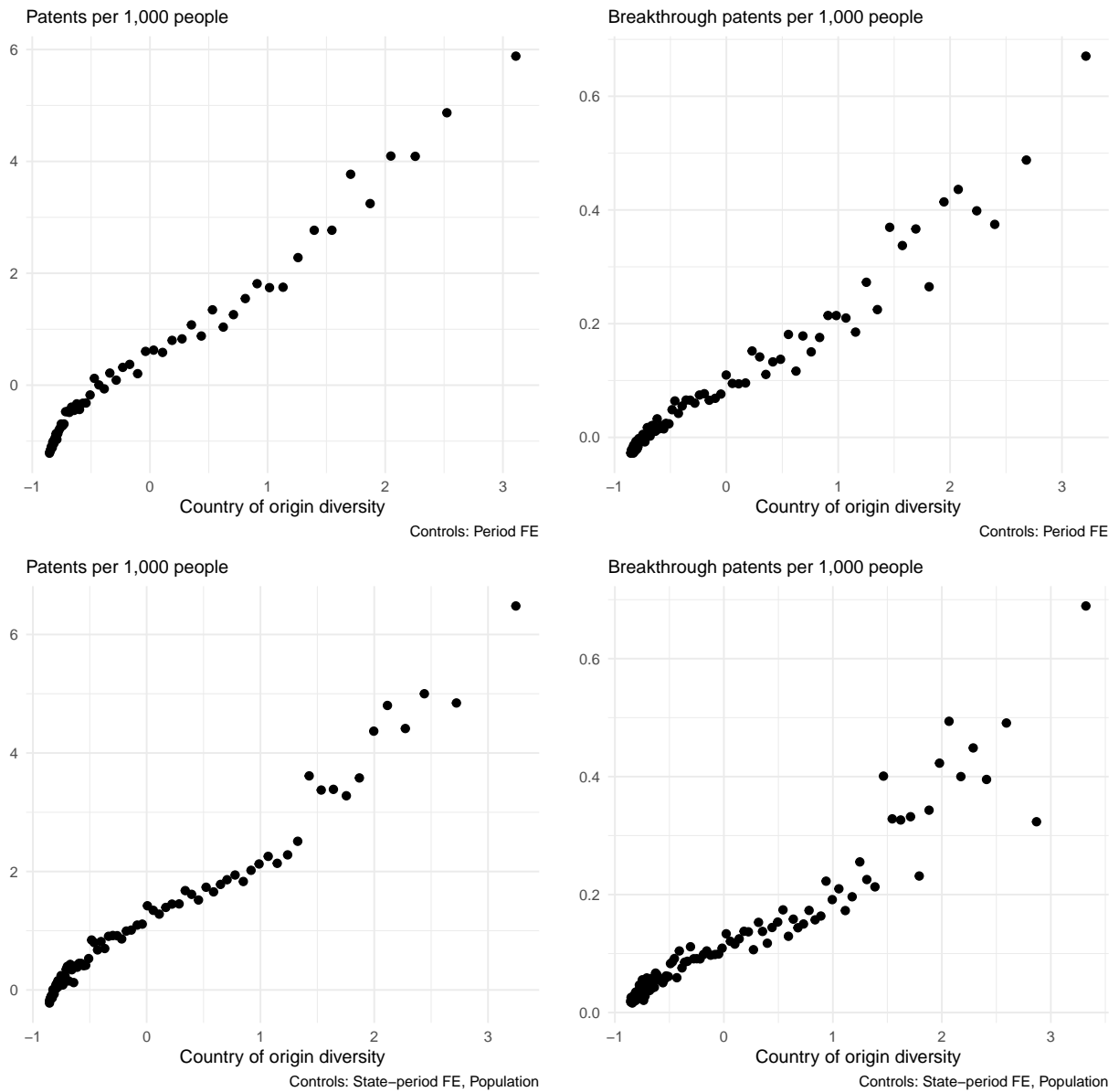


Figure B3: Bivariate relationships between country of origin diversity and (breakthrough) patents

Notes: County-level data from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850 and residualized by census year fixed effects. Bottom graphs: observations are additionally residualized by state-period fixed effects and county population size. Binscatter plot created using the R package written by Cattaneo et al. (2019).

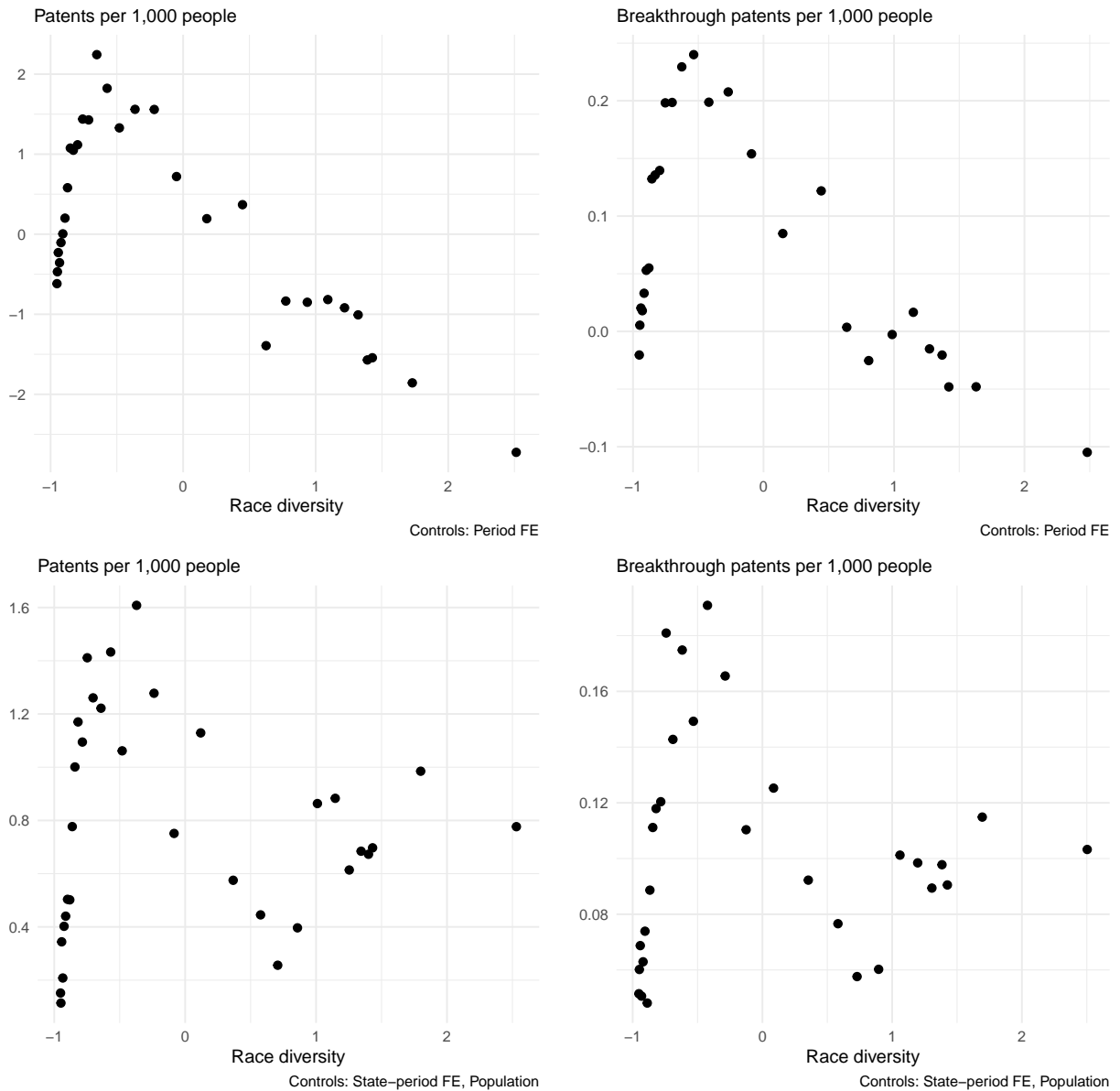


Figure B4: Bivariate relationships between race diversity and (breakthrough) patents
Notes: County-level data from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850 and residualized by census year fixed effects. Bottom graphs: observations are additionally residualized by state-period fixed effects and county population size. Binscatter plot created using the R package written by [Cattaneo et al. \(2019\)](#).

