

Selection into Identification in Fixed Effects Models, with Application to Head Start

Douglas L. Miller*
Cornell University
and NBER

Na'ama Shenhav†
Dartmouth College
and NBER

Michel Z. Grosz‡
Federal Trade Commission

May 22, 2020[§]

Abstract

Many papers use fixed effects (FE) to identify causal impacts of an intervention. When treatment status only varies within some FE groups (e.g., families, for family fixed effects), FE can induce non-random selection of groups into the identifying sample, which we term *selection into identification* (SI). This paper empirically documents SI in the context of several family fixed effects (FFE) applications with a binary treatment. We show that the characteristics of the FFE identifying sample are different than the overall sample (and the policy-relevant population), including having larger families. The main implication of this is that when treatment effects are heterogeneous, the FE estimate may not be representative of the average treatment effect (ATE). We show that a reweighting-on-observables estimator can help recover the ATE for policy-relevant populations from FE estimates, and recommend its use either as a primary estimator or as a diagnostic tool to assess the importance of SI. We apply these insights to re-examine the long-term effects of Head Start in the PSID and the CNLSY using FFE. When we reweight the FFE estimates, we find that Head Start leads to a 2.6 percentage point (p.p.) increase (s.e. = 6.2 p.p.) in the likelihood of attending some college for white Head Start participants in the PSID. This participants' ATE is 78% smaller than the traditional FFE estimate (12 p.p.). We also find that the CNLSY Head Start participants' ATE is smaller than the FE estimates. This suggests that in many contexts SI may distort the external validity of FE estimates.

*Doug Miller, Policy Analysis and Management, Cornell University, Email: dml336@cornell.edu

†Na'ama Shenhav, Department of Economics, Dartmouth College, E-mail: naama.shenhav@dartmouth.edu

‡Michel Grosz, Bureau of Economics, Federal Trade Commission, E-mail: mgrosz@ftc.gov

[§]We would like to thank Colin Cameron, Liz Cascio, Janet Currie, Hilary Hoynes, Pat Kline, Erzo F.P. Luttmer, Jordan Matsudaira, Zhuan Pei, Maya Rossin-Slater, Doug Staiger, Dmitry Taubinsky, Chris Walters, and participants at the AEA Meetings, Cornell, Dartmouth, CSWEP CEMENT Workshop, Hebrew University, McGill University, NBER Labor/Children's Summer Institute, Northwestern, SEA Meetings, SOLE, Syracuse/Cornell Summer Workshop in Education and Social Policy, UC Merced, and the War on Poverty Conference at the University of Michigan. We are also grateful to Alex Magnuson, Wenran Li, and Wenrui Huang for excellent research assistance. The views expressed in this article are not necessarily those of the Federal Trade Commission nor any of its commissioners.

1 Introduction

Fixed Effects (FE) are frequently used to obtain identification of the causal impact of an attribute, intervention, or policy. These models have been used to identify the impacts of academic peers (school-grade FE; Hoxby, 2000; Carrell and Hoekstra, 2010); criminal peers (facility-offense FE; Bayer, Hjalmarsson and Pozen, 2009); the local health care environment (individual FE; Finkelstein, Gentzkow and Williams, 2016); participation in means-tested programs (family FE; Currie and Thomas, 1995; Garces, Thomas and Currie, 2002; Deming, 2009; Rossin-Slater, 2013); neighborhood quality (family FE; Chetty and Hendren, 2018a); and minimum wage laws (county-pair-year FE; Dube, Lester and Reich, 2010), to give a few examples. Many of the estimates in these studies are targeting the average treatment effect for a policy-relevant population (e.g., participants). However, in contrast with other common estimators, there is not yet a comprehensive framework for considering the *external validity* of FE estimates.

In this paper, we document that FE can induce a special type of systematic selection in estimation, which we term “*selection into identification*” (SI). SI occurs because FE estimates are only identified from FE groups (e.g. families, for family FE) that have variation in treatment (“switchers”), which may exclude some groups.¹ This is a distinct problem from whether within-group comparisons are internally valid, which has been the typical subject of debate for FE estimators (e.g., Bound and Solon, 1999), and which is not the focus of this paper.

While researchers may be aware of this issue in an abstract sense, our review of the applied literature indicates that it is rarely discussed and almost never empirically examined. This may be because researchers assume that SI is a minor concern, or because there is not a user-friendly tool to address it. We find across multiple previous applications of FE that switchers are typically (i) a small subset of the sample and (ii) systematically different than the overall population. This highlights a new and important concern about the external validity of FE estimates.² We also show that reweighting-on-observables methods can be used to help “undo” SI to recover the average treatment effect (ATE) for policy-relevant populations. We apply these insights to re-examine the external validity of prior family FE (FFE) estimates of the long-run impact of Head Start, a federally-funded pre-school program.

We begin by illustrating SI in a well-known application of FE: an FFE model where the treatment of interest is an indicator for whether an individual attended Head Start (see, e.g., Garces, Thomas and Currie, 2002, henceforth GTC; Deming, 2009). We use the Panel Study of Income Dynamics (PSID), as in GTC.³ First, we show that FFE uses substantially fewer identifying groups

¹ In the presence of control variables that vary within a group, then there may be variation among non-switchers “net of controls.” We focus on cases where this phenomenon is small in magnitude, and formalize this extension in Section 5.2.

² SI is distinct from issues related to conditional variance weighting in Gibbons, Suarez and Urbancic (2018), which focuses on cases where all groups are switchers.

³ Similar FFE models have been used to evaluate many other treatments; for public housing, see Andersson et al. (2016); for WIC, see Chorniy, Currie and Sonchak (2018); Currie and Rajani (2015); for health, see Almond, Chay

relative to an estimation model without FE. Among the 5,355 children that have siblings in our sample, only 1,098 children reside in switcher households. Second, we find that the loss of sample variation is systematically related to observables. In particular, families that have a very low likelihood of attending Head Start or a very high likelihood of attending Head Start are both unlikely to be switchers. Additionally, families with more children are significantly more likely to be switchers. As a result, switchers are not representative of the overall sample or Head Start participant families along many dimensions. Third, SI varies across subgroups: the FFE identifying sample misses 93% of the sibling sample for white children, but only 62% of the sample for black children.

The main implication of SI is that under heterogeneous treatment effects, the FE estimate is no longer representative of the ATE for the sample. We show that the bias from SI is not addressed by “undoing” conditional variance weighting among switchers (Gibbons, Suarez and Urbancic, 2018), and that in our applications the bias from SI is larger than that from conditional variance weighting.⁴ A second result is that the difference between the OLS estimate and FE estimate in part reflects the disparity in the identifying sample (i.e., the causal treatment effect for non-switchers), and should not be interpreted as solely reflecting OLS bias.

To help recover the ATE for policy-relevant populations (“targets”), we develop a method for reweighting FE estimates that builds on extrapolation techniques from the experimental and IV literatures.⁵ In particular, we show that the ATE for a given target can be computed by weighting group-level FE estimates by the ratio of two propensity scores: (i) the propensity to be in the target (e.g., a program participant) and (ii) the propensity to be in the switcher population. This reweighting is valid under two key assumptions: (1) that these propensity scores can be estimated using observable covariates; and (2) that treatment effects are uncorrelated with being a switcher conditional on covariates (i.e., selection into switching on observables). We discuss testable implications of these assumptions and show empirical support for them in our applications. These assumptions provide a weaker alternative to the typical homogeneous treatment effects assumption required for the FE estimate to deliver the ATE. However, in some cases a researcher might prefer the less-externally-valid FE estimate rather than rely on these assumptions.⁶

We demonstrate the performance of our reweighting using Monte Carlo simulations. We find that reweighting reduces or eliminates bias relative to FE in the presence of covariate-based treatment heterogeneity. In Section 5, we discuss several extensions of our basic setup, such as how the inclusion of covariates that vary within a group can create additional “residual switchers,” and show how reweighting can be applied to a non-linear model.

and Lee (2005); Figlio et al. (2014); Abrevaya (2006); Black, Devereux and Salvanes (2007); Xie, Chou and Liu (2016), among others. We summarize the prevalence of this design in Section 2.

⁴This is consistent with Gibbons, Suarez and Urbancic (2018), whose findings suggest that the bias from conditional variance weighting is less than 5% for 75% of estimates.

⁵See Angrist and Fernandez-Val (2013) and Aronow and Carnegie (2013) for extrapolation from IV, and Stuart et al. (2011) and Andrews and Oster (2019) for extrapolation from experiments.

⁶In these cases, researchers can note that FE estimates are specific to the population of switchers, and characterize that population.

Based on these findings we propose new standards for practice when presenting results using FE research designs: (i) clearly show the sample size when limited to switcher families and quantify the contribution of “residual switchers”; (ii) show the balance of covariates across switcher and non-switcher families (e.g. Table 2); (iii) reweight FFE estimates for a representative population (e.g. Table 6). Reweighted estimates can be presented either as an additional diagnostic tool or as a primary measure of treatment effects. We are not the first to use the more rigorous reporting standards in (i) and (ii), but in our survey of the FFE literature the vast majority do not discuss either of these issues. Out of 35 papers with binary treatments, only one-third reported the sample sizes we recommend in (i), and just one paper included (ii).⁷

In Section 6, we apply these insights to quantify the importance of selection into identification for FFE estimates of the long-run impact of Head Start. Head Start has a budget of \$8.6 billion dollars and annually enrolls roughly 60% of the number of 3 and 4 year old children in poverty, which makes it a quantitatively important intervention for this population (Carneiro and Ginja, 2014).⁸ FFE have been used to identify the long term impacts of Head Start in many of the foundational studies of this program (Currie and Thomas, 1995; Deming, 2009, GTC), which find positive impacts on economic and non-cognitive outcomes of participants measured in adulthood.

First, using data from the PSID and the Children of the National Longitudinal Study of Youth (CNLSY) (as in GTC and Deming, 2009),⁹ we newly document that Head Start has greater returns in larger families across multiple human capital measures. This might result from the fact that parental time investment in children’s human capital is spread more thinly in larger families, which in turn could lead to greater returns to alternative investments, such as Head Start.¹⁰

Second, we show that the FFE estimate of the impact of Head Start is typically too large relative to the reweighted ATE for participants, in part due to over-representation of larger families. Using the PSID, we obtain an FFE estimate that suggests that Head Start leads to a statistically significant 12 p.p. increase in attendance of some college for white individuals. However, when we reweight the FFE estimates to obtain the ATE for white *Head Start participants*, we find that Head Start leads to a 2.6 percentage point (p.p.) increase in the likelihood of attending some college (s.e. = 6.2 p.p.). This reweighted ATE for participants is 78% smaller than the FFE estimate (p<0.05 for the difference.) The participants’ ATE is also 91% smaller than the estimated effects

⁷Important exceptions include Finkelstein, Gentzkow and Williams (2016) and Wiswall (2013), who include a substantive discussion of external validity concerns, and Currie and Rossin-Slater (2013), who show the change in summary statistics when moving from the sample of mothers with multiple births to the sample that had at least one hurricane during pregnancy, and are therefore likely to be “switchers.” GTC report the number of identifying observations used to identify Head Start for the entire sample (not for subsamples), and Deming (2009) reports the aggregate number of identifying observations used to identify the pre-school, Head Start, and no-formal-care coefficients, but not for each coefficient.

⁸See Gibbs, Ludwig and Miller (2013) for an overview of the Head Start program.

⁹Although we rely on the PSID as in GTC, we use a sample of siblings three times as large as in GTC due to the additional data that has been collected since the publication of that paper and somewhat different inclusion criteria.

¹⁰In Section 6, we show that this heterogeneity by family size is not explained by other covariates or by larger families having longer sibling cohort spans. Instead it appears that there is something important about family size per se.

on college-going in GTC; 45% to 91% smaller than unadjusted estimates for all participants from other FFE studies (Bauer and Schanzenbach, 2016; Deming, 2009); and 51% smaller than estimates from the county roll-out of Head Start (Bailey, Sun and Timpe, 2018).

Similarly, in the CNLSY, the FFE estimate suggests that Head Start leads to an 8.5 p.p. increase in high school completion, while the reweighted estimate for Head Start participants is 44% smaller and not statistically significant ($p < 0.10$ for the difference.) Reweighting also attenuates the previously-estimated impact of Head Start on idleness and having a learning disability, relative to the FFE estimates. In contrast, reweighting has little to no impact on the estimates for the poor health outcome.

In Section 7, we discuss the broader relevance of SI for panel fixed effects models, particularly for short panels and “lumpy” treatment variables (e.g. binary treatments). In particular, we walk through how the guidelines that we have put forth could apply to three prior FE applications.

This paper makes three main contributions. First, we empirically document across multiple FE settings that many groups are not switchers, and that switchers are observably different than non-switchers.¹¹ This stands in contrast to the assumption that all FE groups are switchers, which underlies several econometric techniques (Bates et al., 2014; Loken, Mogstad and Wiswall, 2012; Gibbons, Suarez and Urbancic, 2018). It also implies that the method in Gibbons, Suarez and Urbancic (2018) for recovering the ATE from FE models by “undoing” the conditional variance weighting – which relies on this assumption – is not sufficient to recover the ATE in many contexts. Finally, it suggests that more researchers using FE models should examine SI and discuss this in their interpretation of FE estimates. We show that this is rarely done in current practice in Section 2, and provide easy rules-of-thumb to guide such investigations.

Second, we develop a general reweighting method to obtain the ATE for policy-relevant populations from FE models.¹² This provides an alternative to defining the population of interest to be switchers, as discussed in Angrist (1998) (for OLS models) and Imai and Kim (2019) (for FE models). In doing so, we add to a small set of papers that provide methods for researchers to extrapolate FE treatment effects beyond switchers. Existing approaches include bounding methods for nonparametric panel models (Section 4 of Chernozhukov et al., 2013), and using a random coefficients model for treatment effect heterogeneity (Wooldridge, 2005, 2019; Bates et al., 2014). Our literature review suggests that practitioners do not pursue either of these current methods for correction. Relative to these other correction methods, we provide easy-to-implement guidelines for using our reweighting approach that is closely fitted to current empirical practice.

Our reweighting solution is most similar in spirit to Angrist and Fernandez-Val (2013), who

¹¹This is distinct from recent works on the validity of difference-in-difference and other two-way FE strategies, where the empirical specification ensures that SI is unlikely to be a concern. See, e.g. Goodman-Bacon (2018); Borusyak and Jaravel (2017); Callaway and Sant’Anna (2018); Chaisemartin and D’Haultfoeille (2019). We discuss SI in two-way FE designs in Section 5.2.

¹²This is distinct from strategies that use reweighting for internal validity, such as traditional propensity score estimation methods.

reweight IV estimates using discrete covariates.¹³ We tailor this idea to the FE context to develop our reweighting method, and show the set of assumptions under which these weights can recover the ATE from FE estimates. Also different than Angrist and Fernandez-Val (2013), we reweight using propensity scores, which allow for greater flexibility in conditioning variables.

Third, our empirical findings contribute to a growing body of work investigating the long term effects of Head Start using quasi-experimental methods (Ludwig and Miller, 2007; Carneiro and Ginja, 2014; Thompson, 2017; Bauer and Schanzenbach, 2016; Johnson and Jackson, 2017; Bailey, Sun and Timpe, 2018; Pages et al., 2019; Barr and Gibbs, 2018; Leuven and de Haan, 2016, in addition to the FFE papers above). These studies typically present LATE or ITT estimates, and find improvements in childhood health, reductions in adolescent behavioral problems and obesity, and increases in adult educational attainment and earnings.¹⁴ Relative to most of these studies, we evaluate the effect of Head Start on longer-run outcomes; show that these effects vary significantly by family size; and also adjust estimates using covariate re-weighting to get closer to the ATE for Head Start participants. We show that incorporating this adjustment lowers the estimated long term effect of Head Start relative to FE estimates.

2 A Survey of FFE Applications

We begin by performing a survey of FE papers in the literature to gauge the prevalence of these methods, particularly among studies that use binary treatments or have small sample sizes (which may be more susceptible to SI.) Since our application focuses on a FFE model, we focus on applications of this particular method in the literature. This focus will lead us to undercount the prevalence of FE more broadly, but provides an unambiguous example of a short-panel setting which is susceptible to SI concerns. We surveyed publications from January 2000 to May 2017 in 11 leading journals that publish applied microeconomics articles. We include all studies that use family fixed effects as a primary or secondary strategy.¹⁵

Our literature review yields 55 papers published from 2002 to 2017. We provide descriptive statistics of these articles in Table 1, including statistics by journal. Overall, these articles account for between 0.5 and 1 percent of the papers published in our sample of journals in any given year,

¹³More broadly, we relate to a number of studies that use reweighting to account for the discrepancy between “what you want” and “what you get” from common estimators. These include Lochner and Moretti (2015), who reweight OLS with IV weights for greater comparability; Sloczynski (2018), who reweights OLS to obtain the ATE; and Stuart et al. (2011) and Andrews and Oster (2019) who propose reweighting experiments to account for selection into participation.

¹⁴ One exception to this is Pages et al. (2019), who suggest that the effect of Head Start may be negative for recent cohorts, although the identifying sample is not discussed.

¹⁵We surveyed: AEJ: Applied Economics, AEJ: Economic Policy, AER, AER P&P, Journal of Health Economics, Journal of Human Resources, Journal of Labor Economics, Journal of Political Economy, Journal of Public Economics, QJE, Review of Economics and Statistics. To identify these articles, we used the search terms “family,” “within-family,” “sibling,” “twin,” “mother,” “father,” “brother,” “sister,” “fixed effect,” “fixed-effect,” and “birthweight” using queries on journal websites. We then searched within articles to see whether FFE was used in the analysis. Finally, we added some additional papers to the list that we are aware of and did not satisfy these search terms. The resulting list is fairly comprehensive, but still likely to be a slight undercount of FFE articles in these journals.

but this varies from 0 to 9 percent of each journal in a given year. The first panel tabulates the frequency of binary treatments and binary outcomes across the sample of papers, the focus of our methodological insights. Nearly two-thirds (35) of the papers have a binary treatment of interest and 23 have a binary outcome. The second and third panels show the varied topics that appear in the sample, spanning health, public, education, and labor fields.

The final panel of the table summarizes the distribution of sample sizes used with FFE. The samples are frequently not limited to families with variation in the treatment variable; therefore, the sample size in the table is an upper bound on the number of observations used for identification. The median number of sibling observations is 6,315, or roughly 85% of the sample in our analysis. We note that there is a high variance in sample size across samples, indicating that there is not a threshold for FFE analyses. The bottom 25% of papers have fewer than 1,200 observations, while the top 25% have over 160,000 sibling observations.

Appendix Figure B.1 illustrates the popularity of this estimation strategy over time. It shows a steady stream of FFE papers over the past 15 years; and that these papers have an impact on the literature, with a mean 233 citations per article (Google Scholar citations as of May 2019). Moreover, since the survey was completed, additional FFE studies have been published (see e.g. Chetty and Hendren (2018*a,b*)). We discuss the prevalence of switcher counts in the next section.

3 Fixed Effects and Selection into Identification

We employ the FE research design to address the concern that Head Start treatment may be correlated with some fixed characteristics of a family that also determine outcomes. For example, the decision to participate in Head Start of siblings is influenced by low parental income — a requirement for eligibility — and availability of an alternative source of care, which may independently influence long-term outcomes. As a result, in our setting – as well as in many other settings – the cross-sectional estimate of the effect of treatment is likely to be biased.

To formalize our setting, let $D_i \in \{0, 1\}$ indicate whether an individual i participates in treatment (e.g. Head Start) and $g(i)$ be the relevant group (e.g. family) for i of the set of groups G in the sample, and let potential outcomes in the untreated and treated states be $Y_i(0)$, $Y_i(1)$, respectively. We observe for each i one outcome, $Y_i = Y_i(D_i)$, treatment, D_i , and group membership, $g(i)$. For brevity, we will frequently write this simply as g . We refer to groups for whom realized $Var(D_i|i \in g(i)) > 0$ as “switchers,” and denote switching status with a binary variable $S_g = 1$, and the set of switchers as $G_S \subseteq G$.

We assume that treatment may be correlated with group characteristics, e.g. mean family income, but is randomly assigned within groups:

Assumption 1 (Group ID Conditional Independence):

$$Y_i(0), Y_i(1) \perp D_i | g(i) = g \tag{1}$$

Assumption 1 encompasses the standard FE specification assumption in linear models. It rules out Roy (1951)-type selection into treatment within groups, in which the probability of receiving treatment is correlated with treatment effects, $Y_i(1) - Y_i(0)$.¹⁶ In the context of Head Start, within-family variation in treatment has been shown to be uncorrelated with most observable characteristics of children (Deming, 2009, GTC, 2002), suggesting the assumption is reasonable.

Under this assumption, estimated treatment effects $\hat{\delta}_{g,FE}$ are an unbiased estimate of group-level treatment effects, $\delta_{g,FE} \equiv \mathbb{E}[Y_i(1) - Y_i(0)|g(i) = g]$. The FE estimate averages $\hat{\delta}_{g,FE}$ for the $g \in G_S$, using weights that are proportional to the within-group variance of D_i and the number of observations in g , n_g (Angrist, 1998; Angrist and Pischke, 2009, eqn. 3.3.7). The FE estimand is:

$$\delta_{FE} = E_g [\delta_{g,FE} \cdot \omega_{g,FE}] \quad (2)$$

where

$$\omega_{g,FE} = \frac{\text{Var}(D_i|g(i) = g, S_g = 1) \cdot n_g \cdot \text{Pr}(S_g = 1|g(i) = g)}{\sum_g [\text{Var}(D_i|g(i) = g, S_g = 1) \cdot n_g \cdot \text{Pr}(S_g = 1|g(i) = g)]}$$

We examine two methodological issues that arise from the FE research design: (i) reduction in identifying variation moving from G to G_S ; and (ii) a change in the composition of the identifying sample due to the fact that $\text{Pr}(S_g = 1|g(i) = g)$ is not likely to be the same for all groups. Issue (i) is well understood in principle, but the degree to which G_S is smaller than G is often underappreciated, and implicitly assumed to be negligible in many theoretical results. Issue (ii) is more novel, and should cause researchers to update the interpretation of the population for which these estimates are relevant. In our literature review, both of these issues are not reported in empirical practice: one-third of papers with a binary independent variable reported the number of switchers, and one paper showed summary statistics by switching status.¹⁷

3.1 Empirical Relevance

To illustrate ideas, we provide an empirical example - for more detail, see Section 6. The sample consists of 2986 white children born in the years 1954-1987. The regression of interest estimates the effect of ever having attended Head Start on a dummy for ever having attended college. The coefficient on Head Start in a cross-section regression is 0.049 (s.e. = 0.044). When mother fixed effects are added, the coefficient becomes 0.120 (s.e. = 0.053). This result indicates that the impact of Head Start participation on college attendance is meaningful in magnitude, and statistically significantly different from zero.

