

# The Long-Run Effect of Public Libraries on Children: Evidence from the Early 1900s

Ezra Karger

July 21, 2021

## Abstract

Between 1890 and 1921, Andrew Carnegie funded the construction of 1,618 public libraries in cities and towns across the United States. I link these library construction grants to census data and measure the effect of childhood public library access on adult outcomes. Library construction grants increased children's educational attainment by 0.10 years, did not affect wage income, and increased non-wage income by 4%. These income effects are driven by occupational choice. Access to a public library caused children to shift away from occupations like manual labor, factory-work, and mining into safer and more prestigious occupations like farm-ownership, clerical, and technical jobs. I show that compulsory schooling laws had parallel effects on children, increasing educational attainment, non-wage income and occupational prestige without affecting wage income. Economists often rely solely on wage income to measure the returns to education. But public libraries and compulsory schooling laws in the early 1900s increased educational attainment and had positive effects on children's adult labor market outcomes without affecting wage income.

---

Department of Economics, University of Chicago. I would like to thank Martha Bailey for providing me with extensive data describing county-level demographic characteristics and public libraries from prior work with Brian Jacob, Michael Kevane, and William Sundstrom. I would also like to thank Dan Aaronson, Enrico Berkes, Dan Black, Celeste Carruthers, Kerwin Charles, Sergio Correia, Daniel Feenberg, Camilo Garcia-Jimeno, Rick Hornbeck, Max Kellogg, Bhash Mazumder, Derek Neal, Peter Nencka, Matt Notowidigdo, Devin Pope, Lauren Sartain, and Marianne Wanamaker for their feedback and assistance.

## Introduction

Andrew Carnegie immigrated with his family from Scotland to Pennsylvania at the age of 13. He immediately began work at a cotton mill as a low-skilled bobbin boy, earning \$1.20 a week.<sup>1</sup> A nearby resident opened up his personal library to the “local working boys,” giving Carnegie access to a large set of books. Carnegie attributed much of his success as a businessman to this private library, and he wanted to provide a similar experience to children everywhere by funding the construction of public libraries in underserved communities.<sup>2</sup>

After a successful career in business, Carnegie followed through on this resolution, funding the construction of thousands of libraries worldwide, including 1,618 in the United States. He first funded the construction of a library in 1880, in his birthplace of Dunfermline, Scotland. In the following years he paid for the construction of a handful of public libraries in his home state of Pennsylvania. He then formalized his grant program, expanding it to other communities across the United States. Carnegie’s requirements were simple: a town had to agree to supply a public plot of land for the new library and fund the library annually at an amount greater than or equal to 10% of the initial grant value.

I use these Carnegie-funded public libraries to provide the first evidence that public libraries have long-run positive effects on children. I measure these effects by linking children from the 1900–1930 census to adult records in the 1940 census. I use an iterative matching procedure, increasing accuracy by relying on 1940 records that match to multiple early census years. I then use within-family variation in the timing of when children got access to a Carnegie Library to measure the causal effect of public library access on adult outcomes. My regression specifications are similar to those used by Aaronson and Mazumder (2011), who measure the causal effect of Rosenwald school construction grants on children’s school attendance during a similar time period. I show

---

<sup>1</sup>The New York Times, “Obituary: Carnegie Started as a Bobbin Boy,” August 12, 1919.

<sup>2</sup>Carnegie wrote in his autobiography: “I resolved, if ever wealth came to me, that it should be used to establish free libraries, that other poor boys might receive opportunities similar to those for which we were indebted to that noble man” (Carnegie, 1901).

that access to a Carnegie grant increased educational attainment, did not affect wage income, and increased non-wage income by causing children to shift into safer, more entrepreneurial, and more prestigious occupations.

To provide additional evidence that my results are not driven by confounders associated with each town's time-varying decision to apply for a library construction grant, I conduct a series of placebo tests and falsification exercises. I first show that the entry of college and university libraries in towns and cities across the United States had no effect on children's educational attainment. This is a falsification test because children likely did not have access to colleges and universities. I then perturb the date of each Carnegie grant by assuming that Carnegie offered the grant 20 years before (or after) it was actually offered. I show that these placebo grants had no effect on children.

Lastly, readers may worry that towns and cities invested in other local institutions that positively affected children at the same time that they constructed a public library. In particular, Goldin (1994, 1998, 1999) discusses the expansion of public high schools across the United States in the early 1900s, showing that this expansion had large effects on children. I construct a complete panel of all public high schools in the United States from 1890–1951. I use this data to show that controlling for the expansion of public high schools across the United States has no effect on my library results.

I connect my results to the findings of Stephens and Yang (2014) who provide evidence that compulsory schooling laws in the early and mid-1900s increased educational attainment but had null or negative effects on wage income. I replicate their results and extend them to additional outcomes, showing that compulsory schooling laws increased non-wage income and caused children to move into safer, more entrepreneurial, and more prestigious occupations. Both public library access and compulsory schooling laws increased non-wage income and made children more likely to enter occupations such as farm-ownership and technical work instead of occupations like machine-operation and manual labor. My results highlight the importance of non-wage income and non-pecuniary outcomes in measuring the returns to education and the causal effects of local institutions on children. This may be especially important in a historical context where a larger

fraction of the population was self-employed and engaged in highly unsafe occupations.

A handful of papers measure the causal effect of public libraries on individuals or communities. Kevane and Sundstrom (2014, 2016) measure the expansion of public libraries across the United States in the early 1900s and use aggregated data to argue that public libraries had no effect on county-level voting behavior. And contemporaneous work by Berkes and Nencka (2019) shows that Carnegie's library construction grants increased patenting activity when compared to towns receiving a grant offer that never materialized. Examining a more recent time period, Bhatt (2010) uses distance to a public library as an instrument for library access and argues that library access increases the amount of time children spend reading, the amount of time parents spend reading to their children, and homework completion rates. Bhatt also finds that access to a public library decreases the amount of time people spend watching television. And Gilpin, Karger, and Nencka (2020) show that sharp increases in public library investment causes increases in library usage and children's test scores. Lastly, Neto (2019) uses data from a modern public library census to argue that public library resources in Appalachia are not correlated with local employment and labor force participation rates.

My work draws on a large literature from the library sciences field, which ethnographically profiled and gathered information about Andrew Carnegie's library construction grants. Most notably, Bobinski (1969) digitized and standardized information from the Carnegie Corporation's microfilm archives, producing a table of all of the grants Carnegie made to towns and cities. In other important work, Daniel (1961) discusses the expansion of public libraries across the United States; Martin (1993) collects detailed descriptions of why dozens of cities turned down Carnegie's promised grants; and Klinenberg (2018) traces the importance of social infrastructure like public libraries in the functioning of society. These sources provided important historical context for my empirical findings.

# Data

## Library Data

Andrew Carnegie did not publicize his interest in funding the construction of new public libraries. Still, after news spread of his first grants to construct public libraries in Southwestern Pennsylvania, politicians and officials in hundreds of cities and towns sent Carnegie unsolicited requests for funds to build new public libraries. Carnegie did not have time to respond to the requests he received for library funding, so after directing and overseeing the first few grants himself, he quickly put his personal secretary, James Bertram, in charge of the application process. Bertram would only seriously consider requests from city and town officials, but many members of the general public sent him letters requesting funds for a library. When he received a letter from a local resident of a town, he asked that they find an elected official who could submit an official request for a library (Bobinski 1969). In most cases, Bertram conducted the entire application and grant process by mail. Bertram also ensured that library construction grants only went to towns that committed to spending 10% of the grant amount on annual upkeep of the library. Carnegie scaled his grant amounts by the population of each town, targeting a grant amount of \$2 per person in most cases. So in letters to cities and towns requesting funds, Bertram requested current population counts, which he often verified using publicly available tabulations from decennial census data (Bobinski 1969).

Carnegie allowed Bertram to give away the vast majority of \$37 million in the early 1900s with little oversight.<sup>3</sup> After Bertram corresponded by mail with city or town officials, he would forward applications satisfying the basic requirements to Andrew Carnegie, who rubberstamped dozens of grant applications in short meetings. Carnegie trusted Bertram completely, and for good reason. Bertram was well-known for being fair in his assessment of applicants. As George Bobinski ex-

---

<sup>3</sup>\$37 million is the nominal dollar amount Carnegie used to fund library construction in the United States between 1890 and 1921. Using a CPI inflator, that is equivalent to \$1 billion in 2019. But U.S. GDP was only \$34 billion in 1910, so as a constant fraction of GDP, \$37 million in 1910 is equivalent to close to \$21 billion today.

plained in his authoritative book about the Carnegie library grant program, “[Bertram’s] aloofness was attributed by one observer to a desire to maintain a strict impersonal and disinterested attitude toward each and every applicant... He judged proposals strictly on their merit. Personal relations or considerations never influenced his judgment. No worthy applicant was to be rejected, and, yet, no unworthy one was to be accepted” (Bobinski 1969, p. 30).<sup>4</sup> After Carnegie approved each grant, it took an average of 2.4 years for the town to build a Carnegie-funded public library. And 90% of the Carnegie-funded public libraries opened within four years of the initial grant.

The Carnegie Corporation preserved the correspondence between Bertram and community leaders. Drawing on these documents, Bobinski painstakingly collected exact dates, locations, and grant amounts for the grants Carnegie gave to cities and towns. In total, Bobinski calculated that Carnegie funded the construction of 1,618 libraries in 1,417 cities in the United States.<sup>5</sup> Almost all places that received a grant got funding to build one library, but a small number of cities received funding for an entire branch system. For example, New York City received funding to construct 66 libraries and Philadelphia received funding to construct 25 libraries. I hand-collect the opening date of each Carnegie library from archival sources.

In Figure 1, I link the digitized data from Bobinski to cities and I mark the modern-day zip codes of cities that received money from Carnegie.<sup>6</sup> Figure 1 shows several important patterns. First, New England was densely populated but received few grants. This is because it already had hundreds of public libraries before Carnegie began his philanthropic endeavors (see Figure 2, described below). Second, grants were concentrated in the Midwest, close to Carnegie’s hometown in southwest Pennsylvania. Third, communities in the South applied for grants at lower rates because

---

<sup>4</sup>While Bertram funded most applications, Bobinski (1969) compiled an additional list of 225 failed grants. Carnegie offered these 225 grants to cities and towns, but the grant offer fell through for various reasons. In many cases the grant fell through because a local philanthropist offered to fund a new public library on his or her own. And even when a town rejected a Carnegie grant, they often built a new public library soon after. So I do not use these failed grants in my analysis.

<sup>5</sup>Carnegie also funded a significant number of libraries outside of the United States. In this paper, I do not consider those grants.

<sup>6</sup>Kevane and Sundstrom (2014) digitized this information from a table in Bobinski’s book.

they worried that Carnegie would require the libraries he funded to racially integrate (Bobinski, 1969).<sup>7</sup> Fourth, as Figure 1 shows, some cities received grants in the late 1800s while other cities received grants in the early 1900s. My econometric specifications will rely on spatial and intertemporal variation in the availability of Carnegie libraries, and Figure 1 confirms that there is a significant amount of variation within states and over time in the availability of Carnegie libraries.

I augment this data with information from library censuses. Between 1875 and 1929, the U.S. Bureau of Education surveyed libraries in the United States to gather information about the name, location, type, and founding year of every library in the United States with a non-negligible number of volumes. In some years, that number was as low as 300 books.<sup>8</sup> The type variable is not consistent across survey years, so it is difficult to differentiate public libraries from subscription libraries. I follow the methodology of Kevane and Sundstrom (2014) to identify public libraries by looking for keywords in the library name like “public library” or “free library.”<sup>9</sup> In Figure , I use this data to map the founding date of the first public library in each town in the United States, through 1929. And in Figure 3, I plot the evolution of public libraries and Carnegie grants over time. Figure 3 shows that by 1920, more than 30% of the towns and cities in the U.S. with a public library had constructed a Carnegie-funded public library.

## Census Microdata

I use complete count census data to measure the causal effect of public library access on children. Researchers at the University of Minnesota’s Institute for Social Research and Data Innovation cleaned and standardized this data, and the 1900–1940 datasets are publicly available at

---

<sup>7</sup>This worry was unfounded. Carnegie did not require that towns integrate the public libraries that he funded.

<sup>8</sup>Kevane and Sundstrom (2014) digitized this data and Bailey, Jacob, Kevane, and Sundstrom (2011) used the data in unpublished work linking the expansion of public libraries across the United States with county-level economic outcomes.

<sup>9</sup>More specifically, I call a library a **public** library if the library name contains one of these phrases: ‘public library,’ ‘county library,’ ‘city library,’ ‘village library,’ or ‘free library.’ The library name also must not contain any of these phrases: ‘school,’ ‘department of public instruction,’ ‘college,’ or ‘university.’ The founding dates in the public library data are unreliable, and often misreported. Still, I show that Carnegie grants are highly correlated with the founding dates of public libraries from this survey data.

[usa.ipums.org](http://usa.ipums.org) (Ruggles et al., 2020). I access these datasets as part of a license from the National Bureau of Economic Research (NBER), which hosts a restricted copy of the datasets that includes the full name of each respondent. I produce all of the results in this paper on NBER's servers. I measure childhood demographic characteristics in the 1900–1930 census data and I measure adult outcomes in the 1940 census, which was the first national census in the U.S. to contain individual-level educational attainment and income information. I restrict my attention to men, because many women's last names changed from childhood to adulthood, making it difficult to match women across decennial census years.

I link the 1900–1930 census records to 1940 census records using an iterative matching technique, drawing from research, programs, and thoughtful descriptions by Ferrie (1996), Abramitzky, Boustan, and Eriksson (2012, 2014), and Bailey et al. (2019). See the Matching Appendix for more details about the procedure I use. To summarize the method, I use 36 combinations of exact, cleaned, and phonetic first and last name, birthplace, race, and age to search for possible matches between 0–25 year old men in each of the four early census years (1900–1930) and 15–65 year-old men in the 1940 census. In each of the 36 iterations, I relax an exact matching requirement and drop any uniquely identified child who matches to multiple possible adults. As an example, in the first pass I look for unique perfect matches between records from the 1900 and 1940 census using phonetic name codes, race, birthplace, and age. In the second iteration, I allow age to be one year higher in the 1940 data than in the 1900–1930 data, and in the third iteration, I allow age to be one year lower in the 1940 data. After checking for matches within one- and two-year age bands, I relax additional restrictions. This method of matching closely follows Abramitzky, Boustan, and Eriksson (2012, 2014), except that after finding high-quality matches, I search the remaining unmatched records to attempt to find lower-quality matches so that I can better analyze adult census records in 1940 that match to multiple early childhood records.

In the set of 39.2 million black, white, and Native-American men in the 1940 census born between 1875 and 1920, 25.7 million (66%) match to a unique record in at least one earlier census

year. 12.3 million adult men in 1940 match to one childhood census record and 13.3 million adult men in 1940 match to unique record in multiple early census years. I use two tests of match accuracy to further refine my sample. First, in the 1940 decennial census, sample-line respondents reported their parents' birthplaces. I test link accuracy by comparing these responses to reported parental birthplaces from the matched childhood census record (following Bailey et al. 2019). Second, I focus on the 13.3 million adult records in 1940 that match to unique records in two early censuses. If a 45 year-old John Smith in the 1940 census matches to the record of a 5 year-old John Smith in 1900 and a 15 year-old John Smith in 1910, then I can compare the reported names of John Smith's parents in 1900 and 1910. If the first two letters of John Smith's mother's and father's first names in 1900 are different from the first two letters of the mother's and father's first names in 1910, then it is likely that one of those early census records is not a correct match.<sup>10</sup> In those cases, I do not include that 1940 record in my matched sample. For 1940 records that match to multiple early census records, I keep the earliest matched census record in my matched sample.

My final matched sample consists of 17.8 million people with both an adult census record in 1940 and a childhood census record in 1900, 1910, 1920, or 1930. In Table 1, I show average characteristics of the full set of 1940 adult men, my matched sample, and a subset of the matched sample consisting of 5.8 million people for whom I additionally match at least one brother to adult records in the 1940 census. I use this sample for my within-family analysis. My overall match rate of 1940 records to childhood records, at 45.5%, is at the high end when compared with other linking strategies (Abramitzky et al. 2020; Bailey et al. 2019).<sup>11</sup> And I estimate that roughly

---

<sup>10</sup>This is unlikely to happen by chance. The most common first two letters of fathers' first name in my matched sample are 'JO' and the most common first two letters of mothers' first name in my matched sample are 'MA'. Each occurs around 15% of the time because of the large number of people named John and Mary in the early 1900s. So assuming independence, the probability that two randomly selected boys will have a father or mother with the same first two letters of their first names is at most  $2 * 0.15^2 = 4.5\%$ . I rely on first letters rather than exact matches of parent names across childhood census records to allow for transcription errors, which are quite common in digitized early census records. And in cases where a child is reported in a household with only one parent, I use only that one parent for this comparison.