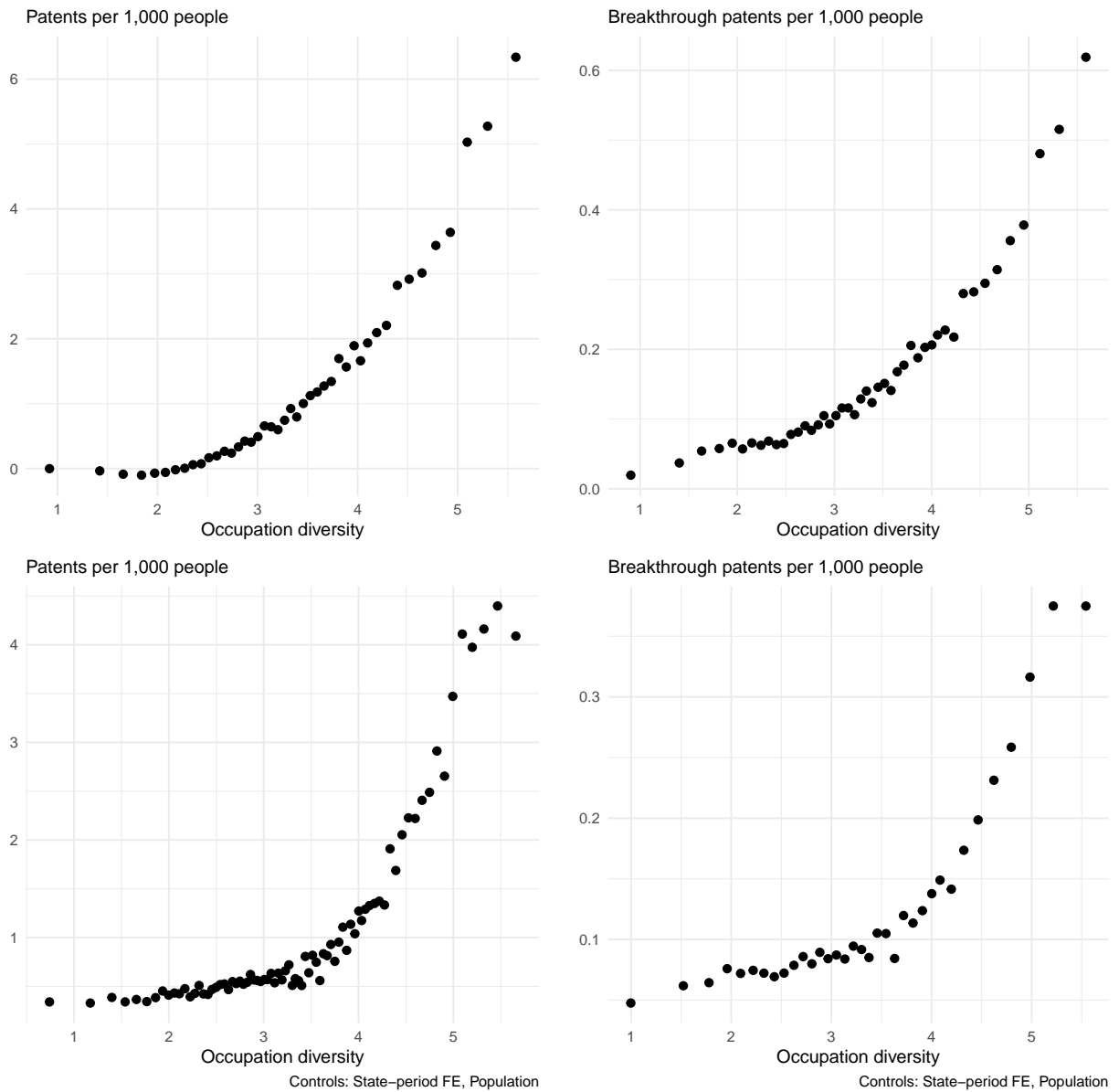


Figure B5: Bivariate relationships between occupational diversity and (breakthrough) patents

Notes: County-level data from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850 and residualized by census year fixed effects. Bottom graphs: observations are additionally residualized by state-period fixed effects and county population size. Binscatter plot created using the R package written by [Cattaneo et al. \(2019\)](#).

Table B5: Least-squares estimates: education

	Patents per 1,000 people (mean = 1.78, sd = 1.88)			Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.21)		
	(1)	(2)	(3)	(4)	(5)	(6)
	Average years of schooling	1.10*** (0.125)	0.072 (0.114)	0.292** (0.140)	0.114*** (0.015)	0.009 (0.009)
Surname diversity		0.589*** (0.145)	0.360** (0.146)		0.056*** (0.013)	0.028* (0.014)
Surname diversity × Average years of schooling		0.573*** (0.105)	0.653*** (0.121)		0.064*** (0.013)	0.069*** (0.016)
Constant	1.28*** (0.096)	0.147** (0.061)		0.129*** (0.011)	0.012* (0.007)	
R ²	0.282	0.563	0.650	0.247	0.495	0.592
State fixed effects			✓			✓
Observations	3,078	3,078	3,078	3,078	3,078	3,078

Notes: The table reports least-squares estimates of regressions of innovation outcomes on surname diversity and individuals' average years of schooling. The unit of observation is a county-period in 1940. Observations are weighted by county population in 1940. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table B6: Least-squares estimates: inverse hyperbolic sine transformed outcomes

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i>						
	IHS Patents per 1,000 people (mean = 1.20, sd = 0.90)					
Surname diversity	0.688** (0.038)		0.435** (0.041)	0.435** (0.045)	0.441** (0.050)	0.123** (0.038)
Country of origin diversity		0.590** (0.027)	0.273** (0.045)	0.284** (0.046)	0.215** (0.054)	0.033 (0.060)
R ²	0.635	0.600	0.671	0.697	0.768	0.901
<i>Panel B:</i>						
	IHS Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.27)					
Surname diversity	0.164** (0.027)		0.063** (0.014)	0.053** (0.014)	0.059** (0.017)	0.063* (0.033)
Country of origin diversity		0.155** (0.020)	0.109** (0.027)	0.120** (0.026)	0.131** (0.036)	0.002 (0.019)
R ²	0.366	0.412	0.428	0.467	0.573	0.760
Immigrant shares by country of origin (59 shares)				✓	✓	✓
Period fixed effects	✓	✓	✓	✓		
Period-State fixed effects					✓	✓
County fixed effects						✓
Observations	22,222	22,222	22,222	22,222	22,222	22,222

Notes: The table reports estimates of least-squares regressions of innovation outcomes on surname diversity and other dimensions of sociocultural diversity. In Panel A (Panel B), the outcome is inverse hyperbolic sine transformed number of (breakthrough) patents issued in a given period per 1,000 people. The unit of observation is a county-period from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

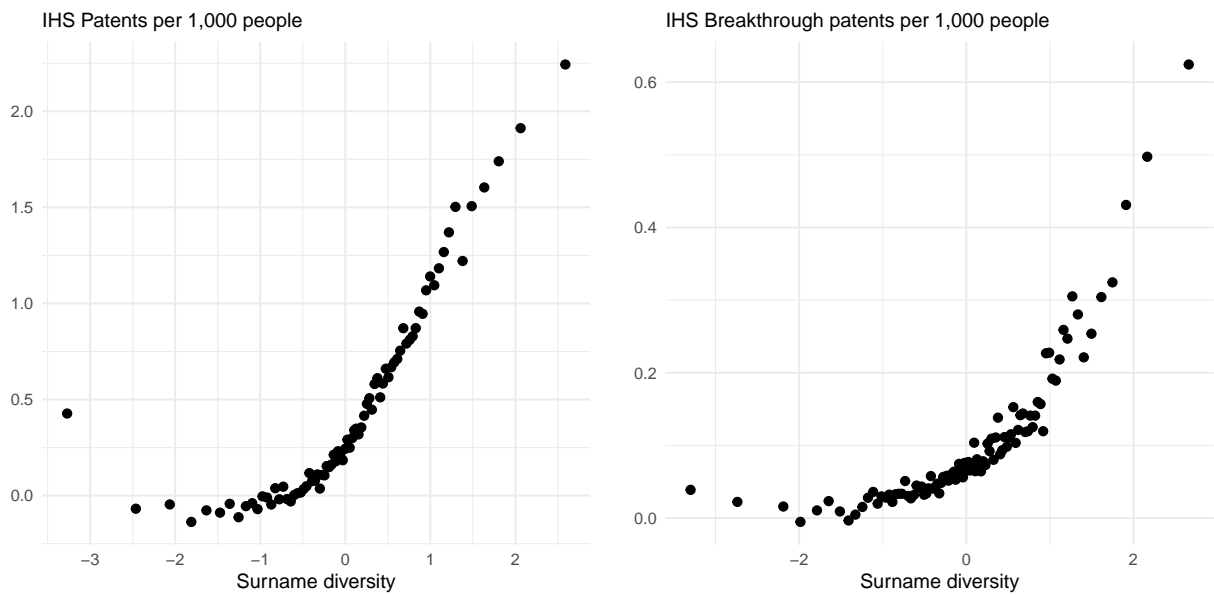


Figure B6: Bivariate relationships between surname diversity and inverse hyperbolic sine transformed innovation outcomes

Notes: County-level data from 1850 to 1940 (excluding the midyears). Observations are weighted by county population in 1850 and residualized by census year fixed effects. Binscatter plot created using the R package written by Cattaneo et al. (2019).

Table B7: Least-squares estimates: Herfindahl

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i>						
	Patents per 1,000 people (mean = 2.26, sd = 2.58)					
Herfindahl surname diversity	2.25*** (0.419)		0.497*** (0.139)	0.361*** (0.132)	0.317*** (0.100)	0.152 (0.091)
Herfindahl country of origin diversity		1.66*** (0.113)	1.55*** (0.138)	1.59*** (0.124)	1.49*** (0.183)	0.155 (0.160)
R ²	0.238	0.528	0.534	0.574	0.660	0.862
<i>Panel B:</i>						
	Breakthrough patents per 1,000 people (mean = 0.18, sd = 0.24)					
Herfindahl surname diversity	0.191*** (0.041)		0.034** (0.013)	0.022* (0.011)	0.019** (0.009)	0.024* (0.012)
Herfindahl country of origin diversity		0.146*** (0.015)	0.139*** (0.018)	0.144*** (0.015)	0.145*** (0.021)	0.011 (0.015)
R ²	0.171	0.445	0.448	0.502	0.597	0.786
Immigrant shares by country of origin (59 shares)				✓	✓	✓
Period fixed effects	✓	✓	✓	✓		
Period-State fixed effects					✓	✓
County fixed effects						✓
Observations	22,299	22,299	22,299	22,299	22,299	22,299

Notes: The table reports least-squares estimates of regressions of innovation outcomes on Herfindahl surname diversity, immigrant diversities and population size. In Panel A (Panel B), the outcome is number of (breakthrough) patents per 1,000 people. The unit of observation is a county-period from 1850 to 1940 (excluding the midyears). Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