¹⁶Some recent FE strategies explore relaxation of this assumption. For example, in a two-period person-level FE design, Lemieux (1998) estimates union wage returns to both observed and unobserved skills. This approach is extended (with application to farmer adoption of HYV seeds) in Suri (2011) and Verdier and Castro (2019).

¹⁷We also note that the goal of the latter comparison was to show that within-family comparisons were more valid than cross-sectional comparisons, not to discuss external validity. Even if we look at the more recent period after 2010, just 5 out of 20 articles reported the number of switchers, suggesting that researchers still do this rarely.

We illustrate the identifying variation for the FFE regression of some college on Head Start attendance in Panel (a) of Figure 1, which shows a scatterplot of the deviation in Head Start attendance for each individual i from the mean attendance in his or her family, $g(i)$, $\overline{HeadStart_i - HeadStart_{g(i)}}$, against the within-family deviation in attainment of some college for the sample, $\overline{AnyCollege_i - AnyCollege_{g(i)}}$.¹⁸ Strikingly, the largest mass of observations is at (0,0): the majority of families have no variation in Head Start participation and no variation in the college attendance of their children. Individuals in families with no variation in Head Start account for 96% of the sample – removing these leaves us with 213 individuals in switching families.

This reduction in identifying observations could result in a selected sample if switching is correlated with family characteristics. To gain intuition about which variables might determine switching we build a simple model of the Head Start participation decision within families. If the probability of attending Head Start is a constant, π , and independent across siblings in a family, then the probability of switching, $Pr(S_g = 1)$ is simply a function of π and family size, n_g :

$$Pr(S_g = 1) = 1 - (1 - \pi)^{n_g} - \pi^{n_g}$$

According to this formula, the probability of switching has an inverse-U-shaped relationship with π , peaking at $\pi = 0.5$. Further, for a given level of π , the likelihood of being in a switching family is increasing with family size. We illustrate these features in Appendix Figure B.2.

The markers in Figure 2 show the actual probability of attending Head Start and of being in a switching family for each family size by black/white race and by whether the mom has some college or not. As in the stylized model, the likelihood of switching is increasing with family size for each of these subgroups.¹⁹ This could reflect the fact that over time, across children, parents are more likely to be exposed to the program, or are more likely to experience a change in family income, which alters eligibility for the program.

We also observe that switching increases with π , following the inverse-U. The probability of Head Start attendance among black families and families with low-educated moms is much higher and closer to 0.5, compared to white families and families with high-educated moms; and the switching probability is correspondingly larger for black and low-educated families. As a result, the sample used for FFE identification is comprised of 7% of the sibling sample for whites, and 38% of the sibling sample for blacks. Note that while we are focusing on race and maternal education, this notion can be generalized to any other family characteristic, such as SES, that determine π .

This pattern is not unique to the PSID or to Head Start. Panels (b) and (c) of Figure 2 show this relationship using data from two other FFE papers, Collins and Wanamaker (2014) and Deming

¹⁸The size of each symbol is weighted by the number of individuals. A value of 0.5 along the horizontal axis, for example, means that a person went to Head Start in a family where half the children attended Head Start. Values other than 0.5 and -0.5 indicate that the share of children that attended Head Start was different than 0.5; e.g. a value of -0.75 means that a person did not go to Head Start in a family where three quarters of the children did.

¹⁹Appendix Table B.1 shows that this pattern is driven by a much larger incidence of no Head Start participation among smaller families. For example, 78% of 2-child families have no Head Start participants, compared with 48% of families with 5 or more children.

(2009). In both papers, the treatment variable of interest is binary; migration to the North and Head Start participation, respectively. In each of these samples, the probability of being a switcher is increasing in family size.

3.1.1 Selection into Identification Driven by Many Variables

Since SI is likely to affect the balance of characteristics other than family size, we now examine a large number of observable characteristics of switcher families and non-switcher families. Panel A of Table 2 indicates that in addition to having a larger family size, children in switcher families tend to have parents with significantly less education than children in non-switcher families (column 3). These differences in parental education are significant even in a regression framework where we control for differences in family size and the other covariates in the table, though only at the 10 percent level (columns 4 and 5). Family income during preschool of children in switcher families is significantly lower than non-switcher families overall (some of which may have incomes too high to ever qualify for Head Start).²⁰ These patterns are consistent with switching increasing with the probability of Head Start participation.

Next, we examine a one-dimensional summary of how much overlap there is in the characteristics of switchers and non-switchers. We do so by comparing propensity-score-type weights across switchers and non-switchers. This is in the spirit of the literature on using propensity score to extrapolate experimental results to broader populations (Stuart et al., 2011). For comparison we use propensity-score-derived weights, $\frac{Pr(S_{g(i)}=1|\mathbf{X}_{ig})}{Pr(HeadStart_i=1|\mathbf{X}_{ig})}$, which we discuss in more detail in Section 4. These weights give a measure of how aligned the characteristics (vector \mathbf{X}_{ig}) of switchers are with the characteristics of Head Start participants, the population of interest. An average value of 1 implies perfect alignment, while a higher value implies that the characteristics of switchers are over-represented relative to the characteristics of Head Start participants. We estimate the elements of this ratio using a multinomial logit.

Panel B of Table 2 shows that this measure is between 1.9 and 2.9 for the switchers sample, which is 0.3 SD larger than for non-switchers. This indicates that the observables of switchers are not aligned with our population of interest, and that this misalignment is worse for switchers than non-switchers²¹.

3.2 Consequences for Estimation: Effective Number of Identifying Observations

A convenient way to summarize the amount of variation used in FE is by the number of individuals in switching families. However, since not all switchers provide the same amount of identifying

²⁰If we limit ourselves to families with Head Start participants, we continue to find many statistically significant differences (see Appendix Table B.2).

²¹For example, Stuart et al. (2011) suggest that a 0.1 to 0.25 SD difference in propensity scores between the experimental and non-experimental population may be too large to rely on extrapolation without further adjustments.

variation, this can be a misleading measure. For example, a 4-sibling family with 1 treated and 3 untreated individuals has an $\omega_{g,FE}$ that is 25% smaller than the $\omega_{g,FE}$ of a family with 2 treated and 2 untreated ($0.25 \cdot 0.75 = 0.1875 < 0.25 = 0.5 \cdot 0.5$).

We develop a formula for the “effective number of observations,” which captures this idea by (i) quantifying the total amount of identifying variation and (ii) converting this into standardized units (person-equivalents).

$$N_{eff} = \frac{\sum_{g \in G_S} \text{Var}(D_i | g(i) = g) \cdot (n_g - 1)}{\text{Var}(D_{i,reference})} \quad (3)$$

The numerator quantifies the “total amount of variation” identifying δ_{FE} by summing over the variation contributed by each group. This expression is similar to the weight given to each group in the FE estimate, except family size is adjusted for the fact that group-level fixed effects remove one degree of information from each family, $(n_g - 1)$. The denominator provides a translation from “total variation” to “person-equivalents” of variation by normalizing by the variation contributed by an individual observation in a fixed, researcher-determined group, $\text{Var}(D_{i,reference})$.

In our application, we report effective observations using two different reference groups. First, we normalize by the variation across units in a cross section regression after controlling for reasonable g-level covariates, $\text{Var}(D_{i,reference})$. This gives the equivalent number of cross-sectional units that would produce the variation used in the FE regression. Second, we normalize by the within-family variation from units in two-unit groups. The difference between this and the number of individuals in switching families informs us about whether the presence of larger groups contributes additional identifying variation.²²

3.3 Consequences for Estimation: OLS weights and FE weights

Under homogeneous treatment effects ($\delta_g = \delta$), SI has no effect on expected bias in estimation of Equation 2, and the FE estimate trivially is unbiased for the ATE for the sample and the population. There is only a loss of precision that accompanies the overall reduction in sample size.

The more interesting case is when treatment effects are heterogeneous. In that case, SI will lead the FE estimate to provide a biased estimate of the ATE, even if one corrects for the conditional variance weighting of FE among switchers. To be concrete, let Z be a discrete covariate that varies at the group level, such as family size, that determines the magnitude of the effect of treatment. We allow for a different treatment effect for each value of Z : $\delta_g = f(z_g) = \delta_z$, and define \mathbb{Z} as the set of values of z_g present in the samples of siblings and switchers. The treatment effect estimated without FE using a sample of groups with $n_g \geq 2$, e.g. siblings, is:

²² $\text{Var}(D_i | g(i) = g)$ is calculated using the population formula for variance, $\text{Var}(D_i | g(i) = g) = \frac{1}{n_g} \sum_{i \in g} \left(D_i - \frac{\sum_{i \in g} \mathbf{1}(D_i=1)}{n_g} \right)^2$, rather than the sample formula (which would divide by $n_g - 1$).

$$\delta_{OLS} = \sum_{z \in \mathbb{Z}} \delta_{z,OLS} \cdot \omega_{z,OLS} \quad (4)$$

where

$$\omega_{z,OLS} = \frac{(\text{Var}(D_i | n_g \geq 2, z_g = z) \cdot \text{Pr}(z_g = z | n_g \geq 2))}{\sum_{z' \in \mathbb{Z}} (\text{Var}(D_i | n_g \geq 2, z_g = z') \cdot \text{Pr}(z_g = z' | n_g \geq 2))}$$

and $\delta_{z,OLS}$ is the OLS estimate without FE of the treatment effect for groups with $z_g = z$, and $\text{Var}(D_i | n_g \geq 2, z_g = z)$ is the conditional variance of treatment among the sample with $n_g \geq 2$ and $z_g = z$.

The FE estimator for the same sample is:

$$\delta_{FE} = \sum_{z \in \mathbb{Z}} \delta_{z,FE} \cdot \omega_{z,FE} \quad (5)$$

where

$$\omega_{z,FE} = \frac{(\text{Var}(D_i | FE, z_g = z) \cdot \text{Pr}(z_g = z | S_g = 1))}{\sum_{z' \in \mathbb{Z}} (\text{Var}(D_i | FE, z_g = z') \cdot \text{Pr}(z_g = z' | S_g = 1))}$$

and $\delta_{z,FE}$ is the FE estimate of the treatment effect for groups with $z_g = z$, $\text{Var}(D_i | FE, z_g = z)$ is the conditional variance of treatment among the sample for groups with $z_g = z$, net of family fixed effects.

Moving from OLS to FE, the δ 's change and also the ω 's change. The change in the δ 's is how we usually interpret the move from OLS to FE: the change is from “between” (bad) variation to “within” (good) variation. But the full change also incorporates the different weightings of different values of z_g . If the OLS sample and the FE sample overlap in the covariates, we can decompose the difference between OLS and FE to identify how much is caused by the change in weights, ω_z , and how much is driven by the change in identification, δ_z , as:

$$\begin{aligned} \delta_{FE} - \delta_{OLS} &= \sum_{z \in \mathbb{Z}} \underbrace{(\omega_{z,FE} - \omega_{z,OLS}) \cdot (\alpha \cdot \delta_{z,FE} + (1 - \alpha) \cdot \delta_{z,OLS})}_{\text{Impact of } \Delta \text{ weighting}} \\ &+ \sum_{z \in \mathbb{Z}} \underbrace{(\delta_{z,FE} - \delta_{z,OLS}) \cdot (\alpha \cdot \omega_{z,OLS} + (1 - \alpha) \cdot \omega_{z,FE})}_{\text{OLS Bias}} \end{aligned} \quad (6)$$

with $\alpha \in [0, 1]$ a researcher-determined weight. The impact of SI is captured in the first summation of Equation 6, which is a function of the disparity in regression weights ω_z between OLS and FE, multiplied by an α -weighted average of the $\delta_{z,OLS}$ and $\delta_{z,FE}$. Setting $\alpha = 0$ in this term uses cross-section coefficients to assess the importance of changing the regression weights from OLS to FE. Setting $\alpha = 1$ uses the FE coefficients to assess this. If there is important heterogeneity among both w_z and δ_z , these two extremes can provide useful benchmarks to compare against the

OLS and FE estimates, as we do in Section 4.3.²³

Equation 6 shows that SI matters more for the difference between OLS and FE when the gap in weights ($\omega_{z,FE} - \omega_{z,OLS}$) is correlated with treatment effects δ_Z (across elements of Z), and there is heterogeneity in δ_Z . This also raises the issue that the FE estimate may not be unbiased for the ATE. When there are heterogeneous effects, unbiasedness would require the weights $\omega_{z,FE}$ to match with the average population shares across values of Z . This will typically not be the case, as shown theoretically in Chernozhukov et al. (2013); Gibbons, Suarez and Urbancic (2018); Imai and Kim (2019), and as we demonstrate in our empirical application below.

SI impacts the probability of each family size appearing in FE, in addition to potentially impacting the conditional variance of treatment. Existing methods to reweight FE estimates (Gibbons, Suarez and Urbancic, 2018; Imai and Kim, 2019) address only the conditional variance of treatment, and so can at best recover the ATE for switchers. Since switchers are typically not a population of interest, this raises concerns for the external validity of the FE estimator.

Illustration of Consequences: Greater Returns to Head Start in Larger Families

We use data from our empirical example to illustrate the change in the components of ω_z across OLS and FE. Panel A of Table 3 shows the proportion of the overall sample, the sibling sample, and the switcher sample comprised by families with 1, 2, 3, 4, or 5 or more kids. Moving from the sibling sample to the switchers sample, the share of the sample comprised of families with 5 or more kids roughly doubles (from 16.9% to 32%), while the share of the sample comprised of 2-child families declines by 40% (from 34.5% to 21%).

Panel B of Table 3 shows the variance in Head Start across the same family sizes and samples, netting out group fixed effects for the switcher sample. Going from the siblings sample to the switcher sample, the variance approximately doubles for all family sizes. This suggests that the change in the conditional variance across OLS and FE plays a relatively minor role in shifting weights across groups in our setting.²⁴

We then calculate $\omega_{z,OLS}$ and $\omega_{z,FE}$ in Panel C. Going from the sibling sample to the switchers sample, $\omega_{2-child}$ declines by over 25% and $\omega_{3-child}$ declines by 15%. Conversely, $\omega_{5-child}$ nearly doubles from 0.134 to 0.243, and $\omega_{4-child}$ families increases by over 25%.

The effect of Head Start also varies by family size in our applications. The first two columns of Panel A of Table 4 shows the estimated effects of Head Start on the likelihood of completing some college by the number of children in a family for our illustrative sample. We show the results with and without family fixed effects. In both specifications, the effect of Head Start is significantly

²³This decomposition is similar in form to Equation 13 in Loken, Mogstad and Wiswall (2012), which uses $\alpha = 1/2$. However, we sum over a group-level covariate that is distinct from the treatment of interest, while Loken, Mogstad and Wiswall (2012) sum over values of an individual covariate (that varies within families), which is also the treatment of interest.

²⁴We provide additional evidence that “undoing” the conditional variance weighting makes little difference in this application in Section 6.

higher among white children in families with 5 or more children and, once fixed effects are added, the effect of Head Start is monotonically increasing with the number of children in a family.

One possible explanation for this heterogeneity is that children with higher initial endowments receive greater parental investments in larger families, and also benefit more from Head Start (Aizer and Cunha, 2012). Another possibility is that Head Start substitutes for parental time, which is more scarce in larger families. Another interpretation is that this heterogeneity reflects the fact that other covariates correlated with family size, such as income, mediate the impacts of Head Start. This final explanation seems less likely, as we find that the heterogeneity in family size survives the inclusion of other interactions, as we discuss in Section 6.

The bottom of Panel A shows the number of Head Start switcher observations and effective observations in terms of cross-sectional and two-sibling switcher individuals.²⁵ It shows that a total of 213 individuals are used to identify these coefficients, less than one tenth of the total sample, and that the variation is equivalent to 236 individuals in 2-person switching families. Hence, by including families with three or more children, on average, each observation is providing more variation than in a similar-sized sample of 2-child families. Further, the variation is equivalent to 731 individuals in a cross-sectional regression. This is because there is relatively little variation in Head Start in the full sample.

Like with SI, the larger Head Start effects we document for big families is not specific to the PSID. Columns (3) to (5) of Table 4 show the CNLSY FFE estimated effects of Head Start by family size for idleness, having a learning disability, and being in poor health.²⁶ For each of these outcomes, the impact of Head Start for 5+ child families is at least twice as large as the impact for 2 or 3 child families. For high school graduation, we also see a large impact for 4-child families, roughly double the impact for 2 and 3 child families. This implies that we should expect an increase in the coefficient going from OLS to fixed effects due to the *change in weighting* across the identifying samples, even without a change in the source of identification.

The number of switchers in the CNLSY sample is 581, less than half of the total number of observations. As in the PSID, the variation in this sample is equivalent to a larger sample of 2-person families (648 individuals.) However, the corresponding cross-sectional observations is smaller (438.7). These two examples illustrate that there are multiple forces driving the effective number of observations calculation: lost information from the group FEs drives down variation; but moving toward larger conditional variance of treatment increases variation. In the PSID example the second effect dominates; in the CNLSY case the first effect dominates.

²⁵For effective cross-sectional individuals, the denominator of Equation 3 is the variance of Head Start, residualized by the family mean of the covariates in the analysis. For the effective number of two-person switcher individuals, the denominator is $[V(D_i|g) \cdot (n_g - 1)] / n_g = [0.5^2 \cdot (2 - 1)] / 2 = 0.125$.

²⁶We focus on these outcomes because individuals that attended Head Start were found to fair significantly better on each of these outcomes in Deming (2009).

4 Extrapolating from Identifying to Target Population

The difference between OLS and FE in implicit weighting of heterogeneous treatment effects leads us to consider translating the FE estimates into an ATE for a (researcher-determined) population of interest. We propose a method to flexibly obtain the ATE for such populations of interest, which we refer to as “target” populations, and denote by an indicator T_g . Commonly, the target population in applied work is the ATE for a nationally representative sample, which may be a reasonable starting place for most researchers. For some treatments, like means-tested programs, one might be interested in the ATE for eligible families, or families with a participating member. Alternatively, one can generate estimates for two target populations, such as urban and rural participants, and compare these to get a sense of differences in the efficacy of the treatment in different settings.

4.1 Assumptions and Proposition

The methods rely on four key assumptions. which are variants of those used for extrapolation from IV (Angrist and Fernandez-Val (2013); Aronow and Carnegie (2013)). First, we assume that Group ID conditional independence (Assumption 1, Equation 1) holds.

Assumption 2 (Conditional Fixed Effect Ignorability (CFEI)):

$$\mathbb{E}[Y_i(1) - Y_i(0)|S_g, P_x, Q_x] = \mathbb{E}[Y_i(1) - Y_i(0)|P_x, Q_x] \quad (7)$$

$$\mathbb{E}[Y_i(1) - Y_i(0)|T_g, P_x, Q_x] = \mathbb{E}[Y_i(1) - Y_i(0)|P_x, Q_x] \quad (8)$$

Second, we assume that conditional on observables, the true treatment effect is independent of a group’s switching or target status, which we refer to as conditional fixed effect ignorability (CFEI).²⁷ We use two propensity scores constructed from the vector of observable group characteristics, \mathbf{X}_g , as the conditioning variables: $P_x := Pr[S_g = 1|\mathbf{X}_g = \mathbf{x}]$ is the propensity to be a switching group, and $Q_x := [T_g = 1|\mathbf{X}_g = \mathbf{x}]$ is the propensity to be in the (researcher-determined) target group.

Assumption 2 is a strong assumption. CFEI eliminates, for example, a second type of Roy (1951)-type selection, whereby switchers have an unobserved quality that increases the effectiveness of treatment compared to observationally equivalent non-switchers. This assumption is required in order to extrapolate treatment effects from switchers to non-switchers. Analogous assumptions are made in the IV literature extrapolating from compliers to a target population (Angrist and Fernandez-Val, 2013; Aronow and Carnegie, 2013). The plausibility of Assumption 2 will depend on the context of the specific application.

In our application this assumption may be reasonable, because observables strongly predict both switching and target status. We have shown that the key determinants of $Pr[S_g = 1]$ are

²⁷We considered using CoFEfe as the acronym for this assumption. This would provide a novel candidate interpretation of US President Donald Trump’s enigmatic tweet of May 31, 2017.

family size and the underlying probability of Head Start participation. Family size is observable, and observable covariates, such as family income, can take us a long way in predicting program participation. Likewise, the family-level determinants of $Pr[T_g = 1]$ for a target such as Head Start participants will be largely tied to observable eligibility requirements for the program, such as income and household size, which together determine the income-to-poverty ratio.

However, there may be unobservables that determine switching or target status, and which are also correlated with treatment effects. To address this concern, we develop two testable implications of Assumption 2 at the end of this subsection. We implement those tests in Section 6.2.1 and Appendix Tables B.10, B.11 and B.12.

In contexts where this assumption may be too strong, the reweighting procedure can still help to reduce bias - balancing on observables will likely reduce the degree of bias that is driven by mismatch on unobservables. Andrews and Oster (2019) develop this idea formally in the context of extrapolating from an experiment to a larger target population. This improvement relies on the assumption that the direction of bias from selection on observables is the same as the direction of bias from selection on unobservables, which may be plausible in many cases.²⁸

Assumption 3 (Correct Propensity Score Specification):

$$Pr(S_g = 1|\mathbf{X}_g) = F(\theta_g; \mathbf{X}_g) \tag{9}$$

$$Pr(T_g = 1|\mathbf{X}_g) = G(\chi_g; \mathbf{X}_g) \tag{10}$$

Third, we assume that the propensity scores that we estimate have the correct functional form, with $F(\cdot)$ and $G(\cdot)$ known, and θ_g and χ_g parameters to be estimated. In our application, we model $F(\cdot)$ and $G(\cdot)$ as a multinomial logit.

Assumption 4 (Overlap in P_x):

$$\text{If } Q_x > 0, \text{ then } P_x > 0, \forall \mathbf{x} \in \text{Support}(\mathbf{X}_g) \tag{11}$$

Fourth, we require a positive probability of being a switcher for each value of \mathbf{X}_g in the target group, which ensures that we can use the switcher sample to recover the distribution of treatment effects in the target sample. Since some covariate values may not be observed in the switcher sample, this assumption implicitly places some restrictions on the relationship between treatment effects and these covariates. For example, since we do not observe 1-unit groups (“singletons”) in the switcher sample – they are “never-switchers” – but they are present in some of our application target populations,²⁹ we cannot allow treatment effects for singletons to be outside the support of

²⁸For example, we show that larger families are likely to have more variation and larger treatment effects. Such families may also be more motivated (or better equipped) to seek out information about public programs, which would also predict greater variation and larger treatment effects.

²⁹Singletons comprise 6% of Head Start participants in the CNLSY, and 18% of Head Start participants in the PSID.

the treatment effects of switchers. This also precludes us from including an indicator for singletons in \mathbf{X}_g .

In our application we preserve overlap by assuming that treatment effects are the same for all groups with 2 or fewer units, which allows us to extrapolate treatment effects for singletons from 2-unit groups.³⁰ An alternative approach is to extrapolate treatment effects for never-switchers using functional form assumptions, which we discuss under extensions below.