<sup>11</sup>This is unsurprising, because for each adult in 1940 I often have two attempts to find their childhood census record in an early census year. Most current approaches to census linking attempt to link records between only two decennial censuses.

80% of these matches are correct,<sup>12</sup> which is a high level of accuracy relative to other matching techniques with a similar overall match rate (see Abramitzky et al. 2020 and Bailey et al. 2019 for accuracy rates using different linking methods.). See the Matching Appendix for details about these calculations.

Each early census record identifies the respondent’s enumeration district, and I use these enumeration districts to link each adult to the city or town in which he lived as a child. Enumeration districts in the early 1900s defined a contiguous geographic area where one census enumerator could knock on every door and collect census information for every household within 2–4 weeks (Census Bureau, 1940). So cities like Chicago contained many enumeration districts, while sparsely populated enumeration districts could contain an entire county. There were around 65,000 enumeration districts in each census year. In the census data, there are also place names associated with most respondents. For example, respondents living in one enumeration district have a place name of “Precinct 16, Cahaba,” referring to the 16th precinct of Cahaba, Alabama: a town with a population of only a few hundred people in 1900. The average enumeration district had a population of roughly 1,200, making enumeration districts a granular measure of geographic place. See the Data Appendix for more details about the raw geographic variables I use in this paper.

I clean and standardize these raw places names and I use the cleaned place names to match census records to Carnegie grants and public libraries. In other words, I link each of Carnegie’s grant to children in the 1900–1930 census living in towns and cities with the exact names corresponding to each grant. Of the 1,417 successful grants, I map all but seven to at least one enumeration district. And many grants match to a large number of enumeration districts (like the grant Andrew Carnegie gave to New York). In Table 2, I list the most populous cities and towns in my sample of matched census records, along with the date Carnegie offered a grant to that city, if a grant was offered. The matched sample is not concentrated in a small number of places. Chicago contained

---

<sup>12</sup>In other words, for every 100 records in 1940 that are in my final matched sample, 80 are matched to a childhood census record corresponding to the same person, and 20 are not.

2.3% of the matched sample, and the boroughs of New York City together make up 3.3% of my matched sample. My results are robust to the exclusion of cities and towns containing more than 0.2% of my matched sample, which excludes all major cities.

Cities and towns provide the best unit of analysis for this project because Carnegie grants targeted individual cities and towns. In a textbook for students interested in working in library circulation occupations, Flexner (1927) differentiates between patrons who could use the public library for free, and patrons who paid a subscription fee. Flexner explains that local students, residents of the local service area, and non-residents with business addresses in the “town or county” would normally receive “free access” to the local public library. But “all other applicants living outside the legal boundaries of the library service area” would be given library access only “on payment of an annual, semi-annual, or monthly fee.” In other words, the public library was freely available only to local residents. And since Carnegie gave library construction grants to cities and towns, which raised revenue to fund the library, in most cases only residents of the city or town who paid for library maintenance would have had free access to the Carnegie library. Another reason to rely on cities and towns as the unit of analysis is distance. I am measuring the effect of library access on children; if children lived in the same county as a Carnegie library, but a different town in that county, getting to the library would have been more difficult. So when I compare children in neighboring towns within a county, I will attribute differences in adult outcomes to this differential library access.

## **Methodology**

For each person in my matched sample, I construct two measures of Carnegie grant exposure. First, in my main specifications I use the binary measure  $1(GrantByAge5)$ —an indicator for whether each child received a Carnegie grant by the age of 5. Almost all Carnegie-funded libraries were constructed within four years of when the grant was given. So, this treatment indicator is an

indicator of whether each child had access to a Carnegie library during their childhood. Second, in robustness checks I analyze *ExposureCarnegie*, a measure of the fraction of years between age 5 and 20 when a child had access to a Carnegie grant. I use the data collected by Bobinski (1969) to calculate these two measures. As an example, if the decennial census records John Smith as an eight year-old living in New York City in 1900, I look in the Carnegie grant data and see that New York City received a Carnegie grant in 1899. So I calculate John Smith's exposure to a Carnegie grant as  $1(\text{GrantByAge5}) = 0$  and  $\text{ExposureCarnegie} = \frac{14}{16} = 0.88$ . Or more generally, for a child born in year  $y$  and living in city  $c$  as a child,

$$\text{ExposureCarnegie}_{y(i),c(i)} = \frac{1}{16} \sum_{t=5}^{20} 1(\text{CarnegieGrant}_{y+t,c})$$

where  $1(\text{CarnegieGrant}_{y+t,c})$  is an indicator for whether city  $c$  received a Carnegie grant during or before year  $y+t$ .

I can now analyze the effect of Carnegie grant exposure on adult outcomes of children. I model that relationship using specifications of the form:

$$\text{Outcome}_i = \beta_0 + \beta_1 * 1(\text{GrantByAge5})_{y(i),c(i)} + \Pi X_i + \varepsilon \quad (1)$$

where  $i$  indexes individuals,  $\text{Outcome}_i$  is an outcome measured in the 1940 census, and  $X_i$  is a matrix of controls. My use of cohort\*city-level variation in exposure to Carnegie grants draws from several papers measuring the effect of school construction on children. Duflo (2001) uses the exposure of different cohorts of students in regions of Indonesia to a school construction program and measures the effect of the program on educational attainment. In related work, Neilson and Zimmerman (2014) link a staggered school construction program in New Haven, Connecticut to neighborhood and child outcomes to measure the effect of the program on children. And Aaronson and Mazumder (2011) link Julius Rosenwald's funding for black schools in the South in the early 1900s to children's school attendance using a methodology that I follow closely.

My choice of specifications draws from Aaronson and Mazumder (2011), who measure the effect of exposure to a constructed Rosenwald school between ages 7 and 13 on school attendance. I follow their strategy of beginning with a model that controls for baseline characteristics, and progressively adding controls until my final specification includes combinations of family fixed effects, county-by-birth year fixed effects, and birth order fixed effects. In my most basic specification,  $X_i$  includes birth year, state of birth, and birth order fixed effects. I then progressively add more control variables, including parental characteristics, enumeration district fixed effects, county-by-birth year fixed effects, and census microfilm page fixed effects. Because enumerators traveled door-to-door collecting data, census microfilm pages often list households in the order of neighboring houses; this makes census microfilm pages a granular measure of place (Logan and Parman, 2017; Magnuson and King, 2010). In my within-family specification,  $X_i$  also includes household fixed effects, allowing me to compare siblings within a household to each other and to similarly situated children in the same county who received no (or less) access to a Carnegie grant. I interact all of these fixed effects, and the fixed effects included in specifications throughout this paper, with race-by-childhood census year fixed effects. The coefficient of interest is  $\beta_1$ , which measures the effect of exposure to a Carnegie grant on children’s adult outcomes.

In my within-family specification, I perform a more complicated version of this thought experiment: Consider two households (A and B), living in different towns in the same county in 1910. Each household has two sons, ages 0 and 16, who I will index  $A_{oldest}$ ,  $A_{youngest}$ ,  $B_{oldest}$ , and  $B_{youngest}$ . Household A lives in a town that receives a Carnegie library grant in 1915. Household B’s town never received any grant. This means that  $A_{youngest}$  has  $1(GrantByAge5) = 1$  and the other three children have  $1(GrantByAge5) = 0$ . Let  $Y_i$  be years of educational attainment in 1940.  $\beta_1$  measures the difference-in-difference of educational attainment between these four children. Or algebraically,  $\beta_1$  would be  $(Y_{A_{youngest}} - Y_{A_{oldest}}) - (Y_{B_{youngest}} - Y_{B_{oldest}})$ . And if the youngest child in household A obtains an abnormally high level of educational attainment, I will attribute this to his exposure to a Carnegie grant. In the regression specifications that I present, there are

more than two million families and  $\beta_1$  represents a more complicated weighted average over the within-family differences in grant exposure and outcomes.

To theoretically motivate the within-family model, consider a framework where unobserved heterogeneity at the family-level is correlated with exposure to a Carnegie grant. Following Todd (2008), let  $Y_{ij}$  be the outcome for person  $i$  from family  $j$  and assume that the data-generating process is:

$$Y_{ij} = \phi(X_{ij}) + D_{ij}\beta + \theta_j + v_{ij}.$$

Where  $D_{ij}$  is an indicator for whether person  $i$  in family  $j$  was treated,  $X_{ij}$  is a series of controls, like birth cohort or birthplace, and  $\beta$  is the coefficient of interest.  $\theta_j$  is a family effect, which could represent genetics, family tradition, or the decision of the family to live in a particular place, all of which may be correlated with access to a Carnegie grant. For example, parents may move to a town because they expect that it will receive a Carnegie grant in the near future. Or, families with a high  $\theta_j$  might be more likely to live in towns with a Carnegie-funded library.  $v_{ij}$  is an unobserved error term that is uncorrelated with  $\theta_j$ ,  $D_{ij}$ , and  $X_{ij}$ . Controlling for family fixed effects allows me to control for any family-level characteristics that are correlated with Carnegie grant exposure. If I did not include  $\theta_j$  in my regression specifications, and it was correlated with  $X_{ij}$  or  $D_{ij}$ , we would have a standard omitted variable problem. Chetty and Hendren (2018) rely on a similar framework to measure the causal effect of neighborhoods on adult outcomes of children. If there is heterogeneity in the treatment effect so that  $\beta_i$  varies with  $i$ , then the within-family specification recovers a weighted average of the individual-level treatment effects, where the weights correspond to within-family variation in treatment.

# Results

## Education

I begin by looking at the effect of Carnegie grant exposure on years of educational attainment, high school graduation rates, and college attendance rates. These results provide strong evidence that access to a Carnegie grant increased educational attainment. In Table 3, I present the main set of regression results (see Equation 1) with years of schooling as an outcome. Access to a Carnegie grant is consistently associated with increased educational attainment across specifications. With a baseline set of controls, children with access to a Carnegie grant by age five had 0.57 more years of educational attainment than children without access to a Carnegie grant. But most of this gap is driven by selection on observables. When I control for neighborhood fixed effects, the difference drops to 0.18 years, and the effect of a Carnegie grant on years of educational attainment settles between 0.08 and 0.13 years when I add additional controls. In my within-family specification, a Carnegie grant received by age five increases educational attainment by 0.10 years.<sup>13</sup>

The difference between the magnitude of the treatment effect in the first five models and the sixth within-family model is not because of the sample's changing composition. I can use only 30% of my sample to estimate the within-family model because the other 10 million children do not have a brother in my matched sample. But when I re-run all models with only the sample of matched brothers from the within-family model, identical patterns emerge (see Appendix Table A2). Exposure to a Carnegie grant had a positive and statistically significant effect on educational attainment in all models using this sample. The magnitude of the treatment effect remains between 0.10 and 0.12 for this sample of brothers in all specifications containing enumeration district fixed effects and parental controls.

---

<sup>13</sup>Throughout this paper, I cluster standard errors at the level of the childhood county of residence, following Aaronson and Mazumder (2011). The United States contains over 3,000 counties. Clustering at the state of residence level (51 clusters) only marginally increases the size of standard errors in some specifications, and does not change the statistical significance of my coefficient of interest (see Appendix Table A1). I calculate large numbers of multi-way fixed effects using the `reghdfe` Stata package, introduced in Correia (2016).

In Appendix Table A3, I present the same baseline models measuring the effect of Carnegie grants on educational attainment, but I use the continuous measure of exposure (*ExposureCarnegie*) instead of the binary measure of exposure that I use throughout the rest of this paper. In the within-family model, exposure to a Carnegie grant throughout childhood increased educational attainment by 0.20 years. This treatment effect estimate is larger than the estimate using a binary indicator because it reflects the causal effect of moving a child from no exposure to a Carnegie grant to complete exposure (from ages 5–20). But the within-family variation in exposure to a Carnegie grant is well below one because very few families have children born 15 or more years apart.

There are several clear concerns about exogeneity in this model. First, the results may be driven by sharply improving conditions for children in the South or by higher-quality public libraries in the Northeast. The county-by-birth year fixed effects should capture most of this variation, but in Appendix Table A4 I show that the education results are unchanged when I subset the sample to only children in the Midwest, the location that received the most Carnegie grants. My results also do not change when I subset the sample to white children (see Appendix Table A5).

Another concern is that parents may have differentially migrated with their children to towns that received a Carnegie library. And the younger children in families that moved may have had particularly high innate ability or expected educational attainment. In Appendix Table A6, I remove from my sample any child whom I see in the census only after a Carnegie grant was given to the town. For example, if a town received a grant in 1905, I remove all children in that town who I observe in the 1910, 1920, and 1930 decennial censuses from the analysis sample. This reduces my sample size significantly, and the within-family treatment effect is slightly less precisely estimated than in my main specification, but still shows that access to a Carnegie grant by age five causes children to obtain a statistically significant 0.10 additional years of educational attainment.

In Tables 4 and 5, I present linear probability models in the same regression framework from Equation 1, with high school graduation and college attendance as outcomes, instead of years of educational attainment. Access to a Carnegie grant increased high school graduation rates across

specifications. In the within-family specification, access to a Carnegie grant increased high school graduation rates by 1.4 percentage points, from a baseline of 26 percent in the matched sample of siblings. Carnegie grants also increased the probability that a child would attend at least some college by 1.5 percentage points from a baseline of 11 percent in the matched sample of siblings used in the regression model. These are large effects, and they imply a positive effect of Carnegie-funded public libraries on the right tail of educational attainment.

I now present two-stage-least-squares estimates of the same model, but instead of measuring the effect of Carnegie grants on educational attainment, I instrument for the availability of an opened Carnegie library by the age of 5 (as an indicator) with the year of each Carnegie grant. As Table 6 shows, the first stage is highly significant, with a coefficient of 0.43 in the within-family specification and an F-statistic of 352. Access to a Carnegie-funded public library induced by a Carnegie grant leads to 0.24 additional years of educational attainment in the within-family specification. I present only the 2SLS estimates for this main outcome to show the strong first stage. In the rest of the paper, I focus on reduced form results that rely on the precise Carnegie grant dates and the fact that the construction of Carnegie-funded public libraries took an average of 2.4 years.

## **Placebo Tests**

There are three clear concerns with my interpretation of the coefficients in Table 3 as causal effects of library exposure on children. First, cities and towns may have invested in public libraries at the same time as they invested in other local institutions. In dozens of state reports on education systems from the early 1900s, officials discuss the expansion of public libraries and public high schools across their state. Goldin (1994, 1998, 1999) presents state-level data on the dramatic expansion of high schools across the United States in the early 1900s, and it is possible that towns constructed public high schools at the same time that they accepted Carnegie's library construction grants. To address this concern, I construct the first complete panel of public high schools in the

United States. I collect information about each high school's founding date and exact location from 8,000 pages of books and reports that I scan, digitize, and standardize. For more details about the sources, see the Data Appendix. In Figure 4 I plot the first year in which each town and city in the United States opened its first public high school, and as with public libraries, we can see that high schools spread across the Northeast and Midwest much earlier than in the South. In Appendix Table A7, I take my main educational attainment result (Table 3) and I include as a control variable county fixed effects interacted with fixed effects indicating the age of each child when a public high school first entered their town. I calculate these ages at high school entry using the complete panel of public high schools. Controlling flexibly for the entry of each town's first public high school has no effect on my coefficient of interest. Access to a Carnegie grant increased educational attainment by 0.10 years.

A second related concern is that when a town received a Carnegie grant, it could function as a proxy for contemporaneous town-level financial success. When a town can support an institution like a library, it may be a sign of time-varying growth in population, financial independence, or tax revenue that would positively affect children in a way that correlates with, but is not a function of, public library access. To address this concern, I use the library censuses described in the Data section to create a dataset of all college and university libraries in the United States, along with their founding date. I link 1,126 of these libraries to census microdata and I construct an indicator for whether each child had access to a college or university library by the age of 5. In Appendix Table A8, I re-run my main specification from Table 3, looking at the effect of exposure to a college or university library on children's educational attainment as adults. Exposure to a college library by the age of 5 is associated with educational attainment until I control for neighborhood (enumeration district) fixed effects and parental characteristics. In addition, when I add county-by-birth year fixed effects, census reel fixed effects, and household fixed effects to my main specification, we can see that exposure to a college library has a precise null effect on the educational attainment of affected children.