B.3 IV Results

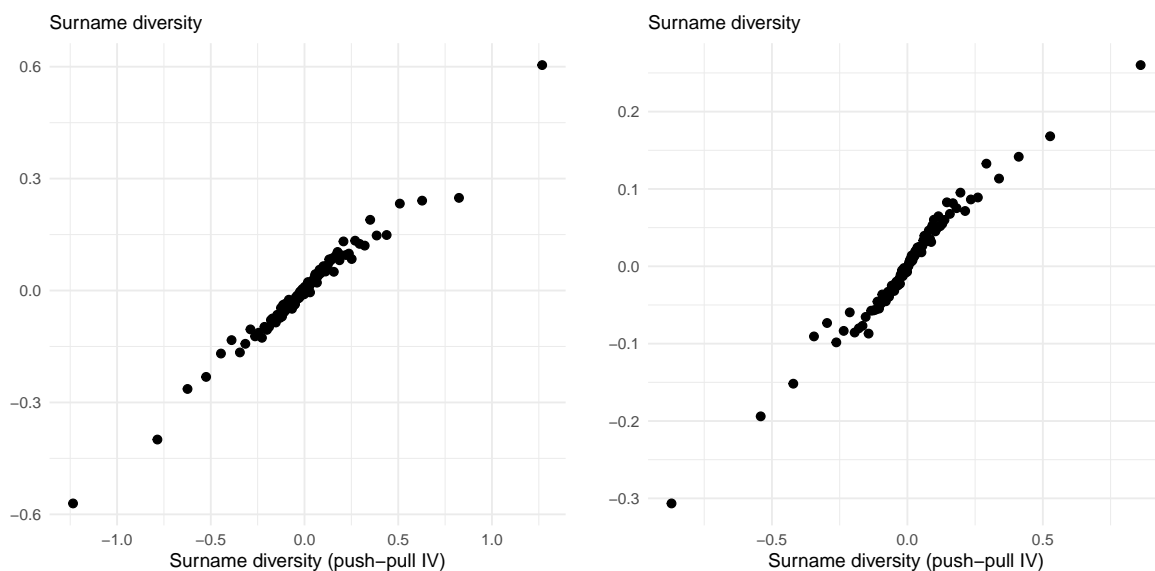


Figure B7: Binned scatter plots of surname diversity (pull-push IV) and actual surname diversity from 1900 to 1940

Notes: County-level data from 1900 to 1940 (including midyears). Observations are weighted by county population in 1900 and residualized by county fixed effects and state-period fixed effects (left plot) and county-specific time trends (right plot).

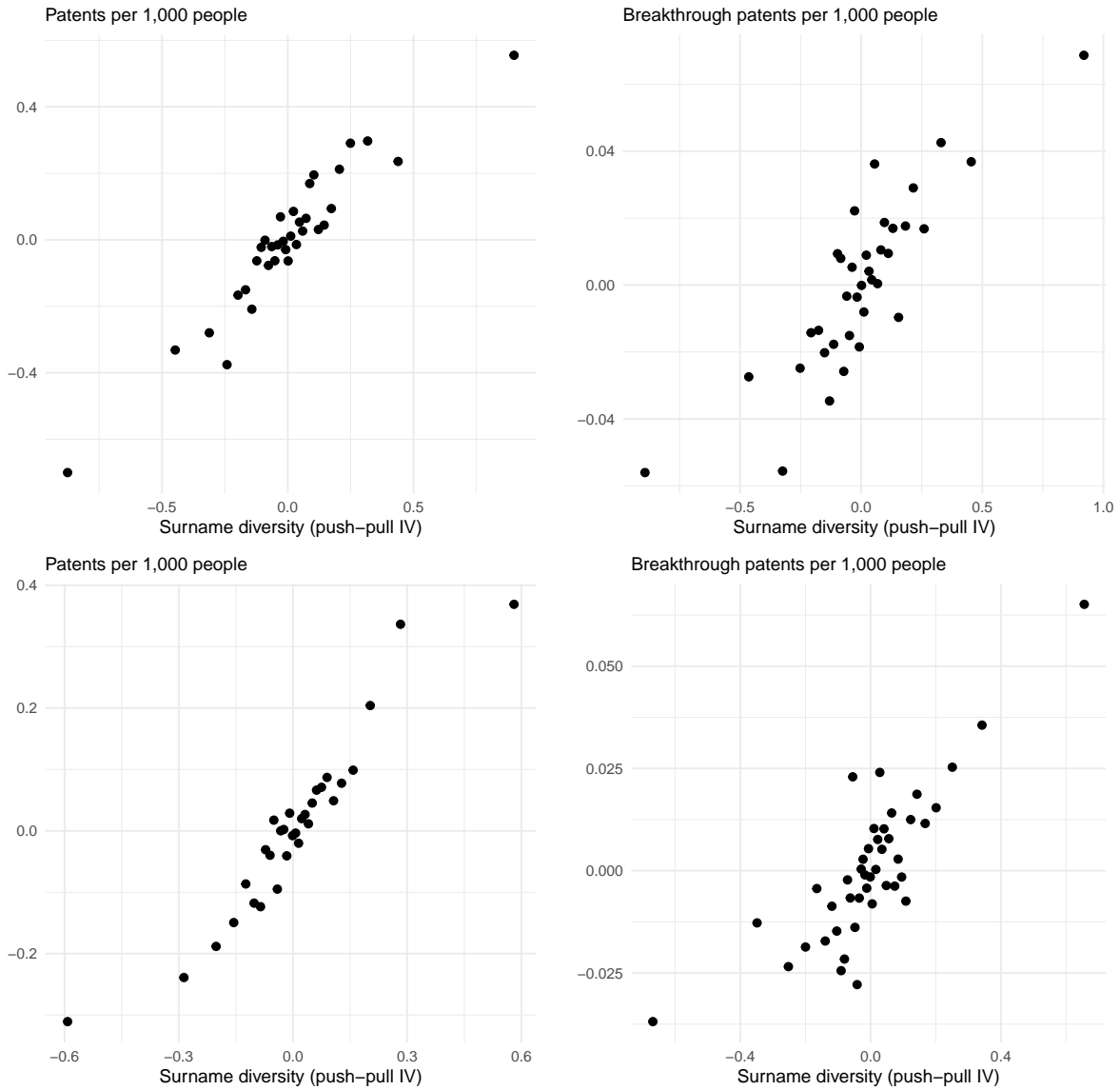


Figure B8: Binned scatter plots of surname diversity (pull-push IV) and innovation outcomes from 1900 to 1940

Notes: County-level data from 1900 to 1940 (including midyears). Observations are weighted by county population in 1900 and residualized by county fixed effects and state-period fixed effects (top plots), and county-specific time trends (bottom plots).

Table B8: Controlling and instrumenting for population size

	Patents per 1,000 people (mean = 2.04, sd = 2.6)			Breakthrough patents per 1,000 people (mean = 0.14, sd = 0.24)		
	(1)	(2)	(3)	(4)	(5)	(6)
	<i>Panel A: Least-squares estimates</i>					
Surname diversity	1.428*** (0.322)	1.425*** (0.339)	1.206*** (0.292)	0.170*** (0.041)	0.136*** (0.041)	0.111*** (0.038)
Population	0.069** (0.030)	0.079** (0.036)	0.118*** (0.042)	0.009*** (0.003)	0.010*** (0.002)	0.021*** (0.007)
<i>Panel B: Reduced-form estimates</i>						
Surname diversity (push-pull IV)	0.651*** (0.184)	0.753*** (0.131)	0.638*** (0.164)	0.084*** (0.018)	0.081*** (0.021)	0.060* (0.034)
Population (push-pull IV)	0.023 (0.043)	0.013 (0.062)	0.082 (0.121)	0.004 (0.004)	0.003 (0.005)	0.017 (0.018)
<i>Panel C: Instrumental-variable estimates</i>						
Surname diversity	1.460*** (0.363)	1.707*** (0.352)	1.523*** (0.522)	0.189*** (0.044)	0.184*** (0.056)	0.137 (0.096)
Population	0.032 (0.043)	0.031 (0.058)	0.163 (0.178)	0.005 (0.003)	0.005 (0.004)	0.031 (0.027)
KP <i>F</i> -statistic, Surname diversity	44.298	39.749	17.979	44.298	39.749	17.979
KP <i>F</i> -statistic, Population	76.245	93.742	295.820	76.245	93.742	295.820
County fixed effects	✓	✓	✓	✓	✓	✓
Period fixed effects	✓			✓		
State-Period fixed effects		✓	✓		✓	✓
County-specific linear time trends			✓			✓
Observations	23,660	23,660	23,660	23,660	23,660	23,660

Notes: The table reports the estimates of the least-squares, reduced-form, and IV estimates for the specification described in equation (5) but additionally controlling or instrumenting for (the instrument for) county population size, as predicted by our estimates for equation (3). An observation is a county-period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table B9: Placebo test and persistence: Instrumenting for population size

	$t-2$	$t-1$	t	$t+1$	$t+2$	$t+3$
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A:</i> Patents per 1,000 people						
Surname diversity	-0.569 (0.515)	0.074 (0.358)	1.523*** (0.522)	1.198*** (0.373)	1.280* (0.671)	0.851 (0.570)
Population	0.166 (0.142)	-0.027 (0.037)	0.163 (0.178)	0.049 (0.098)	-0.143*** (0.040)	-0.226 (0.144)
<i>Panel B:</i> Breakthrough patents per 1,000 people						
Surname diversity	-0.115 (0.101)	-0.005 (0.085)	0.137 (0.096)	0.261** (0.100)	0.198** (0.095)	0.071 (0.097)
Population	0.023* (0.013)	-0.008 (0.015)	0.031 (0.027)	-0.003 (0.021)	-0.010 (0.012)	-0.050 (0.035)
KP F -statistic, Surname diversity	15.161	15.802	17.979	26.215	21.463	30.603
KP F -statistic, Population	52.137	11.091	295.820	54.947	156.464	55.557
Observations	17,743	17,746	23,660	17,746	17,743	14,785
County fixed effects	✓	✓	✓	✓	✓	✓
State-Period fixed effects	✓	✓	✓	✓	✓	✓
Surname-Period fixed effects	✓	✓	✓	✓	✓	✓
County-specific linear time trends	✓	✓	✓	✓	✓	✓