Proposition 1. *Define the re-weighted FE estimator for target population t as*

$$\widehat{\delta}^t := \frac{1}{\sum_i \mathbf{1}(S_{g(i)} = 1)} \sum_{i|S_{g(i)}=1} \widehat{w}_{g(i)}^t \cdot \widehat{\delta}_{g,FE}, \quad (12)$$

with $\widehat{w}_{g(i)}^t$ our estimate of $w_{g(i)}^t$,

$$w_{g(i)}^t := \frac{Q_x \cdot Pr[S_{g(i)} = 1]}{P_x \cdot Pr[T_{g(i)} = 1]} \quad (13)$$

Under Assumptions 1 through 4, $\widehat{\delta}^t$ is consistent for the ATE of the target population, $\mathbb{E}[Y(1) - Y(0)|T_g = 1]$.

The proof is in Appendix A. Intuitively, the weights are increasing in Q_x and decreasing in P_x , such that we upweight observations that are more similar to the target, and downweight observations that are overrepresented in the switching population. The treatment estimate for each switcher group g is weighted proportionately to match the share of the target population with observable characteristics matching g , which gives the ATE under the assumptions above.³¹

Testable Implications

CFEI requires that treatment effects should be balanced across T_g , conditional on P_x and Q_x (Equation 8). Hence, switchers in the target population should on average have the same treatment effects as switchers in the non-target population once observables are balanced across the two groups. This is potentially testable if some switchers are not in the target population (i.e. if the target is not “everyone,” “multi-unit groups,” or otherwise contains the set of switchers G_S .) For instance, if the target population is families that participate in a safety net program, groups that live in rural areas, or firms that are in a particular industry, we could test whether treatment effects are different for participants and non-participants, rural and urban individuals, etc.

CFEI also requires that treatment effects should not vary across switchers and non-switchers (Equation 7). This could be violated, for example, if families with a lower return to Head Start

³⁰We implement this by including an indicator for “1 or 2 child families” in \mathbf{X}_g together with indicators for other family sizes. Alternatively, the target group can be defined to only include families that are ever switchers, such as “siblings” or “multi-child Head Start families.”

³¹See Appendix A for a simple derivation of the weights.

make less of an effort to have any of their children attend, or if families with a higher return to Head Start make a greater effort to have all of their children attend.

Because we never observe the treatment effects of non-switchers, we can not directly test for this type of violation. However, this mechanism generates a testable implication: if families systematically choose the number of children to attend Head Start based on private information, we should also expect to see this pattern within switchers. In particular, switcher families that have a higher fraction of children that attend Head Start should be expected to experience systematically larger (or smaller) gains. Thus, a useful diagnostic test is whether treatment effects vary with the number or fraction of treated units.

We implement both of these tests for our setting in Section 6.2.1.

4.2 Reweighting Methodology

In order to implement this reweighting strategy, we first need to obtain estimates of P_x and Q_x . When the identifying sample is a subset of the target population, $Q_x = 1$, and P_x can be estimated by a logit or probit model. Otherwise, these two elements can be calculated by using a multinomial logit model to estimate the probability of each of the four possible combinations of having $S_g = 1/S_g = 0$ and $T_g = 1/T_g = 0$. \hat{Q}_x is then constructed as the sum of the predicted $Pr(T_g = 1, S_g = 0)$ and the predicted $Pr(T_g = 1, S_g = 1)$ for each unit; and \hat{P}_x is constructed as the sum of the predicted $Pr(T_g = 1, S_g = 1)$ and the predicted $Pr(T_g = 0, S_g = 1)$ for each unit.³²

With these weights in place, the ATE for the target population can be estimated in one of two ways. The first is a two-step “post-regression weighting” of $\hat{\delta}_g$, where $\hat{\delta}_g$ is estimated from a regression of the outcome on interactions between D_i and group-specific dummies. Then aggregate $\widehat{w_{g(i)}^t}$ to the group-level and perform a normalization to obtain the final estimation weights, $\hat{s}_g^t = \frac{\widehat{w_{g(i)}^t \cdot n_g}}{\sum_{g' \in G_S} \widehat{w_{g(i)'}^t \cdot n_{g'}}} = \frac{\hat{Q}_x \cdot n_g}{\sum_{g' \in G_S} \hat{P}_x \cdot n_{g'}}$. The 2-step ATE combines these using:

$$\widehat{\delta_{2step}^t} = \sum_{g \in G_S} \hat{s}_g^t \cdot \hat{\delta}_g \quad (14)$$

A second approach is to obtain the ATE in a single step using “in-regression weights.” For this, we need to adjust for the fact that the FE estimator uses weights ω_{FE} rather than population shares. We address this by incorporating inverse conditional variance weights, as $v_g = (Var(D_i | g(i) = g))^{-1}$ (Gibbons, Suarez and Urbancic, 2018).³³ Then, the ATE can be estimated by $\widehat{\delta_{1step}^t}$ from a one-step regression using $\widehat{w_{g(i)}^t} \cdot v_g$ as regression weights.

For both of these methods, we obtain standard errors by bootstrapping, which accounts for the uncertainty in the estimation of the propensity scores.

³² $\widehat{w_{g(i)}^t}$ can also be multiplied by survey weights, as we do in our PSID example.

³³This variance is computed using the population formula, (dividing by n_g), rather than the sample formula (dividing by $n_g - 1$). As with the two-step estimator, these weights can also be multiplied by sample weights.

4.3 Special Case: Univariate Heterogeneity

If the source of heterogeneity in estimates is a single, discrete covariate, we can obtain further insight from performing the decomposition captured in Equation 6. Taking the OLS family-size-specific coefficients from column (1) of Table 4 and reweighting by the fixed-effects regression weights ($\alpha = 0$ in Eq. 6), we obtain a weighted coefficient of 0.069, shown in the bottom row of Table 4. This implies that approximately 1/3 of the change from OLS to FE ($\frac{0.069-0.049}{0.12-0.049}$) is driven by the change in family size weights; with the other 2/3 driven by change in identifying variation. Further, reweighting the FE estimates using the OLS weights ($\alpha = 1$ in Eq. 6) produces a coefficient is 0.083. This implies that the imbalance in family size alone causes the FFE estimate to be 50% higher than the estimates without FE.

4.4 Monte Carlo Experiments

We perform a Monte Carlo analysis to examine the properties of our proposed reweighting estimators. We use naturally occurring selection into identification from our PSID application and model treatment effects for three settings, allowing the true ATE to be known. Each setting has a different model of heterogeneity in treatment (assumed to be known to the researcher), which determines the covariates that we use to generate the propensity score.

We generate the data for the Monte Carlo as follows. We construct an untreated outcome for each individual by running a linear probability model of attainment of “some college or more” on demographic variables, income during childhood, and parental education. We then use the estimates from this model to generate a prediction of a continuous probability that an individual completes some college, X_{ig} .³⁴ All simulations start with this constructed variable X_{ig} and the variable $HeadStart_{ig}$ from the original data.

We then construct latent outcomes inclusive of treatment as $Y_{ig}^* = X_{ig} + \beta_{ig}HeadStart_{ig}$, where β_{ig} is the treatment effect of Head Start and varies across different models of heterogeneity. We always scale Y_{ig}^* to ensure that these probabilities lie within the range $[0, 1]$. We then randomly generate the binary outcome variable as $Pr(Y_{ig} = 1) = Y_{ig}^*$.

We consider three models of heterogeneity in treatment effects. First, $\beta_{ig} = 0.08$. We use the variable X_{ig} to generate propensity scores. Second, $\beta_{ig} = 0.192$ for large families (with 4 or more siblings) and $\beta_{ig} = 0$ for small families (3 or fewer children). We use a dummy variable for “large family” to generate propensity scores. Third, we allow the treatment effect heterogeneity to vary smoothly: $\beta_{ig} = 0.08 \cdot \left(1 - \frac{X_{ig} - \bar{X}_{ig}}{s.d.(X_{ig})}\right) \cdot \frac{1}{3}$, with \bar{X}_{ig} and $s.d.(X_{ig})$ the sample mean and standard deviation of X_{ig} . This produces a treatment effect that is larger for lower-baseline-probability individuals and ranges from 0.01 to 0.15 for most of the population. For this more complex treatment effect, we generate propensity scores in two ways: using X_{ig} and, more flexibly, using a spline in X_{ig} , with knots at the 5th, 20th, 50th, 80th, and 95th percentiles of X_{ig} . The latter

³⁴For simplicity, we restrict the sample to those with $X_{ig} \in [0, 1]$ at baseline.

model presumes that the researcher has some intuition that the treatment effect or selection into identification may vary non-linearly with baseline outcomes.

We run 3,000 replications of our Monte Carlo simulation. In each replication, we keep track of the true ATE for each target population of interest, the FE estimate of the treatment effect, and the reweighted regression estimate of the treatment effect for each target population.³⁵ The FE estimate is the same for all target populations. We consider four target populations: (i) individuals in Head Start switching families;³⁶ (ii) all siblings; (iii) all individuals in the sample (including singletons); and (iv) all Head Start participants. We multiply all estimates by 1,000 for easier readability.

Panel A of Table 5 presents results for the model with constant treatment effects. In this setting, the average treatment effect is the same for all target populations, all estimators are unbiased, and the FE model is the minimum variance estimator. The reweighting estimators have mean squared errors 3 to 20% larger than for OLS.

Panel B of Table 5 presents results for the model with zero treatment effect for small families, and large treatment effects for large (4+ children) families. It shows that for every target population, FE is biased, while the reweighting estimator is always unbiased. This improvement in bias over FE leads to much better mean squared error results for the reweighting estimator.³⁷

Panels C and D of Table 5 examine the third model with heterogeneous treatment effect that varies with X_{ig} . Here the FE model has relatively little bias for the switcher and Head Start participant targets (-0.2 p.p. and -0.08 p.p. on a base of 9 p.p.), but has much larger bias for the remaining targets. Panel C shows that the regression reweighting estimator which uses X_{ig} in the propensity score estimation has less bias than FE for all target populations, with no detectable bias for the switcher, or Head Start populations. The small bias for the reweighting estimator for the other target populations results from an imperfect balance in the X_{ig} variable, even after reweighting.³⁸

Panel D shows that when we re-estimate the model including a spline in X_{ig} in the propensity score estimation, the reweighting estimator has no detectable bias for any of the target groups. This suggests that allowing for greater flexibility in the functional form relationship between covariates and the propensity score can achieve greater reductions in bias.

Overall, the results of this exercise show that that the reweighted estimator has significantly less bias than FE for the types of treatment effect heterogeneity we consider, and can be successfully targeted toward different target populations. Consistent with the conditioning on observables

³⁵Both post-regression and in-regression reweighting produce the same results.

³⁶This will not necessarily be the same as the FE estimate because of differences in the conditional variance across families.

³⁷In results not reported, we have examined adding X_{ig} as a covariate to the propensity score estimation stage in this model. This introduces a small amount of bias in the reweighting estimator (-0.1 p.p., relative to the 2 to 3 p.p. bias in FE) for the “siblings” and “all” target groups.

³⁸This is because $Pr(S_g = 1)$ is misspecified as a linear function of X_{ig} , which causes us to misassign the weight for each treatment effect.

requirements of this estimator, its performance is best when it is given the appropriate covariates for the particular type of heterogeneity, and when the model for the probability of switching is correctly specified.

5 Extensions

5.1 Projecting Treatment Effects for “Never-Switchers”

As noted above, the reweighting estimator in Proposition 1 only recovers the ATE for the target population if (i) the target does not include never-switchers or (ii) if the treatment effects for never-switchers in the target can be assumed to be the same as some other target groups with $P_x > 0$. Otherwise, the reweighting estimator only obtains the ATE for the subset of the target with $P_x > 0$, for whom treatment effects are identified.

A slight variant of (ii) above which could also enable recovery of the full target ATE is to extrapolate treatment effects for never-switchers. This requires a stronger form of CFEI: that treatment effects are not only a function of observable characteristics, but that the researcher can correctly specify the functional form of this relationship.³⁹ This assumption may not be warranted if heterogeneity is primarily driven by unobservable characteristics; in cases where there is support in the data for such a relationship with observed covariates (e.g. increasing effects with group size), this may be a reasonable way to proceed. The weighted average of estimated treatment effects for $P_x > 0$ and extrapolated effects for $P_x = 0$ gives the ATE for the target group.

5.2 Unit i Covariates

We now consider FE models that include covariates C_i that vary across i units within a group. Researchers may want to include C_i in their models in order to (i) make Assumption 1 more reasonable; (ii) improve precision of estimates (iii) allow extrapolation to target groups defined at the unit level.

Once these covariates are included, the typical intuition that “groups with variation in treatment” provide identification breaks down. This is because for some groups, who we refer to as “residual switchers,” there can be variation in the treatment residualized of C_i , even if there is no within-group variation in D_i .⁴⁰ Thus, treatment effects can also be estimated for residual switchers; however, identifying variation comes from within-family variation in C_i , not D_i .

How much do residual switchers matter for estimates? We can quantify this by calculating the share of variation in D_i coming from residual switchers, using a formula similar to the calculation of the effective number of observations.⁴¹ In our PSID application, residual switchers provide 3% of

³⁹See Appendix A.1.1 for a formalization of this assumption and an extension of Proposition 1 using extrapolation.

⁴⁰See Appendix A.3 for a formalization of this.

⁴¹In particular, the share of identification from residual switchers is equal to $1 - \frac{\sum_g \text{Var}(D_i | C_i, g(i)=g) \cdot (n_g - 1) \cdot \mathbf{1}(S_g=1)}{\sum_g \text{Var}(D_i | C_i, g(i)=g) \cdot (n_g - 1)}$. Alternatively, calculating the effective number of observations using Equation 3 (altering the variance to condition

the variation used for identification of the Head Start FFE coefficient. Therefore, this contributes minimally to the FE estimate.

We can also consider incorporating residual switchers into the reweighting methods. For a general discussion of how our key assumptions and proposition can be extended to accommodate C_i , see Appendix A.3. The decision to include residual switchers can vary across contexts, and should depend on the extent to which variation from residual switchers is valid for identifying treatment effects. For example, in our application, residual switchers are primarily families where no children attended Head Start. As a result, we believe that variation from these families is not aligned with our desired thought experiment, which leads us to ignore these families in reweighting. In contrast, in a setting like difference-in-difference, untreated “residual switchers” can provide equal identifying variation as the switchers, which makes it appropriate to include variation from all groups.

5.3 Nonlinear Functional Form

Next, we relax the linear functional form assumption used to demonstrate SI in our Monte Carlo simulations. One reason this may make a difference is that conditional or fixed effect logit models use only “double switchers,” families with variation in both the outcome variable and the treatment variable, rather than “switchers”. In Appendix E, we show that the biases from SI are similar in the linear probability model and conditional logit, and that the reweighting we propose is equally effective at reducing bias in both cases.

5.4 Continuous D_i

Finally, while we have focused on the case where D_i is binary, it is worth noting that SI can also be present when D_i is continuous (since $\hat{\delta}_{g,FE}$ is still only estimated for switching families.) It is not clear how frequently this will manifest in practice, however, since groups are more likely to have variation in a continuous covariate. Even so, it may still be worthwhile to verify the number of switchers, since there may be persistent bunching at one value of D_i , such as at zero maternal income or at zero instances of an uncommon event.

6 Effects of Head Start

6.1 Data and Replication of GTC and Deming (2009)

We now turn to examining the impact of Head Start on long run outcomes using the PSID and CNLSY, which were used to analyze this question in GTC and Deming (2009).

on C_i) would produce similar results

6.1.1 PSID

The PSID sample includes the sample of individuals surveyed in the PSID by 2011. The PSID began in 1968 as a survey of roughly 5,000 households and has followed the members of these founding households and their children longitudinally. The longitudinal nature of the study allows sibling comparisons during early adulthood as well as later in life.

We begin our analysis with a replication of GTC. The sample includes all black or white individuals born between 1966 and 1977, and excludes Hispanic individuals. We provide a detailed description of our replication of GTC in Appendix D. Despite some minor differences, the two PSID samples are qualitatively similar. The summary statistics are often within a third of a standard deviation of each other. Moreover, the estimated effects of Head Start in this sample are similar to those estimated in GTC. We find large (23 p.p.) and significant effects of Head Start on the probability that whites attain some college, and large point estimates (9.3 p.p.) for high school graduation, though in our case these are not statistically significant. We do not find that Head Start meaningfully reduces the probability of committing a crime.⁴²

For the remaining analyses from here, we use a sample that substantially expands and modifies the GTC sample. First, we expand the sample to include individuals born between 1978 and 1987. The individuals in these cohorts were too young when the analysis in GTC was performed to observe their education and early career outcomes. Second, we include older siblings of all individuals, including those born prior to 1966. These early cohorts were typically too old to benefit from the introduction of Head Start, and serve as a plausible control group for the early cohorts.

In addition to modifications of the sample, we also expand the number of outcomes under analysis in order to gain a more extensive understanding of the channels by which Head Start affects children’s lives. We follow the established practice of distilling the measures to summary indices to lessen problems with multiple hypothesis testing (see, e.g., Anderson, 2008; Kling, Liebman and Katz, 2007; Hoynes, Schanzenbach and Almond, 2016). We create four indices to capture economic and health outcomes observed for individuals at age 30 and 40. The “economic sufficiency index” includes measures of educational attainment, receipt of AFDC/TANF, food stamps, mean earnings, mean family income relative to the poverty threshold, the fraction of years with positive earnings, the fraction of years that the individual did not report an unemployment spell, and homeownership. The “good health index” summarizes the following component measures: non-smoking, report of good health, and negative of mean BMI.⁴³

The process of creating each index follows the procedure described in Kling, Liebman and Katz (2007). In particular, we standardize each component of the index by subtracting the mean outcome for non-treated children, defined as children that did not attend any form of preschool,

⁴²In some subsamples, we even find an effect in the opposite direction. We believe these cases are driven by situations where there are rather few observations identifying the coefficients, and that the lack of correspondence may be driven by very minor (and un-diagnosable) differences in specification and/or dataset construction.

⁴³ See Appendix Table B.5 for descriptive statistics of the inputs to the indices.

and then dividing the result by the standard deviation of the outcome for non-treated children. The summary index takes a mean of these standardized measures.⁴⁴ We also extract the first principal component of the standardized variables for “economic sufficiency” and for “good health”. Later we use these as alternative outcome variables.

Appendix Table B.3 reports sample descriptive statistics for the expanded sample we construct. For ease of comparison with our earlier replication, we include means for the entire sample, the subsamples of Head Start participants/non-participants, and for the sample of individuals with siblings. We present the means of the analyzed outcomes in Appendix Table B.4.⁴⁵

6.1.2 CNLSY

We obtain the CNLSY sample from the Deming (2009) replication files, which ensures that the samples are identical. The CNLSY is a longitudinal survey that follows the children born to the roughly 6,000 women that took part in the NLSY79 survey. The sample we use includes all children who were at least 4 years old by 1990.

6.2 Head Start Estimation

The empirical strategy takes advantage of within-family variation in participation in Head Start to identify the long term impact of the program. Following GTC and Deming (2009), we estimate:

$$Y_{ig} = \alpha + \beta_1 \text{HeadStart}_{ig} + \beta_2 \text{OtherPreSchool}_{ig} + \mathbf{X}_{ig}\gamma + \delta_g + \varepsilon_{ig} \quad (15)$$

where Y_{ig} represents a long-term outcome for individual i with mother g . HeadStart_{ig} indicates whether a child reports participation in the program, and $\text{OtherPreSchool}_{ig}$ indicates participation in other preschool (and no participation in Head Start). These two variables are in this way defined so as to be mutually exclusive, with “neither Head Start nor other preschool” as the omitted category.⁴⁶ δ_g is a mother fixed effect which enables comparisons across siblings with a shared mother. The vector \mathbf{X}_{ig} includes a large number of controls for individual and family characteristics to absorb differences in personal and household characteristics which may be correlated with one’s participation in Head Start and long term outcomes. These controls vary due to data availability across sources and specification used in earlier work, but fall into three broad categories: demographics, family background, and family economic circumstances during early childhood.⁴⁷

⁴⁴Consistent with Kling, Liebman and Katz (2007), we generate a summary index for any individual for whom we observe a response for one component of the index. Missing components of the index are imputed as the mean of the outcome conditional on treatment status. For example, if a former Head Start participant is missing an outcome, it is imputed as the mean outcome of other Head Start participants. Likewise for other preschool, or non-preschool participants.

⁴⁵Appendix Table B.5 includes summary statistics for the inputs to the summary indices. Appendix Tables B.6, B.7, and B.8 contain the number of observations for each outcome and control variable in the analysis .

⁴⁶Since Head Start only became available in 1965, we recode Head Start attendance to be “other preschool” for the 1961 and older cohorts.

⁴⁷For the PSID, these include: individual’s year of birth, sex, race, and an indicator for being low birth weight,

Missing control variables are imputed at the mean, and we include an indicator variable for these imputed observations. We cluster standard errors on mother id, and use population-representative weights where appropriate.⁴⁸ When Y_{ig} is binary, we estimate linear probability models as a main specification and check the sensitivity of our results to alternative models.

The coefficient of interest is β_1 , the impact of Head Start on long term outcomes compared to no preschool. We generate propensity score weights to obtain the ATE for three target populations: (1) Head-Start-eligible individuals, based on family income between ages 2 and 5;⁴⁹ (2) all Head Start participants; and (3) all siblings. For parsimony, we use a subset of the variables in Table 2 to generate the propensity score for each race: year of birth, gender, mother’s years of education, income at age 3, and income at age 4, and indicators for family size (grouping together 1 and 2 child families).⁵⁰ We include results for the post-regression weighting method; results are qualitatively similar when we use in-regression weighting.

6.2.1 Evidence on Model Assumptions: Identifying and Conditional Ignorability

The standard test of the identifying assumption (Assumption 1) is to look for balance in observables across siblings within families. Deming (2009) finds little evidence that Head Start attendance is correlated with observable differences across siblings, which suggests that the magnitude of selection may be small. In Appendix Table B.9, we examine the plausibility of the identifying assumption in the PSID by testing the correlation between participation in Head Start and observable pre-Head Start individual and family characteristics. For the white sample which forms our focus, there are few statistically significant correlations, which suggest that the assumption may be reasonable.⁵¹

As a first test of CFEI (Assumption 2), Appendix Table B.10 analyzes whether treatment effects vary across observationally comparable families in the target and non-target populations. Recall that CFEI implies that conditional on P_X and Q_x , treatment effects should be independent of observable characteristics, including membership in the target population. To implement this test, we regress estimated family-specific treatment effects on an indicator for whether an individual is a member of the target population, employing traditional inverse propensity score weights to

mother and father’s years of education, an indicator for having a single mother at age 4, 4-knot splines in annual family income for each age 0, 1, and 2, a fourth spline based on average family income between ages 3 and 6, indicators for mother’s employment status at ages 0, 1, and 2, and household size at age 4. For the CNLSY, these include: health conditions before age 5, PPVT test score at age 3, measures of birth weight, measures of mother’s health and health behaviors, mother’s working behavior and income prior to age 4, indicator for being first born, participation in Medicaid, relative care, and indicators for early care types.