Finally, Andrew Carnegie may have given library construction grants to towns that were regularly investing in local institutions more than neighboring towns. If this were the case, then younger children living in the Carnegie-funded towns might always achieve higher levels of educational attainment than their older brothers, even if Carnegie-funded libraries had no effect on children. To test this, I perturb the Carnegie grant dates. In Appendix Table A9, I assume that Carnegie gave grants to towns 20 years after he actually did. In the baseline specification, these placebo Carnegie grants are still associated with children’s educational attainment. But once I add parental characteristics (column 3), exposure to these placebo grants has no effect on the final educational attainment of children. In Appendix Table A10, I perform the same analysis, but I assume that Carnegie gave grants to towns 20 years before he actually did. In the within-family specification (column 6), access to these lagged Carnegie grants has a precise null effect on children’s educational attainment.

## **Income**

Now that I have established that Andrew Carnegie’s library construction grants increased children’s educational attainment, I analyze the effect of Carnegie grants on income. In Table 7, I measure the effect of Carnegie grants on log annual wage income. I treat anyone who reports zero wage income as having zero log wage income and I use the same set of specifications discussed above.<sup>14</sup> The first column shows that Carnegie gave grants to communities where children would grow up to have more wage income than children from communities which did not receive a grant. But once I control for neighborhood fixed effects (in column 2), exposure to a Carnegie grant is associated with a 2 log point decline in adult wage income. Adding additional controls implies that Carnegie grants had a negative effect on wage income, and in the within-family specification, Carnegie grants had a precise null effect on wage income.

---

<sup>14</sup>While these results use the full sample of 20–65 year old men in the 1940 census, all results are unchanged if I subset to prime-age 25–55 year old men.

The wage effects may be puzzling, because economists often find that human capital interventions increase wage income. But wage income is only one component of total income. Wage income includes all money a worker earned as an employee, but wage income excludes “the earnings of businessmen, farmers, or professional persons derived from business profits, sale of crops, or fees” (Census Bureau, 1940). In 1940, the Census Bureau did not ask respondents to report their non-wage income. However, the census did ask all non-institutionalized respondents to state whether or not they had received at least \$50 of non-wage income in the reference year. 33% of men in my matched sample of siblings reported having at least than \$50 of non-wage income. In Table 8, I show that Carnegie grants had large positive effects on non-wage income. The baseline specification (column 1) shows that Carnegie gave grants to communities where children would grow up to have less non-wage income. But once I add granular geographic controls, we see that Carnegie grants increased the probability that children had at least \$50 of non-wage income by 0.4–0.7 percentage points. And in my within-family specification (column 6), we see that Carnegie grants increased the probability of having at least \$50 of non-wage income by 0.7 percentage point relative to a baseline of 33 percent in this sample (see Table 1).

To summarize the income effects, Andrew Carnegie’s library construction grants had no effect on wage income and increased non-wage income. But because of how the Census Bureau collected income information in the 1940 census, it is impossible to directly measure the effect of Carnegie-funded public libraries on total income.

## **Imputed Income Measures**

In 1950, the Census Bureau asked respondents to state the dollar value of their non-wage income in the decennial census for the first time. I use the non-wage income distribution from 1950 to impute a measure of total income in 1940. Because the Census Bureau did not collect income information before 1940 and collected only some income information in 1940, many economists

impute wage and non-wage income in early census years using the wage distribution in 1950.<sup>15</sup>

I begin my imputation procedure with a sample of 20–65 year old white, black, and Native-American men in the 1950 census who were asked to report wage, non-wage, and total annual income. I then regress each of those income measures on an indicator for whether the respondent earned more than \$120 of non-wage income,<sup>16</sup> indicators for 20 quantiles of wage income, and fixed effects for age, adult state of residence, industry, occupation, years of educational attainment, and annual hours worked. I interact these fixed effects with race. I use the coefficients from this regression in 1950 to predict wage income, non-wage income, and total income for each respondent in the 1940 census. When imputed values are missing because fixed effect cells contained only one individual, I use predicted values from a simpler secondary regression that only includes these fixed effects as regressors: state, race, industry, occupation, age, weekly hours worked, and annual weeks worked.<sup>17</sup>

In Tables 9–11, I measure the effect of Carnegie grants on imputed total income, imputed wage income, and imputed non-wage income. First, we see in Table 9 that access to a Carnegie grant by age five had a positive effect on imputed total income in the within-family specification. This coefficient is a precisely estimated 2.5 log points. In Table 10, we can see that exposure to a Carnegie grant had no effect on wage income once I include county-by-birth year fixed effects as controls. This is consistent with the actual effect of Carnegie exposure on wage income in 1940 (see Table 7), which is unsurprising because my imputation procedure conditions on reported wage income percentiles. Lastly, in Table 11 we see that Carnegie grants had large positive effects on

---

<sup>15</sup>For examples of this type of imputation, see Bailey and Collins (2006), Feigenbaum (2015), and Bayer and Charles (2018). A recent paper by Saavedra and Twinam (2020) discusses many more examples, and proposes a LASSO-based technique to construct an income measure for early censuses using income measures from the 1950 census. Although Saavedra and Twinam’s method is ideal for imputing income measures when no income information exists, in 1940 the Census Bureau collected exact wage income information and an indicator for whether each respondent earned more than \$50 of non-wage income. I therefore use these two variables as inputs into my imputation procedure. This improves accuracy, and is feasible because I can condition on those variables in the 1940 and 1950 census microdata.

<sup>16</sup>\$120 in 1950 is equivalent to \$50 in 1940, using the CPI to inflate 1940 dollars to 1950 dollars.

<sup>17</sup>There are roughly 92,000 observations in these regressions of income on demographic characteristics in the 1950 census, 4,000 indicator variables in the first regression, and 700 indicator variables in the second regression.

imputed non-wage income once I include a basic set of controls. The point estimate ranges from 4–5 log points across models 3-6. In column 6, the within-family estimate implies that access to a Carnegie grant by age five increased non-wage income by 4 log points. To summarize, Carnegie grants had no effect on wage income, increased non-wage income by 4 log points, and either had no effect on total income or increased total income by 2.5 log points (in the within-family model). These imputed income measures are based on the 1950 wage distribution.

## Occupational Prestige

In Table 12, I measure the effect of Carnegie grants on occupational prestige. As my measure of occupational prestige, I use a coding by Siegel (1971) of 412 occupations, which is the only direct measure of occupational standing harmonized by IPUMS to match occupation codes in the census.<sup>18</sup> Siegel constructed these measures using a 1963 study that asked respondents to rate the prestige of different occupations. While this question is amorphous, Siegel (1971) discusses a study of 490 respondents that found a correlation of  $\geq 0.9$  in their ranking of 30 occupations along any of these dimensions: “opportunity for advancement,” “security,” “influence over others,” “responsibility for supervising others,” “work calls for originality and creativity,” “education required,” “training required,” “interesting and challenging work,” and “regarded as desirable to associate with” (Garbin and Bates, 1966). Although Siegel performed his study well after the 1940 census, a large literature shows that occupational prestige remained stable between the 1940s and the 1960s. Prestige rankings from a 1947 study and a 1963 replication had a correlation of 0.99 (Hodge, Siegel, and Rossi 1964).

The main problem with occupational prestige metrics is that the units are not easily inter-

---

<sup>18</sup>Many other measures of occupational standing rely on the average or median income or educational attainment from a set of respondents in an occupation or industry in one census year. These occupational scores can then be used to impute some measure of income or education, based on occupation, in analyses using pre-1940 census data. But my 1940 data contains individual income and educational attainment for respondents, so those measures are duplicative. See Tables 3, 7, and 8 for the direct effect of Carnegie’s grants on reported education, wage income, and non-wage income.

pretable. Siegel put his study participants in front of a numbered nine-step cardboard ladder. He then gave them small paper cards, each with an occupation. Siegel told respondents to put occupation cards on the top position in the ladder if it “has the highest possible social standing,” at the bottom of the ladder if the respondent thought that the occupation had “the lowest possible social standing,” and otherwise, “somewhere in between.” Siegel gave each rung an equally spaced numeric score, ranging from 0 for the lowest rung through 100 for the highest rung of the ladder. The prestige score for an occupation was then the weighted average of individual respondent ratings for each occupation. As Table 1 shows, the occupation-weighted average score in 1940 was around a 34 for my matched sample, and the standard deviation of occupational prestige was 18. But the meaning of the effect size here is difficult to interpret.

I standardize the occupational prestige measure to be mean zero and standard deviation one in my sample. As I show in Table 12, exposure to a Carnegie grant increased occupational prestige across specifications. In my baseline model, Carnegie grants increased occupational prestige by 0.126 standard deviations, but most of this effect is driven by selection on observables. When I add neighborhood fixed effects, the point estimate attenuates to 0.050. And in my within-family specification (column 6), exposure to a Carnegie grant increases occupational prestige by 0.020 standard deviations in the matched sample of siblings.

## **Compulsory Schooling Laws**

I have shown that access to a public library in the early 1900s increased educational attainment, had no effect on wage income, and increased non-wage income. These results mirror the findings of Stephens and Yang (2014), who showed that compulsory schooling laws (CSLs) increased educational attainment but did not affect the wages of adults in the 1960–1980 censuses.

Stephens and Yang wrote their paper in response to a large literature measuring the effect of CSLs on children. For example, Angrist and Krueger (1991) and Oreopoulos and Salvanes (2011)

show that CSLs increase educational attainment and increase wage income. But Stephens and Yang demonstrate that the positive effect of CSLs on wages is driven by differential changes in income by census region, over time. They add region\*birth year fixed effects to the standard regression models linking education to income, instrumenting for education with compulsory schooling laws. The addition of these fixed effects attenuates (to zero) the effect of education on income in their models.

In Table 13, I use the main specification from Stephens and Yang’s paper to measure the effect of educational attainment on respondents’ weekly wage, annual wage income, annual non-wage, and annual total income. I also measure the effect of educational attainment on occupational prestige, and housing values. The specification, taken from their paper, is:

$$Outcome_{st,i} = \alpha Educ_{st,i} + \chi_s + \delta_t + \beta x_{st,i} + \varepsilon_{st,i} \quad (2)$$

where  $i$  indexes individuals,  $s$  indexes state of birth, and  $t$  indexes year of birth. I follow their paper and instrument for years of educational attainment ( $Educ_{st,i}$ ), using three indicators (RS7, RS8 and RS9) for whether each adult in the 1960–1980 census was required to attend school for seven, eight, or nine years. In other words, the first stage is:

$$Educ_{st,i} = \pi_1 RS7_{st} + \pi_2 RS8_{st} + \pi_3 RS9_{st} + \lambda_s + \theta_t + \mu x_{st,i} + v_{st,i} \quad (3)$$

I use Stephens’ and Yang’s dataset and replication programs to produce these results, but I merge on additional outcomes from IPUMS microdata (Ruggles et al., 2020). In Panel A of Table 13, I present results for all adults. In Panel B, I present results for only men. The first column of Table 13 replicates the 2SLS result from Table 1 of Stephens and Yang (2014). Stephens and Yang use weekly wage income as their main outcome, but the weeks worked variable is quite noisy, so I also present results using annual income measures. Similar to my library results, educational attainment had negative (albeit noisy) effect on annual wage income, decreasing the annual wage

income of men by 24 log points. Education also dramatically increased non-wage income by 52 log points, but this outcome is zero around 70% of the time, making the point estimate quite sensitive to how I treat zeroes. In the fourth column, I show that an additional year of education increases the probability a respondent had a significant amount of non-wage income by 5 percentage points, relative to a baseline of 34%. Educational attainment also increased occupational prestige by 0.09 standard deviations. These effects on wage and non-wage income combine to imply that educational attainment had a noisy but positive effect on total income of 7 log points for men and women (and 4 log points for men).

The 1960–1980 censuses also contain housing price values, and in the seventh column of Table 13, I show that an additional year of education, induced by compulsory schooling laws, increases respondents' home value by 9 log points, although this outcome is only measured for home-owners. Column 8 shows that education decreased home-ownership rates by around 2 percentage points. To summarize these results, Carnegie-funded public libraries and compulsory school laws had null or negative effects on wage income, positive effects on non-wage income, and positive effects on occupational prestige.

## **Occupational Choice**

I now show that occupational choices drove this large increase in non-wage income in both the Carnegie library and compulsory schooling law settings. In Table 14, I use a series of linear probability models to measure the effect of a Carnegie grant on occupational choice. The sample is the set of workers who report a valid occupation in 1940.

Census enumerators in 1940 asked respondents to state their occupation. IPUMS aggregated these tens of thousands of occupations into roughly 200 granular categories and nine broad categories. I examine the effect of access to a Carnegie grant on the probability that children chose occupations in one of these nine broad categories as an adult. The broad categories, along with some granular examples from each broad category, are:

1. Professional and Technical: accountants, dentists, professors, nurses, and engineers
2. Farm Owner: farm owners and tenant farmers<sup>19</sup>
3. Managers, Officials, and Proprietors: buyers for stores and building managers
4. Clerical: banktellers, clerks, and telephone operators
5. Sales: advertising agents, insurance brokers, and salesmen
6. Craftsmen: bakers, mechanics, plumbers, and tailors
7. Operatives: electricians, machinists, furnacemen, miners, and taxi drivers
8. Service: hospital attendants, barbers, housekeepers, and waiters
9. Laborers: farm laborers, fisherman, and lumbermen

In Table 14, we see that access to a Carnegie grant increased the probability of becoming a professional worker, a farm owner, a clerical worker, and a salesman by a combined 2.5 percentage points. This is relative to 31% of the matched sample of siblings in those four occupational groups. Access to a Carnegie grant decreased the probability of becoming a manager, a craftsman, an operative, or a laborer by 2.5 percentage points from a baseline of 62 percentage points in the matched sample of siblings. The effect on managers may be surprising, but most managers in 1940 were small-scale proprietors of stores and businesses. While these were not dangerous or low-paid occupations, clerical and technical occupations were better-paid and more prestigious.

---

<sup>19</sup>It is difficult to differentiate farm owners and tenant farmers in the census data, since both responses were often enumerated as “farmers” in the 1940 decennial census. But these two occupations were more similar than readers might expect. For example, Hodge, Siegel, and Rossi (1964) discuss a survey from 1947 which asked people for their opinions about the occupational standing of 90 occupations. The most prestigious occupations were U.S. Supreme Court justice, physician, and nuclear scientist. The least prestigious were garbage collector, street sweeper, and shoe shiner. ‘Farm owner’ ranked as the 43rd most prestigious occupation, and ‘tenant farmer’ ranked as the 51st most prestigious occupation. Farm owners ranked directly below trained machinists and above undertakers. Tenant farmers ranked directly below bookkeeper and above insurance agent. So while farm owners did have higher occupational standing than tenant farmers in the 1940s, they had similar occupational standing and were both considered to be decent middle-class occupations. Also of note is that in 1940, 38% of farm operators were tenant farmers, and not land-owners (Census Bureau, 1950).

In Table 15, I measure the effect of educational attainment on occupational choice, instrumenting for educational attainment using compulsory schooling laws and the main regression model described above (Equations 2–3). As in Table 14, the outcomes are indicators for whether an adult worked in the given broad occupational category. These effects are quite similar to the Carnegie grant effects from Table 14. Increases in educational attainment induced by compulsory schooling laws led to large shifts in occupational choice. An additional year of education increased the probability a person engaged in professional or technical work, or became a farm-owner. An additional year of education also caused fewer children to work as craftsmen or operatives, although the effect on entering professional/technical work and moving away from craftsmanship are not statistically significant when I focus on prime-age men in my sample.

## Life Expectancy

It is difficult to place a direct monetary value on increases in educational attainment, changes in occupational choice, or increases in occupational prestige, but in this section of the paper I use the Social Security Death Master File (SSDMF) to measure the effect of occupational choice on mortality. I then use these occupational safety measures to show that Carnegie grants moved people into safer occupations.