Notes: The table reports IV estimates of the leads and lags of innovation outcomes on surname diversity for the specifications described in equation (5), but instrumenting county population using the push-pull instrument for county population size, as predicted by our estimates for equation (3). Columns 1 and 2 use the two-period and one-period lag of the dependent variables, respectively. Column 3 repeats the baseline specification (contemporaneous values of the dependent variables). Columns 4 to 6 use the one-period, two-period and three-period lead of the dependent variables, respectively. Observations are county-periods and weighted by county population in 1900. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table B10: Regional heterogeneity in the effect of surname diversity on innovation

	Patents per 1,000 people (mean = 2.04, sd = 2.60)			Breakthrough patents per 1,000 people (mean = 0.14, sd = 0.24)		
	(1)	(2)	(3)	(4)	(5)	(6)
	<i>Panel A: Least-squares estimates</i>					
Surname diversity × Region = Midwest	2.927*** (0.433)	2.807*** (0.446)	3.069*** (0.509)	0.251*** (0.038)	0.241*** (0.038)	0.286*** (0.068)
Surname diversity × Region = Northeast	1.596** (0.697)	3.112** (1.171)	1.983*** (0.721)	0.264*** (0.084)	0.378** (0.156)	0.358* (0.178)
Surname diversity × Region = South	0.826*** (0.150)	0.433*** (0.140)	0.367** (0.153)	0.071*** (0.017)	0.037*** (0.008)	0.015** (0.008)
Surname diversity × Region = West	1.697*** (0.564)	0.941** (0.380)	1.325*** (0.380)	0.201** (0.090)	0.066*** (0.023)	0.071*** (0.021)
<i>Panel B: Reduced-form estimates</i>						
Surname diversity (push-pull IV) × Region = Midwest	1.685*** (0.362)	1.919*** (0.328)	1.648*** (0.419)	0.151*** (0.035)	0.175*** (0.039)	0.168*** (0.062)
Surname diversity (push-pull IV) × Region = Northeast	0.430 (0.305)	0.646*** (0.170)	0.547*** (0.120)	0.099*** (0.036)	0.114*** (0.018)	0.098*** (0.026)
Surname diversity (push-pull IV) × Region = South	0.540*** (0.155)	0.357** (0.170)	0.287*** (0.106)	0.056*** (0.014)	0.035*** (0.011)	0.012** (0.005)
Surname diversity (push-pull IV) × Region = West	0.729*** (0.197)	0.369*** (0.085)	0.660*** (0.128)	0.061*** (0.022)	0.002 (0.010)	0.031* (0.019)
<i>Panel C: Instrumental-variable estimates</i>						
Surname diversity × Region = Midwest	3.300*** (0.653)	3.721*** (0.670)	3.978*** (1.093)	0.313*** (0.065)	0.339*** (0.079)	0.405** (0.171)
Surname diversity × Region = Northeast	1.237* (0.631)	2.875*** (0.640)	2.889*** (0.893)	0.254*** (0.077)	0.509*** (0.139)	0.517** (0.246)
Surname diversity × Region = South	1.098*** (0.244)	0.606** (0.265)	0.464*** (0.169)	0.122*** (0.025)	0.060*** (0.017)	0.019* (0.010)
Surname diversity × Region = West	1.637** (0.697)	0.928* (0.461)	1.753* (0.910)	0.142* (0.073)	0.005 (0.023)	0.083 (0.071)
Kleibergen-Paap <i>F</i> -statistic 1st coefficient	200.951	105.238	27.024	200.951	105.238	27.024
Kleibergen-Paap <i>F</i> -statistic 2nd coefficient	21.441	3.995	2.791	21.441	3.995	2.791
Kleibergen-Paap <i>F</i> -statistic 3rd coefficient	106.312	29.253	21.172	106.312	29.253	21.172
Kleibergen-Paap <i>F</i> -statistic 4th coefficient	30.594	2.363	1.981	30.594	2.363	1.981
County fixed effects	✓	✓	✓	✓	✓	✓
Period fixed effects	✓			✓		
State-Period fixed effects		✓	✓		✓	✓
County-specific linear time trends			✓			✓
Observations	23,660	23,660	23,660	23,660	23,660	23,660

Notes: The table reports regional heterogeneity in the least-squares, reduced-form, and instrumental-variable (IV) estimates for the specifications described in equation (5). An observation is a county in a period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table B11: Inverse hyperbolic sine transformed outcomes

	IHS Patents per 100,000 people (mean = 1.08, sd = 0.91)			IHS Breakthrough patents per 100,000 people (mean = 0.15, sd = 0.97)		
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Least-squares estimates</i>						
Surname diversity	0.479*** (0.081)	0.458*** (0.075)	0.381*** (0.068)	0.211*** (0.057)	0.159*** (0.053)	0.122*** (0.041)
<i>Panel B: Reduced-form estimates</i>						
Surname diversity (push-pull IV)	0.204*** (0.053)	0.228*** (0.042)	0.200*** (0.026)	0.107*** (0.026)	0.098*** (0.021)	0.086*** (0.027)
<i>Panel C: Instrumental-variable estimates</i>						
Surname diversity	0.457*** (0.086)	0.528*** (0.082)	0.517*** (0.119)	0.240*** (0.060)	0.226*** (0.065)	0.224** (0.091)
Kleibergen-Paap <i>F</i> -statistic	63.280	51.050	28.341	63.280	51.050	28.341
<i>Panel D: First-stage estimates</i>						
Surname diversity (push-pull IV)	0.445*** (0.056)	0.431*** (0.060)	0.386*** (0.073)	0.445*** (0.056)	0.431*** (0.060)	0.386*** (0.073)
County fixed effects	✓	✓	✓	✓	✓	✓
Period fixed effects	✓			✓		
State-Period fixed effects		✓	✓		✓	✓
County-specific linear time trends			✓			✓
Observations	23,660	23,660	23,660	23,660	23,660	23,660

Notes: The table reports least-squares, reduced-form, and instrumental-variable (IV) estimates for the specifications described in equation (5) and first-stage estimates for equation (4). An observation is a county in a period from 1900 to 1940. Observations are weighted by county population in 1900. The endogenous variable is county-level surname diversity in t . In columns 1 to 3, the dependent variable is inverse hyperbole sine (IHS) transformed number of patents filed in the county in the five-year period starting in t divided by county population size in 1900. In columns 4 to 6, the dependent variable is IHS transformed number of breakthrough patents filed in the county in the five-year period starting in t divided by county population size in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. The sources and construction of all variables are explained in Appendix A. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

Table B12: Herfindahl surname diversity

	Patents per 1,000 people (mean = 2.04, sd = 2.60)			Breakthrough patents per 1,000 people (mean = 0.14, sd = 0.24)		
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Least-squares estimates</i>						
Herfindahl surname diversity	0.476*** (0.125)	0.374*** (0.104)	0.389*** (0.099)	0.046*** (0.017)	0.029** (0.012)	0.024** (0.011)
<i>Panel B: Reduced-form estimates</i>						
Herfindahl surname diversity (push-pull IV)	-0.013 (0.138)	0.094 (0.155)	0.402*** (0.128)	0.027*** (0.010)	0.036* (0.019)	0.064*** (0.015)
<i>Panel C: Instrumental-variable estimates</i>						
Herfindahl surname diversity	-0.102 (1.092)	0.668 (1.232)	3.568* (1.927)	0.214* (0.122)	0.253 (0.208)	0.571 (0.375)
Kleibergen-Paap <i>F</i> -statistic	7.434	6.131	3.020	7.434	6.131	3.020
<i>Panel D: First-stage estimates</i>						
Herfindahl surname diversity (push-pull IV)	0.126*** (0.046)	0.141** (0.057)	0.113* (0.065)	0.126*** (0.046)	0.141** (0.057)	0.113* (0.065)
County fixed effects	✓	✓	✓	✓	✓	✓
Period fixed effects	✓			✓		
State-Period fixed effects		✓	✓		✓	✓
County-specific linear time trends			✓			✓
Observations	23,660	23,660	23,660	23,660	23,660	23,660

Notes: The table reports the estimates of the least-squares, reduced-form, and IV estimates for the specification described in equation (5) but using Herfindahl surname diversity as endogenous variable. An observation is a county-period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

B.4 Alternative Shift-share Instruments

The conventional shift-share approach in the immigration literature rests on the observations that migrants locate near people from the same country of origin (Altonji and Card, 1991; Card, 2001). We can adapt this approach to our context because migrants also settle near family members. That is, other than in the push-pull approach, we allocate newly arriving migrants (i.e., between t and $t - 1$) according to the preexisting share of people in a county (in year $t - 1$) with the same surname. This procedure allows us to calculate the predicted inflow of migrants by surnames in this period.

The construction of the instrument for surname diversity based on this method requires counties' previous-period stocks (and not just inflow) of each surname. To get the current stocks (i.e., in t), we add the last-period inflow (predicted via shift-share) to counties' previous-period stocks of each surname (i.e., in $t - 1$). Then, we apply the entropy formula to obtain the instrument for diversity.¹⁷

A concern with this calculation is that previous-period surname stocks are endogenous. The inclusion of county fixed effects in the estimation mitigates this concern somewhat because it shifts the focus from levels to changes in diversity; the previous period surname stocks are hence less important as a source of bias. Nevertheless, we follow the approach in Burchardi et al. (2021) and construct an additional shift-share instrument of surname diversity, which relies on the *predicted* stock of surnames in the previous period. These predicted previous-period stocks are calculated based on the historical push-pull approach. We add the predicted stock of surname in $t - 1$ (calculated based on the push-pull approach) to the predicted inflow between t and $t - 1$ (calculated based on the shift-share approach) to arrive at the predicted surname stocks in t .