⁴⁸We follow our predecessors’ weighting practices: for the PSID, we generate representative population weights from the 1995 March CPS, and for the CNLSY do not use weights.

⁴⁹An individual is considered Head-Start-eligible if at any point between the ages of 2 and 5 her family income was below 150% of the poverty level, to account for our imperfect ability to observe reportable income.

⁵⁰Results are similar when we substitute family size indicators with linear and quadratic terms in family size.

⁵¹For the black sample, participation in Head Start is correlated with a greater likelihood of having higher income at age 1, and lower income at age 2, which may raise concerns that black families may tend to send their children to Head Start after a rupture in the family or after an income shock. However, given the many hypotheses being tested in this table, these significant findings might be spurious; and these results are somewhat sensitive, becoming insignificant when we drop observations with imputed controls.

balance observables (and therefore treatment effects, if CFEI holds) between target and non-target switchers.⁵² When target status varies within family, another way of viewing this test is as asking: is the treatment effect related to the share of individuals in the family that are in the target? This test passes with no sign of systematic differences across target and non-target individuals across all outcomes, two datasets, and two different target populations (Head-Start-eligible and Head-Start-participants.)⁵³

As a second, related, test, we examine whether the impacts of Head Start vary with the fraction of children that attend Head Start in a family. Appendix Tables B.11 and B.12 show the results for the PSID and CNLSY samples, respectively. In the PSID, we find no significant interaction with the fraction of children that attend Head Start. In the CNLSY, we find that 5-child families that have more children attend Head Start have smaller improvements in high school and larger improvements in learning disabilities. Averaging across all family sizes, however, we find no significant relationship between the fraction that attend Head Start and treatment effects. This lends further support to our assumption that it is appropriate to extrapolate from switchers to non-switchers.⁵⁴

Our reweighting procedure also relies on adequate overlap of P_x across switchers and individuals in the target population in the non-switching sample. In Appendix Figure B.3 we show the density of the estimated probabilities of being a Head Start participant for the switching sample and the non-switching Head Start participant sample. This figure shows that there is a good deal of overlap across the two groups, but also that there are a few Head Start participants whose p-scores lie outside the range of the switchers. These observations represent 4 individuals, 5% of the Head Start non-switcher observations, and 2% of all Head Start participants. We interpret this magnitude of violation of the overlap assumption as mild enough to disregard in our subsequent analysis.⁵⁵

6.3 Head Start Results

6.3.1 Reweighted Estimates

We begin by presenting results for our illustrative outcome, attainment of some college for whites in the PSID, in Panel A of Table 6. Column (1) of the table presents the estimated impact of Head Start on some college in GTC, column (2) presents the results using our expanded sample, and columns (3) to (5) present reweighted estimates for the three target populations. As reported earlier, we estimate that Head Start increases the likelihood of attaining some college by a statistically

⁵²For target individuals the weights are $1/Pr[T_i = 1, S_g = 1|X_{ig}]$, and for non-target individuals the weights are $1/Pr[T_i = 0, S_g = 1|X_{ig}]$.

⁵³We have run analogous models at the family level, which give qualitatively similar results.

⁵⁴If we take the negative slope on the fraction attending at face value (and disregard the noise in the estimates), it suggests that for the target population of HS participants, our treatment effects might be too large. In contrast, for a target population that has many non-switcher non-participants, our estimated treatment effect might be too small. If there are systematic differences, one approach to correct for this would be to include fraction participating as an additional covariate, with dummies that span bins that cover both switching and non-switching families; although in that case, it may also be prudent to place less emphasis on the reweighted estimates.

⁵⁵We provide the equivalent figure for the “Head-Start-eligible” target population in Appendix Figure B.4. For this target group, the range of switching sample estimated p-scores encompasses that for non-switching target observations.

significant 12 p.p. (se: 0.053) using the baseline FFE model. This estimate is 57% smaller than the estimate reported in GTC, 0.281 (se: 0.108).⁵⁶ The standard errors are also roughly 50% smaller, corresponding to the roughly tripling of sample size (2,986 compared with 1,036).

As we foreshadowed earlier, these estimates are unlikely to represent the ATE for policy relevant populations, such as the Head Start eligible population and Head Start participants. Figure 3 shows a scatter of the FFE weights and the Head-Start-representative weights for each family in the white sample, divided by 2 to 3 child families (Panel A) and 4 or more child families (Panel B). The larger (smaller) markers signify that the estimated effect of Head Start on some college for the family is above (below) median. We also include a 45 degree line for reference. The figure shows that, in general, the Head-Start-representative weights are higher than the FFE weights for small families that experience smaller impacts of Head Start. Conversely, the representative weights are lower relative to the FFE weights for large families that experience larger impacts of Head Start. Hence, we should expect the reweighted estimates to show a reduced impact of Head Start relative to FFE.

The reweighted estimate of the impact of Head Start for the eligible, participant, and sibling populations is between 0.068, 0.026, and 0.079, respectively, and are all statistically insignificant. Setting aside the lack of precision in the estimates, these represent moderately large impacts relative to the 43.7% average rate of college going among Head Start eligible children. But comparing to the FFE coefficient, these effects imply a 34% to 78% smaller impact on college attendance. Putting these estimates in broader perspective, they are 45 to 91% smaller than the unadjusted estimates for *all* participants from other FFE studies (Bauer and Schanzenbach, 2016; Deming, 2009) and 51% smaller than the estimate from the county roll-out of Head Start (Bailey, Sun and Timpe, 2018), although the lower end of the confidence intervals for these estimates include our ATE.

Panel B of Table 6 presents results for the Economic Sufficiency Index in the PSID. Our FFE estimate shows a statistically insignificant 0.023 SD decline in this index associated with Head Start. When we reweight the effects, we find slightly larger negative effects for Head Start eligible children and Head Start participants, and a positive effect (0.03 SD) for siblings. It bears emphasizing, though, that the results are not precisely estimated, such that the 95% confidence intervals allow for a sizeable positive impact of Head Start in spite of the small or negative point estimate. For example, the confidence interval for the economic index for whites allows for a Head-Start-induced improvement of 0.16 SD or a reduction of 0.22 SD for Head Start participants. This limits our ability to make firm conclusions about Head Start's impact on this outcome.

The following four panels of Table 6 show the CNLSY FFE estimates, those reported in Deming (2009) and our replication, and our reweighted estimates. The panels report effects for high school graduation, idleness (not in school or at work), diagnosis of a learning disability, and poor health (based on self-reported health status). The FFE estimates indicate that Head Start leads to an 8.5 p.p. increase in high school graduation ($p < 0.01$), a 7.2 p.p. decline in idleness ($p < 0.10$), a

⁵⁶We show in the appendix that this discrepancy is not due to faulty replication of the GTC estimates in a smaller sample. We estimate a coefficient of 0.232 (se: 0.094) for this sample and outcome in our replication.

5.9 p.p. decline in having a learning disability ($p < 0.01$), and a 6.9 p.p. decline in reporting poor health ($p < 0.01$). The reweighted estimate for participants for high school is 44% smaller, and not statistically significant. We also see substantial 24% and 28% declines in the estimated impact on idleness and having a learning disability, respectively, when we consider the impact on participants. The poor health estimates are relatively more stable; the reweighted impacts on participants are just 3% smaller than the FFE estimate.

In the final column of the table, we test whether the difference between the reweighted estimate for participants and the FFE estimate is statistically significant. We bootstrap the standard errors for this difference by taking draws with replacement from the sample and performing the FFE estimation and reweighting again. We do this 1,000 times and obtain the standard error of our difference as the standard deviation of the 1,000 estimated FFE-reweighted differences. We find that the reweighted estimates for some college (PSID) and high school graduation (CNLSY) are statistically different from the FFE estimate at the 5% and 10% levels, respectively. The remainder of the outcomes are more imprecisely estimated, and therefore we can not reject that the reweighted estimate is the same as the FFE estimate.

Returning to the PSID, Appendix Tables B.13 and B.14 show the PSID FFE estimates and reweighted results for high school and the good health index for whites, and the corresponding results for blacks. Overall, the results suggest little support for a positive long term effect of Head Start. This is true for the FFE estimates and the reweighted estimates. Nonetheless, the magnitude of the estimates can vary importantly with reweighting, particularly for whites. This makes sense since the identifying sample is a much smaller share of the overall sample for whites relative to blacks. For example, the FFE estimate for the good health index for whites is -0.265 SD, but reweighting for the Head Start participant population changes this estimate to -0.439. In contrast, the coefficients are relatively stable for blacks.⁵⁷

We explore other reweighting strategies in Appendix Tables B.15 and B.16. Reweighting using linear extrapolation of treatment effects to singletons in Table B.15 produces qualitatively similar results to the baseline reweighting, although the point estimates for Head Start participants are often smaller.⁵⁸ Table B.16 presents the results when we reweight the FFE estimates using sample shares instead of propensity score weights. Across all outcomes, these estimates are quite similar to the FFE estimates, underscoring that the conditional-variance-weighting plays a minor role in this setting.

⁵⁷For the black sample, most estimates are also statistically insignificant. However, for the age 30 Economic Sufficiency Index, the reweighted estimates indicate statistically significant negative impacts of Head Start. For example, for a target population of participants the reweighted coefficient on Head Start is -0.211 (s.e. = 0.073).

⁵⁸We have also explored excluding singletons altogether from the target. The estimates for non-singleton Head Start participants and non-singleton Head-Start-eligible children typically lie between the reweighted estimates for siblings and Head Start participants.

6.3.2 More Evidence on the Role of Family Size

One key pattern in our findings is that larger families appear to have larger returns to Head Start than smaller families. We believe this to be a new finding in the Head Start literature. We note that this was not a pattern we initially set out to test in this study, so there is some chance of this finding being inadvertently driven by chance and our limited sample sizes. However we think that this may provide an interesting hypothesis for future studies. Also, we first observed this pattern in the PSID data, and so our CNLSY results (see e.g. Columns 3, 4, and 5 of Table 4) are to some degree an out of sample confirmation of this pattern.

We have examined whether the larger coefficients for larger family sizes in Table 4 are driven by family size standing in for other covariates. In Appendix Table B.17 we perform a “horse race” analysis, comparing whether heterogeneous coefficients load on to family size, or other covariates. This table shows that the heterogeneity with family size is robust to also allowing for heterogeneity along other covariates. We have also experimented with specifications that test for whether larger family size is merely proxying for “longer sibling cohort span,” and do not find evidence that this is the case.

6.3.3 Additional FFE Estimates

Continuing our analysis of the PSID, we also investigate effects of Head Start on a variety of additional short-term outcomes, outcomes at age 40, as well as heterogeneity by race, gender by cohort in Appendix C. We do not find any systematic evidence of effects on any of these outcomes, or important heterogeneity along these dimensions.

7 Other Applications

We have shown empirical evidence for selection into identification for three FFE applications relating to the returns to human capital investment and returns to domestic migration. In each of these contexts, there appears to be a mechanical relationship between $Pr(S_g = 1)$ and group size. In the Head Start setting, heterogeneity along these lines creates an upward-bias in the FE estimate. Since returns to migration may also be heterogeneous by family characteristics, it may be useful to also reweight the estimates from Collins and Wanamaker (2014) to obtain the ATE for a representative set of migrants.

We now discuss three additional FE designs present in the education, labor, and environmental literatures that illustrate settings where the tools that we have developed may apply. First, a number of studies examine the effect of peers in the classroom within a school-grade (or school) using school-grade FE (or school FE). For example, Carrell and Hoekstra (2010) and Carrell, Hoekstra and Kuka (2018) examine the effect of having a peer exposed to domestic violence (DV) using this strategy, finding large negative impacts on contemporaneous achievement that persist to

reduce long-term earnings. While the DV measure in these studies is continuous, it is reasonable to think that this may be a “lumpy” variable in the sense that some schools (or school-grades), which have a low probability of DV, will never have a student exposed to DV during the 8 year window of observation, and some school-grades, which have a high probability of DV, may always have a student exposed to DV. Given the likely correlation between $Pr(DV_g = 1)$ and $Pr(S_g = 1)$, non-switcher schools probably also have a different set of school resources (e.g. share of highly experienced teachers) and student composition (e.g. mean family income) than switchers, which could either exacerbate or mitigate the effects of DV. As a result, the effects estimated from switching schools may not generalize to low-probability-DV non-switchers or high-probability-DV non-switchers.

Second, a set of influential papers by Dube, Lester and Reich (2010, 2015) identify the impact of minimum wage laws within border-county pairs (using border-pair-by-year FE). This strategy produces bounds on minimum wage elasticities which are less negative than those estimated with other strategies. The authors report that 91% of the county pairs in the data have variation in the minimum wage at some point during the analysis, however states with more border counties and who have more frequent changes to the minimum wage relative to neighboring states will contribute more variation to the design. Hence, in practice, identification may be concentrated among a subset of the 91%. At the same time, the characteristics of switching border counties are likely to be different from interior counties, in terms of the education distribution, population density, or industry composition, which could influence the response to minimum wage increases (Cengiz et al., 2019). Thus, reweighting the estimates of switching border pairs to account for these characteristics could yield a different estimate for the impact of the minimum wage.

Third, it has become common to estimate the effect of environmental shocks on health and human capital using variation in temperature or rainfall within a local area (e.g. district FE). For example, Shah and Steinberg (2017) employ this strategy and find that a positive rain shock (top 20% rainfall) reduces the likelihood that students attend school, and vice versa for droughts. Since shocks are by definition infrequent events, it is likely that some districts that have more moderate climates will have no shocks over the 4 years of analysis. These non-switching districts may be located in a different geography, have distinct industrial composition, or population characteristics, which could in turn affect the elasticity of school attendance. Hence, extrapolating from switcher to non-switcher districts may require reweighting strategies such as those we propose.

Some recent work notes empirical issues related to SI. In the context of estimating returns to experience for teachers, Wiswall (2013) observes that teacher-FE models, paired with modest panel lengths and experience categories such as “4+ years,” imply that the coefficients on these categories will be solely identified from relatively new teachers, and lead to misestimation of the returns to experience. Spears, Coffey and Behrman (2019) raise SI concerns for FFE models used by Jayachandran and Pande (2017) to examine the relationship between birth order and child height in the Demographic and Health Surveys. Since these outcomes are only observed for those under 5,

Spears, Coffey and Behrman (2019) argue that a FE design selects for families with shorter birth spacing, which in turn may have more negative impacts of birth order.

As researchers using FE consider the relevance of SI for a particular setting, we also note that greater degrees of SI may, on the margin, favor a design that uses a more coarse level of fixed effects in order to increase the number of switchers. For example, Geruso and Spears (2018) examine the impacts of heat on infant mortality. They use district-by-calendar-month FEs in one specification, and village-by-calendar-month FEs as an alternative. The latter specification, with smaller units, may have greater SI concerns. It also has larger standard errors, and the authors suggest that the coarser aggregation is their preferred model. The reweighting methods we propose could be employed to mitigate the SI concerns in applications that use more narrow FEs.

Overall, these applications highlight the fact that selection into identification is likely to be relevant across the numerous domains where FE are applied. We leave it to future researchers to quantify the role of this selection and apply reweighting techniques to test the sensitivity of the conclusions.

8 Conclusion

Fixed effects can provide a useful approach for treatment effect estimation. The *internal* validity of this strategy, which has been the subject of much debate, relies on the assumption that treatment is randomly assigned to units in a group. In this article, we show that an additional assumption is needed for the *external* validity of results: that groups with variation (switchers) have comparable treatment effects to groups without variation (non-switchers). In other words, fixed effects estimates are generalizable only if there is no *selection into identification*.

We show that this assumption is not trivial in the context of family fixed effects. We document across multiple settings that switching families are systematically larger and show that this can induce bias in estimation. We develop a novel approach to recover ATE's for representative populations, which upweights observations that are under-represented in the identifying sample relative to the population of interest. We demonstrate that this reweighting approach performs well using Monte Carlo simulations.

We apply these lessons to an analysis of the long term effects of Head Start in the PSID and CNLSY using family fixed effects. Relative to prior evaluations of Head Start using FFE in the PSID, we use a sample three times as large in size, include longer run (up to age 40) outcomes, and expand the set of outcomes under consideration. Echoing prior findings, we find using FFE that Head Start significantly increases the likelihood of completing some college and graduating from high school, and decreases the likelihood of being idle, having a disability, or reporting poor health.

Using our reweighting methods, we estimate that Head Start leads to a 2.6 p.p. increase in the likelihood of attending some college for Head Start participants, and a 6.8 p.p. increase for Head Start eligible. The ATE estimate for participants is 78% smaller than the FFE estimate, a difference

which is statistically significant at the 5% level. We examine several other outcomes and find few statistically significant results. In sum, the FFE results in the PSID indicate that Head Start has little effect on many long term outcomes on average, with the exception of completing some college, and perhaps even detrimental effects for men. In the CNLSY, for high school graduation we find that the reweighted estimate for participants (4.8 p.p.) is 44% smaller than the FFE estimate, a difference which is statistically different at the 10% level. We find relatively less change associated with reweighting for other outcomes.

Overall, we interpret our findings as pointing primarily toward “increased uncertainty” and to a limited degree toward “zero effects” of the Head Start program. This suggests that there is some discordance between the long-term results from the FFE design, and new estimates using other designs, which generally produce larger and more robust effects of this intervention. Reconciling these findings is beyond the scope of this paper, but would be a productive avenue for future work.

Based on our findings, we propose new standards for practice when using FE or similar research designs to diagnose, and potentially correct for, the role of changes in sample composition in explaining the gap between OLS and FE estimates.

1. First, analyses should report the switching sample size in addition to the total sample size, including for relevant subsamples of the data (e.g. whites and blacks). It may also be useful to calculate the effective number of observations and share of identifying variation from true switchers to increase transparency into the variation among switchers.
2. Second, we suggest that researchers show a balance of observables across switching status to complement evidence of within-sample balance across treatment status. These covariates should include the number of units in a group (if there is imbalance) and correlates of treatment. For example, in the case of movers, one might consider testing for balance of urbanicity, age, and occupations. If there are differences in these covariates, researchers should examine heterogeneity along these dimensions. These tests are likely to have limited power to detect issues if there are interactions between covariates, but are a useful bellweather for important external validity concerns.
3. As a subsequent step, we recommend using propensity-score reweighting of the FE estimates to obtain estimates for a representative population or a policy-relevant population, such as program participants. Since these methods can perform unevenly under some models of heterogeneity, we suggest testing for sensitivity of results and reporting a range of estimates where applicable.

References

- Abrevaya, Jason.** 2006. “Estimating the effect of smoking on birth outcomes using a matched panel data approach.” *Journal of Applied Econometrics*, 21(4): 489–519.
- Aizer, Anna, and Flavio Cunha.** 2012. “The Production of Human Capital: Endowments, Investments and Fertility.” National Bureau of Economic Research Working Paper 18429. DOI: 10.3386/w18429.
- Almond, Douglas, Kenneth Y. Chay, and David S. Lee.** 2005. “The Costs of Low Birth Weight.” *The Quarterly Journal of Economics*, 120(3): 1031–1083.
- Anderson, Michael L.** 2008. “Multiple Inference and Gender Differences in the Effects of Early Intervention: A Reevaluation of the Abecedarian, Perry Preschool, and Early Training Projects.” *Journal of the American Statistical Association*, 103.
- Andersson, Fredrik, John C. Haltiwanger, Mark J. Kutzbach, Giordano E. Palloni, Henry O. Pollakowski, and Daniel H. Weinberg.** 2016. “Childhood Housing and Adult Earnings: A Between-Siblings Analysis of Housing Vouchers and Public Housing.” National Bureau of Economic Research Working Paper 22721.
- Andrews, Isaiah, and Emily Oster.** 2019. “A Simple Approximation for Evaluating External Validity Bias.” *Economics Letters*, 178: 58–62. Working Paper.
- Angrist, Joshua, and Jorn-Steffen Pischke.** 2009. *Mostly Harmless Econometrics*. Princeton University Press.
- Angrist, Joshua D.** 1998. “Estimating the Labor Market Impact of Voluntary Military Service Using Social Security Data on Military Applicants.” *Econometrica*, 66(2): 249–288.
- Angrist, Joshua D., and Ivan Fernandez-Val.** 2013. “ExtrapoLATE-ing: External Validity and Overidentification in the LATE Framework.” *Advances in Economics and Econometrics: Tenth World Congress*, , ed. Daron Acemoglu, Manuel Arellano and Eddie Dekel Vol. 3 of *Econometric Society Monographs*, 401 – 434. Cambridge University Press.
- Aronow, Peter M., and Allison Carnegie.** 2013. “Beyond LATE: Estimation of the Average Treatment Effect with an Instrumental Variable.” *Political Analysis*, 21(4): 492–506.
- Bailey, Martha J., Shuqiao Sun, and Brenden Timpe.** 2018. “Prep School for Poor Kids: The Long-Run Impacts of Head Start on Human Capital and Economic Self-Sufficiency.” Working Paper.
- Barr, Andrew, and Chloe R. Gibbs.** 2018. “Breaking the Cycle? Intergenerational Effects of an Anti-Poverty Program in Early Childhood.” *mimeo*.

- Bates, Michael David, Katherine E. Castellano, Sophia Rabe-Hesketh, and Anders Skrondal.** 2014. “Handling Correlations Between Covariates and Random Slopes in Multilevel Models.” *Journal of Educational and Behavioral Statistics*, 39(6): 524–549.
- Bauer, Lauren, and Diane Whitmore Schanzenbach.** 2016. “The Long-Term Impact of the Head Start Program.” *The Hamilton Project*.
- Bayer, Patrick, Randi Hjalmarsson, and David Pozen.** 2009. “Building Criminal Capital behind Bars: Peer Effects in Juvenile Corrections*.” *The Quarterly Journal of Economics*, 124(1): 105–147.
- Beck, Nathaniel.** 2015. “Estimating Grouped Data Models with a Binary Dependent Variable and Fixed Effects: What are the Issues? Comments Prepared for Delivery at the Annual Meeting of the Society for Political Methodology.” *mimeo*.
- Black, Sandra E., Paul J. Devereux, and Kjell G. Salvanes.** 2007. “From the Cradle to the Labor Market? The Effect of Birth Weight on Adult Outcomes*.” *The Quarterly Journal of Economics*, 122(1): 409–439.
- Borusyak, Kirill, and Xavier Jaravel.** 2017. “Revisiting Event Study Designs, with an Application to the Estimation of the Marginal Propensity to Consume.” Working Paper.
- Bound, John, and Gary Solon.** 1999. “Double Trouble: On the Value of Twins-based Estimation of the Return to Schooling.” *Economics of Education Review*, 18(2): 169–182.
- Callaway, Brantly, and Pedro H. C. Sant’Anna.** 2018. “Difference-in-Differences with Multiple Time Periods and an Application on the Minimum Wage and Employment.” Working Paper.
- Cameron, A. Colin, and Pravin K. Trivedi.** 2005. *Microeconometrics: Methods and Applications*. Cambridge University Press.
- Carneiro, Pedro, and Rita Ginja.** 2014. “Long-Term Impacts of Compensatory Preschool on Health and Behavior: Evidence from Head Start.” *American Economic Journal: Economic Policy*, 6(4): 135–173.
- Carrell, Scott E., and Mark L. Hoekstra.** 2010. “Externalities in the Classroom: How Children Exposed to Domestic Violence Affect Everyone’s Kids.” *American Economic Journal: Applied Economics*, 2(1): 211–228.
- Carrell, Scott E., Mark Hoekstra, and Elira Kuka.** 2018. “The Long-Run Effects of Disruptive Peers.” *American Economic Review*, 108(11): 3377–3415.
- Cengiz, Doruk, Arindrajit Dube, Attila Lindner, and Ben Zipperer.** 2019. “The Effect of Minimum Wages on Low-Wage Jobs*.” *The Quarterly Journal of Economics*.