The 2013 SSDMF is a dataset containing the birth date, death date, and social security number for 55 million people who died between 1899 and 2013.<sup>20</sup> The SSDMF is not comprehensive, but contains more than 70% of deaths among 65+ year-olds who died between 1967 and 1972, and more than 95% of deaths among 65+ year-olds who died after 1972 (Hill and Rosenwaike 2001). Coverage is significantly lower for deaths at ages below 65.

I link the SSDMF to my matched census dataset, searching for exact matches using first name, last name, gender, birth year, and birth month when possible. Of the 12.9 million records in my

---

<sup>20</sup>The dataset is no longer publicly available, but was previously purchasable through the Social Security Administration for around \$2,000. A blogger made his copy of the 2013 dataset public for use by researchers. I downloaded the data from <http://ssdmf.info/download.html>

final matched sample, 4.0 million match to a record in the SSDMF with an age at death of 65 or older. I do not have enough statistical power to directly measure the causal effect of Carnegie grants on individual life expectancy, but I use these linked records to measure the effect of entering different occupations and industries on life expectancy. I then show that Carnegie grants and compulsory schooling laws moved children into physically safer occupations and industries.

I begin with a hedonic regression, decomposing respondents' age at death according to this linear model:

$$AgeAtDeath_i = \beta_0 + \alpha_{O(i)}Occupation_i + \delta_{I(i)}Industry_i + \beta X_i + u_i \quad (4)$$

where  $\alpha_O$  is a fixed effect for occupation  $O$  and  $\delta_I$  is a fixed effect for industry  $I$ .  $X_i$  is a matrix of controls including household fixed effects, county\*birth year fixed effects, birth order, and parental characteristics.<sup>21</sup>

I run this regression on a sample of roughly 600,000 people who lived in households as children where two or more siblings match to the SSDMF data, and who matched to records in the 1940 census. I extract these occupation and industry fixed effects (there are around 200 of each), and I use these fixed effects as measures of occupation and industry safety. In this within-family specification, each fixed effect measures the effect of occupation or industry on life expectancy, conditional on living to age 65. This strategy is similar to the literature calculating the value of a statistical life. In that literature, researchers begin with a worker-level dataset and regress wages on demographic characteristics and occupational fatality risk. The coefficient on fatality risk is used to determine the difference in wages that similar workers expect to receive when choosing jobs with higher or lower fatality risk (Gentry and Viscusi, 2016). The household and county\*birth year fixed effects in my regression are particularly granular demographic characteristics which are

---

<sup>21</sup>The parental characteristics are mother and father's age, industry, occupation, and birthplace. These are not perfectly collinear with the household fixed effects because some households contain children with different mothers or fathers. All fixed effects, except for the occupation and industry fixed effects, are interacted with childhood census year and race.

not available in the standard studies using hedonic analysis to estimate the value of a statistical life. In recent work, in the same spirit, Aldy (2019) uses within-couple variation in occupational risk choices to measure the value of a statistical life.

In Table 16, I list the 24 occupations in my 1940 matched sample with more than 100,000 workers. 13.5 million of the 17.8 million men (76%) in my matched sample worked in one of these occupations. I rank these jobs by my measure of occupational safety ( $\alpha_O$ ). The safest occupation is teaching, which increased life expectancy by 1.45 years relative to the group of people with no listed occupation. The riskiest is mining, which decreased life expectancy by 0.99 years relative to teaching. In Table 16, I also present two external measures of occupational fatality risk. First, I present a measure of age-adjusted mortality risk for these occupations using occupations reported on death certificates for men in 1950 (Guralnick, 1963). Second, I present the 2017 occupational fatality rate for the most similar occupation listed in the summary tables of the 2017 Census of Fatal Occupational Injuries (CFOI). The correlation between my measure of each occupation's causal effect on life expectancy and the 1950 measure of excess deaths is -0.61, and this variation is mainly driven by the extremes: both my life expectancy measure and the 1950 excess mortality measure identify teaching as a safe occupation and laborers, mine operatives, and painters as quite unsafe occupations.

For each respondent  $i$  in the 1940 matched sample, I now have an estimate of the value-added life expectancy of the industry and occupation they chose to work in. That measure is the sum of the industry fixed effect and occupation fixed effect ( $\alpha_O(i) + \delta_I(i)$ ) from Equation 4. In Table 17, I regress this measure of occupation and industry life expectancy for each respondent on my measure of exposure to a Carnegie grant, removing any measure of occupational safety in the top or bottom 1% of the distribution so that I can increase precision by avoiding outliers due to noisy estimates of  $\alpha_O(i)$  and  $\delta_I(i)$ . After adding neighborhood controls, access to a Carnegie grant moved people into occupations and industries that were safer by between 0.006 years (two days) and 0.027 years (ten days). Though this effect may seem small, if we interpret the estimates from Equation 4 as causal

effects of occupational choice on life expectancy, an additional two days of life, summed over the millions of children who received access to a Carnegie grant at some point, implies a significant increase in total expected life years.

In the third-to-last column of Table 13, I use the main specification from Stephens and Yang (2014) to measure the effect of a year of educational attainment on occupational safety, instrumenting for educational attainment with CSLs. In a sample of both men and women, education had a positive effect on occupational safety, shifting people into occupations which increased life expectancy by 0.025–0.033 years (9–12 days).

## **Comparison of Libraries and Compulsory Schooling Laws**

I can now quantitatively compare the effect of Carnegie library access and the effect of compulsory schooling laws on children. I do this by comparing two sets of coefficients: the effect of education on outcomes measured in the 1960–1980 census, instrumenting for education with CSLs; and the effect of Carnegie library exposure on those same outcomes for adults in the 1940 census. The two sets of coefficients are quite similar. In Figure 5, I plot the treatment effects from the Carnegie library regressions discussed above against the treatment effects from my analysis of the effect of compulsory schooling laws on men.

Twelve outcomes from these regression specifications are directly comparable: the effect of each intervention on wage income, the probability of having  $\geq$  \$50 of non-wage income, occupational prestige, and the probability of working in each of nine broad occupational categories. These outcomes are not all independent, but the strong correlation is still stark. For completely different cohorts, in disjoint datasets, subject to different interventions—public library access and compulsory schooling laws—we see remarkably similar effects on children. The raw correlation between these T-statistics is 0.52, with a rank-rank correlation of 0.81. Both Carnegie library grants and compulsory schooling laws increased educational attainment and increased non-wage income by

giving children the opportunity to enter safer, more entrepreneurial, and more prestigious occupations.

## Conclusion

Many eminent politicians, scientists, and writers attribute their professional success to public libraries. Supreme Court Justice Sonia Sotomayor stumbled across her local public library while out shopping with her mother. Though her English language skills were poor at the time, the librarian told her about library cards and Sotomayor proceeded to use the library to develop a love of reading and excel in school.<sup>22</sup> And Justice Clarence Thomas similarly attributed his love of books and learning to his time as a child in a Carnegie-funded Library.<sup>23</sup> But economists have largely ignored the institution.

In this paper, I link the rollout of Andrew Carnegie's public library construction grants to complete count census data. I show that access to a public library increased educational attainment, had no effect on wage income, and increased non-wage income. Children exposed to Carnegie grants shifted into safer, more entrepreneurial, and more prestigious occupations. The same patterns explain the puzzling results of Stephens and Yang (2014), who showed that compulsory schooling laws increased educational attainment and had no effect on children's wage income. Like Carnegie's library construction grants, compulsory schooling laws increased educational at-

---

<sup>22</sup>Sotomayor, in a 2018 interview with David Axelrod, said: "I needed to escape from home and I fortuitously found the local library... The local library was in a shopping center [near to where I lived]... It was on the second floor of a building that housed the Macy's... I went to that Macy's [with my mom]... I saw the library... [My mom] walked in with me and I asked what it was. The [librarian] came to us and told us you could borrow books. My mom asked how and [the librarian] said: get a library card... I started to read... it became my escape... [I only excelled in school] after I found reading... I was a marginal student in school... I was having difficulty understanding English... It took a very long time for me to understand the different meanings and usages of words that sound the same..."

<sup>23</sup>"U.S. Supreme Court Justice Clarence Thomas and Pulitzer Prize winning author James Alan McPherson both spent their time engrossed in the stacks of books [at the Carnegie Library]... When Thomas was growing up, he spent most of his free time in the Carnegie Library, on the black side of town. It wasn't until he was a teenager that integration gave him access to the Big Library, as he called it, but once access was granted Thomas took advantage of it... When he was a kid, Thomas told them, the library was how he expanded his world, using books to visit places that were beyond his reach" (The Washington Post, "Supreme Discomfort," by Kevin Merida and Michael A. Fletcher. August 4, 2002).

tainment and increased non-wage income by shifting people into safer, more entrepreneurial, and more prestigious occupations. Economists often use wage income to measure the returns to education and the causal effects of local institutions on children, but wage income is only one component of the returns to investment in human capital: non-wage income, occupational prestige, and occupational safety are other important outcomes to consider and in this paper I show that both Carnegie-funded public libraries and compulsory schooling laws had positive effects on adult outcomes that operated through these measures.

## References

- [Aaronson and Mazumder, 2011] Daniel Aaronson and Bhashkar Mazumder. “The Impact of Rosenwald Schools on Black Achievement.” *Journal of Political Economy* Vol. 119, No. 5, pp. 821-888, 2011.
- [Abramitzky, Boustan, and Eriksson, 2012] Ran Abramitzky, Leah Platt Boustan, and Katherine Eriksson. “Europe’s Tired, Poor, and Huddled Masses: Self-Selection and Economic Outcomes in the Age of Mass Migration.” *American Economic Review* Vol. 102, No. 5, pp. 1832-56, 2012.
- [Abramitzky, Boustan, and Eriksson, 2014] Ran Abramitzky, Leah Platt Boustan, and Katherine Eriksson. “A Nation of Immigrants: Assimilation and Economic Outcomes in the Age of Mass Migration.” *Journal of Political Economy* Vol. 122, No. 3, pp. 467-506, 2014.
- [Abramitzky et al., 2020] Ran Abramitzky, Leah Platt Boustan, Katherine Eriksson, James J. Feigenbaum, and Santiago Pérez. “Automated Linking of Historical Data.” *Journal of Economic Literature* , forthcoming.
- [Aldy, 2019] Joseph E. Aldy. “Birds of a Feather: Estimating the Value of Statistical Life from Dual-Earner Families” *Journal of Risk and Uncertainty*, Vo. 58, pp. 187-205, 2019.
- [Angrist and Krueger, 1991] Joshua D. Angrist and Alan B. Krueger. “Does Compulsory School Attendance Affect Schooling and Earnings?.” *Quarterly Journal of Economics* Vol. 106, No. 4, pp. 979-1014, 1991.
- [Bailey and Collins, 2006] Martha J. Bailey and William J. Collins. “The Wage Gains of African-American Women in the 1940s.” *Journal of Economic History* Vol. 66, No. 3, pp. 737-777, 2006.

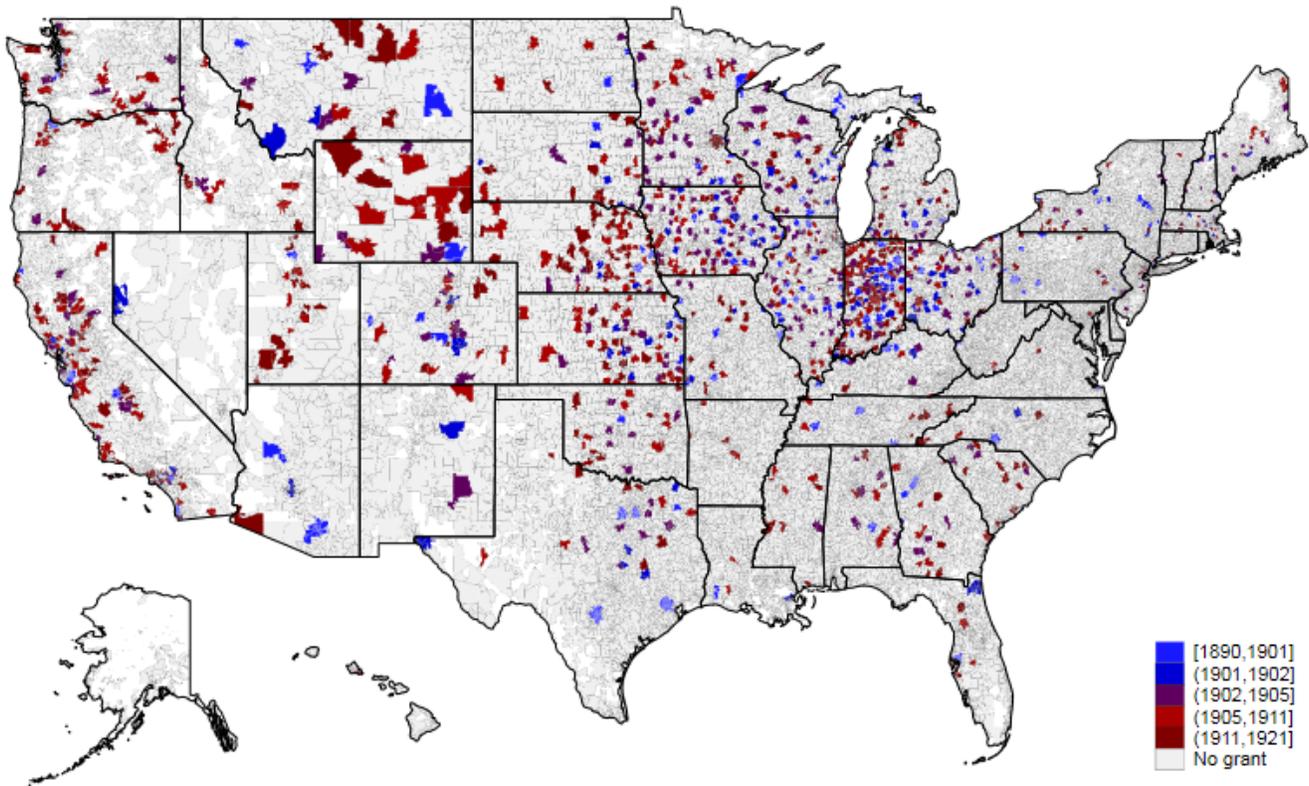
- [Bailey et al., 2011] Martha Bailey, Brian Jacob, Michael Kevane, and William Sundstrom. “Carnegie’s Legacy and the Growth of American Cities: Did Public Libraries Have Measurable Effects?” *Unpublished Data*, 2011.
- [Bailey et al., 2019] Martha Bailey, Connor Cole, Morgan Henderson, and Catherine Massey. “How Well do Automated Linking Methods Perform? Lessons from U.S. Historical Data” Vol. 58, No. 4, pp. 997-1044, 2020.
- [Bayer and Charles, 2018] Patrick Bayer and Kerwin Kofi Charles. “Divergent Paths: A New Perspective on Earnings Differences between Black and White Men since 1940.” *Quarterly Journal of Economics* Vol. 133, No. 3, pp. 1459-501, 2018.
- [Berkes and Nencka, 2019] Enrico Berkes and Peter Nencka. “Knowledge access: The effects of Carnegie libraries on innovation.” *Working Paper*, 2019.
- [Bhatt, 2010] Rachana Bhatt. “The Impact of Public Library Use on Reading, Television, and Academic Outcomes.” *Journal of Urban Economics* Vol. 68, No. 2, pp. 148-166, 2010.
- [Bobinski, 1969] George Sylvan Bobinski. “Carnegie Libraries: Their History and Impact on American Public Library Development.” American Library Association, 1969.
- [Carnegie, 1901] Andrew Carnegie. “The gospel of wealth and other timely essays.” Century, 1901.
- [Census Bureau, 1940] Census Bureau. “Instructions to Enumerators, Population and Agriculture, 1940.” *1940 Census Publications* Form PA-1.
- [Census Bureau, 1950] Census Bureau. “Agriculture 1950 - A Graphic Summary.” *1950 Census Publications* Vol. 5, No. 6, pp. 69-102, 1950.
- [Chetty and Hendren, 2018] Raj Chetty and Nathaniel Hendren. “The impacts of neighborhoods on intergenerational mobility I: Childhood exposure effects.” *Quarterly Journal of Economics* Vol. 133, No. 3, pp. 1107-62, 2018.