Appendix Table B13 reports the estimates for both the IV specification that rests on the shift-share instrument alone (Panel A) and on the one that combines the shift-share and historical push-pull approach (Panel B). Consistent with our baseline results, the estimates are positive and highly significant for both patents and breakthrough patents per capita. Their point estimates are larger, though, they also estimated with more noise. In summary, our results are robust to the use of more conventional shift-share IV strategies.

¹⁷Since the 1940 census does not provide information on the immigration year, we cannot calculate a shift-share instrument for this period. Therefore, the sample in this robustness check is restricted to 1900 to 1930.

Table B13: Robustness: Alternative shift-share instruments

	Patents per 1,000 people (mean = 2.10, sd = 2.59)			Breakthrough patents per 1,000 people (mean = 0.14, sd = 0.25)		
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Shift-share IV using realized surname shares</i>						
Surname diversity	2.216*** (0.449)	2.197*** (0.488)	3.624*** (0.706)	0.301*** (0.061)	0.222*** (0.063)	0.474*** (0.140)
Kleibergen-Paap <i>F</i> -statistic	1,697	1,845	87	1,697	1,845	87
<i>Panel B: Shift-share IV using push-pull predicted surname shares</i>						
Surname diversity	2.539*** (0.670)	3.223*** (0.697)	4.917*** (1.142)	0.394*** (0.075)	0.379*** (0.090)	0.816** (0.373)
Kleibergen-Paap <i>F</i> -statistic	85	60	16	85	60	16
County fixed effects	✓	✓	✓	✓	✓	✓
Period fixed effects	✓			✓		
State-Period fixed effects		✓	✓		✓	✓
County-specific linear time trends			✓			✓
Observations	20,704	20,704	20,704	20,704	20,704	20,704

Notes: The table reports IV estimates for the specifications described in equation (5), but based on alternative shift-share procedures to construct the instrument for surname diversity. Panel A reports estimates for a shift-share instrument using realized surname shares, akin to [Card \(2001\)](#). Panel B reports estimates for a shift-share instrument using predicted surname shares based on the push-pull approach described in equation (3), akin to [Burchardi et al. \(2021\)](#). An observation is a county-period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered on states and reported in parentheses. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

B.5 Patent Technology Class Fixed Effects

Another potential concern with the interpretation of our findings is that patenting practices vary across industries and technologies (Moser, 2013), and these differences might affect our results.

Using the fact that the USPTO assigns a technology class to each granted patent, we assess this concern by estimating specifications that include patent class fixed effects to absorb any technology-specific traits. Similar to the surname fixed effects specifications in our main analysis, this requires us to change the unit of observation from county-period to patent class-county-period. The estimating equations are given by equations (8) and (9), where equation (8) is the first stage and equation (9) is the second stage.

$$\text{Surname diversity}_i^t = \gamma \widehat{\text{Surname diversity}_i^t} + \mu_{t,s(i)} + \mu_i + \mu_{t,c} + v_{i,c}^t \quad (8)$$

$$Y_{i,c}^t = \beta \text{Surname diversity}_i^t + \alpha_{t,s(i)} + \alpha_i + \alpha_{t,c} + \varepsilon_{i,c}^t \quad (9)$$

where i indexes counties, s states, t census years (including the midyears), and c patent class. There are 408 patent classes in our sample from 1900 to 1944. Examples of the patent class level are “Geometrical Instruments”, “Stoves and Furnace”, and “Chemistry: Electrical and Wave Energy”. As before, $\widehat{\text{Surname diversity}_i^t}$ is county i 's surname diversity in t , and $\text{Surname diversity}_i^t$ is county i 's predicted surname diversity in t . $Y_{i,c}^t$ now is the number of (breakthrough) patents (per 1,000 residents) in patent class c , filed in county i in the five-year period starting in t . Therefore, the innovation outcomes vary at the patent class-county-period level, while surname diversity remains defined at the county-period level. Importantly, we can now include patent class-period fixed effects, denoted by the parameter $\alpha_{t,c}$, which implies we non-parametrically control for patent class-specific confounders across periods, including differences in patenting practices across industries. The coefficient of interest is β . Observations are weighted by the number of people in a county in the year 1900. Standard errors are clustered in two ways, on states and patent class.

The results are reported in Table B14 and show that estimates are virtually unaffected by the inclusion of patent class fixed effects (in columns 2 to 4 and 6 to 8). All the estimates are highly significant in all specifications. Thus, we conclude that differences across technological categories do not affect our results.

Table B14: Patent technology class fixed effects

	Patents per 1,000 people (mean = 0.01, sd = 0.03)				Breakthrough patents per 1,000 people (mean = <0.01, sd = 0.01)			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
<i>Panel A: Least-squares estimates</i>								
Surname diversity	0.006*** (0.001)	0.006*** (0.001)	0.006*** (0.001)	0.004*** (0.001)	0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)	0.000*** (0.000)
<i>Panel B: Reduced-form estimates</i>								
Surname diversity (push-pull IV)	0.002*** (0.001)	0.002*** (0.001)	0.002*** (0.001)	0.002** (0.001)	0.000*** (0.000)	0.000*** (0.000)	0.000*** (0.000)	0.000** (0.000)
<i>Panel C: Instrumental-variable estimates</i>								
Surname diversity	0.005*** (0.001)	0.005*** (0.001)	0.005*** (0.001)	0.007** (0.003)	0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)	0.001** (0.000)
Kleibergen-Paap <i>F</i> -statistic	57.355	57.355	46.053	30.024	57.355	57.355	46.053	30.024
<i>Panel D: First-stage estimates</i>								
	Surname diversity							
Surname diversity (push-pull IV)	0.423*** (0.056)	0.423*** (0.056)	0.407*** (0.060)	0.356*** (0.065)	0.423*** (0.056)	0.423*** (0.056)	0.407*** (0.060)	0.356*** (0.065)
County fixed effects	✓	✓	✓	✓	✓	✓	✓	✓
Period fixed effects	✓				✓			
Patent class-Period fixed effects		✓	✓	✓		✓	✓	✓
State-Period fixed effects			✓	✓			✓	✓
County-specific linear time trends				✓				✓
Observations	8,264,856	8,264,856	8,264,856	8,264,856	8,264,856	8,264,856	8,264,856	8,264,856

Notes: The table reports least-squares, reduced-form, and instrumental-variable (IV) estimates for the specifications described in equation 9 and first-stage estimates for equation 8. An observation is a patent class in a given county in a period from 1900 to 1940. Observations are weighted by the population in a given county in the year 1900. In columns 1 to 3, the dependent variable is number of patents with *c* as the main technological category and filed by individuals in county *i* in the five-year period starting in *t* divided by population size in county *i* in 1900. The dependent variable in columns 4 to 6 is the corresponding number of breakthrough patents. Standard errors are two-way clustered on states and technological category and reported in parentheses. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

C Additional Results

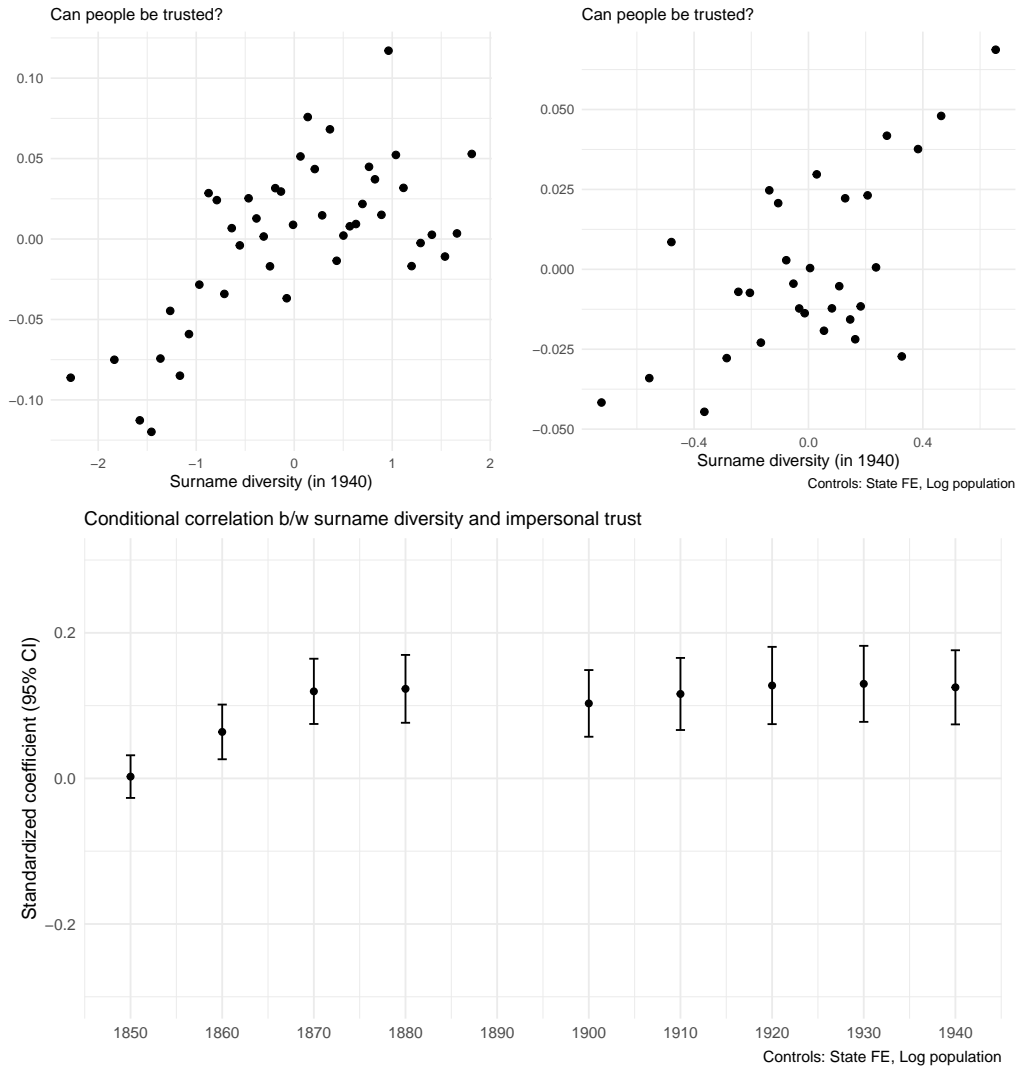


Figure C1: Relationship between surname diversity and impersonal trust

Notes: An observation is an individual. Top left: Bivariate relationship in 1940. Top right: Variables residualized by state fixed effects and log county population in 1940. Bottom: Coefficients of regressions of impersonal trust today on surname diversity conditional on state fixed effects and log county population by census year (1850-1940) and survey year, sex, age, and race fixed effects. The trust question is taken from the General Social Survey, waves 1972 to 2016.