- Chaisemartin, Clement, and Xavier D’Haultfoeille.** 2019. “Two-way Fixed Effects Estimators with Heterogeneous Treatment Effects.” *mimeo*.
- Chamberlain, Gary.** 1980. “Analysis of Covariance with Qualitative Data.” *The Review of Economic Studies*, 47(1): 225–238.
- Chernozhukov, Victor, Iván Fernández-Val, Jinyong Hahn, and Whitney Newey.** 2013. “Average and Quantile Effects in Nonseparable Panel Models.” *Econometrica*, 81(2): 535–580.
- Chetty, Raj, and Nathaniel Hendren.** 2018*a*. “The Impacts of Neighborhoods on Intergenerational Mobility I: Childhood Exposure Effects*.” *The Quarterly Journal of Economics*, 133(3): 1107–1162.
- Chetty, Raj, and Nathaniel Hendren.** 2018*b*. “The Impacts of Neighborhoods on Intergenerational Mobility II: County-Level Estimates.” *The Quarterly Journal of Economics*, 133(3): 1163–1228.
- Chorniy, Anna V, Janet Currie, and Lyudmyla Sonchak.** 2018. “Does Prenatal WIC Participation Improve Child Outcomes?” National Bureau of Economic Research Working Paper 24691.
- Collins, William J., and Marianne H. Wanamaker.** 2014. “Selection and Economic Gains in the Great Migration of African Americans: New Evidence from Linked Census Data.” *American Economic Journal: Applied Economics*, 6(1): 220–252.
- Currie, Janet, and Duncan Thomas.** 1995. “Does Head Start Make a Difference?” *American Economic Review*, 85(3): 341–364.
- Currie, Janet, and Ishita Rajani.** 2015. “Within-Mother Estimates of the Effects of WIC on Birth Outcomes in New York City.” *Economic Inquiry*, 53(4): 1691–1701.
- Currie, Janet, and Maya Rossin-Slater.** 2013. “Weathering the storm: Hurricanes and birth outcomes.” *Journal of Health Economics*, 32(3): 487 – 503.
- Deming, David.** 2009. “Early Childhood Intervention and Life-Cycle Skill Development: Evidence from Head Start.” *American Economic Journal: Applied Economics*, 1(3): 111–134.
- Dube, Arindrajit, T. William Lester, and Michael Reich.** 2010. “Minimum Wage Effects Across State Borders: Estimates Using Contiguous Counties.” *The Review of Economics and Statistics*, 92(4): 945–964.
- Dube, Arindrajit, T. William Lester, and Michael Reich.** 2015. “Minimum Wage Shocks, Employment Flows, and Labor Market Frictions.” *Journal of Labor Economics*, 34(3): 663–704.

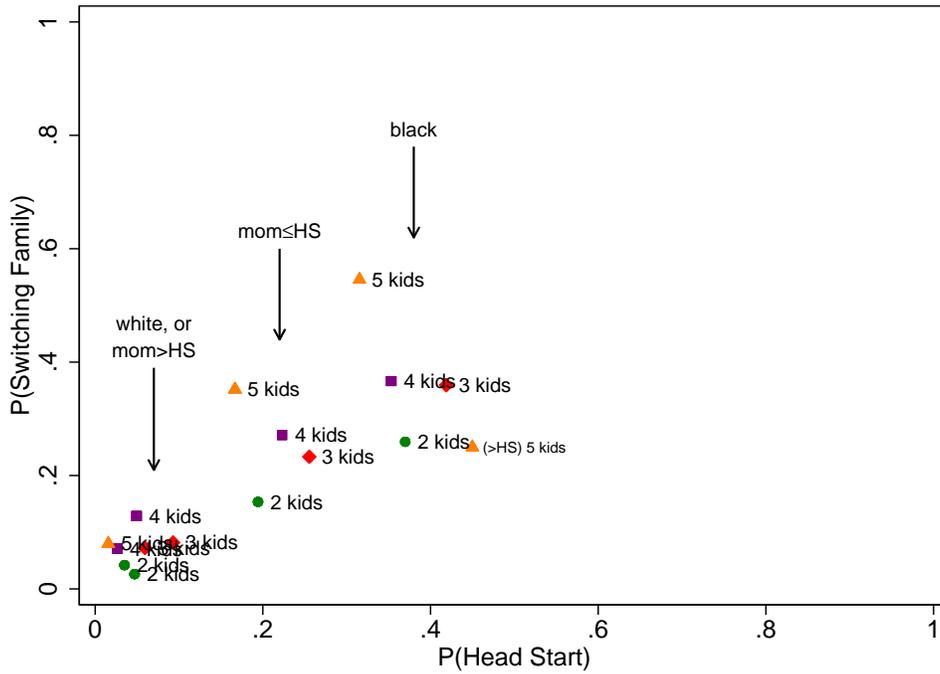
- Fernandez-Val, Ivan.** 2009. “Fixed Effects Estimation of Structural Parameters and Marginal Effects in Panel Probit Models.” *Journal of Econometrics*, 150(1): 71 – 85.
- Figlio, David, Jonathan Guryan, Krzysztof Karbownik, and Jeffrey Roth.** 2014. “The Effects of Poor Neonatal Health on Children’s Cognitive Development.” *American Economic Review*, 104(12): 3921–3955.
- Finkelstein, Amy, Matthew Gentzkow, and Heidi Williams.** 2016. “Sources of Geographic Variation in Health Care: Evidence From Patient Migration*.” *The Quarterly Journal of Economics*, 131(4): 1681–1726.
- Garces, Eliana, Duncan Thomas, and Janet Currie.** 2002. “Longer-Term Effects of Head Start.” *American Economic Review*, 92(4): 999–1012.
- Geruso, Michael, and Dean Spears.** 2018. “Heat, Humidity, and Infant Mortality in the Developing World.” National Bureau of Economic Research Working Paper 24870.
- Gibbons, Charles E., Serrato Juan Carlos Suarez, and Michael B. Urbancic.** 2018. “Broken or Fixed Effects?” *Journal of Econometric Methods*, 0(0).
- Gibbs, Chloe, Jens Ludwig, and Douglas L Miller.** 2013. “Does Head Start Do Any Lasting Good?” *Legacies of the War on Poverty*, ed. Martha J. Bailey and Sheldon Danziger. Russell Sage Foundation.
- Goodman-Bacon, Andrew.** 2018. “Difference-in-Differences with Variation in Treatment Timing.” National Bureau of Economic Research Working Paper 25018.
- Hoxby, Caroline M.** 2000. “The Effects of Class Size on Student Achievement: New Evidence from Population Variation.” *Quarterly Journal of Economics*, 115(4): 1239–1285.
- Hoynes, Hilary, Diane Whitmore Schanzenbach, and Douglas Almond.** 2016. “Long-Run Impacts of Childhood Access to the Safety Net.” *American Economic Review*, 106(4): 903–934.
- Imai, Kosuke, and In Song Kim.** 2019. “When Should We Use Unit Fixed Effects Regression Models for Causal Inference with Longitudinal Data?” *American Journal of Political Science*, 63(2): 467–490.
- Imbens, Guido W., and Donald B. Rubin.** 2015. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press.
- Jayachandran, Seema, and Rohini Pande.** 2017. “Why Are Indian Children So Short? The Role of Birth Order and Son Preference.” *American Economic Review*, 107(9): 2600–2629.

- Johnson, Rucker C., and C. Kirabo Jackson.** 2017. “Reducing Inequality Through Dynamic Complementarity: Evidence from Head Start and Public School Spending.” National Bureau of Economic Research Working Paper 23489. DOI: 10.3386/w23489.
- Kling, Jeffrey R, Jeffrey B Liebman, and Lawrence F Katz.** 2007. “Experimental Analysis of Neighborhood Effects.” *Econometrica*, 75(1): 83–119.
- Lemieux, Thomas.** 1998. “Estimating the Effects of Unions on Wage Inequality in a Panel Data Model with Comparative Advantage and Nonrandom Selection.” *Journal of Labor Economics*, 16(2): 261–291.
- Leuven, Edwin, and Monique de Haan.** 2016. “Head Start and the Distribution of Long Term Education and Labor Market Outcomes.” IZA Institute of Labor Economics Working Paper.
- Lochner, Lance, and Enrico Moretti.** 2015. “Estimating and Testing Models with Many Treatment Levels and Limited Instruments.” *The Review of Economics and Statistics*, 97(2): 387–397.
- Loken, Katrine V., Magne Mogstad, and Matthew Wiswall.** 2012. “What Linear Estimators Miss: The Effects of Family Income on Child Outcomes.” *American Economic Journal: Applied Economics*, 4(2): 1–35.
- Ludwig, Jens, and Douglas L. Miller.** 2007. “Does Head Start Improve Children’s Life Chances? Evidence from a Regression Discontinuity Design.” *The Quarterly Journal of Economics*, 122(1): 159–208.
- Mundlak, Yair.** 1978. “On the Pooling of Time Series and Cross Section Data.” *Econometrica*, 46(1): 69–85.
- Pages, Remy J.-C., Dylan J. Lukes, Drew H. Bailey, and Greg J. Duncan.** 2019. “Elusive Longer-Run Impacts of Head Start: Replications Within and Across Cohorts.” *arXiv:1903.01954 [econ, q-fin]*. arXiv: 1903.01954.
- Robins, James M., Andrea Rotnitzky, and Lue Ping Zhao.** 1995. “Analysis of Semiparametric Regression Models for Repeated Outcomes in the Presence of Missing Data.” *Journal of the American Statistical Association*, 90(429): 106–121.
- Rossin-Slater, Maya.** 2013. “WIC in your neighborhood: New evidence on the impacts of geographic access to clinics.” *Journal of Public Economics*, 102: 51–69.
- Roy, A. D.** 1951. “Some Thoughts on the Distribution of Earnings.” *Oxford Economic Papers*, 3(2): 135–146.
- Shah, Manisha, and Bryce Millett Steinberg.** 2017. “Drought of Opportunities: Contemporaneous and Long-Term Impacts of Rainfall Shocks on Human Capital.” *Journal of Political Economy*, 125(2): 527–561.

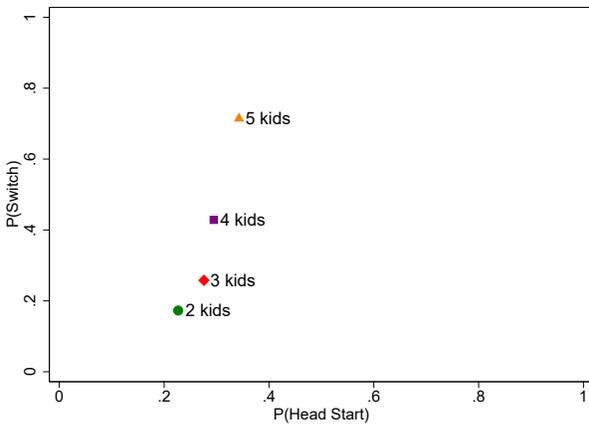
- Sloczynski, Tymon.** 2018. “A General Weighted Average Representation of the Ordinary and Two-Stage Least Squares Estimands.” Working Paper.
- Spears, Dean, Diane Coffey, and Jere Behrman.** 2019. “Birth Order, Fertility, and Child Height in India and Africa.” National Bureau of Economic Research Working Paper 12289.
- Stuart, Elizabeth A., Stephen R. Cole, Catherine P. Bradshaw, and Philip J. Leaf.** 2011. “The Use of Propensity Scores to Assess the Generalizability of Results from Randomized Trials.” *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 174(2): 369–386.
- Suri, Tavneet.** 2011. “Selection and Comparative Advantage in Technology Adoption.” *Econometrica*, 79(1): 159–209.
- Thompson, Owen.** 2017. “Head Start’s Long-Run Impact: Evidence from the Program’s Introduction.” *Journal of Human Resources*.
- Verdier, Valentin, and Andrew Castro.** 2019. “Average Treatment Effects for Stayers with Correlated Random Coefficient Models of Panel Data.” *mimeo*.
- Wiswall, Matthew.** 2013. “The dynamics of teacher quality.” *Journal of Public Economics*, 100: 61–78.
- Wooldridge, Jeffrey M.** 2005. “Fixed-Effects and Related Estimators for Correlated Random-Coefficient and Treatment-Effect Panel Data Models.” *The Review of Economics and Statistics*, 87(2): 385–390.
- Wooldridge, Jeffrey M.** 2010. *Econometric Analysis of Cross Section and Panel Data*. MIT Press.
- Wooldridge, Jeffrey M.** 2019. “Correlated random effects models with unbalanced panels.” *Journal of Econometrics*, 211(1): 137 – 150. Annals Issue in Honor of Jerry A. Hausman.
- Xie, Zong-Xian, Shin-Yi Chou, and Jin-Tan Liu.** 2016. “The Short-Run and Long-Run Effects of Birth Weight: Evidence from Large Samples of Siblings and Twins in Taiwan.” *Health Economics*, 26(7): 910–921.

Figure 2: Likelihood of Being a Switcher Family Increases with Family Size and P(treatment)

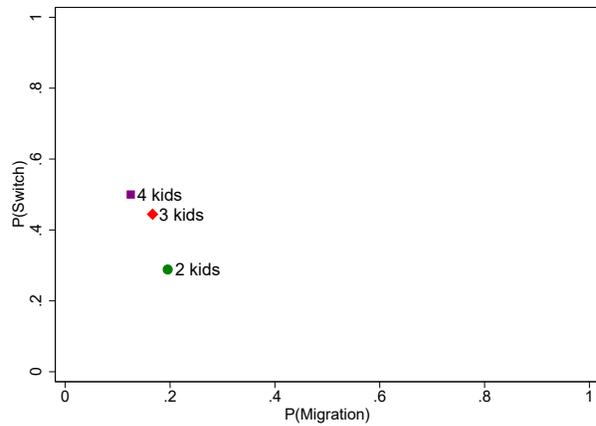
(a) Head Start in PSID



(b) Head Start in CNLSY



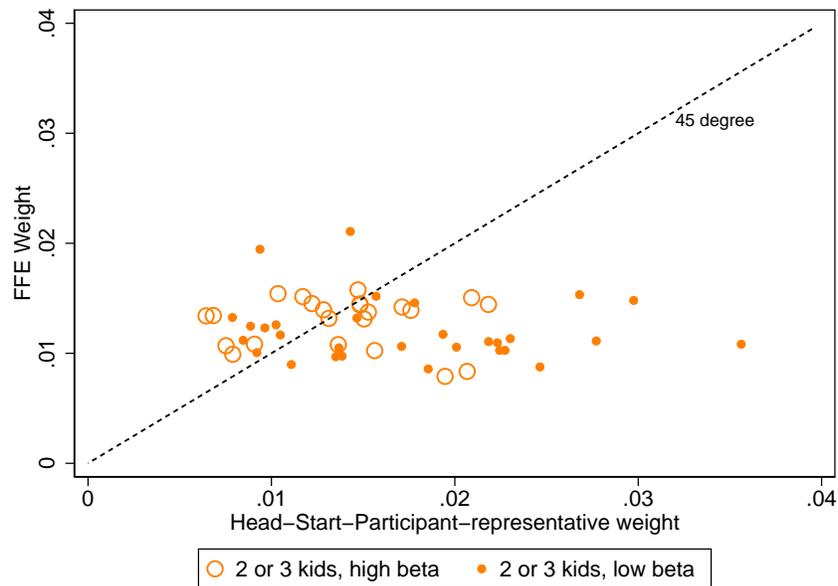
(c) Migration in Census



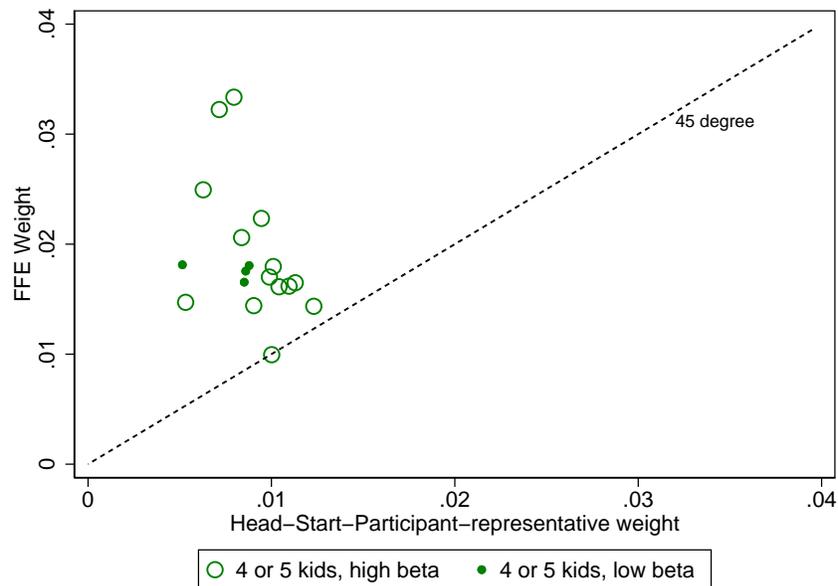
Notes: This figure shows the probability of being in a switching family and the probability of “treatment” by family size using three datasets and varying treatments. Panel (a) plots the probability of being in a switching family and of attending Head Start by family size for the following groups in the PSID: Whites, Blacks, children of mothers with at most a high school degree, and children of mothers with at least some college. Figure (b) is a simplified version of (a) using data on Head Start participation and family size from the CNLSY. Figure (c) shows the probability of being in a switching family and the probability of migrating to the northern US, using a linking of the 1910 to 1930 censuses used in Collins and Wanamaker (2014).

Figure 3: FFE Weights and Head-Start-Participant-Representative Weights by Family Size and Some College β (PSID White Sample)

(a) Families with 2-3 Children



(b) Families with 4+ Children



Notes: Each marker in this figure indicates the FFE weights and Head-Start-participant-representative (post-regression) weight for one white switching family. The color of the marker indicates whether the family has 2-3 children or 4 or more children. The size of the marker indicates the estimated family-specific beta from a regression of attainment of some college on interactions between Head Start and family id fixed effects. A larger marker indicates an above median beta, while a smaller marker indicates a below-median beta. The 45 degree line is included for reference. Observations above (below) the line are overweighted (underweighted) in the FFE sample relative to a representative Head Start sample. Source: Panel Study of Income Dynamics, 1968-2011 waves.

10 Tables

Table 1: Family FE Articles in Top Applied Journals, 2002 to 2017

	Binary Indep.	Binary Dep.	Both Binary	Total
AEJ: Applied	6	4	3	8
AEJ: Economic Policy	1	1	1	1
AER	3	1	1	5
AER Papers and Proceedings	2	2	1	3
Journal of Health Economics	5	3	2	7
Journal of Human Resources	7	2	2	12
Journal of Labor Economics	2	1	1	5
Journal of Political Economy	2	1	1	2
Journal of Public Economics	4	4	4	5
QJE	1	4	1	4
Review of Economics and Statistics	2	0	0	3
Total	35	23	17	55
<i>Common Dependent Variables</i>				
Schooling/Attainment	23			
Test Score	17			
Employment/Earnings	15			
Birth Weight	6			
Health	6			
Behavioral Issues/Crime	5			
<i>Common Independent Variables</i>				
Schooling	8			
Birth Weight	5			
Health	5			
Parental Traits	4			
Employment	3			
Birth order	3			
Means-Tested Public Program	2			
Death of Family Member	2			
Bombing/Radiation	2			
<i>Observations by Sample</i>				
	Siblings N	Total N		
p10	469	1,212		
p25	1,167	2,142		
p50	6,315	17,501		
p75	160,122	551,630		
p90	750,697	1,582,142		
Year Publication Min/Max	2002	2017		
Articles with Balance Table if Binary Ind.	1			

Notes: This table presents a summary of FFE articles published between January 2000 and May 2017 in 11 top applied journals, which are listed in the first panel of the table. For reference, between 2002 and 2017 the number of articles published in AEJ: Applied was 310; AEJ: Policy was 313; AER was 1722; AER P&P was 1676; JoLE was 434; Journal of Political Economy was 548; QJE was 639; JHR was 543; JPubE was 1688; REStat was 1033; JHE was 1017. Articles were initially identified using the search terms “family,” “within family,” “sibling,” “twin,” “mother,” “father,” “brother,” “sister,” “fixed effect,” “fixed-effect,” and “birthweight” using queries on journal websites. Siblings N is the number of observations reported for the sample of siblings, while Total N represents the number of total observations reported. See text for details.