- [Correia, 2016] Sergio Correia. “A Feasible Estimator for Linear Models with Multi-Way Fixed Effects.” *Working Paper*, 2016.
- [Daniel, 1961] Hawthorne Daniel. “Public Libraries for Everyone: The Growth and Development of Library Services in the United States, Especially Since the Passage of the Library Services Act.” Doubleday, 1961.
- [Duflo, 2001] Esther Duflo. “Schooling and Labor Market Consequences of School Construction in Indonesia: Evidence from an Unusual Policy.” *American Economic Review* Vol. 91, No. 4, pp. 795-813, 2001.
- [Feigenbaum, 2015] James Feigenbaum. “Intergenerational Mobility during the Great Depression” *Working Paper*, 2015.
- [Ferrie, 1996] Joseph P. Ferrie. “A New Sample of Males Linked from the 1850 Public Use Micro Sample of the Federal Census of Population to the 1860 Federal Census Manuscript Schedules.” *Historical Methods* Vol. 29, No. 4, pp. 141-156, 1996.
- [Flexner, 1927] Jennie M. Flexner. “Circulation Work in Public Libraries.” American Library Association, 1927.
- [Garbin and Bates, 1966] A.P. Garbin and Frederick L. Bates. “Occupational Prestige and Its Correlates: A Re-examination.” *Social Forces* Vol. 44, No. 3, pp. 295-302, 1966.
- [Gilpin, Karger, and Nencka, 2020] Gregory Gilpin, Ezra Karger, and Peter Nencka. “The Returns to Public Library Investment.” *Working Paper*, 2020.
- [Goldin, 1994] Claudia Goldin. “How America Graduated from High School: 1910 to 1960.” *NBER Working Paper* June, 1994.
- [Goldin, 1998] Claudia Goldin. “America’s Graduation from High School: The Evolution and Spread of Secondary Schooling in the Twentieth Century.” *Journal of Economic History* Vol. 58, No. 2, pp. 345-374, 1998.

- [Goldin, 1999] Claudia Goldin. “Egalitarianism and the Returns to Education during the Great Transformation of American Education.” *Journal of Political Economy* Vol. 107, No. S6, pp. S65-S94, 1999.
- [Guralnick, 1963] Lillian Guralnick. “Mortality by Occupation and Cause of Death among Men 20 to 64 Years of Age: United States, 1950.” *Vital Statistics Special Reports* Vol. 53, No. 3, September 1963
- [Hill and Rosenwaive, 2001] Mark E. Hill and Ira Rosenwaive. “The Social Security Administration’s Death Master File: the completeness of death reporting at older ages.” *Social Security Bulletin* Vol. 64, No. 1, pp. 45-51, 2001.
- [Hodge, Siegel, and Rossi, 1964] Robert W. Hodge, Paul M. Siegel and Peter H. Rossi. “Occupational Prestige in the United States, 1925-63.” *American Journal of Sociology* Vol. 70, No. 3, pp. 286-302, 1964.
- [Kevane and Sundstrom, 2014] Michael J. Kevane and William A. Sundstrom. “The Development of Public Libraries in the United States, 1870–1930: A Quantitative Assessment.” *Information & Culture* Vol. 49, No. 2, pp. 117-144, 2014.
- [Kevane and Sundstrom, 2016] Michael J. Kevane and William A. Sundstrom. “Public Libraries and Political Participation: 1870-1940.” *Working Paper*, 2016.
- [Klinenberg, 2018] Eric Klinenberg. “Palaces for the people: How social infrastructure can help fight inequality, polarization, and the decline of civic life.” Broadway Books, 2018.
- [Logan and Parman, 2017] Trevon D. Logan and John M. Parman. “The National Rise in Residential Segregation.” *Journal of Economic History* Vol. 77, No. 1, pp. 127-170, 2017.
- [Magnuson and King, 1995] Diana L. Magnuson and Miriam L. King. “Comparability of the Public Use Microdata Samples: Enumeration Procedures.” *Historical Methods* Vol. 28, No. 1, pp. 27-32, 1995.

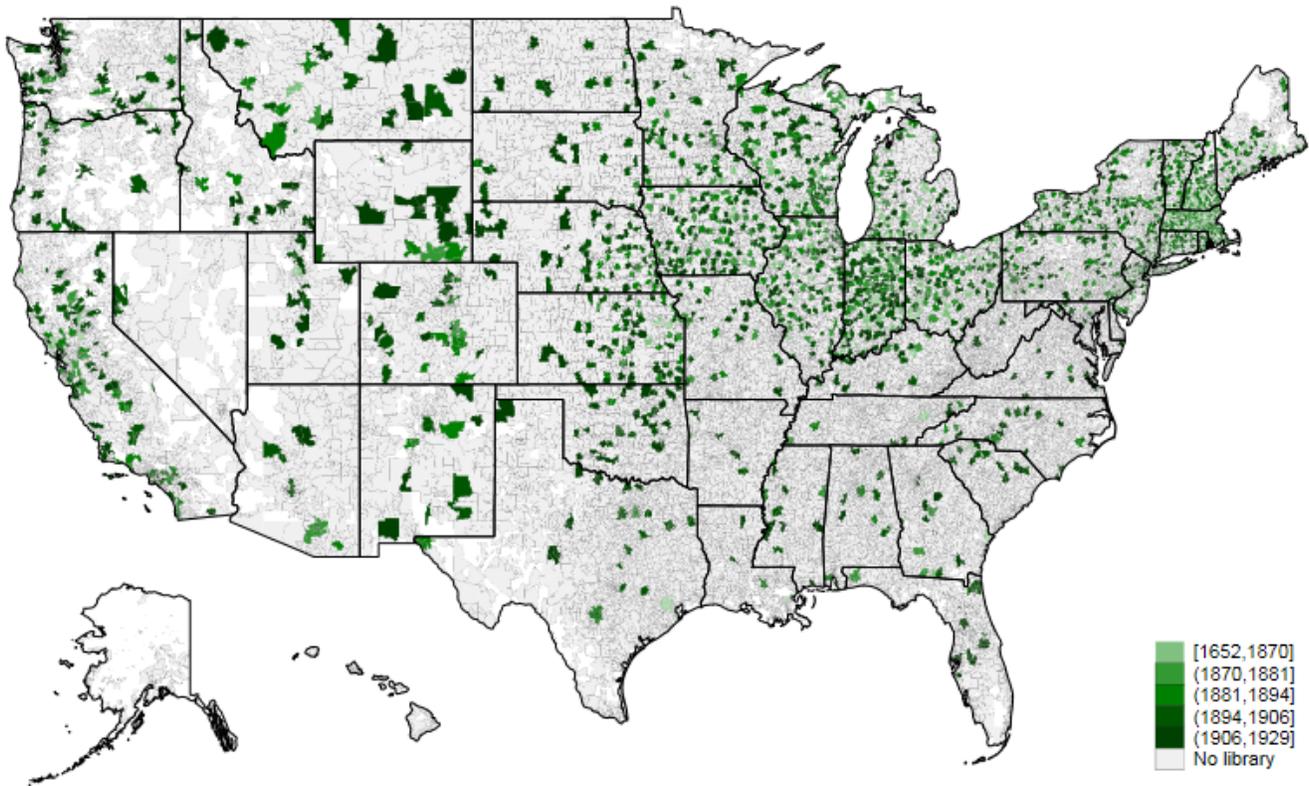
- [Martin, 1993] Robert Sidney Martin. “Carnegie Denied: Communities Rejecting Carnegie Library Construction Grants, 1898-1925.” Greenwood Press, 1993.
- [Neilson and Zimmerman, 2014] Christopher A. Neilson and Seth D. Zimmerman. “The Effect of School Construction on Test Scores, School Enrollment, and Home Prices.” *Journal of Public Economics* Vol. 120, pp. 18-31, 2014.
- [Neto, 2019] Amir B Ferreira Neto. “Do Public Libraries Impact Local Labor Markets? Evidence from Appalachia.” *Working Paper*, 2019.
- [Oreopoulos and Salvanes, 2011] Philip Oreopoulos and Kjell G. Salvanes. “Priceless: The Non-pecuniary Benefits of Schooling.” *Journal of Economic Perspectives* Vol. 25, No. 1, pp. 159-184, 2011.
- [Pisati, 2018] Maurizio Pisati. “SPMAP: Stata module to visualize spatial data.” *Statistical Software Package*, 2018.
- [Ruggles et al., 2020] Steven Ruggles, Sarah Flood, Ronald Goeken, Josiah Grover, Erin Meyer, Jose Pacas, and Matthew Sobek. ”IPUMS USA: Version 10.0 [dataset].” Minneapolis, MN: IPUMS, 2020. <https://doi.org/10.18128/D010.V10.0>
- [Saavedra and Twinam, 2020] Martin Saavedra and Tate Twinam. “A Machine Learning Approach to Improving Occupational Income Scores.” *Explorations in Economic History* Vol. 75, 2020.
- [Siegel, 1971] Paul Mathew Siegel. “Prestige in the American Occupational Structure.” *Ph.D. dissertation*, University of Chicago, 2018.
- [Stephens and Yang, 2014] Melvin Stephens Jr and Dou-Yan Yang. “Compulsory Education and the Benefits of Schooling.” *American Economic Review* Vol. 104, No. 6, pp. 1777-92, 2011.
- [Todd, 2008] Petra Todd. “Evaluating Social Programs with Endogenous Program Placement and Selection of the Treated.” *Handbook of Development Economics* Vol. 4, pp. 3848-94, 2008.

Figure 1:  
Carnegie Library Grants in United States  
by year Carnegie agreed to grant



Note: I match Carnegie grants to city names in Census microdata. I then use a modern-day mapping of cities to zip codes to identify the zip codes associated with each Carnegie grant. Color indicates the year that Carnegie agreed to the grant. There is no complete list of Carnegie library opening dates, but anecdotal evidence suggests that the vast majority of communities receiving a grant constructed their library within four years of the grant announcement. Map constructed using the spmap package (Pisati 2018).

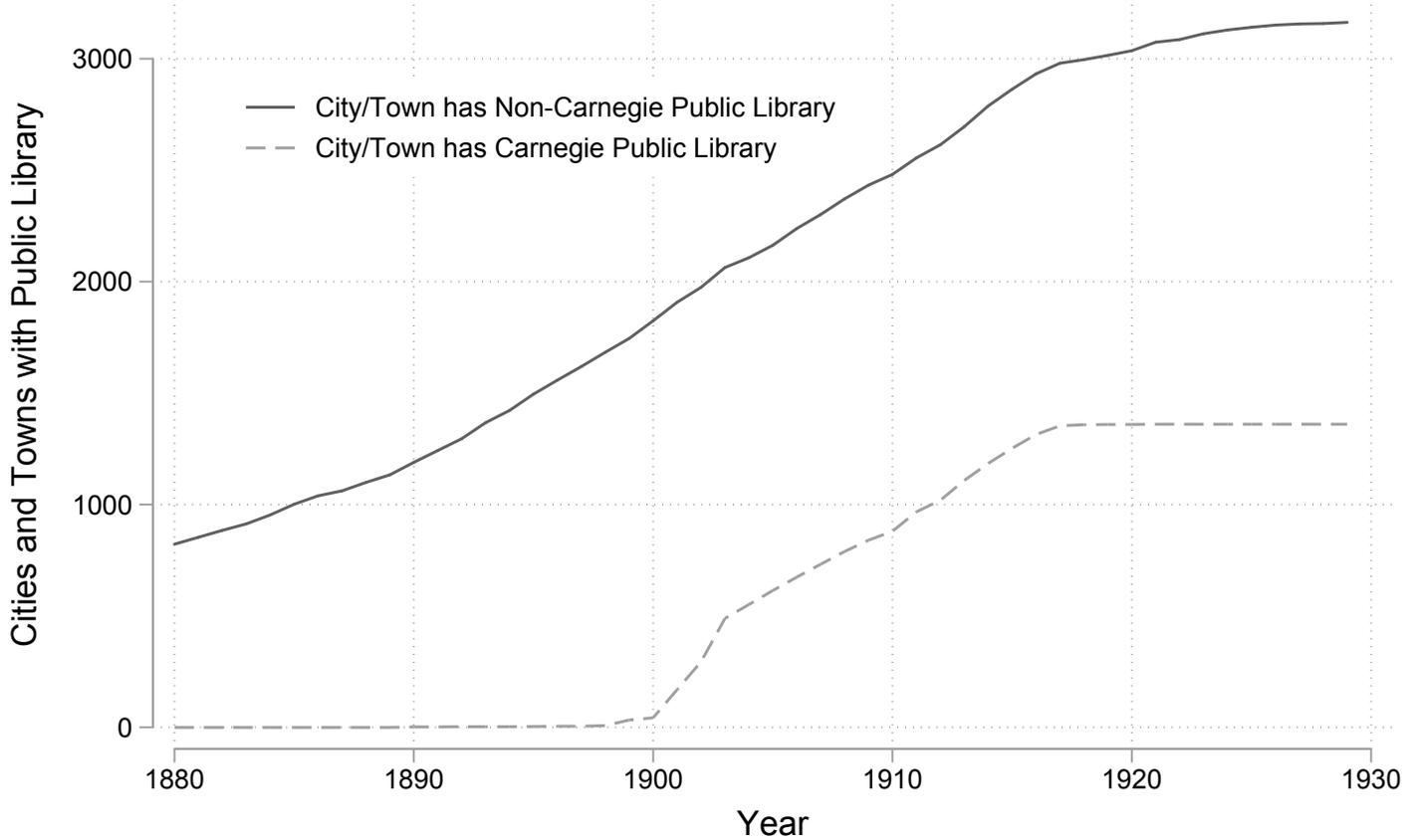
Figure 2:  
Public Libraries in United States as of 1929  
by year of founding



Note: I match public libraries to city names in Census microdata. I then use a modern-day mapping of cities to zip codes to identify the zip codes associated with each public library. Color indicates the year that the public library was founded. Data come from the public library censuses produced by the Bureau of Education between 1875 and 1929. See Data section for more details. Map constructed using the `spmap` package (Pisati 2018).

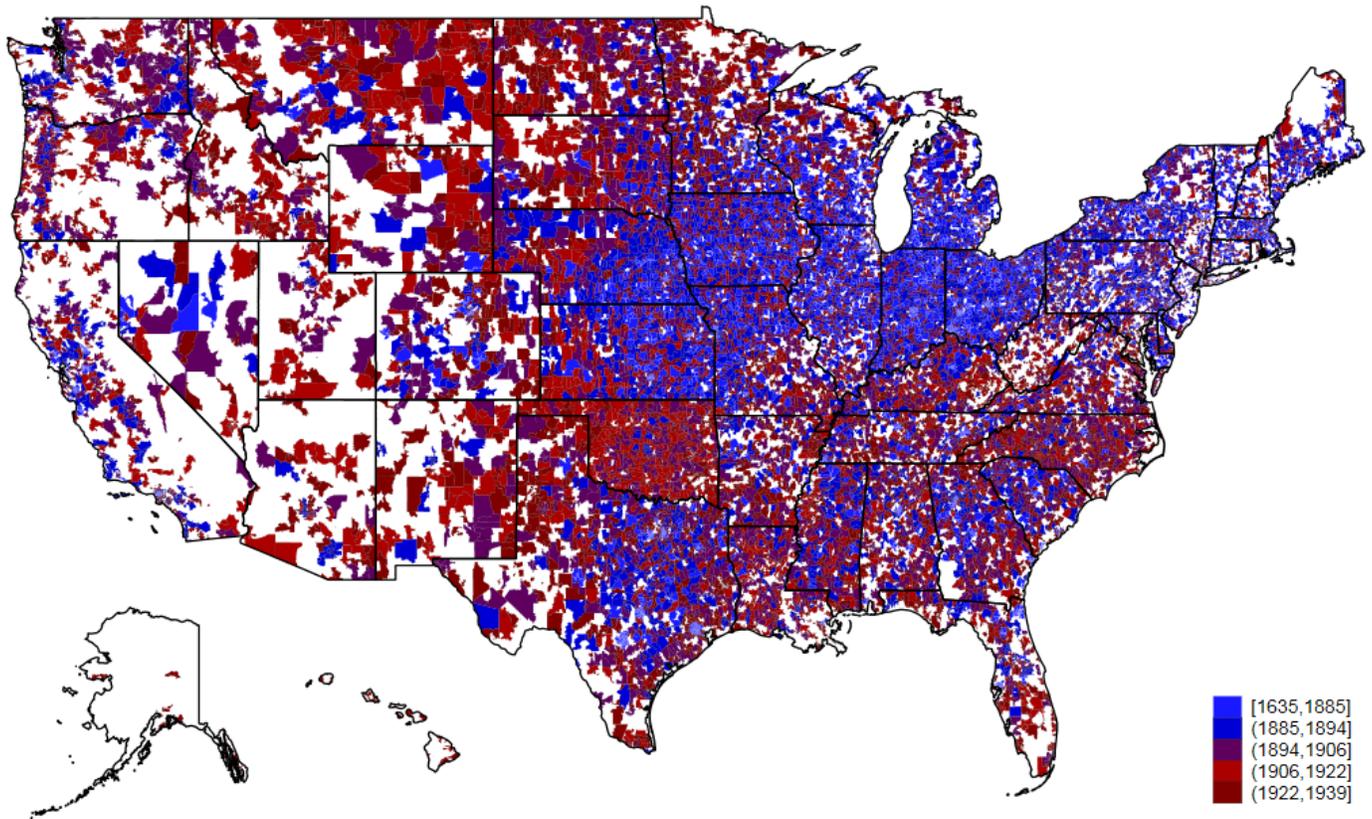
Figure 3:

### Public Libraries in United States, 1880-1929



Note: I match public libraries and Carnegie grants to city names in Census microdata. I then use a modern-day mapping of cities to zip codes to identify the zip codes associated with each public library and each Carnegie grant. This figure shows the number of cities and towns with a non-Carnegie grant as of each year and the number of cities and towns with a Carnegie grant promised as of each year. Carnegie grant data was originally collected by Bobinski (1969) and public library data come from the public library censuses produced by the Bureau of Education between 1875 and 1929. See Data section for more details.