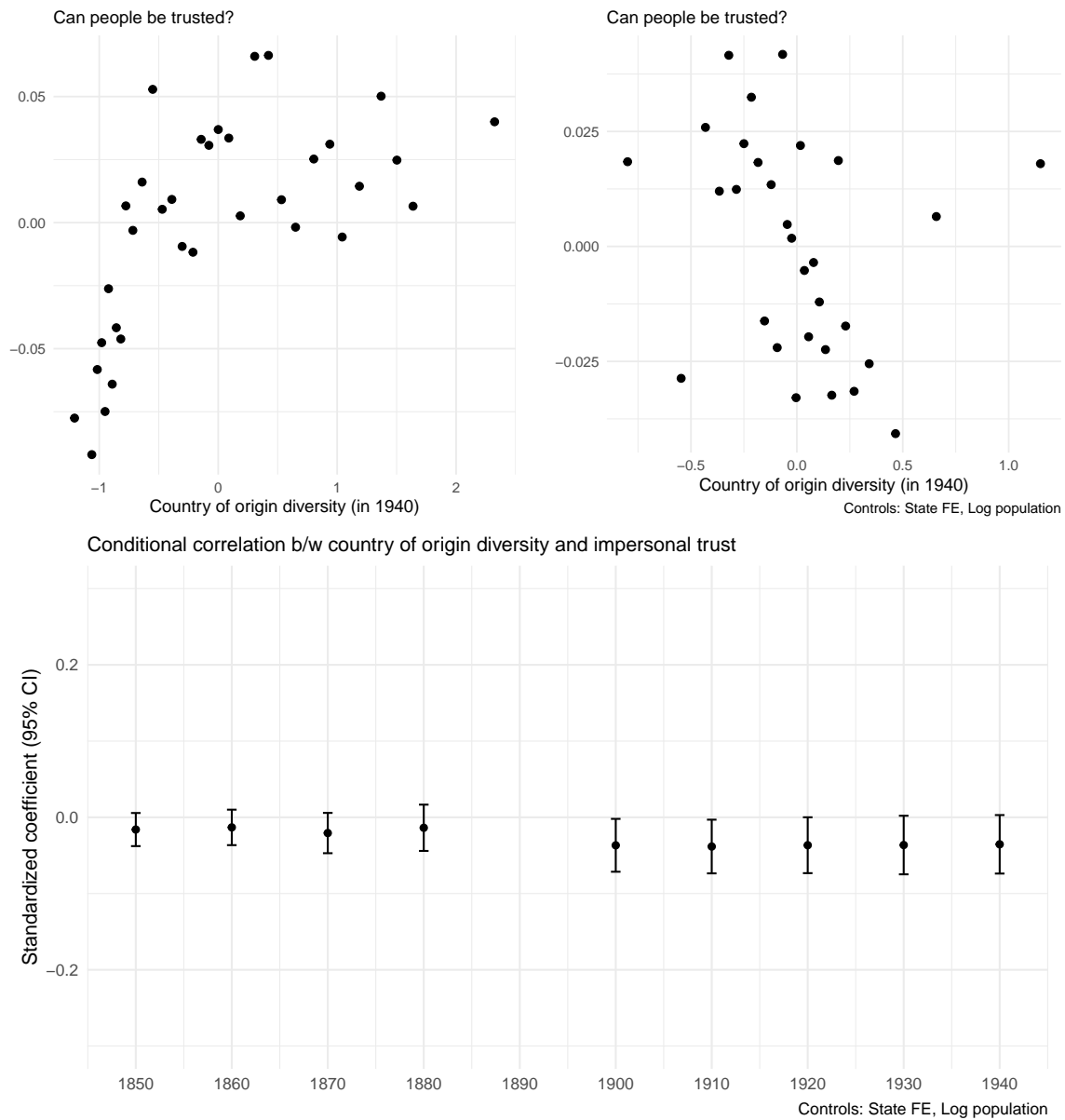


Figure C3: Relationship between country of origin diversity and impersonal trust

Notes: An observation is an individual. Top left: Bivariate relationship in 1940. Top right: Variables residualized by state fixed effects and log county population in 1940. Bottom: Coefficients of regressions of impersonal trust today on country of origin diversity conditional on state fixed effects and log county population by census year (1850-1940) and survey year, sex, age, and race fixed effects. The trust question is taken from the General Social Survey, waves 1972 to 2016.

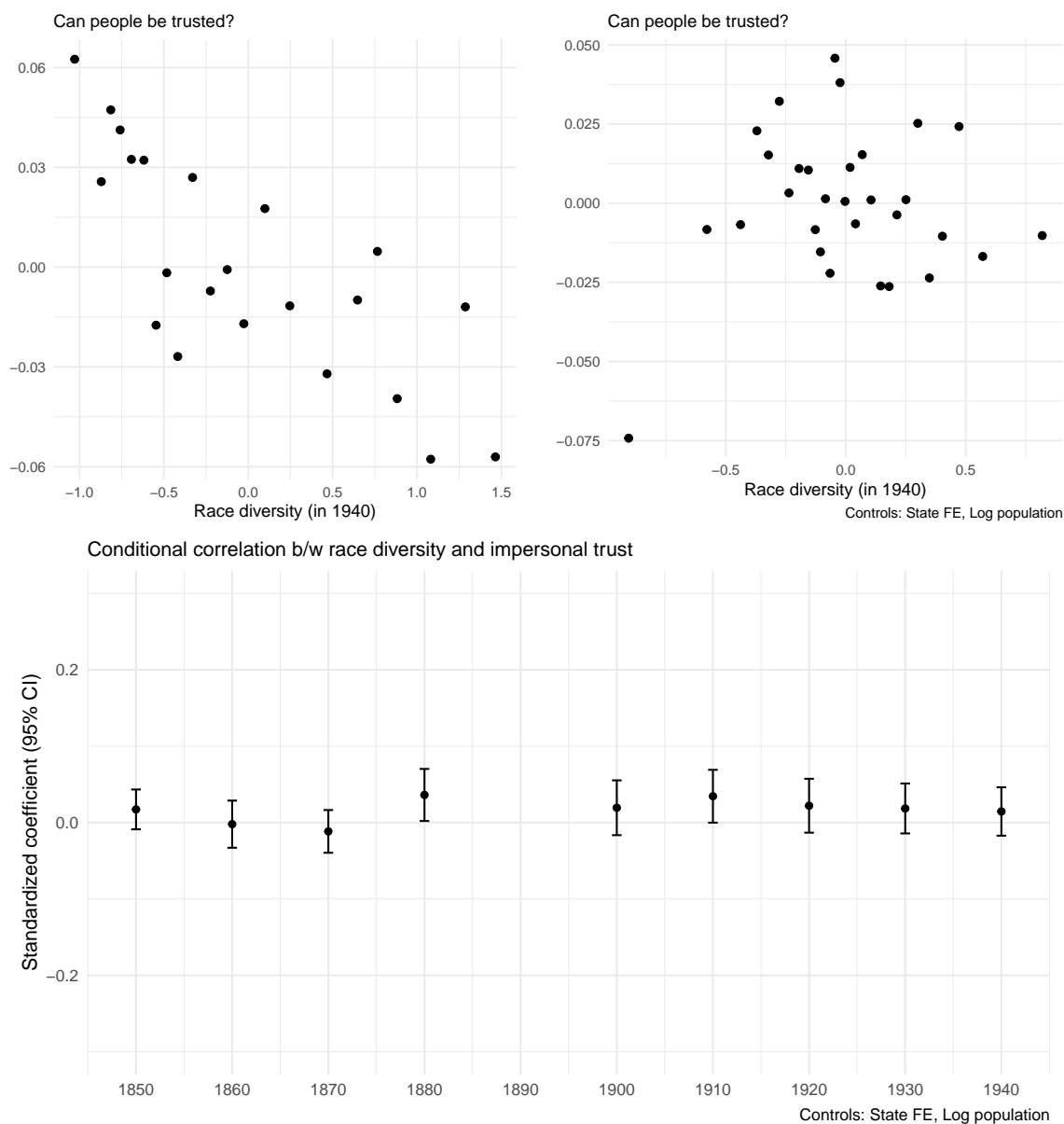


Figure C4: Relationship between race diversity and impersonal trust

Notes: An observation is an individual. Top left: Bivariate relationship in 1940. Top right: Variables residualized by state fixed effects and log county population in 1940. Bottom: Coefficients of regressions of impersonal trust today on race diversity conditional on state fixed effects and log county population by census year (1850-1940) and survey year, sex, age, and race fixed effects. The trust question is taken from the General Social Survey, waves 1972 to 2016.