Table 2: Switchers and Non-Switchers Vary Along Dimensions Other Than Family Size

	(1)	(2)	(3)	(4)	(5)
	Switch	Non-Switch	T-Stat. (1)=(2)	Beta Switch	T-Stat (4)
<i>A. Individual Covariates</i>					
Fraction female	0.562	0.495	4.067	0.024	0.719
Fraction African-American	0.516	0.111	25.877	0.249	5.640
Mother's yrs education	9.283	11.230	-21.590	-0.140	-0.751
Father's yrs education	9.190	11.371	-19.594	-0.389	-1.784
Had a single mother at age 4	0.252	0.099	10.049	0.055	2.543
Family income (age 3-6) (CPI adjusted)	31809	52574	-24.735	-4759	-5.719
Mother employed, age 0	0.508	0.570	-3.099	0.055	2.339
Mother employed, age 1	0.517	0.543	-1.342	0.058	2.359
Mother employed, age 2	0.536	0.554	-0.951	0.118	3.565
Household size at age 4	5.487	4.451	12.343	0.755	4.936
Fraction low birth weight	0.077	0.058	1.971	0.010	0.702
Observations	1103	5500	6603	7372	7372
<i>B. Inverse Selection into Identification Wts.</i>					
Pr(switch)/Pr(Head Start), Whites	2.976	2.318			
	(1.99)	(1.98)			
Pr(switch)/Pr(Head Start), Blacks	1.987	1.148			
	(1.21)	(1.10)			

Notes: Panel A of this table presents comparisons of the characteristics of individuals in switching families and non-switching families. Columns 1, 2, and 3, respectively, show the mean characteristics of individuals in families that are switchers; individuals in families that are not switchers; and individuals that attended Head Start (HS) in non-switcher families. Column 3 presents the t-statistic for the test that columns 1 and 2 are equal. Column 4 shows the estimates from a regression of each row heading on an indicator for being in a switcher family, with the corresponding t-statistic shown in Column 5, with standard errors clustered on id1968. All controls from the main specification are included excluding the variable shown in the row heading. All estimates are weighted to be representative of 1995 population; see text for details. Panel B shows the mean and standard deviation (in parenthesis) of the inverse of the post-regression propensity score weights when the target is Head Start participants. This gives a measure of how aligned the characteristics of switchers are with the characteristics of Head Start participants, the population of interest. An average value of 1 implies perfect alignment, while a higher value implies that the characteristics of switchers are over-represented relative to the characteristics of Head Start participants. Pr(switch) and Pr(Head Start) are estimated from a multinomial logit model of these outcomes on family size and other covariates described in the text. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table 3: Change in Weighting of Regression Estimates Across Sibling and Switcher Samples (PSID)

	Number of Children in Family:				
	1	2	3	4	5 +
<i>A. Share of Sample</i>					
All Sample	0.123	0.273	0.238	0.147	0.134
Siblings Sample	0.000	0.345	0.300	0.186	0.169
Switchers Sample	0.000	0.210	0.271	0.197	0.322
<i>B. Variance in Head Start</i>					
All Sample	0.089	0.104	0.121	0.127	0.132
Siblings Sample	0.000	0.024	0.050	0.059	0.068
Switchers Sample	0.000	0.045	0.098	0.131	0.174
<i>C. Regression weights</i>					
All Sample	0.171	0.257	0.284	0.117	0.101
Siblings Sample	0.000	0.338	0.374	0.154	0.134
Switchers Sample	0.000	0.256	0.307	0.190	0.248

Notes: This table shows the change in the composition of the PSID sample moving from all individuals (“All Sample”) to individuals that have at least one other sibling in the sample (“Siblings Sample”) to individuals in families that have variation in Head Start attendance (“Switchers sample.”) Panel A shows the share of individuals in each sample that come from a family with 1 child (zero siblings), 2 children, etc. Panel B shows the variance in Head Start for each family size and sample. For the switcher sample, this is calculated net of family fixed effects. Panel C shows the “regression weight” given to each family size in a given sample, denoted as ω_z and defined formally in Section 3. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table 4: Returns to Head Start by Family Size,
and Implications for Regression Estimates

	PSID		CNLSY		
	Some College		HS Grad	Idle	Ln. Disab.
	CX (1)	FE (2)	FE (3)	FE (4)	FE (5)
<i>A. Effects by Family Size</i>					
Head Start x 1 child family	0.169* (0.091)				
Head Start x 2 child family	0.038 (0.079)	-0.126 (0.099)	0.033 (0.042)	-0.067 (0.052)	-0.028 (0.025)
Head Start x 3 child family	-0.030 (0.087)	0.152** (0.075)	0.061 (0.060)	-0.038 (0.068)	-0.070 (0.043)
Head Start x 4 child family	-0.053 (0.100)	0.251*** (0.091)	0.156* (0.086)	-0.002 (0.111)	-0.064 (0.049)
Head Start x 5+ child family	0.572*** (0.119)	0.348*** (0.126)	0.277*** (0.097)	-0.306** (0.139)	-0.157* (0.081)
Head Start x Unknown child family	-0.099 (0.108)				
Observations	4258	2986	1251	1251	1247
Head Start Switchers		213	581	581	581
Effective Obs. (Indivs. 2-Person Fams)		235.9	647.9	647.9	647.9
Effective Obs. (CX Indivs.)		731.8	438.7	438.7	438.7
<i>B. Simulated Estimates across Samples using Family-Size Regression Weights</i>					
All	0.046				
Siblings	0.037	0.083	0.074	-0.068	-0.053
Switchers	0.069	0.123	0.088	-0.073	-0.060

Notes: Panel A of this table shows the coefficients from a regression of some college on a series of indicators for whether an individual attended Head Start interacted with an indicator for the number of children in one's family. The sample is composed of white individuals. Columns 1 include controls, but not mother f.e., and standard errors are clustered at 1968 family id. Column 2 includes mother fixed effects, and standard errors clustered by mother id. The number of Head Start switchers is equal to the number of individuals in families that have variation in Head Start. "Effective Obs. (CX Indivs.)" is the equivalent number of cross-sectional units that provide the same amount of variation as switchers. "Effective Obs. (Indivs. 2-Person Fams)" is the equivalent number of individuals in 2-person switching families that provide the same amount of variation as switchers. Both of these are calculated using Equation 3, where the denominator is the variance of Head Start, residualized by the family mean of the covariates in the analysis, or 0.125, respectively. Panel B shows the weighted average of the coefficients when using regression weights, ω_z (defined in Section 3), determined by the overall distribution of families ("All"), the distribution of 2+ child families ("Siblings"), and the distribution of 2+ child families that have variation in Head Start attendance ("Switchers"). * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves and Children of the National Longitudinal Study of Youth.

Table 5: Monte Carlo Experiments: Bias of Reweighting and FFE Relative to True ATE, and Efficiency of Reweighting Relative to FFE

	True ATE	Bias:		MSE of Reweight
		FE	Reweight	MSE of FE
<i>A. Constant TE; p-score: X_{ig}</i>				
Switchers	80	-0.3	-0.2	1.03
Siblings	80	-0.3	-0.5	1.19
All	80	-0.3	-0.5	1.20
HS Participants	80	-0.3	-0.3	1.04
<i>B. Large family TE; p-score: large family</i>				
Switchers	83.0	-11.1*	-0.6	0.92
Siblings	49.6	22.2*	-0.1	0.70
All	40.3	31.6*	0.1	0.54
HS Participants	41.1	30.7*	0.1	0.55
<i>C. TE linear in X_{ig}; p-score: X_{ig}</i>				
Switchers	94.2	-2.0*	-0.6	1.03
Siblings	80.1	12.2*	1.6*	0.99
All	80.0	12.2*	1.7*	1.00
HS Participants	91.5	0.8	-0.2	1.03
<i>D. TE linear in X_{ig}; p-score: X_{ig} spline</i>				
Switchers	94.2	-1.5*	-0.3	1.04
Siblings	80.1	12.7*	-0.4	1.08
All	80.0	12.8*	-0.4	1.09
HS Participants	91.5	1.3	-0.2	1.09

Notes: This table shows the results from 3,000 Monte Carlo simulations. Each panel of the table shows results from a different DGP and/or different covariates used in the p-score, and each row within panel is for a different target population. The true DGP is linear, and is discussed in Section 4.4. The first panel shows results where Head Start has a constant treatment effect (TE) for all individuals; the second shows results where Head Start (HS) has no effect on individuals from small families (3 or fewer children) and a large effect for families with many children (4 or more children); and the third and fourth panels show results where treatment effects that are linear in X_{ig} . Column 1, “True Beta,” presents the true average increase in the probability of completing some college for participants in Head Start in the sample, which is a function of the DGP and sample composition. Columns 2 and 3 present the bias of various estimation strategies, defined as the difference between the estimated effects of Head Start and the true beta. The estimated effects come from a LPM, propensity-score weighted LPM, respectively. Column 4 presents the ratio of the mean squared error (MSE) of the reweighting estimators relative to LPM. Reweighted estimates are obtained using in-regression weighting, with weights adjusting for the representativeness of switchers (using the variable(s) indicated in each of the panel headings as predictors in the multinomial logit step) and the conditional variance of Head Start within families. All betas are multiplied by 1,000. * $p < .01$.

Table 6: Head Start Impact for Representative Eligible Children, Participants, and Siblings

Using Reweighting

	FFE		Reweighted ATE, Target =			Diff. b/w
	GTC/Deming	Expand Sample/ Replicate	HS Eligible	Participants	Siblings	FFE and Participant ATE
<i>A. Some College (PSID)</i>						
Head Start	0.281** (0.108)	0.120** (0.053)	0.068 (0.063)	0.026 (0.065)	0.079 (0.060)	0.094** (0.042)
Y Mean in Target	–	0.556	0.387	0.437	0.556	
<i>B. Economic Sufficiency Index, Age 30 (PSID)</i>						
Head Start	–	-0.023 (0.102)	-0.038 (0.096)	-0.032 (0.106)	0.021 (0.097)	0.009 (0.090)
Y Mean in Target	–	0.213	-0.198	-0.485	0.213	
<i>C. High School Graduation (CNLSY)</i>						
Head Start	0.086*** (0.031)	0.085*** (0.030)	0.033 (0.042)	0.048 (0.037)	0.020 (0.044)	0.037* (0.023)
Y Mean in Target	–	0.776	0.734	0.766	0.776	
<i>D. Idle (CNLSY)</i>						
Head Start	-0.071* (0.038)	-0.072* (0.037)	-0.061 (0.045)	-0.055 (0.042)	-0.067 (0.050)	-0.017 (0.026)
Y Mean in Target	–	0.197	0.221	0.201	0.197	
<i>E. Learning Disability (CNLSY)</i>						
Head Start	-0.059*** (0.020)	-0.059*** (0.021)	-0.031 (0.026)	-0.042* (0.022)	-0.040 (0.026)	0.017 (0.015)
Y Mean in Target	–	0.051	0.055	0.041	0.051	
<i>F. Poor Health (CNLSY)</i>						
Head Start	-0.070*** (0.026)	-0.069*** (0.026)	-0.063* (0.037)	-0.067** (0.034)	-0.050 (0.038)	-0.003 (0.020)
Y Mean in Target	–	0.103	0.098	0.074	0.103	

Notes: Column 1 of this table shows the FFE estimated impacts of Head Start for whites from GTC or for the whole sample from Deming (2009). Column 2 shows the FFE estimate using our expanded sample for PSID outcomes and using our replication sample for CNLSY outcomes. The outcomes in Panels A and B are taken from the PSID white sample, and the outcomes in Panels C to F are taken from the CNLSY sample. Columns 3 to 5 present reweighted estimates of the effect of Head Start for three target populations (shown in the column header) using the post-regression reweighting procedure, in which we multiply group-level estimates of the impact of Head Start by the representative weight for the target population of interest. Column 6 presents the difference in the estimate in column 2 (FFE) and column 4 (reweighted for participants), with the standard error obtained from a bootstrap procedure described in the text. "–" is used to indicate that the information is not available. Sample size is N=2,986 for the expanded sample, and 1,036 for GTC. Standard errors obtained by bootstrapping. * p < .10, ** p < .05, *** p < .01.

Appendix

Contents

- A. Derivations and Proofs
- B. Supplementary Figures and Tables
- C. Supplementary PSID FFE Results
- D. Replication of GTC (2002)
- E. Functional Form Choices with Binary Treatment and Binary Outcome

A Derivations and Proofs

A.1 Proof of Proposition 1

The proof of Proposition 1 closely follows the proofs of Theorem 2 in Angrist and Fernandez-Val (2013) and Theorem 1 in Aronow and Carnegie (2013). There are two key differences. First, we rely on Group ID Conditional Independence (Assumption 1), instead of the IV exclusion restriction. Second, we condition on two propensity scores, unlike Aronow and Carnegie (2013), who condition on $Pr(D_i = 1)$, and Angrist and Fernandez-Val (2013), who condition on discrete covariates.

Recall that we define $\hat{\delta}^t := \frac{1}{\sum_i \mathbf{1}(S_{g(i)=1})} \sum_i \widehat{w_{g(i)}^t} \cdot \hat{\delta}_{g,FE}$ and $w_{g(i)}^t := \frac{Q_x \cdot Pr[S_{g(i)=1}]}{P_x \cdot Pr[T_{g(i)=1}]}$.

By Assumptions 1 and 3,

$$\hat{\delta}^t \rightarrow_p \mathbb{E} \left[w_{g(i)}^t \cdot \hat{\delta}_g | S_{g(i)} = 1 \right] \quad (16)$$

By Assumption 1 and the law of iterated expectations,

$$\mathbb{E} \left[w_{g(i)}^t \cdot \hat{\delta} | S_{g(i)} = 1 \right] = \mathbb{E} \left[w_{g(i)}^t (Y_i(1) - Y_i(0)) | S_{g(i)} = 1 \right] \quad (17)$$

By the law of iterated expectations,

$$\mathbb{E}[w_{g(i)}^t \cdot (Y_i(1) - Y_i(0)) | S = 1] = \mathbb{E}[\mathbb{E}[w_{g(i)}^t (Y(1) - Y(0)) | S = 1, P_x, Q_x] | S_{g(i)} = 1] \quad (18)$$

$$= \mathbb{E}[\mathbb{E}[w_{g(i)}^t (Y(1) - Y(0)) | P_x, Q_x] | S_{g(i)} = 1] \quad (19)$$

$$= \mathbb{E}[w_{g(i)}^t \mathbb{E}[(Y(1) - Y(0)) | P_x, Q_x] | S_{g(i)} = 1] \quad (20)$$

$$= \mathbb{E}[w_{g(i)}^t \Delta(P_x, Q_x) | S_{g(i)} = 1]$$

Line (19) follows from (18) from the CFIE assumption that $\mathbb{E}[Y(1) - Y(0) | S, P_x, Q_x] = \mathbb{E}[Y(1) - Y(0) | P_x, Q_x]$. Line (20) follows from line (19) because $w_{g(i)}^t$ is a function of P_x and Q_x only, by definition.

Now let F be the distribution of (P_x, Q_x) , and let $F(\cdot | S = 1)$ be the distribution conditional on $S = 1$. By Bayes rule,

$$\begin{aligned} \int w_{g(i)}^t \Delta(P_x, Q_x) dF(P_x, Q_x | S = 1) &= \int w_{g(i)}^t \cdot \Delta(P_x, Q_x) \frac{Pr(S = 1 | P_x, Q_x)}{Pr(S = 1)} dF(P_x, Q_x) \\ &= \int w_{g(i)}^t \cdot \Delta(P_x, Q_x) \frac{P_x}{Pr(S = 1)} dF(P_x, Q_x) \\ &= \int \frac{Q_x P(S = 1)}{P_x P(T = 1)} \Delta(P_x, Q_x) \frac{P_x}{Pr(S = 1)} dF(P_x, Q_x) \\ &= \int \Delta(P_x, Q_x) \frac{Q_x}{Pr(T = 1)} dF(P_x, Q_x) \\ &= \int \Delta(P_x, Q_x) \frac{Pr(T = 1 | P_x, Q_x)}{Pr(T = 1)} dF(P_x, Q_x) \\ &= \int \Delta(P_x, Q_x) dF(P_x, Q_x | T = 1) \end{aligned}$$

By CFEL, $\int \Delta(P_x, Q_x) dF(P_x, Q_x|T = 1) = \int \mathbb{E}[Y(1) - Y(0)|T = 1, P_x, Q_x] dF(P_x, Q_x|T = 1) = \mathbb{E}[Y(1) - Y(0)|T = 1]$.

A.1.1 Extrapolating treatment effects to never-switchers

Note that the ATE for the target population can be written as a weighted average of the ATE for switchers, $\delta^{t, P_x > 0}$, and the ATE for never-switchers, $\delta^{t, P_x = 0}$:

$$\delta^t := Pr(P_x > 0|T_g = 1) \cdot \delta^{t, P_x > 0} + (1 - Pr(P_x > 0|T_g = 1)) \cdot \delta^{t, P_x = 0} \quad (21)$$

where $\delta^{t, P_x > 0} = \mathbb{E}[Y(1) - Y(0)|T = 1, P_x > 0]$ and $\delta^{t, P_x = 0} = \mathbb{E}[Y(1) - Y(0)|T = 1, P_x = 0]$.

Since we can not identify treatment effects for $P_x = 0$ from the FFE design, we must impose an additional assumption that allows extrapolation of treatment effects for this group.

Assumption 5 (Treatment Effect Functional Form):

$$\mathbb{E}[Y(1) - Y(0)|\mathbf{X}_g] = H(\Phi; \mathbf{X}_g)$$

with $H(\cdot)$ known, and Φ parameters that can be consistently estimated. Under Assumption 5, $\hat{\Phi}$ can be estimated using e.g. the regression $\hat{\delta}_g = \Phi' \mathbf{X}_g + u_g$.

Proposition 2. *Under Assumptions 1, 2, 3, and 5, and assuming $Pr(P_x > 0|T_g = 1)$ can be consistently estimated, the ATE for the target t can be consistently estimated by*

$$\hat{\delta}^t = Pr(P_x > 0|T_g = 1) \cdot \widehat{\delta^{t, P_x > 0}} + (1 - Pr(P_x > 0|T_g = 1)) \cdot \widehat{\delta^{t, P_x = 0}}$$

where $\delta^{t, P_x > 0}$ comes from Equation 12 and $\delta^{t, P_x = 0}$ is estimated from a projection of $\hat{\delta}_g$ on X_g .

Proof:

Define $\widehat{\delta^{t, P_x = 0}} := \frac{1}{\sum \mathbf{1}(T_g = 1, P_x = 0)} \sum_{T_g = 1, P_x = 0} H(\hat{\Phi}; \mathbf{X}_g)$.

Assumptions 1, 3, and 5 imply that

$$plim \widehat{\delta^{t, P_x = 0}} = \mathbb{E}[Y(1) - Y(0)|T = 1, P_x = 0] \quad (22)$$

From the proof of Proposition 1, we have that

$$plim \widehat{\delta^{t, P_x > 0}} = \mathbb{E}[Y(1) - Y(0)|T = 1, P_x > 0] \quad (23)$$

Then,

$$plim \hat{\delta}^t = Pr(P_x > 0|T_g = 1) \cdot \mathbb{E}[Y(1) - Y(0)|T = 1, P_x > 0] \quad (24)$$

$$\begin{aligned} &+ (1 - Pr(P_x > 0|T_g = 1)) \cdot \mathbb{E}[Y(1) - Y(0)|T = 1, P_x = 0] \\ &= \mathbb{E}[Y(1) - Y(0)|T = 1] \end{aligned} \quad (25)$$

A speculative alternative approach would be to take a “double robust” estimation approach. This is modeled after the double robust approach for estimating causal effects, as in Robins, Rotnitzky and Zhao (1995). The discussion in Chapter 17 of Imbens and Rubin (2015, pp. 399-400)

notes that in the traditional setting, consistency of the estimated treatment effect requires either correct specification of the propensity score, or of the regression model. In our setting, instead of trying to model potential outcomes, we are estimating average treatment effects. Implementation would proceed by (i) defining weights $w_{g(i)}^t$ as in equation (13); (ii) estimating $\hat{\delta}_g = H(\hat{\Phi}; \mathbf{X}_g)$ using these weights, and then predicting $\frac{1}{\sum \mathbf{1}(T_g=1)} \sum_{T_g=1} H(\hat{\Phi}; \mathbf{X}_g)$ over the full target population.

A.2 Intuition for Propensity Score Weighting

In this section, we provide a simple derivation of the weighting scheme that we propose to obtain the ATE from the switchers sample by introducing a concrete example in which the treatment effect is determined by one discrete covariate, x , and in which there are only few groups in the switcher sample. For ease of exposition, we refer to groups as families and units within groups as kids.

A.2.1 Thought Experiment

Suppose that the target population is comprised of 75% black individuals and 25% white individuals. The switchers sample has 1 white family with 3 kids and 2 black families with 3 and 5 kids, respectively. Thus, to be representative of the target population, the white family should be given a weight of 25%. The share for each black family is proportional to the number of individuals in the family, normalized so that the total share across the two families is 75%. Thus, the first family should be given a weight of $0.75 \cdot \frac{3}{8}$, and the second family should be given a weight of $0.75 \cdot \frac{5}{8}$.

A.2.2 Notation

Under the setup above, the weight that should be given to a switcher family g where all individuals have race $x_i = x$, can be written as:

$$s_{gx} = \frac{\sum(\mathbf{1}(x_i = x)|T_{g(i)} = 1)}{\sum T_{g(i)}} \cdot \frac{\sum \mathbf{1}(g(i) = g)}{\sum(S_{g(i)}|x_i = x)} \quad (26)$$

The first term, $\frac{\sum(\mathbf{1}(x_i = x)|T_{g(i)} = 1)}{\sum T_{g(i)}}$, gives the share of individuals in the target population with race x . The second term, $\frac{\sum \mathbf{1}(g(i) = g)}{\sum(S_{g(i)}|x_i = x)}$ gives the size of family g as a proportion of the switcher sample with race x .

Equivalently,

$$s_{gx} = Pr(x_i = x|T_{g(i)} = 1) \cdot \frac{\sum \mathbf{1}(g(i) = g)}{Pr(x_i = x|S_{g(i)} = 1) \cdot \sum S_{g(i)}} \quad (27)$$

$$= \frac{Pr(x_i = x|T_{g(i)} = 1)}{Pr(x_i = x|S_{g(i)} = 1)} \cdot Pr(g(i) = g|S_{g(i)} = 1) \quad (28)$$

A.2.3 Estimation

1. We obtain an estimate of $\hat{Q}_x = Pr(T_{g(i)} = 1|x_i = x)$ as fitted values from a regression of T on X .

This is equal to $\frac{Pr(x_i = x|T_{g(i)} = 1) \cdot Pr(T_{g(i)} = 1)}{Pr(x_i = x)}$ by Bayes rule.

2. We obtain an estimate of $\hat{P}_x = Pr(S_{g(i)} = 1|x_i = x)$ as fitted values from a regression of S on X .

This is equal to $\frac{Pr(x_i=x|S_g=1) \cdot Pr(S_{g(i)}=1)}{Pr(x_i=x)}$ by Bayes rule. The ratio of (1) and (2) is $\frac{Pr(x_i=x|T_{g(i)}=1)}{Pr(x_i=x|S_{g(i)}=1)} \cdot \frac{Pr(T_{g(i)}=1)}{Pr(S_{g(i)}=1)}$.

3. To get s_{gx} , we need to multiply this ratio by $Pr(g(i) = g|S_{g(i)} = 1)$ and divide by $\frac{Pr(T_{g(i)}=1)}{Pr(S_{g(i)}=1)}$.

We then normalize the weights, which gives $s_{gx} = \frac{\frac{Q_x \cdot n_g}{P_x}}{\sum_{g \in G_S} \frac{Q_x}{P_x} \cdot n_g}$

A.3 Extension to Unit i Covariates

A.3.1 Modified Assumptions

We begin with a simplification of the model, in which outcomes are a linear function of treatment, individual covariates, and additively separable individual error terms: $Y_{ig} = \delta_g \cdot D_i + \beta \cdot C_i + \alpha_g + (u_{ig} \cdot D_i + \epsilon_{ig})$. We assume constant coefficients on C_i , and require the systematic part of the treatment effect to be constant within a group. There can also be an idiosyncratic component of the treatment effect, denoted by u_{ig} . We also now allow for individual covariates to enter into the propensity to be in the switching or target populations, respectively, as: $P_{X,C} = Pr[S_{g(i)} = 1|X_g, C_i]$ and $Q_{X,C} = Pr[T_i = 1|X_g, C_i]$. The IPW weights therefore vary at the individual level: $w_i^t = \frac{Q_{X,C}}{P_{X,C}} \cdot \frac{Pr(S=1)}{Pr(T=1)}$.