Figure 4:  
High School Entry by 1944



Note: I match all public high schools from the public high school censuses described in the Data section to city names in Census microdata. I then use a modern-day mapping of cities to zip codes to identify the zip codes associated with each public high school. Color indicates the year that the first public high school was founded in each town. Map constructed using the `smap` package (Pisati 2018).

Figure 5:

Comparison of Treatment Effects from Carnegie and CSL Regressions



Notes:

1. Treatment effects from Compulsory Schooling Law (CSL) regressions are drawn from the models presented in Tables 13 and 15.
2. Treatment effects from the library regressions are from the within-family models in Tables 3, 7–9, and 14.
3. Data for the CSL regressions are the publicly available 1960–1980 IPUMS long-form census samples.
4. Data for the library regressions are the restricted 1900–1940 complete count census data, available through NBER.
5. Standard errors for the CSL regressions are clustered at the state of birth-by-year of birth level using the methods described in Stephens and Yang (2014).
6. Standard errors for the library regressions are clustered at the childhood county of residence level.
7. I suppress the wage income effect to ensure that the other coefficients can be easily viewed on the same figure.

Table 1: Comparison of All 1940 Men and Matched Sample

Variable	All 1940 Records	Matched Samples		All Matched, by Age of Carnegie Access			
		All	Siblings	By Age 5	Age 6-15	Age>15	Never
N	39,175,147	17,818,121	5,844,571	3,171,477	912,487	730,274	13,003,883
Age in Childhood Census	-	10.37	9.65	9.01	12.78	17.93	10.12
Age in 1940	39.20	38.65	39.38	31.75	45.51	54.45	38.99
Fraction White	0.91	0.91	0.93	0.95	0.95	0.95	0.90
Fraction Black	0.09	0.09	0.07	0.05	0.05	0.05	0.10
Born in U.S.	0.87	0.94	0.94	0.93	0.84	0.84	0.95
Years of Education	8.66	8.95	8.79	10.16	9.01	8.57	8.68
High School Graduation Rate	0.26	0.28	0.25	0.40	0.27	0.23	0.25
College Attendance Rate	0.11	0.12	0.10	0.17	0.13	0.11	0.10
Prestige	34.43	34.91	34.85	35.37	37.81	37.57	34.45
Weeks Worked	38.20	39.16	39.56	38.92	40.59	37.26	39.22
Hours Worked	33.78	34.73	35.22	34.24	35.13	31.94	34.97
Exposure to Carnegie Grant	-	0.18	0.14	1	0	0	0
Exposure to Public Library	-	0.38	0.34	0.93	0.77	0.61	0.21
Occupation: Professional/Technical	0.06	0.06	0.05	0.08	0.08	0.07	0.05
Occupation: Farm owner	0.13	0.14	0.15	0.04	0.06	0.10	0.16
Occupation: Manager	0.09	0.09	0.09	0.09	0.15	0.15	0.08
Occupation: Clerical	0.06	0.07	0.06	0.11	0.08	0.07	0.06
Occupation: Sales	0.06	0.07	0.06	0.11	0.08	0.07	0.06
Occupation: Craftsman	0.16	0.16	0.16	0.17	0.21	0.20	0.15
Occupation: Operative	0.18	0.18	0.18	0.21	0.16	0.13	0.17
Occupation: Service	0.06	0.05	0.05	0.06	0.06	0.07	0.05
Occupation: Laborer	0.19	0.19	0.19	0.15	0.11	0.13	0.20
Wage	851	867	854	1,016	1,221	1,012	797
Wage (median)	600	600	600	900	1,000	644	520
Probability $\geq$ \$50 Non-Wage Income	0.31	0.31	0.33	0.21	0.33	0.40	0.33

Note: See Data Appendix for definition of each variable and description of processing.

Table 2: Fraction of Matched Sample in Largest Cities and Towns

Place Name	State	Count	% of Matched Sample	Cumulative Percentage	Carnegie Grant	Public Library
Chicago	IL	395,768	2.2%	2.3%	1902	1872
Manhattan	NY	303,022	1.7%	3.9%	1899	1652
Philadelphia	PA	237,505	1.3%	5.3%	1903	1820
Brooklyn	NY	227,985	1.3%	6.5%	1899	1823
Detroit	MI	134,446	0.8%	7.3%	1901	1865
St. Louis	MO	121,565	0.7%	8.0%	1901	1860
Cleveland	OH	115,359	0.6%	8.6%	1903	1868
Baltimore	MD	108,224	0.6%	9.2%	1906	1874
Boston	MA	107,998	0.6%	9.8%		1848
Milwaukee	WI	82,650	0.5%	10.2%		1847
Buffalo	NY	79,618	0.4%	10.7%		1836
Pittsburgh	PA	73,498	0.4%	11.1%	1890	1872
Los Angeles	CA	70,593	0.4%	11.6%	1911	1872
Newark	NJ	66,317	0.4%	11.9%		1888
San Francisco	CA	63,506	0.5%	12.3%	1901	1866
Cincinnati	OH	63,474	0.4%	12.6%	1902	1802
New Orleans	LA	60,568	0.3%	13.0%	1902	1843
Bronx	NY	57,004	0.3%	13.3%	1899	1652
Minneapolis	MN	53,921	0.3%	13.6%	1912	1885
Jersey City	NJ	53,201	0.3%	13.9%		1889
Washington	DC	49,314	0.3%	14.2%	1899	1865
Providence	RI	44,884	0.3%	14.4%		1874
Rochester	NY	41,978	0.2%	14.7%		1912
Indianapolis	IN	39,280	0.2%	14.9%	1909	1831
Louisville	KY	39,104	0.2%	15.1%	1899	1871

Note: See Data Appendix for definition of each variable and description of processing.

Table 3: Effect of Carnegie Grant on Years of Schooling

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.572 (0.034)	0.179 (0.019)	0.076 (0.018)	0.090 (0.009)	0.134 (0.010)	0.101 (0.017)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,412,080	17,395,212	17,393,991	17,265,943	14,867,729	5,375,550
R <sup>2</sup>	0.20	0.26	0.29	0.32	0.55	0.73
Mean dep. var.	8.99	8.99	8.99	9.01	9.00	8.88

## Notes:

1. Data are the restricted 1900–1940 complete count census data, available through NBER.
2. Standard errors clustered at the childhood county of residence level.
3. Child controls are birthplace, birth order, and birth year. All are interacted with race (white, black, and Native American)-by-childhood census year (indicators for whether I see each 1940 adult in the 1900, 1910, 1920, or 1930 census).
4. Enumeration district fixed effects are enumeration district indicators (roughly 65,000 in each census year) interacted with race-by-childhood census year. Enumeration districts are occasionally not unique within a county. In these rare cases, the fixed effect combines individuals in potentially geographically discontinuous neighborhoods.
5. Parent controls are mother's birthplace, father's birthplace, mother's age, father's age, mother's occupation, father's occupation, mother's industry, and father's industry. All are interacted with race-by-childhood census year fixed effects. Missing values are coded as their own category.
6. County-by-birth year fixed effects include county fixed effects interacted with birth year fixed effects, interacted with race-by-childhood census year fixed effects.
7. Census microfilm page numbers are numbered in the data, but occasionally repeat within microfilm reel. In those cases, my Census microfilm page fixed effects will group together
8. In the model including household fixed effects, I exclude roughly 1,500 children who lived in a household with more than 8 matched siblings.
9. I use the reghdfe Stata package to iteratively drop fixed effects that uniquely separate one observation from all other observations. This is why the sample size decreases from model to model.
10. There are 17.8 million observations in my analysis sample, but only 17.5 million have valid educational attainment information, so the model in column 1 contains 17.5 million records (after dropping singletons).

Table 4: Effect of Carnegie Grant on High School Graduation Rate

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.064 (0.004)	0.029 (0.003)	0.017 (0.003)	0.015 (0.001)	0.020 (0.001)	0.014 (0.003)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,412,080	17,395,212	17,393,991	17,265,943	14,867,729	5,375,550
R <sup>2</sup>	0.10	0.15	0.19	0.21	0.47	0.68
Mean dep. var.	0.28	0.28	0.28	0.29	0.28	0.26

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 5: Effect of Carnegie Grant on College Attendance

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.043 (0.002)	0.024 (0.002)	0.017 (0.002)	0.018 (0.001)	0.021 (0.001)	0.015 (0.002)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,412,080	17,395,212	17,393,991	17,265,943	14,867,729	5,375,550
R <sup>2</sup>	0.03	0.08	0.12	0.14	0.42	0.64
Mean dep. var.	0.12	0.12	0.12	0.12	0.12	0.11

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 6: IV Effect of Carnegie Grant on Years of Schooling

	1	2	3	4	5	6
1(Carnegie Library by Age 5)	0.75 (0.05)	0.33 (0.05)	0.14 (0.04)	0.17 (0.02)	0.27 (0.02)	0.24 (0.04)
First stage:						
Coefficient	0.77	0.54	0.53	0.53	0.51	0.43
F-statistic	756	272	266	508	462	352
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
N	17,412,080	17,395,212	17,393,991	17,265,943	14,867,729	5,375,550

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 7: Effect of Carnegie Grant on Log Annual Wage Income

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.381 (0.032)	-0.022 (0.013)	-0.058 (0.011)	-0.045 (0.008)	-0.044 (0.009)	-0.003 (0.018)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	16,766,377	16,749,418	16,748,204	16,618,602	14,182,462	5,027,116
R <sup>2</sup>	0.07	0.11	0.12	0.15	0.41	0.61
Mean dep. var.	4.91	4.91	4.91	4.91	4.90	4.87

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 8: Effect of Carnegie Grant on Probability  $\geq$  50 Dollars of Non-Wage Income

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	-0.049 (0.004)	-0.001 (0.002)	0.004 (0.001)	0.004 (0.001)	0.005 (0.001)	0.007 (0.003)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,240,954	17,223,995	17,222,752	17,094,494	14,687,385	5,292,713
R <sup>2</sup>	0.07	0.10	0.11	0.14	0.40	0.60
Mean dep. var.	0.31	0.31	0.31	0.31	0.32	0.33

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 9: Effect of Carnegie Grant on Log Imputed Total Income from 1950 Data

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.166 (0.013)	0.029 (0.004)	0.008 (0.004)	-0.004 (0.003)	-0.002 (0.004)	0.025 (0.011)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,342,421	17,325,555	17,324,336	17,196,140	14,793,112	5,339,327
R <sup>2</sup>	0.09	0.12	0.13	0.16	0.42	0.60
Mean dep. var.	7.21	7.21	7.21	7.21	7.24	7.25

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 10: Effect of Carnegie Grant on Log Imputed Wage Income from 1950 Data

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.425 (0.035)	0.006 (0.015)	-0.038 (0.012)	-0.037 (0.007)	-0.033 (0.008)	0.016 (0.018)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,342,421	17,325,555	17,324,336	17,196,140	14,793,112	5,339,327
R <sup>2</sup>	0.07	0.11	0.11	0.15	0.41	0.61
Mean dep. var.	5.85	5.85	5.85	5.86	5.85	5.81

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 11: Effect of Carnegie Grant on Log Imputed Non-Wage Income from 1950 Data

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	-0.296 (0.025)	0.023 (0.013)	0.048 (0.011)	0.041 (0.007)	0.040 (0.008)	0.039 (0.019)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,342,421	17,325,555	17,324,336	17,196,140	14,793,112	5,339,327
R <sup>2</sup>	0.07	0.10	0.11	0.14	0.40	0.60
Mean dep. var.	3.21	3.21	3.21	3.21	3.26	3.32

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Table 12: Effect of Carnegie Grant on Occupational Prestige

	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>
1(Carnegie Grant by Age 5)	0.126 (0.007)	0.050 (0.003)	0.030 (0.003)	0.021 (0.003)	0.027 (0.003)	0.020 (0.007)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	16,127,337	16,110,400	16,109,139	15,978,774	13,517,806	4,746,742
R <sup>2</sup>	0.08	0.12	0.14	0.17	0.43	0.63
Mean dep. var.	0.00	0.00	0.00	0.00	0.02	0.01

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. Occupational prestige is standardized so that the average in the full matched sample is zero and the standard deviation is 1.

Table 13: Effect of Education on Various Outcomes (Instrumented with Compulsory Schooling)

	<b>Weekly Wage</b>	<b>Annual Wage</b>	<b>Annual Non-Wage</b>	<b>Annual Non-Wage <math>\geq</math> \$275</b>	<b>Annual Total</b>	<b>Occ. Prestige</b>	<b>Occ. Safety</b>	<b>Housing Price</b>	<b>Owns House</b>
<b>Panel A: Men and Women</b>									
Years of Education	-0.003 (0.016)	-0.09 (0.07)	0.33 (0.06)	0.03 (0.01)	0.07 (0.06)	0.09 (0.02)	0.025 (0.011)	0.089 (0.019)	-0.018 (0.010)
First stage F-statistic:	39	.	.	.	.	.	.	.	.
N	3,680,223	5,066,060	.	.	.	.	3,021,522	5,015,842	.
Average in Sample	5.43	6.68	2.47	0.27	7.60	0.00	0.14	10.70	0.72
<b>Panel B: Men</b>									
Years of Education	-0.014 (0.021)	-0.24 (0.10)	0.52 (0.12)	0.05 (0.01)	0.04 (0.05)	0.08 (0.03)	0.033 (0.015)	0.099 (0.025)	-0.024 (0.013)
First stage F-statistic:	21	.	.	.	.	.	.	.	.
N	2,166,387	2,502,089	.	.	.	.	.	1,467,960	2,466,891
Average in Sample	5.80	8.26	3.20	0.34	9.34	0.00	0.06	10.71	0.72
State of Birth (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Year of Birth (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Region*Birth Year (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Age <sup>4</sup> , census year, gender	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

1. See notes to Table 3.
2. Column 1 replicates the 2SLS estimates from Columns 4 and 6 of Table 1 in Stephens and Yang (2014) by subsetting to the sample of respondents with non-zero wage income.
3. In Columns 2-6, respondents with zero reported income are treated as zeroes and I log annual income measures.
4. In Column 7, the universe is the set of respondents who report owning a house.
5. Standard errors clustered at the state of birth-by-year of birth level.
6. Occupational prestige is standardized so that the average in the sample is zero and the standard deviation is 1.