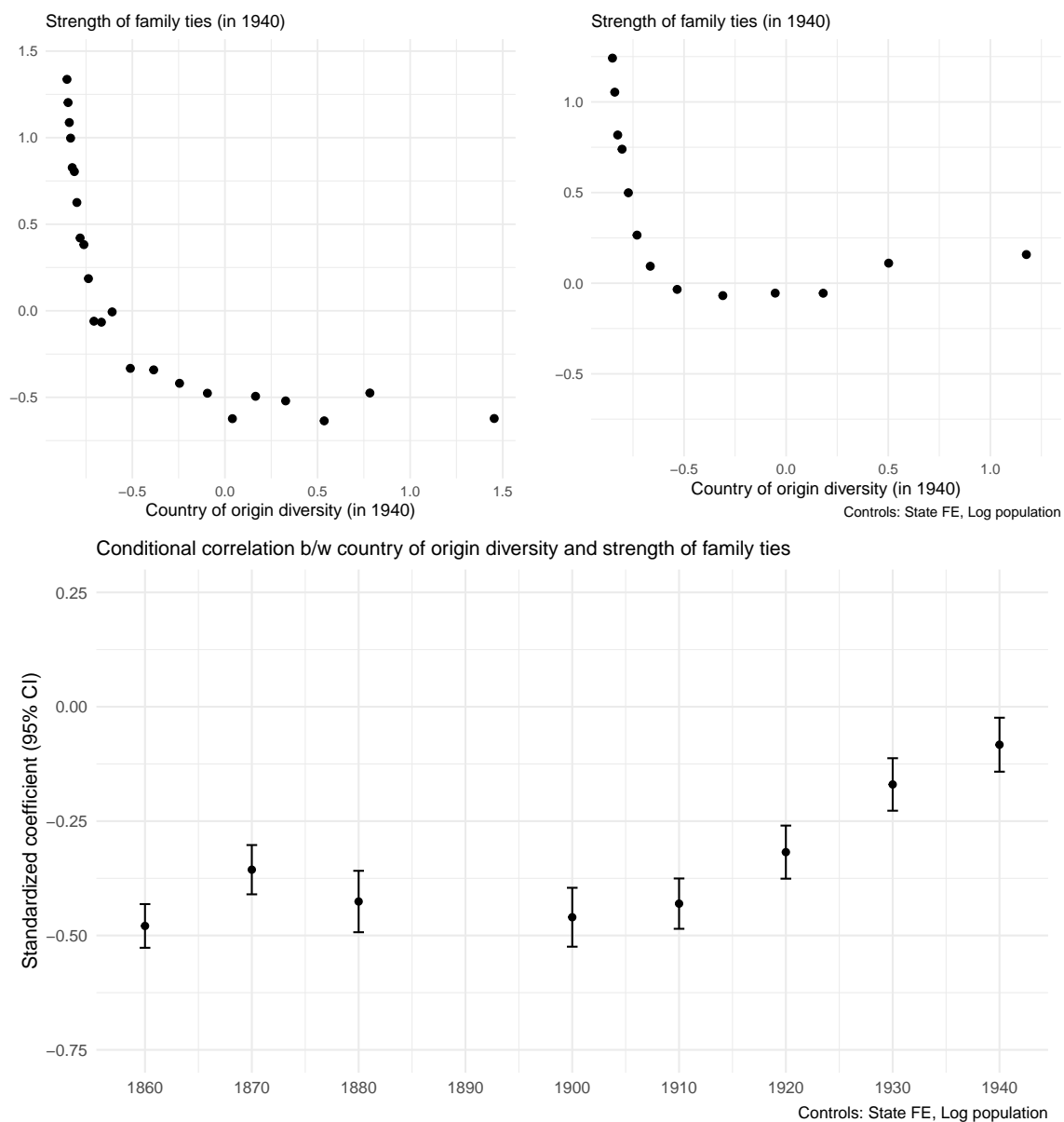


Figure C5: Relationship between country of origin diversity and strength of family ties
Notes: An observation is a county. Top left: Bivariate relationship in 1940. Top right: Variables residualized by state fixed effects and log county population in 1940. Bottom: Coefficients of regressions of strength of family ties on country of origin diversity conditional on state fixed effects and log county population by census year (1860-1940). The strength of family ties data is constructed following [Raz \(2023\)](#).

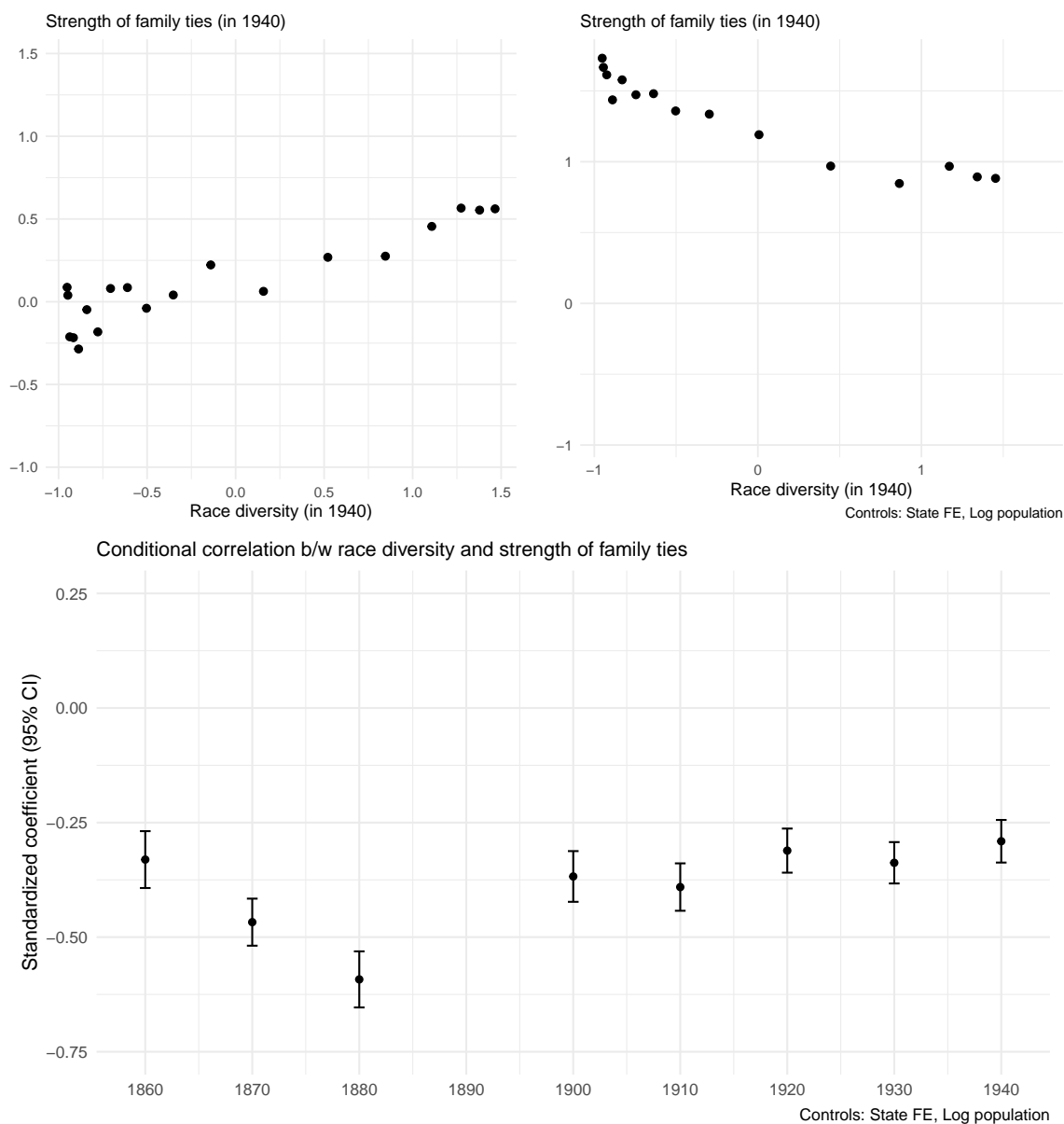


Figure C6: Relationship between race diversity and strength of family ties

Notes: An observation is a county. Top left: Bivariate relationship in 1940. Top right: Variables residualized by state fixed effects and log county population in 1940. Bottom: Coefficients of regressions of strength of family ties on race diversity conditional on state fixed effects and log county population by census year (1860-1940). The strength of family ties data is constructed following Raz (2023).

Table C1: Number of technologies on patents

	Average number of technologies on patents (mean = 2.34, sd = 0.67)		
	(1)	(2)	(3)
<i>Panel A: Least-squares estimates</i>			
Surname diversity	0.131*** (0.035)	0.071** (0.033)	0.040 (0.062)
<i>Panel B: Reduced-form estimates</i>			
Surname diversity (push-pull IV)	0.090*** (0.029)	0.071** (0.035)	0.045 (0.055)
<i>Panel C: Instrumental-variable estimates</i>			
Surname diversity	0.214** (0.085)	0.175* (0.104)	0.128 (0.172)
Kleibergen-Paap <i>F</i> -statistic	57.192	45.938	25.609
County fixed effects	✓	✓	✓
Period fixed effects	✓		
State-Period fixed effects		✓	✓
County-specific linear time trends			✓
Observations	20,257	20,257	20,257

Notes: The table reports the estimates of the least-squares, reduced-form, and IV estimates for the specification described in equation (5) with average number of technologies on patents as dependent variable. An observation is a county-period from 1900 to 1940. Observations are weighted by county population in 1900. Standard errors are clustered at the state level. All independent variables are standardized to mean zero and unit variance. ***, **, and * indicate significance at the 1%, 5%, and 10% levels.