Table 14: Effect of Carnegie Grant on Occupational Choice

	Professional and Technical	Farm Owner	Manager	Clerical	Sales	Craftsmen	Operative	Service	Laborer
1(Carnegie Grant by Age 5)	0.005 (0.001)	0.007 (0.001)	-0.007 (0.002)	0.008 (0.001)	0.005 (0.002)	-0.006 (0.002)	-0.007 (0.002)	-0.001 (0.001)	-0.005 (0.002)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Parent controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Household (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
% in matched sibling sample	0.06	0.13	0.09	0.06	0.06	0.16	0.18	0.06	0.19
N	5,590,275	same	same	same	same	same	same	same	same
R <sup>2</sup>	0.56	0.66	0.55	0.54	0.54	0.55	0.57	0.56	0.60

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. Occupations defined using the major categories from the IPUMS occ1950 variable. I classify a small number of farm managers as managers, so that the farm owner category contains only farm owners and tenant-farmers. A complete mapping of granular occupations to these 9 broader occupations can be found here: [https://usa.ipums.org/usa-action/variables/occ1950#codes\\_section](https://usa.ipums.org/usa-action/variables/occ1950#codes_section)
4. These regressions are subset to individuals reporting a valid occupation.

Table 15: Effect of Education on Occupational Choice (Instrumented with Compulsory Schooling)

	<b>Professional and Technical</b>	<b>Farm Owner</b>	<b>Manager</b>	<b>Clerical</b>	<b>Sales</b>	<b>Craftsmen</b>	<b>Operative</b>	<b>Service</b>	<b>Laborer</b>
<b>Panel A: Men and Women</b>									
Years of Education	0.018 (0.006)	0.021 (0.004)	0.001 (0.005)	-0.001 (0.008)	0.005 (0.005)	-0.012 (0.006)	-0.042 (0.009)	0.004 (0.005)	0.005 (0.004)
First stage F-statistic:	39	.	.	.	.	.	.	.	.
N	4,374,405	.	.	.	.	.	.	.	.
Fraction in Sample	0.186	0.014	0.109	0.190	0.073	0.142	0.150	0.097	0.038
<b>Panel B: Men</b>									
Years of Education	0.007 (0.009)	0.033 (0.008)	-0.006 (0.009)	0.003 (0.007)	0.006 (0.006)	-0.011 (0.011)	-0.036 (0.011)	0.004 (0.006)	-0.001 (0.006)
First stage F-statistic:	21	.	.	.	.	.	.	.	.
N	2,442,498	.	.	.	.	.	.	.	.
Fraction in Sample	0.178	0.023	0.148	0.062	0.070	0.237	0.171	0.057	0.055
State of Birth (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Year of Birth (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Region*Birth Year (FEs)	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Age <sup>4</sup> , census year, gender	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

1. See notes to Table 13.
2. Standard errors clustered at the state of birth-by-year of birth level.
3. Models are linear probability models where the outcome in each column is an indicator for whether a respondent was in that occupation. Sample restricted to those with valid occupations.

Table 16: Effect of Occupation on Life Expectancy, Conditional on Reaching Age 65+

Occupation	Broad Category	Count in Matched Sample	Effect on Life Expectancy	1950 Excess Mortality	2017 Mortality Rate
Teachers	Professional	142,231	1.45	58	0.4
Farm owner	Farm owner	2,251,993	0.87	96	24.0
Professional and technical workers (general)	Professional	100,467	0.46	85	3.0
Electricians	Craftsmen	106,774	0.45	113	8.3
Janitors and sextons	Service	126,975	0.44	117	2.8
Insurance agents and brokers	Sales	107,309	0.15	96	0.5
Farm laborers, unpaid family	Laborers	797,318	0.13	96	16.4
Farm laborers, wage workers	Laborers	260,365	0.12	96	16.4
Mechanics (automobile)	Craftsmen	199,788	0.12	96	7.2
Manager, officials, and proprietors	Managers	1,273,996	0.07	86	1.1
Bookkeepers	Clerical	212,975	0.03	94	0.6
No listed occupation	—	1,184,701	0.00		
Carpenters	Craftsmen	326,215	-0.00	85	7.5
Clerical workers (general)	Clerical	599,441	-0.08	83	0.6
Salesmen and sales clerks (general)	Sales	803,248	-0.09	86	1.6
Foremen	Craftsmen	261,173	-0.16	96	1.0
Mechanics (general)	Craftsmen	190,796	-0.16	78	8.1
Operatives (general)	Operatives	1,225,092	-0.32	97	9.4
Machinists	Craftsmen	225,744	-0.39	120	5.4
Deliverymen and routemen	Operatives	167,220	-0.42	107	26.8
Truck and tractor drivers	Operatives	614,081	-0.43	107	26.8
Laborers (general)	Laborers	1,827,286	-0.60	169	9.7
Mine operatives and laborers	Operatives	327,067	-0.64	173	11.7
Painters, Construction, and Maintenance	Craftsmen	182,390	-0.99	130	14.2
Total		13,514,642			

Notes:

1. This table only includes occupations with more than 100,000 workers in my 1940 matched sample.
2. Broad Category defined using the major categories from the IPUMS occ1950 variable. A complete mapping of granular occupations to these 9 broader occupations can be found here: [https://usa.ipums.org/usa-action/variables/occ1950#codes\\_section](https://usa.ipums.org/usa-action/variables/occ1950#codes_section)

3. Effect on Life Expectancy is the value of the occupation fixed effect from the regression in Equation 4.
4. 1950 Excess Mortality is defined in Guralnick (1963) as the ratio of the number of deaths from death certificates indicating a “usual occupation” equal to the given occupation relative to the expected number of deaths for people of that age. So for example, teachers have a 1950 Excess Mortality of 58, indicating that only 58 teachers died relative to the 100 who you would expect to have died given the age distribution of teachers who died.
5. 2017 data from the “Hours-based fatal injury rates by industry, occupation, and selected demographic characteristics, 2017” table of the Census of Fatal Occupational Injuries (CFOI), from: <https://www.bls.gov/iif/oshcfoi1.htm#2017>. 2017 Mortality Rate is number of on-the-job deaths per 100,000 full-time equivalent workers. I include it in the table when I can find a comparable occupation in the 2017 CFOI data for a given occupation from 1940.

Table 17: Projected Effect of Carnegie Grant on Life Expectancy, Conditional on Reaching Age 65

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.009 (0.003)	0.027 (0.002)	0.025 (0.002)	0.017 (0.001)	0.019 (0.002)	0.006 (0.003)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,448,037	17,431,086	17,429,864	17,302,197	14,905,019	5,384,479
R <sup>2</sup>	0.03	0.05	0.06	0.09	0.37	0.57
Mean dep. var.	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. For each respondent, the outcome is the sum of the industry and occupation coefficients from Equation 4.

# Appendix Tables

Appendix Table A1: Effect of Carnegie Grant on Educational Attainment, State Cluster

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.572 (0.061)	0.179 (0.024)	0.076 (0.021)	0.090 (0.008)	0.134 (0.007)	0.101 (0.016)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,412,080	17,395,212	17,393,991	17,265,943	14,867,729	5,375,550
R <sup>2</sup>	0.20	0.26	0.29	0.32	0.55	0.73
Mean dep. var.	8.99	8.99	8.99	9.01	9.00	8.88

1. See notes to Table 3.
2. Standard errors clustered at the childhood state of residence level.

Appendix Table A2: Effect of Carnegie Grant on Educational Attainment, only Siblings

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.621 (0.034)	0.171 (0.019)	0.099 (0.018)	0.098 (0.014)	0.121 (0.015)	0.101 (0.017)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	5,716,477	5,713,985	5,713,550	5,569,796	5,421,502	5,375,550
R <sup>2</sup>	0.17	0.27	0.30	0.35	0.66	0.73
Mean dep. var.	8.83	8.83	8.83	8.86	8.88	8.88

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A3: Effect of Carnegie Grant on Educational Attainment, Continuous Measure

	1	2	3	4	5	6
%Years Age 5–20 with Grant	0.626 (0.035)	0.273 (0.033)	0.098 (0.030)	0.136 (0.014)	0.230 (0.015)	0.199 (0.027)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,412,080	17,395,212	17,393,991	17,265,943	14,867,729	5,375,550
R <sup>2</sup>	0.20	0.26	0.29	0.32	0.55	0.73
Mean dep. var.	8.99	8.99	8.99	9.01	9.00	8.88

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A4: Effect of Carnegie Grant on Educational Attainment, Midwest

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.507 (0.048)	0.196 (0.024)	0.111 (0.022)	0.104 (0.013)	0.149 (0.013)	0.073 (0.023)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	6,153,830	6,148,172	6,147,000	6,109,797	5,463,382	2,071,528
R <sup>2</sup>	0.12	0.19	0.23	0.26	0.49	0.71
Mean dep. var.	9.54	9.54	9.54	9.55	9.49	9.28

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A5: Effect of Carnegie Grant on Educational Attainment, Whites

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.559 (0.033)	0.171 (0.019)	0.071 (0.018)	0.088 (0.010)	0.133 (0.010)	0.104 (0.017)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	15,878,910	15,878,220	15,877,839	15,814,796	13,874,495	5,064,397
R <sup>2</sup>	0.15	0.21	0.25	0.28	0.52	0.72
Mean dep. var.	9.28	9.28	9.28	9.28	9.23	9.06

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A6: Effect of Carnegie Grant on Educational Attainment, Only Pre-Exposure

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.619 (0.043)	0.125 (0.021)	0.084 (0.019)	0.076 (0.013)	0.116 (0.015)	0.102 (0.025)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	13,825,663	13,813,393	13,812,192	13,693,528	11,998,357	4,451,156
R <sup>2</sup>	0.19	0.25	0.28	0.31	0.54	0.73
Mean dep. var.	8.76	8.76	8.76	8.77	8.79	8.69

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A7: Effect of Carnegie Grant on Educational Attainment, High School Control

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.453 (0.030)	0.167 (0.018)	0.070 (0.017)	0.089 (0.009)	0.130 (0.010)	0.097 (0.017)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Age at Public HS Entry (FEs)	Yes	Yes	Yes			
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year- by-Age at Public HS Entry (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,411,951	17,395,059	17,393,831	17,265,883	14,867,646	5,375,276
R <sup>2</sup>	0.20	0.26	0.29	0.32	0.55	0.73
Mean dep. var.	8.99	8.99	8.99	9.01	9.00	8.88

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A8: Effect of College Libraries on Educational Attainment

	1	2	3	4	5	6
1(College Library by Age 5)	0.635 (0.027)	0.108 (0.029)	0.023 (0.030)	-0.005 (0.026)	0.036 (0.026)	0.019 (0.048)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	17,343,845	17,327,003	17,325,785	17,197,694	14,799,176	5,349,331
R <sup>2</sup>	0.20	0.26	0.29	0.32	0.55	0.73
Mean dep. var.	8.99	9.00	9.00	9.01	9.00	8.88

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.

Appendix Table A9: Effect of Carnegie Grant on Educational Attainment, Plus 20 Years

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.489 (0.030)	0.061 (0.014)	-0.012 (0.015)	-0.006 (0.013)	0.014 (0.019)	0.053 (0.038)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	13,117,519	13,104,271	13,103,075	12,984,218	11,222,540	4,126,383
R <sup>2</sup>	0.19	0.25	0.28	0.32	0.54	0.73
Mean dep. var.	8.71	8.71	8.71	8.72	8.74	8.65

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. Unlike in the baseline specification, where I use the entire matched sample, this placebo test subsets to children living in places that did not receive a Carnegie grant and children who received a Carnegie grant when they were 21 years old or older. This ensures that none of the children in this sample had access to a Carnegie grant before they made their human capital investment decisions.

Appendix Table A10: Effect of Carnegie Grant on Educational Attainment, Lagged 20 Years

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.434 (0.035)	0.055 (0.012)	-0.016 (0.014)	0.029 (0.016)	0.038 (0.016)	0.016 (0.021)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	15,812,214	15,795,879	15,794,656	15,666,064	13,242,318	4,758,184
R <sup>2</sup>	0.20	0.27	0.30	0.33	0.56	0.74
Mean dep. var.	9.01	9.01	9.01	9.02	9.00	8.87

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. Unlike in the baseline specification, where I use the entire matched sample, this placebo test subsets to children living in places that did not receive a Carnegie grant and children who received a Carnegie grant when before the age of 5. This ensures that all of the children in this sample had access to a Carnegie grant before they made their human capital investment decisions or never had access to a Carnegie grant.

Appendix Table A11: Effect of Carnegie Grant on Educational Attainment, Perfect Matches

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.659 (0.037)	0.103 (0.015)	0.072 (0.014)	0.077 (0.013)	0.119 (0.016)	0.060 (0.027)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Age at Public HS Entry (FEs)	Yes	Yes	Yes			
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year- by-Age at Public HS Entry (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	9,058,899	9,041,858	9,040,595	8,904,648	6,467,682	2,248,509
R <sup>2</sup>	0.17	0.26	0.31	0.34	0.61	0.78
Mean dep. var.	9.67	9.67	9.67	9.70	9.74	9.61

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. This sample consists of matches that are accurate according to the two match accuracy measures discussed in the Matching Appendix.

Appendix Table A12: Effect of Carnegie Grant on Educational Attainment, All Matches

	1	2	3	4	5	6
1(Carnegie Grant by Age 5)	0.530 (0.031)	0.152 (0.017)	0.063 (0.015)	0.070 (0.008)	0.116 (0.009)	0.087 (0.012)
Child controls	Yes	Yes	Yes	Yes	Yes	Yes
Age at Public HS Entry (FEs)	Yes	Yes	Yes			
Enumeration District (FEs)		Yes	Yes	Yes	Yes	Yes
Parent controls			Yes	Yes	Yes	Yes
County-by-Birth Year- by-Age at Public HS Entry (FEs)				Yes	Yes	Yes
Census Microfilm Page (FEs)					Yes	
Household (FEs)						Yes
Observations	21,755,034	21,738,814	21,737,626	21,620,669	19,503,367	7,850,353
R <sup>2</sup>	0.19	0.25	0.28	0.30	0.51	0.70
Mean dep. var.	9.01	9.01	9.01	9.01	9.01	8.83

1. See notes to Table 3.
2. Standard errors clustered at the childhood county of residence level.
3. This sample consists of all adults in 1940 who match to at least one early childhood census record, independent of match quality.

See the Matching Appendix for more details.

# Data Appendix

## 0.1 Decennial Census Data

Below, I describe the construction of each variable I use from 1900–1940 Census microdata. While the data are accessed on an NBER server, they are in the process of being cleaned and standardized by IPUMS (Ruggles et al., 2020), so many of the variables have the same format as the variables in public IPUMS census microdata.

- Race: Constructed from the RACE variable. I rely on race as it was reported in the childhood census year. I treat anyone enumerated as “Mulatto” as black.
- Years of educational attainment: Constructed from the EDUC variable, represents level of educational attainment in 1940. I map adults who have not progressed past kindergarten and those without any schooling to 0. First grade is 1 year, second grade is 2 years, ..., and a 4-year degree is 16. Lastly, I map 5+ years of college to 18.
- High school graduation: Constructed from the EDUC variable, is an indicator equal to 1 if respondent was a high school graduate or had attended some college in 1940.
- College attendance: Constructed from the EDUC variable, is an indicator equal to 1 if respondent attended at least one year of college by 1940.
- Occupational prestige: Constructed from the PRESGL variable. Assigns an occupational standing score (based on Siegel, 1971) to each occupation in the 1940 census microdata.
- Weeks worked: Constructed from the WKSWORK1 variable. Reports the number of weeks worked, ranging from 0 to 52, for respondents in the 1940 census.
- Hours worked: Constructed from the HRSWORK1 variable. Reports the number of hours worked during the week of March 24-30, 1940. Includes unpaid family work (eg. on a farm). Topcoded at 98 hours.

- Enumeration district: Constructed from the ENUMDIST variable. Not comparable across years.
- City name: Constructed by extensively cleaning and manipulating these variables:
  1. 1900: us1900m\_0045 (ward), us1900m\_0052 (city)
  2. 1910: us1910m\_0052 (city) us1910m\_0053 (city) us1910m\_0063 (ward)
  3. 1920: us1920c\_0057 (city) us1920c\_0058 (city) us1920c\_0068 (city) us1920c\_0069 (ward) stdmcd (city)
  4. 1930: stdmcd (city)

I remove common prefixes (eg. precinct) and suffixes (eg. district). I trim white space, remove special characters, standardize words like “mount” and “mt.”, and enforce that respondents in each enumeration district must map to the same city. When this does not happen, I replace city name with the modal city name reported within that enumeration district. In some cases, respondents in enumeration districts will consistently be linked to multiple city names. For example, respondents from Brooklyn, New York would report living in Brooklyn and New York City. In these cases, I use both city names in an attempt to link respondents to Carnegie grants and public libraries.

- Occupation: Constructed from the OCC1950 variable. Represents Ruggles et al. (2020) work to standardize reported occupations across census microdata. I use the 200 granular occupations to estimate Equation 4. And I aggregate the granular occupations into nine broader groups based on the classification reported here, for use in the linear probability models in Tables 14 and 15: [usa.ipums.org/usa-action/variables/occ1950#codes\\_section](https://usa.ipums.org/usa-action/variables/occ1950#codes_section).
- Industry: Constructed from the IND1950 variable. Represents Ruggles et al. (2020) work to standardize reported industries across census microdata.
- Birthplace: Constructed from the BPL variable, standardized by Ruggles et al. (2020).

- Wage: Constructed from the INCWAGE variable. Topcoded at \$5,001 dollars in 1940.
- Probability  $\geq$  \$50 of Non-Wage Income: Constructed from INCNONWG in 1940.
- Non-wage income: Available in the 1950 census, but not the 1940 census. constructed as the difference between the INCTOT variable and the INCWAGE variable. I use this measure to impute non-wage income in 1940, as described in the text.

## High School Data

I collect data about the founding dates of all high schools in the United States from 1890–1951 using four sources:

1. Annual and Biennial Bureau of Education reports from 1890–1940 listing the number of high school students and teachers in towns and cities with population greater than a threshold level of 2,500, 4,000, 5,000, or 10,000 depending on the year.
2. Censuses of all high schools in the United States, collected by the Bureau of Education every 1-2 years from 1890–1905 and in 1912. After 1912, the next census of high schools was published in 1951. These censuses contain information about the number of students, the number of teachers, the length of study, and the founding year of each high school; the founding year was only collected in a few of the censuses, including in 1903.
3. Lists of all accredited high schools in the United States, published by the Bureau of Education every 2-6 years from 1915 through 1944. Accreditation standards varied by state, and public colleges used these published bulletins to offer admission to local students with a diploma from a high school that met a set of criteria. These criteria often included (1) a requirement that the high school offer four years of study, and (2) that the high school offer at least a minimum number of math, English, and history credits. This source only contains the name of each high school, the accrediting body, and the location of the school. There is no information about enrollment or school size in these lists of accredited high schools.

4. Patterson's College and School directory, available from 1904-1924 annually. These directories list the cities and towns that had high schools in each year, with almost no additional information beyond the location and the name of each principal.

I used a grant from the Institute of Education Sciences to digitize and standardize the documents described above. The data come from around 8,000 pages of tables that I scanned or extracted from hundreds of books that several university libraries have recently scanned and made available online. I worked with a group of data-enterers to transcribe these tables. I then standardized and validated the transcriptions, linking city and town names across years and sources.

For each of the 23,000 cities and towns listed in these tables, I identify the first year when they had a public public high school. I match 76% of these cities and towns with at least one high school before 1940 to places in the 1900–1930 complete count censuses. In Figure 4, I plot the founding date of the earliest high school in each modern-day town in the United States, as of 1944.

## Matching Appendix

Below, I describe the matching process that I use to match children and young adults from the 1900–1930 decennial censuses to adults in the 1940 decennial census. of each variable I use from 1900–1940 Census microdata. I follow the standard matching algorithm from Abramitzky, Bousttan, and Eriksson (2012), first independently matching each childhood census to the 1940 census using an iterative matching procedure. After each step, I collect all matches and remove them from the data—I then attempt to match the remaining records. If any childhood record matched to multiple adult records from 1940, I remove that childhood record and those multiple adult records from the search for matches because it is impossible to ascertain which of the adult records should match to the childhood record.

**Step 1:** I begin by matching 0–25 year old males in the 1900 census to adult men aged 20–65 in the 1940 census exactly on:

- The NYSIIS<sup>24</sup> of the first name
- The NYSIIS of the last name
- Exact age
- State or country of birth
- Race

I then iterate on age, allowing for differences of one year and then two years in either direction. In Step 1, I search for matches in five iterations because of this age mismatch allowance: I search for exact matches on age, matches where the adult record’s reported age is one year before the childhood record, matches where the adult record’s reported age is one year after the childhood record, matches where the adult record’s reported age is two years before the childhood record, and matches where the adult record’s reported age is two years after the childhood record.

---

<sup>24</sup>NYSIIS codes are New York State Identification and Intelligence System Phonetic codes: a phonetic coding system for names.

**Step 2:** After this first set of matches, I repeat Step 1, but instead of looking for NYSIIS codes corresponding to each respondent's name, I match exactly using standardized first and last names, where I standardize nicknames and remove any remaining middle initials. Lastly, I repeat Step 1 using the raw name strings in the census.

**Step 3:** I then repeat Steps 1 and 2 without the restriction that records match exactly on race. This allows for changes in reported race between the childhood and adult census year.

**Step 4:** I then perform three more searches for matches: looking for exact matches on (1) NYSIIS first/last name, then (2) cleaned name, and then (3) raw name, along with birth state/country, and race (no age restriction, to allow for transcription errors in the age variable).

**Step 5:** Lastly, I repeat Step 4, relaxing the restriction that the childhood record match to the adult record on race.

I then append together the four sets of matches from 1900 to 1940, 1910 to 1940, 1920 to 1940, and 1930 to 1940. For each of the adult men in 1940, I keep the two highest-quality matches, based on the ordering defined in Steps 1–5 above. As a tiebreaker, I always choose the earliest record that matched to an adult in 1940. So for example, if an adult male in 1940 matches exactly on NYSIIS first/last name, age, state of birth, and race to a record in each of 1900, 1910, 1920, and 1930, I will keep the 1900 and 1910 records as potential matches.

For each adult record that matched to at least one childhood record, I calculate two measures of match accuracy:

1. **Parental Name Accuracy:** For adult records in 1940 that matched to two childhood records (A and B), I compare the reported parents of A and B in the early censuses. If the first two letters of A's mother and B's mother are the same or the first two letters of A's father and B's father are the same, I call this a high-quality match, because it is very unlikely to happen by chance unless A and B represent the same person.
2. **Birthplace Accuracy:** In 1940, a subset of adults were asked to report the birthplace of their mother and father. So, for all adult records in 1940 that match to at least one childhood record, I compare the reported birthplace of the adult in 1940 with the childhood record that

matched to the adult record. If the reported birthplaces of the records' mother or father match across censuses, I call this a high-quality match.

Below, I present a table showing the fraction of potential matches from the first quality check, for records that matched to two childhood census records. These 21,364,440 records represent 10,682,220 adult males who each matched to two childhood census records (from two separate decennial census years between 1900–1930) such that both childhood census records have at least one parent in the household with a non-missing name. For each of the 36 iterations discussed above, I report the fraction of these matched childhood-adult record tuples that have at least one parent (mother or father) whose first two letters of their first names are the same across the two childhood records in the tuple.

Matching Table 1: Parental Name Accuracy Measure by Loop Iteration

Iteration	# Low-Quality	% Low-Quality	# High-Quality	% High-Quality	Total
1	2,647,659	26	7,529,692	74	10,177,351
2	855,068	35	1,601,099	65	2,456,167
3	637,356	28	1,622,583	72	2,259,939
4	388,385	51	368,776	49	757,161
5	260,936	52	244,288	48	505,224
6	355,820	41	514,657	59	870,477
7	108,645	53	94,549	47	203,194
8	79,033	47	87,882	53	166,915
9	55,407	70	23,626	30	79,033
10	39,788	71	16,122	29	55,910
11	194,585	52	177,261	48	371,846
12	61,608	63	36,915	37	98,523
13	47,120	59	32,473	41	79,593
14	33,660	74	11,796	26	45,456
15	24,723	75	8,215	25	32,938
16	148,035	64	82,435	36	230,470
17	107,070	72	41,018	28	148,088
18	80,409	71	32,711	29	113,120
19	79,217	78	21,792	22	101,009
20	56,350	78	15,611	22	71,961
21	12,621	76	4,063	24	16,684
22	7,786	81	1,799	19	9,585
23	5,757	80	1,438	20	7,195
24	5,442	85	985	15	6,427
25	4,345	86	726	14	5,071
26	4,187	83	884	17	5,071
27	3,672	84	696	16	4,368
28	3,013	84	590	16	3,603
29	3,210	86	516	14	3,726
30	2,454	86	399	14	2,853
31	1,013,371	67	508,752	33	1,522,123
32	315,203	79	82,932	21	398,135
33	236,460	83	49,605	17	286,065
34	166,011	80	40,809	20	206,820
35	34,586	87	5,086	13	39,672
36	20,004	88	2,663	12	22,667
Total	8,098,996	38	13,265,444	62	21,364,440

I now present the same table for the second measure of accuracy, below.

Matching Table 2: Birthplace Accuracy Measure by Loop Iteration

Iteration	# Low-Quality	% Low-Quality	# High-Quality	% High-Quality	Total
1	643,931	11	5,203,196	89	5,847,127
2	231,149	20	926,495	80	1,157,644
3	175,153	16	954,147	84	1,129,300
4	112,087	30	262,979	70	375,066
5	77,548	31	176,493	69	254,041
6	76,934	13	505,112	87	582,046
7	26,405	24	82,484	76	108,889
8	18,962	20	77,911	80	96,873
9	15,287	34	30,031	66	45,318
10	11,327	34	21,748	66	33,075
11	44,061	16	224,098	84	268,159
12	14,549	26	41,389	74	55,938
13	10,936	22	37,786	78	48,722
14	9,222	33	18,816	67	28,038
15	7,095	33	14,109	67	21,204
16	50,073	28	126,433	72	176,506
17	29,779	33	59,342	67	89,121
18	21,548	32	45,054	68	66,602
19	20,336	35	37,542	65	57,878
20	14,963	36	26,791	64	41,754
21	4,347	29	10,877	71	15,224
22	2,358	33	4,778	67	7,136
23	1,688	32	3,649	68	5,337
24	1,531	33	3,055	67	4,586
25	1,260	34	2,448	66	3,708
26	935	27	2,473	73	3,408
27	868	30	2,043	70	2,911
28	663	27	1,790	73	2,453
29	679	27	1,859	73	2,538
30	635	30	1,517	70	2,152
31	423,746	39	655,240	61	1,078,986
32	120,088	40	181,679	60	301,767
33	79,616	36	139,186	64	218,802
34	58,668	35	108,112	65	166,780
35	10,928	35	20,452	65	31,380
36	5,104	28	13,043	72	18,147
Total	2,324,459	19	10,024,157	81	12,348,616

I now present the same table for the first measure of accuracy, but instead of examining all potential matches, I focus on pairs of childhood records in two decennial censuses from 1900–1930 that each matched to the same adult record in 1940 where each of the two childhood records independently matched to the adult record in the same iteration of the loop described above.

Matching Table 3: Parental Name Accuracy Measure by Loop Iteration for Matches from the Same Loop

Iteration	# Low-Quality	% Low-Quality	# High-Quality	% High-Quality	Total
1	919,616	14	5,664,724	86	6,584,340
2	156,606	17	770,426	83	927,032
3	95,132	13	639,040	87	734,172
4	36,646	18	170,702	82	207,348
5	19,410	16	98,548	84	117,958
6	77,994	21	297,204	79	375,198
7	10,232	25	30,668	75	40,900
8	5,160	19	22,292	81	27,452
9	2,572	26	7,308	74	9,880
10	1,512	25	4,506	75	6,018
11	46,146	34	90,660	66	136,806
12	6,804	38	11,206	62	18,010
13	3,886	33	7,736	67	11,622
14	2,090	35	3,862	65	5,952
15	1,178	32	2,554	68	3,732
16	15,758	35	29,294	65	45,052
17	5,674	31	12,520	69	18,194
18	3,616	29	8,982	71	12,598
19	3,270	29	8,024	71	11,294
20	1,992	26	5,714	74	7,706
21	710	48	782	52	1,492
22	188	47	216	53	404
23	96	35	182	65	278
24	72	39	112	61	184
25	64	46	74	54	138
26	192	47	216	53	408
27	122	46	144	54	266
28	82	39	130	61	212
29	100	49	106	51	206
30	74	50	74	50	148
31	288,988	45	356,608	55	645,596
32	39,188	55	31,798	45	70,986
33	35,844	66	18,262	34	54,106
34	25,976	51	24,964	49	50,940
35	1,910	67	934	33	2,844
36	1,698	70	730	30	2,428
Total	1,810,598	18	8,321,302	82	10,131,900

To interpret this table, consider the first row, which shows that 6,584,340 potential childhood records (for 3,292,170 adult records) matched in the first iteration of the loop described above. 86% of these tuples have two childhood records whose fathers or mothers have the same first two letters of their first names in each of the childhood census records. If false positives in this set of matches are distributed independently across potential matches for all potential matches within a loop iteration, and if this measure of accuracy correctly identifies valid matches, then  $0.86^{0.5} = 93\%$  of matches in the first iteration of the loop are accurate. I use the tabulation from Matching Table 3 to calculate accuracy measures for each iteration of the looped matching algorithm described above. Accuracy ranges from 93% for the first iteration of the loop down to  $0.30^{0.5} = 55\%$  for the last (36th) iteration of the loop.

Lastly, I select one childhood census record for each adult census record in 1940 that matched to two childhood census records. I prioritize childhood records in three steps:

1. I first prioritize 0–15 year old childhood census records over 16–25 childhood census records
2. , I then prioritize the earliest childhood census record if a pair of census records matched according to the parental name accuracy measure described above
3. I then prioritize matches from the first 20 iterations of the loop described above, because those matches are on average more accuracy than matches from loop iterations 21–36
4. Lastly, I prioritize the earliest match (prioritizing matches from 1900 over matches from 1910 and so on)

Below, I show which matches I keep for the final analysis sample using this prioritization of matches where there were multiple childhood records to choose from.

Matching Table 4: Potential Matches for the Final Analysis Sample

Iteration	# Drop	% Drop	# Keep	% Keep	Total
1	5,846,785	38	9,406,353	62	15,253,138
2	1,349,945	34	2,612,177	66	3,962,122
3	1,403,608	39	2,231,742	61	3,635,350
4	393,326	30	912,405	70	1,305,731
5	365,490	38	605,397	62	970,887
6	513,599	34	994,743	66	1,508,342
7	112,056	31	246,907	69	358,963
8	109,368	36	194,999	64	304,367
9	38,142	27	103,068	73	141,210
10	43,492	37	74,093	63	117,585
11	232,207	29	568,947	71	801,154
12	56,420	28	148,073	72	204,493
13	56,214	33	115,342	67	171,556
14	22,442	24	70,060	76	92,502
15	26,630	34	52,860	66	79,490
16	135,609	27	366,944	73	502,553
17	81,372	29	195,423	71	276,795
18	78,107	35	147,692	65	225,799
19	48,899	27	131,797	73	180,696
20	55,674	37	95,431	63	151,105
21	13,261	29	32,964	71	46,225
22	7,622	34	14,719	66	22,341
23	6,433	37	10,947	63	17,380
24	4,868	36	8,477	64	13,345
25	5,016	40	7,414	60	12,430
26	4,437	35	8,135	65	12,572
27	3,701	35	6,910	65	10,611
28	3,345	36	6,037	64	9,382
29	3,045	35	5,548	65	8,593
30	2,910	38	4,807	62	7,717
31	988,681	36	1,770,324	64	2,759,005
32	267,271	38	437,845	62	705,116
33	195,571	36	344,787	64	540,358
34	140,320	35	255,775	65	396,095
35	27,268	38	45,025	62	72,293
36	15,366	37	26,449	63	41,815
Total	12,658,500	36	22,260,616	64	34,919,116

From the 22,260,616 childhood records at this, I drop an additional three records that do not have a valid state of residence in the childhood census data. And I drop 4,442,495 matches that were low-quality according to either of the two accuracy measures defined above.<sup>25</sup> This leaves me with 17,818,118 records in my final analysis sample out of 39,175,147 of the 20–65 year old men in 1940 who I attempted to match to 0–25 year old males in the 1900–1930 decennial censuses.

<sup>25</sup>In Appendix Table A11 I show that my main educational attainment results are unchanged when I subset to matches that are accurate, according to one of the two measures described above. And in Appendix Table A12 I show that including these 4.4 million matches that are low-quality has only a small effect on my main results.

This represents a match rate of 45.5%. And based on the accuracy estimates from Matching Table 3, 80% of these matches are true matches and 20% of these matches are false positives. This is towards the frontier of the methods analyzed in Abramitsky et al. (2020), although my calculation of accuracy measures is new in the literature, making comparisons difficult.