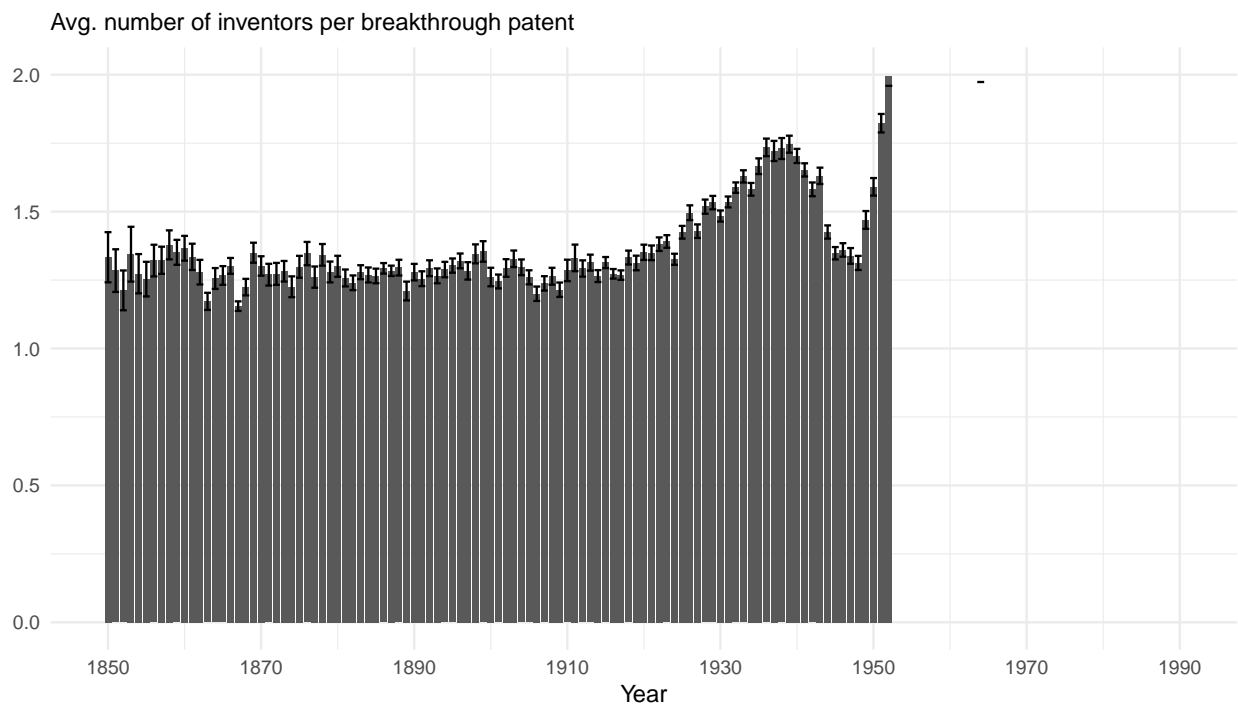
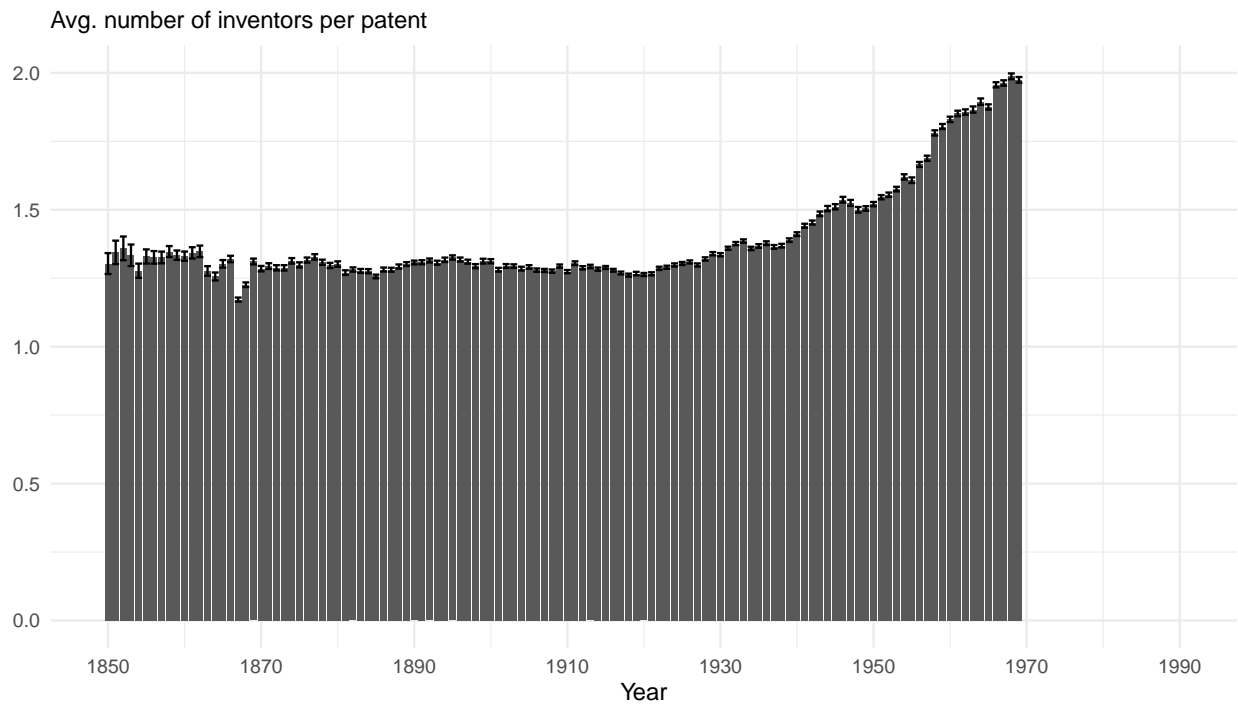


Figure C7: Average number of inventors per (breakthrough) patent over time
Notes: The figures show the average number of distinct inventors per (breakthrough) patent from 1850 to 1990. Error bars indicate standard errors.

Table C2: Do same-surname inventors produce lower-quality patents?

	Breakthrough patent indicator		
	(1)	(2)	(3)
Constant	0.107*** (0.001)		
Same-surname indicator	-0.048*** (0.002)	-0.030*** (0.002)	-0.011*** (0.002)
R ²	0.002	0.033	0.184
Observations	200,818	200,818	200,818
Year fixed effects		✓	✓
Patent technology class fixed effects			✓

Notes: An observation is a patent from 1850 to 1949 with at least two or more distinct inventors. The same-surname indicator takes value one if all inventors on the patent have the same surname. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

D Additional Descriptives

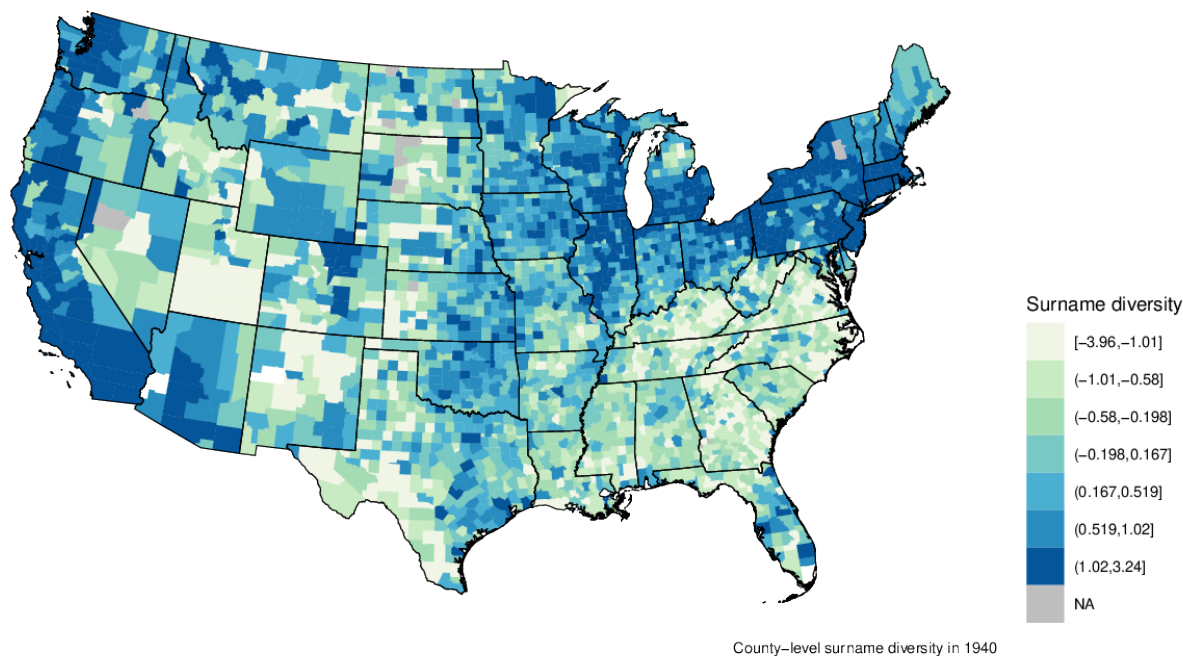


Figure D1: The figures show the geographic variation in surname diversity in 1940.

Table D1: Correlations between surname diversity and other diversities

	Country of origin diversity	Share immigrants	Race diversity	Occupational diversity
Raw Corr.	0.39	0.27	-0.24	0.60
Partial Corr. (Log Population)	0.60	0.50	-0.29	0.57
Partial Corr. (State FE, Log Population)	0.40	0.27	0.08	0.48

Notes: This table reports standardized coefficients of regressions of county-level surname diversity on other dimensions of sociocultural diversity from 1850 to 1940. The first row reports the relationship conditional on year fixed effects. The second row reports the coefficients of regressions additionally controlling for log county population. The third row reports the correlations additionally controlling for state fixed effects. An observation is a county from 1850 to 1940 (excluding the midyears). The sources and construction of all variables are explained in Appendix Section A.