The assumptions from earlier now must now be modified slightly to recover the ATE. Assumption 1, the conditional fixed effects assumption now requires $\epsilon_{ig}, u_{ig} \perp D_i|C_i, \alpha_g$. This gives that $E[\hat{\delta}_g] = \delta_g$.

Assumptions 2--4 and Proposition 1 will carry forward with these redefined terms.

This model can be adapted to allow treatment effects to vary with individually-varying covariates, such as gender of the treated individual. One approach is to re-write the treatment effects as a function of a group-level measure of these covariates (e.g. the share of Head Start participants in the group that are female). Then our main estimator framework applies. These group-level measures are included in the prediction of P_x and Q_x which should be sufficient to recover the ATE.

An alternative approach is to add a $\beta_C \cdot (C_i D_i)$ term to the estimating equation. Our principal re-weighting method can be used to estimate the average of the δ_g for the target population; and this average could then be added to $\hat{\beta}_C$ times the average covariate C_i value for the target population to arrive at an estimate for the ATE for this group.

A.3.2 Defining residual switchers

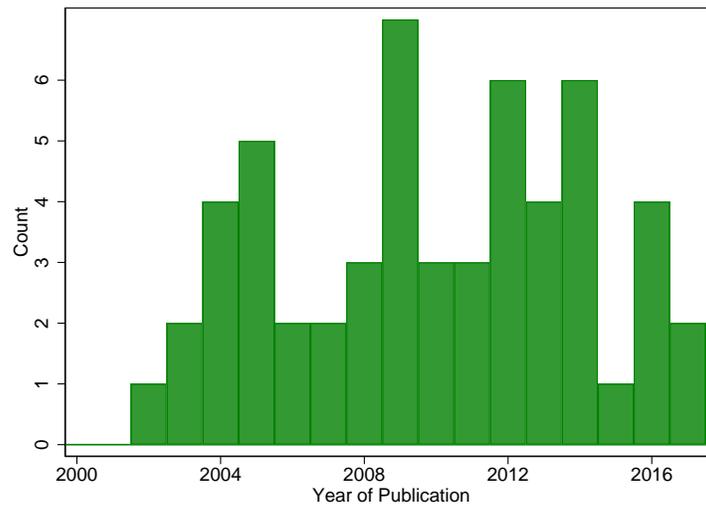
Consider the “deviations from group means” projection matrix, $M = I - H(H'H)^{-1}H$, with H a matrix of dummy variables for group membership: $H[i, j] = 1$ if unit i is a member of group j , and 0 otherwise. Let $\tilde{D} = M \cdot D$ be deviations in treatment from group means. Basic switcher groups (“true switchers”) are defined by having within-group variation in treatment: $V_g := Var(\tilde{D}_i|g(i) = g) > 0$. Next consider the residual-maker matrix projecting on covariates C after taking deviations from group means, $L = I - (M \cdot C)(C' \cdot M \cdot C)(C' \cdot M)$, and let $\tilde{\tilde{D}} = L \cdot M \cdot D$. With unit-varying covariates C_i , the variation that identifies the treatment effects is $V_{g,C} :=$

$Var(\ddot{D}|g(i) = g)$. For some groups g , it could be the case that $V_g = 0$, and also $V_{g,C} > 0$. We call these groups residual switchers.

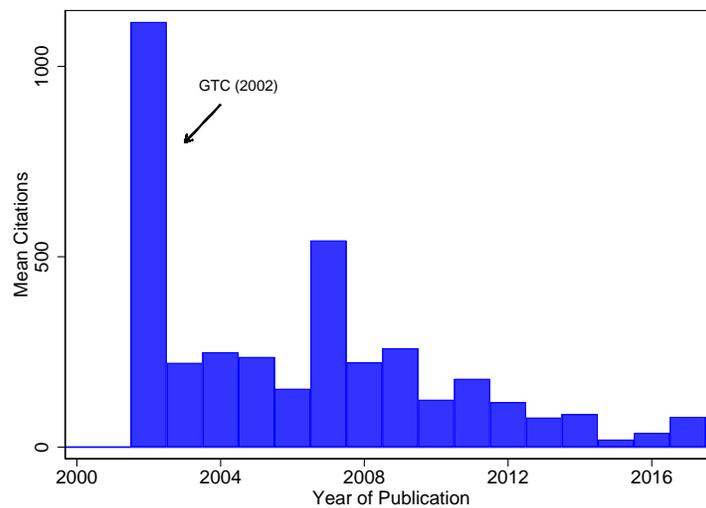
B Supplementary Figures and Tables

Figure B.1: Popularity of Family Fixed Effects Articles

(a) Publications by Year



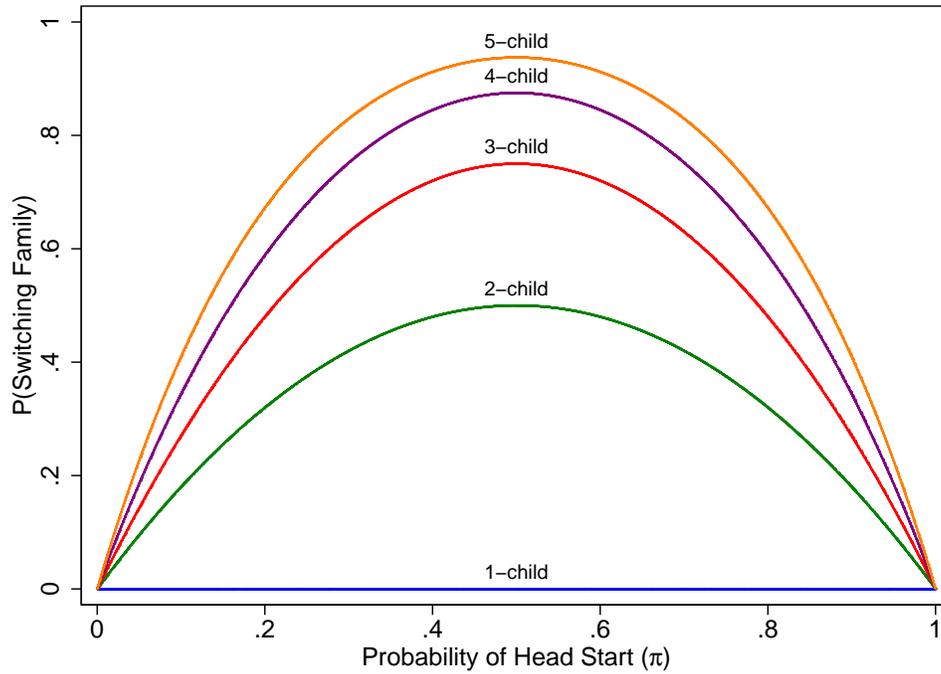
(b) Average Citations by Year of Publication



Notes: These figures display the data from our survey of FFE papers published from January 2000 to May 2017 in 11 leading journals that publish applied microeconomics articles. Figure (a) plots the number of FFE articles published in each year, and Figure (b) plots the average number of Google Scholar citations, as of May 2019, among the articles published in a given year.

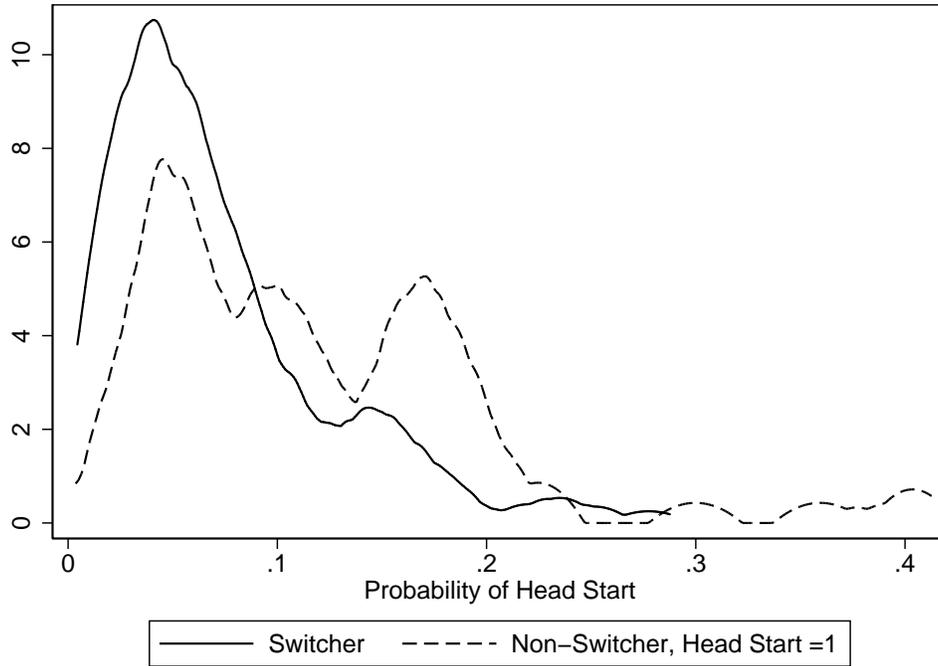
Figure B.2: Illustrative Model of the Role of Family Size in Switching

$$Pr(S_g = 1) = 1 - (1 - \pi)^{n_g} - \pi^{n_g}$$



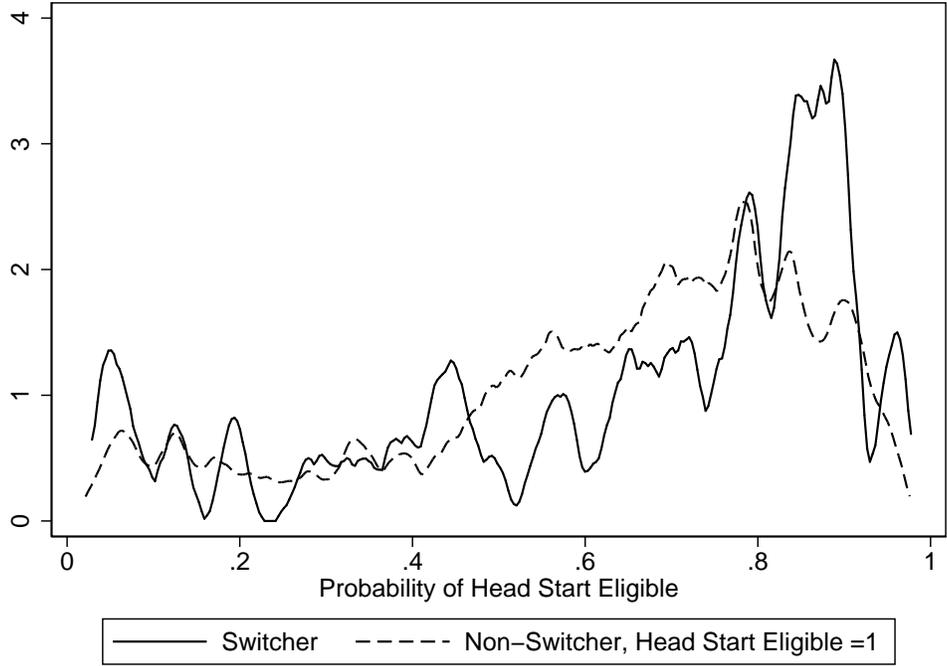
Notes: This figure plots the theoretical function: $P(HSSwitchingFamily) = 1 - (1 - \pi)^{n_g} - \pi^{n_g}$, where n_g is the number of children in a family and π is the probability of attending Head Start, for 2-, 3-, 4-, and 5 (plus)- child families.

Figure B.3: Examining P-Score Overlap: Predicted Probability of Being in Head Start (PSID White Sample)



Notes: This figure shows kernel density plots (bandwidth = 0.01) of the predicted probability of being a Head Start participant for switchers and non-switchers that are Head Start participants. The sample consists of white individuals in the PSID. There are 4 non-switchers (5%) who have a probability of being in Head Start that is outside of the support of the switcher sample.

Figure B.4: Examining P-Score Overlap: Predicted Probability of Being Head-Start-Eligible (PSID White Sample)



Notes: This figure shows kernel density plots (bandwidth = 0.01) of the predicted probability of being Head Start eligible for switchers and non-switchers that are Head-Start-eligible. The sample consists of white individuals in the PSID.

Table B.1: Head Start Attendance and Within-Family Variation in Attendance by Family Size (PSID)

	Number of Children in Family:				
	2	3	4	5+	Total
Share of Family in Head Start (π)	0.157	0.222	0.195	0.206	0.182
Share with Switching	0.121	0.202	0.242	0.471	0.174
All Participants in HS in Family	0.096	0.125	0.093	0.049	0.102
No Participants in HS in Family	0.783	0.672	0.665	0.480	0.724

Notes: This table shows the sources of switching by family size. The first two rows show the likelihood of attending Head Start by family size and the likelihood of having variation in Head Start within a family (switching). The final two rows examines whether differences in rates of switching across family sizes are attributable to variation across family sizes in having all children attend Head Start (row 3) or variation in having no children attend Head Start (row 4).

Table B.2: Switchers are Different than Head-Start-Participating Non-Switchers

	(1)	(2)
	Switch	Non-Switch, HS=1
Fraction female	0.562	0.530
Fraction African-American	0.516	0.714***
Mother's yrs education	9.283	10.304***
Father's yrs education	9.190	10.445***
Had a single mother at age 4	0.252	0.332***
Family income (age 3-6) (CPI adjusted)	31809	25915***
Mother employed, age 0	0.508	0.542
Mother employed, age 1	0.517	0.541
Mother employed, age 2	0.536	0.596**
Household size at age 4	5.487	4.353***
Fraction low birth weight	0.077	0.135***

Notes: This table presents summary statistics for Head-Start-switching families (column 1) and non-switching families that attended Head Start (column 2) together with stars indicating the p-value for the test of equality of the clums. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.3: Demographic Characteristics of Head Start Sample (PSID)

	All	Head Start	No Head Start	Sibling Sample
Head Start	0.076	1.000	0.000	0.073
Other preschool	0.282	0.000	0.305	0.259
Fraction African-American	0.150	0.618	0.111	0.154
Fraction female	0.504	0.548	0.501	0.501
Fraction low birth weight	0.060	0.114	0.056	0.061
Had a single mother at age 4	0.112	0.296	0.091	0.103
Fraction whose mother completed hs	0.717	0.632	0.724	0.689
Fraction whose father completed hs	0.683	0.557	0.692	0.654
Fraction eldest child in family	0.368	0.341	0.371	0.339
Age in 1995	23.830 (9.84)	18.605 (7.76)	24.262 (9.87)	25.063 (10.06)
Mother's yrs education	11.116 (2.76)	10.208 (2.32)	11.190 (2.78)	10.942 (2.81)
Father's yrs education	11.238 (3.23)	10.159 (2.70)	11.314 (3.25)	11.076 (3.35)
Family income (age 3-6) (CPI adjusted)	50339 (35814.01)	28553 (17212.32)	52719 (36509.36)	50973 (37315.99)
Household size at age 4	4.535 (1.68)	4.814 (2.06)	4.504 (1.63)	4.778 (1.64)
Observations	7363	1345	6018	5355

Notes: This table shows the mean demographic characteristics of the sample, weighted to be representative of 1995 population; see text for details. Standard deviations, shown in parentheses, are omitted for binary variables. CPI-adjusted income reported in 1999 dollars. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.4: Outcomes of Interest for Head Start Sample (PSID)

	All	Head Start	No Head Start	Sibling Sample
Fraction completed hs	0.913	0.878	0.916	0.912
Fraction attended some college	0.531	0.428	0.539	0.532
Fraction not booked/charged with crime	0.899	0.889	0.900	0.898
Avg. Earnings age 23-25 (CPI adjusted)	20410 (24927)	14391 (12000)	20818 (25517)	20633 (26547)
Economic Sufficiency Index at 30	0.094 (1.03)	-0.601 (1.05)	0.151 (1.01)	0.096 (1.03)
Economic Sufficiency Index at 40	0.020 (1.01)	-0.532 (0.95)	0.053 (1.01)	0.025 (1.04)
Good Health Index at 30	0.004 (1.03)	-0.558 (1.26)	0.050 (0.99)	0.017 (0.99)
Good Health Index at 40	0.011 (1.01)	-0.486 (1.25)	0.033 (1.00)	0.015 (0.96)
Observations	7363	1345	6018	5355

Notes: This table shows the means for the main outcomes of interest, weighted to be representative of 1995 population; see text for details. Note that the fraction not booked/charged with a crime restricted to individuals that responded to the PSID in 1995 who were between the ages of 16 and 50 in that year. CPI-adjusted income reported in 1999 dollars. Standard deviations, shown in parentheses, are omitted for binary variables. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.5: Summary Statistics for Inputs to Summary Indices (PSID)

	All	Head Start	No Head Start	Sibling Sample
<i>Inputs to Economic Sufficiency Index, 30</i>				
Ever on AFDC/TANF by age 30	0.062	0.220	0.049	0.060
Fraction of last 5 yrs on Food Stamps/SNAP, age 30	0.064 (0.20)	0.151 (0.30)	0.056 (0.19)	0.071 (0.22)
ln(mean earnings in last 5 years), age 30	9.661 (1.06)	9.415 (0.91)	9.676 (1.07)	9.659 (1.07)
Fraction of last 5 yrs with positive earnings, age 30	0.895 (0.25)	0.887 (0.26)	0.896 (0.25)	0.898 (0.25)
Fraction of last 5 yrs ever unemployed, age 30	0.146 (0.24)	0.173 (0.27)	0.144 (0.23)	0.150 (0.24)
Mean Inc. Rel. Pov. in last 5 years, age 30	385.831 (305.98)	233.796 (155.44)	396.729 (311.18)	385.933 (291.36)
Fraction completed college	0.209	0.073	0.220	0.220
<i>Inputs to Economic Sufficiency Index, 40</i>				
Ever on AFDC/TANF by age 40	0.068	0.163	0.062	0.067
Fraction of last 5 yrs on Food Stamps/SNAP, age 40	0.043 (0.16)	0.098 (0.25)	0.040 (0.16)	0.043 (0.16)
ln(mean earnings in last 5 years), age 40	9.962 (1.15)	9.779 (0.90)	9.968 (1.16)	9.957 (1.15)
Fraction of last 5 yrs with positive earnings, age 40	0.850 (0.31)	0.867 (0.29)	0.849 (0.31)	0.849 (0.31)
Fraction of last 5 yrs ever unemployed, age 40	0.094 (0.20)	0.122 (0.24)	0.093 (0.19)	0.098 (0.20)
Mean Inc. Rel. Pov. in last 5 years, age 40	436.769 (366.03)	281.489 (183.89)	443.338 (370.36)	434.280 (361.58)
Fraction of last 5 yrs owned home, age 40	0.500 (0.44)	0.287 (0.42)	0.510 (0.44)	0.522 (0.44)
<i>Inputs to Good Health Index, 30</i>				
Fraction of last 5 yrs smoked less than 1 cigarette/day, age 30	0.745 (0.41)	0.668 (0.45)	0.753 (0.41)	0.755 (0.40)
Fraction of last 5 yrs reported good or better health, age 30	0.948 (0.17)	0.903 (0.24)	0.951 (0.17)	0.950 (0.17)
Mean BMI in last 5 years, age 30	26.569 (6.68)	28.766 (6.74)	26.333 (6.63)	26.615 (6.85)
<i>Inputs to Good Health Index, 40</i>				
Fraction of last 5 yrs smoked less than 1 cigarette/day, age 40	0.738 (0.42)	0.714 (0.44)	0.739 (0.42)	0.728 (0.42)
Fraction of last 5 yrs reported good or better health, age 40	0.919 (0.22)	0.871 (0.29)	0.921 (0.22)	0.922 (0.22)
Mean BMI in last 5 years, age 40	27.504 (5.92)	30.191 (7.42)	27.327 (5.77)	27.433 (5.85)
Observations	7363	1345	6018	5355

Notes: Weighted to be representative of 1995 population; see text for details. SD, in parentheses, are omitted for binary variables. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.6: N's for Control Covariates (PSID)

	All	Head Start	No Head Start	Sibling Sample
Head Start	7372	1354	6018	5361
Other preschool	7372	1354	6018	5361
Fraction African-American	7372	1354	6018	5361
Fraction female	7372	1354	6018	5361
Fraction low birth weight	5366	970	4396	4555
Had a single mother at age 4	6678	1285	5393	4672
Fraction whose mother completed hs	7231	1332	5899	5360
Fraction whose father completed hs	6596	1034	5562	4875
Fraction eldest child in family	7372	1354	6018	5361
Age in 1995	7372	1354	6018	5361
Mother's yrs education	7223	1331	5892	5356
Father's yrs education	6596	1034	5562	4875
Family income (age 3-6) (CPI adjusted)	6086	1145	4941	4338
Household size at age 4	6251	1187	5064	4420
Observations	7372	1354	6018	5361

Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.7: N's for Main Outcomes (PSID)

	All	Head Start	No Head Start	Sibling Sample
Fraction completed hs	7372	1354	6018	5361
Fraction attended some college	7372	1354	6018	5361
Fraction not booked/charged with crime	5005	802	4203	3591
Avg. Earnings age 23-25 (CPI adjusted)	4866	783	4083	3675
Economic Sufficiency Index at 30	7372	1354	6018	5361
Economic Sufficiency Index at 40	4085	613	3472	2845
Good Health Index at 30	4749	791	3958	3600
Good Health Index at 40	2228	312	1916	1673
Observations	7372	1354	6018	5361

Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.8: N's for Auxiliary Outcomes (PSID)

	All	Head Start	No Head Start	Sibling Sample
<i>Inputs to Economic Sufficiency Index, 30</i>				
Ever on AFDC/TANF by age 30	7372	1354	6018	5361
Fraction of last 5 yrs on Food Stamps/SNAP, age 30	4186	713	3473	2805
ln(mean earnings in last 5 years), age 30	4202	620	3582	3159
Fraction of last 5 yrs with positive earnings, age 30	4378	656	3722	3295
Fraction of last 5 yrs ever unemployed, age 30	4259	634	3625	3184
Mean Inc. Rel. Pov. in last 5 years, age 30	5293	891	4402	4068
Fraction completed college	7372	1354	6018	5361
<i>Inputs to Economic Sufficiency Index, 40</i>				
Ever on AFDC/TANF by age 40	4085	613	3472	2845
Fraction of last 5 yrs on Food Stamps/SNAP, age 40	1972	250	1722	1423
ln(mean earnings in last 5 years), age 40	1695	221	1474	1266
Fraction of last 5 yrs with positive earnings, age 40	1829	236	1593	1369
Fraction of last 5 yrs ever unemployed, age 40	1825	236	1589	1365
Mean Inc. Rel. Pov. in last 5 years, age 40	2152	296	1856	1613
Fraction of last 5 yrs owned home, age 40	2292	290	2002	1625
<i>Inputs to Good Health Index, 30</i>				
Fraction of last 5 yrs smoked less than 1 cigarette/day, age 30	2267	385	1882	1742
Fraction of last 5 yrs reported good or better health, age 30	3763	579	3184	2806
Mean BMI in last 5 years, age 30	3248	587	2661	2528
<i>Inputs to Good Health Index, 40</i>				
Fraction of last 5 yrs smoked less than 1 cigarette/day, age 40	1280	182	1098	930
Fraction of last 5 yrs reported good or better health, age 40	1463	182	1281	1116
Mean BMI in last 5 years, age 40	2037	307	1730	1486
Observations	7372	1354	6018	5361

Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.9: Effect of Head Start on Pre-Head-Start Outcomes (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>Low birth weight</i>					
Head Start	0.040*	0.045*	-0.016	-0.018	-0.029
	(0.021)	(0.023)	(0.026)	(0.033)	(0.042)
Other preschool	0.003	0.003	-0.012	-0.056**	-0.003
	(0.012)	(0.013)	(0.023)	(0.027)	(0.027)
Observations	5366	4555	4500	1872	2622
<i>Disabled</i>					
Head Start	-0.006	-0.017	-0.010	-0.016	-0.006
	(0.027)	(0.030)	(0.030)	(0.036)	(0.051)
Other preschool	0.018	0.018	0.021	0.032	0.017
	(0.019)	(0.022)	(0.028)	(0.049)	(0.032)
Observations	3516	2955	2661	1102	1555
<i>Single mom at age 4</i>					
Head Start	0.020	0.025	0.027	-0.007	0.051
	(0.015)	(0.020)	(0.024)	(0.022)	(0.040)
Other preschool	0.022**	0.020*	0.008	0.006	0.011
	(0.009)	(0.011)	(0.017)	(0.031)	(0.018)
Observations	6678	4672	4467	1939	2522
<i>Family income (age 1) (CPI adjusted)</i>					
Head Start	0.000**	-0.000***	0.000	-0.000	-0.000
	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)
Other preschool	-0.000***	-0.000***	-0.000	0.000	-0.000
	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)
Observations	6219	4313	4023	1719	2298
<i>Family income (age 2) (CPI adjusted)</i>					
Head Start	0.000	-0.000	-0.000	0.000	-0.000
	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)
Other preschool	-0.000	-0.000	-0.000	-0.000	-0.000
	(0.000)	(0.000)	(0.000)	(0.000)	(0.000)
Observations	6274	4391	4151	1757	2388
<i>Mom working at age 1</i>					
Head Start	0.001	0.011	0.049	0.002	0.080
	(0.018)	(0.022)	(0.039)	(0.033)	(0.073)
Other preschool	-0.001	-0.002	-0.017	-0.078*	-0.014
	(0.013)	(0.016)	(0.030)	(0.043)	(0.034)
Observations	6219	4313	4023	1719	2298
<i>Mom working at age 2</i>					
Head Start	0.025	0.028	-0.041	-0.008	-0.077
	(0.021)	(0.023)	(0.040)	(0.036)	(0.073)
Other preschool	0.026*	0.032*	0.015	-0.013	0.017
	(0.015)	(0.018)	(0.031)	(0.044)	(0.036)
Observations	6274	4391	4151	1757	2388

Notes: Weighted to be representative of 1995 population; see text for details. SE clustered at 1968 family id in columns 1 and 2 and at mother id level otherwise. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.10: Test of Conditional Ignorability Assumption:

Do Individuals in the Target Population Have Differential Treatment Effects?

	Target= Eligible	Target= Participants
<i>Some College (Whites, PSID)</i>		
In Target	–	-0.084 (0.124)
Observations	–	213
<i>Economic Sufficiency Index (Whites, PSID)</i>		
In Target	–	0.145 (0.199)
Observations	–	213
<i>High School Graduation (NLSY)</i>		
In Target	0.014 (0.069)	-0.101 (0.062)
Observations	467	581
<i>Idle (NLSY)</i>		
In Target	0.036 (0.085)	0.033 (0.097)
Observations	467	581
<i>Learning Disability (NLSY)</i>		
In Target	-0.025 (0.046)	-0.095 (0.063)
Observations	467	581
<i>Poor Health (NLSY)</i>		
In Target	0.015 (0.068)	-0.026 (0.058)
Observations	467	581

Notes: Each cell of this table shows an estimate from a regression of the family-specific impact of Head Start on an indicator for whether an individual is in the target population. Regressions are weighted for balance on observables: target individuals are assigned a weight of $Pr(T_g = 1, S_g = 1 | X_{ig})$ and non-target individuals are assigned a weight of $Pr(T_g = 0, S_g = 1 | X_{ig})$. The first two panels use data from the PSID white sample, and the final four panels use data from the CNLSY. There are no individuals in the PSID white switcher sample that are ineligible for Head Start, which causes us to have missing estimates (“–”).

Table B.11: Test of Conditional Ignorability Assumption:

Do Gains depend on the Fraction of Children attending Head Start? (PSID, College)

Head Start x 2 child family	-0.126 (0.099)	0.103 (0.191)
Head Start x 3 child family	0.395* (0.221)	0.368** (0.174)
Head Start x 4 child family	0.377 (0.271)	0.441** (0.186)
Head Start x 5plus child family	0.522** (0.260)	0.429*** (0.139)
Head Start x 3 child family x Fraction	-0.515 (0.431)	
Head Start x 4 child family x Fraction	-0.304 (0.524)	
Head Start x 5plus child family x Fraction	-0.979 (1.010)	
Head Start x Fraction attend		-0.457 (0.328)
Observations	2986	2986

Notes: This table shows estimates from regressions of some college on an indicator for Head Start, an interaction of Head Start with indicators for family size, and an interaction of Head Start with the fraction of children in the family that attended (column 2) or with the fraction of children in the family that attended and indicators for family size (column 1). The data come from the PSID white sample. Standard errors are clustered on mother id.

Table B.12: Test of Conditional Ignorability Assumption:
Do Gains depend on the Fraction of Children attending Head Start? (CNLSY)

	HS	Idle	LD	Poor Hlth	HS	Idle	LD	Poor Hlth
Head Start x 2 child family	0.059 (0.050)	-0.075 (0.060)	-0.019 (0.025)	-0.098** (0.048)	0.201* (0.117)	-0.078 (0.145)	0.033 (0.090)	0.008 (0.116)
Head Start x 3 child family	0.040 (0.191)	-0.050 (0.208)	-0.049 (0.137)	0.106 (0.183)	0.186 (0.127)	-0.004 (0.151)	-0.021 (0.102)	0.048 (0.122)
Head Start x 4 child family	0.394** (0.182)	-0.059 (0.268)	-0.215* (0.112)	-0.147 (0.164)	0.280** (0.124)	-0.066 (0.169)	0.010 (0.094)	0.111 (0.122)
Head Start x 5plus child family	0.666*** (0.171)	-0.332 (0.312)	0.266** (0.114)	0.116 (0.236)	0.409*** (0.127)	-0.319** (0.154)	-0.124 (0.083)	0.037 (0.129)
Head Start x 3 child family x Fraction	0.000 (0.340)	0.097 (0.387)	-0.049 (0.243)	-0.329 (0.345)				
Head Start x 4 child family x Fraction	-0.516 (0.396)	-0.010 (0.516)	0.344 (0.241)	0.301 (0.284)				
Head Start x 5plus child family x Fraction	-0.995*** (0.369)	0.043 (0.839)	-1.196*** (0.411)	-0.444 (0.564)				
Head Start x Fraction					-0.286 (0.209)	0.006 (0.264)	-0.103 (0.170)	-0.214 (0.209)
Observations	1251	1251	1247	1251	1251	1251	1247	1251

Notes: This table shows estimates from regressions of outcomes (shown in the headers) on an indicator for Head Start, an interaction of Head Start with indicators for family size, and an interaction of Head Start with the fraction of children in the family that attended (columns 5-8) or with the fraction of children in the family that attended and indicators for family size (columns 1-4). The data come from the CNLSY sample. Standard errors are clustered on mother id.

Table B.13: Additional Estimates for Representative White Populations (PSID)

Using Post-Regression Reweighting Method

	FFE		Reweighted, Target =		
	GTC	Expand Sample	HS Eligible	Participants	Siblings
<i>A. High School Graduation</i>					
Head Start	0.203** (0.098)	-0.015 (0.045)	-0.034 (0.046)	-0.030 (0.052)	-0.029 (0.050)
Y Mean in Target	–	0.921	0.852	0.848	0.921
<i>B. Good Health Index, Age 30</i>					
Head Start	–	-0.265 (0.249)	-0.253 (0.225)	-0.439 (0.270)	-0.162 (0.254)
Y Mean in Target	–	0.074	-0.061	-0.583	0.074

Notes: Columns 1 and 2 of this table show the FFE estimated impacts of Head Start from GTC (2002) and using our expanded sample for completion of high school (panel A) and the Good Health Index at age 30 (panel B). The remaining columns present reweighted estimates of the effect of Head Start for three target populations (shown in the column header) using the post-regression reweighting procedure described in the text. "–" is used to indicate that the information is not available. Sample size is N=2,986 for the expanded sample in panel A, and 1,959 for the expanded sample in panel B, and 1,036 for GTC. Standard errors obtained by bootstrapping. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.14: Head Start Impact for Representative Black Eligible, Participants, and Siblings (PSID)
Using Post-Regression Reweighting Method

	FFE		Reweighted, Target =		
	GTC	Expand Sample	HS Eligible	Participants	Siblings
<i>A. High School Graduation</i>					
Head Start	-0.025 (0.065)	-0.024 (0.031)	-0.017 (0.030)	-0.014 (0.031)	-0.015 (0.029)
Y Mean in Target	–	0.862	0.854	0.896	0.862
<i>B. Some College</i>					
Head Start	0.023 (0.066)	-0.016 (0.036)	-0.031 (0.038)	-0.028 (0.041)	-0.032 (0.038)
Y Mean in Target	–	0.396	0.376	0.423	0.396
<i>C. Economic Sufficiency Index, Age 30</i>					
Head Start	–	-0.023 (0.102)	-0.190** (0.078)	-0.211*** (0.082)	-0.167** (0.077)
Y Mean in Target	–	-0.552	-0.626	-0.674	-0.552
<i>D. Good Health Index, Age 30</i>					
Head Start	–	-0.265 (0.249)	0.052 (0.129)	0.062 (0.136)	0.033 (0.124)
Y Mean in Target	–	-0.357	-0.381	-0.539	-0.357

Notes: Columns 1 and 2 of this table show the FFE estimated impacts of Head Start from GTC (2002) and using our expanded sample for completion of high school (panel A) and the Good Health Index at age 30 (panel B). The remaining columns present reweighted estimates of the effect of Head Start for three target populations (shown in the column header) using the post-regression reweighting procedure described in the text. "–" is used to indicate that the information is not available. Sample size is N=2,369 for the expanded sample in panels A, B, and C, and 1,150 for the expanded sample in Panel D, and 762 for GTC. Standard errors obtained by bootstrapping. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table B.15: Head Start Impact for Representative Eligible Children, Participants, and Siblings

Using Reweighting with Regression Extrapolation to Singletons

	Reweighted ATE, Target =	
	HS Eligible	Participants
<i>A. Some College (PSID)</i>		
Head Start	0.083 (0.102)	0.009 (0.109)
Y Mean in Target	0.387	0.437
<i>B. Economic Sufficiency Index, Age 30 (PSID)</i>		
Head Start	-0.056 (0.139)	0.003 (0.229)
Y Mean in Target	-0.198	-0.485
<i>C. High School Graduation (CNLSY)</i>		
Head Start	0.041 (0.032)	0.047 (0.030)
Y Mean in Target	0.734	0.766
<i>D. Idle (CNLSY)</i>		
Head Start	-0.059 (0.038)	-0.057 (0.037)
Y Mean in Target	0.221	0.201
<i>E. Learning Disability (CNLSY)</i>		
Head Start	-0.034* (0.020)	-0.041** (0.017)
Y Mean in Target	0.055	0.041
<i>F. Poor Health (CNLSY)</i>		
Head Start	-0.059* (0.028)	-0.065** (0.030)
Y Mean in Target	0.098	0.074

Notes: Columns 1 and 2 present reweighted estimates of the effect of Head Start for the Head-Start-eligible and Head-Start-participant target populations, where the treatment effects for singletons are extrapolated from switchers using OLS. Sample size is N=2,986 for the expanded sample, and 1,036 for GTC. Standard errors obtained using bootstrap. * p < .10, ** p < .05, *** p < .01.

Table B.16: FFE Estimates Reweighted using Gibbons, Suarez, Urbancic (2018) Method

	FFE	GSU (2018) Reweight
	Baseline	Switchers
<i>A. Some College (PSID)</i>		
Head Start	0.120** (0.053)	0.134*** (0.053)
<i>B. Economic Sufficiency Index, Age 30 (PSID)</i>		
Head Start	-0.023 (0.102)	-0.081 (0.094)
<i>C. High School Graduation (CNLSY)</i>		
Head Start	0.085*** (0.030)	0.084*** (0.027)
<i>D. Idle (CNLSY)</i>		
Head Start	-0.072* (0.037)	-0.068** (0.034)
<i>E. Learning Disability (CNLSY)</i>		
Head Start	-0.059*** (0.020)	-0.053*** (0.019)
<i>F. Poor Health (CNLSY)</i>		
Head Start	-0.069*** (0.026)	-0.059** (0.025)

Notes: Column 1 reprints the FFE estimate using our expanded sample for PSID outcomes and using our replication sample for CNLSY outcomes. Column 2 presents the estimate weighting family-level estimates by the sample share, as suggested in Gibbons, Urbancic, Suarez Serrato (2018). This “undoes” the conditional variance weighting of FFE, and produces an estimate that is interpretable as the ATE for switchers. Sample size is N=2,986 for the expanded sample. Standard errors are clustered on mother id. * p < .10, ** p < .05, *** p < .01.

Table B.17: Horse Race between Family Size and Index of Non-Family-Size Covariates (PSID White Sample)

	x Fam Size	x Index	Horse Race
<i>Index = Predicted Head Start</i>			
Head Start	0.025 (0.063)	0.073 (0.069)	0.008 (0.072)
Head Start x 4plus child family	0.281** (0.112)		0.250** (0.105)
Head Start x Tercile 1 Predicted Head Start		-0.049 (0.094)	-0.116 (0.101)
Head Start x Tercile 2 Predicted Head Start		0.212* (0.113)	0.125 (0.111)
Observations	2986	2986	2986
<i>Index = Predicted Finish College</i>			
Head Start	0.025 (0.063)	-0.088 (0.083)	-0.130 (0.100)
Head Start x 4plus child family	0.281** (0.112)		0.266** (0.112)
Head Start x Tercile 1 Predicted Finish College		0.237** (0.112)	0.155 (0.121)
Head Start x Tercile 2 Predicted Finish College		0.260** (0.131)	0.207 (0.142)
Observations	2986	2986	2986

This table shows estimates from a FFE regression of attainment of some college on an indicator for attendance of Head Start, and an indicator for having a family with 4 or more children (Column 1), dummies for terciles of an index of predicted Head Start attendance (Column 2, Panel A), dummies for terciles of an index of the predicted likelihood of finishing college (Column 2, Panel B), and the combination of family size indicator and terciles of the index (Column 3). The predicted Head Start (finish college) index is created by regressing Head Start attendance (finish college) on all of the control variables in the PSID analysis, except for the household size variable.

C Supplementary PSID FFE Results

In this section, we discuss additional FFE results obtained using our expanded PSID sample.

We present the FFE results for the economic and health indices measured at age 40, together with the indices at age 30 for comparison, in Table C.1. Overall, the results suggest little support for a positive long term effect of Head Start. We come to the same conclusions when we aggregate the inputs using principal components analysis (see Table C.2). Our overall conclusions are not changed importantly by looking at specific outcomes or subsamples. We have also estimated regressions for each of the inputs to the economic and health indices, which we include in Tables C.3, and C.4, C.5, C.6. Table C.7 shows the regression results for the additional outcomes analyzed in GTC, earnings between ages 23 to 25, and not having committed a crime. Across these tables, there is no systematic evidence that Head Start impacts long term outcomes.⁵⁹

Motivated by the prior findings of differential effects by gender in Carneiro and Ginja (2014); Deming (2009), in Table C.8 we look to see whether our mean results are obscuring this form of heterogeneity in our setting. Curiously, we find some evidence of significant negative effects of Head Start among men, in particular for health and economic outcomes at age 40. On the other hand, we find a positive and significant effect of Head Start on the probability that men attain some college. The effects estimated for women are never individually significant, but also not statistically different from men for many outcomes as indicated by the p-value of the difference in the table. The one exception is for economic outcomes observed at age 40, where women are found to have significantly better returns to Head Start participation than observed for men.

Another source of heterogeneity which could generate a discrepancy between our results and GTC is the fact that our sample includes later (younger) cohorts, whose Head Start experience may differ from earlier participants. In Table C.9, we find some support for a decreasing impact of Head Start across cohorts for the age 40 indices, but also find a larger improvement in the health index at age 30 for more recent cohorts. Thus, this does not appear to reconcile our findings.⁶⁰

⁵⁹Moreover, while we find a significant increase in attainment of some college, when we examine the outcome of college completion, we obtain insignificant negative point estimates for the pooled sample (beta = -0.033, se = 0.023), for black children (beta = -0.014, se = 0.018), and for white children (beta = -0.058, se = 0.043).

⁶⁰Moreover, when we instead use a binary indicator for more recent cohorts, we do not find a statistically significant difference in the impacts of Head Start, indicating that these results are sensitive to functional form assumptions.

Table C.1: Impact of Head Start on Economic Sufficiency Index and Good Health Index (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<u><i>Economic Sufficiency Index, age 30</i></u>					
Head Start	-0.147*** (0.043)	-0.117** (0.050)	-0.090 (0.064)	-0.117 (0.081)	-0.023 (0.102)
Other preschool	0.184*** (0.035)	0.181*** (0.040)	0.091 (0.062)	0.050 (0.109)	0.099 (0.072)
Mean Y	0.094	0.096	0.096	-0.552	0.213
Observations	7372	5361	5361	2369	2986
<u><i>Economic Sufficiency Index, age 40</i></u>					
Head Start	-0.080 (0.066)	-0.071 (0.077)	-0.059 (0.100)	-0.170 (0.134)	-0.081 (0.125)
Other preschool	0.112* (0.059)	0.085 (0.077)	0.043 (0.107)	-0.270 (0.223)	0.118 (0.122)
Mean Y	0.020	0.025	0.025	-0.670	0.142
Observations	4085	2845	2503	1065	1435
<u><i>Good Health Index, Age 30</i></u>					
Head Start	-0.349*** (0.058)	-0.320*** (0.064)	-0.148 (0.143)	0.024 (0.149)	-0.265 (0.249)
Other preschool	0.087** (0.038)	0.096** (0.045)	0.081 (0.076)	0.040 (0.159)	0.106 (0.084)
Mean Y	0.004	0.017	0.017	-0.357	0.074
Observations	4749	3600	3114	1150	1959
<u><i>Good Health Index, Age 40</i></u>					
Head Start	-0.201* (0.118)	-0.175 (0.141)	-0.147 (0.202)	0.031 (0.201)	-0.146 (0.393)
Other preschool	0.117 (0.094)	0.095 (0.115)	0.119 (0.130)	0.382* (0.210)	0.038 (0.150)
Mean Y	0.011	0.015	0.015	-0.290	0.062
Observations	2228	1673	1306	511	795

Notes: This table shows the estimates from regressions of either the Economic Sufficiency Index at age 30 (panel A), the Economic Sufficiency Index at age 40 (panel B), the Good Health Index at age 30 (panel C), or the Good Health Index at age 40 (panel D) on an indicator for participation in Head Start and control variables described in the text. Regressions are run on the whole sample (column 1), siblings (columns 2 and 3), black siblings (column 4) and white siblings (column 5). All columns include control variables, and columns 3, 4, and 5 include mother fixed effects. The Good Health Index includes measures of not smoking cigarettes, good self reported health and BMI, averaged over the previous 5 years. The Economic Sufficiency Index includes measures of high school graduation, attendance of some college, no receipt of Food Stamps/SNAP, no receipt of AFDC/TANF, average earnings, employment, and unemployment, averaged over the previous 5 years. Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and at mother id level otherwise. * $p < .10$, ** $p < .05$, *** $p < .01$. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.2: Effect of Head Start on Economic and Health Principal Components (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>Economic Sufficiency Principal Component, age 30</i>					
Head Start	-0.174*** (0.058)	-0.140** (0.068)	-0.100 (0.084)	-0.138 (0.109)	-0.031 (0.128)
Other preschool	0.295*** (0.051)	0.285*** (0.057)	0.150* (0.087)	0.071 (0.150)	0.166 (0.101)
Mean Y	0.154	0.160	0.160	-0.731	0.321
Observations	7372	5361	5361	2369	2986
<i>Economic Sufficiency Principal Component, age 40</i>					
Head Start	-0.113 (0.090)	-0.093 (0.106)	-0.082 (0.131)	-0.219 (0.180)	-0.127 (0.155)
Other preschool	0.209** (0.086)	0.173 (0.113)	0.091 (0.145)	-0.291 (0.296)	0.183 (0.167)
Mean Y	0.026	0.032	0.032	-0.968	0.199
Observations	4085	2845	2503	1065	1435
<i>Good Health Principal Component, Age 30</i>					
Head Start	-0.248*** (0.047)	-0.228*** (0.052)	-0.073 (0.121)	0.057 (0.131)	-0.159 (0.208)
Other preschool	0.070** (0.031)	0.069* (0.037)	0.063 (0.063)	0.033 (0.137)	0.083 (0.069)
Mean Y	0.003	0.013	0.013	-0.309	0.062
Observations	4749	3600	3114	1150	1959
<i>Good Health Principal Component, Age 40</i>					
Head Start	-0.143 (0.107)	-0.126 (0.128)	-0.101 (0.200)	0.044 (0.200)	-0.174 (0.400)
Other preschool	0.101 (0.089)	0.077 (0.110)	0.121 (0.104)	0.288 (0.221)	0.062 (0.117)
Mean Y	0.009	0.015	0.015	-0.259	0.056
Observations	2228	1673	1306	511	795

Notes:

Table C.3: Effect of Head Start on Inputs to Economic Sufficiency Index at age 30 (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>High School Graduate</i>					
Head Start	0.007 (0.018)	-0.002 (0.021)	-0.011 (0.026)	-0.024 (0.031)	-0.015 (0.045)
Other preschool	-0.002 (0.011)	-0.008 (0.014)	0.036* (0.021)	-0.012 (0.048)	0.046* (0.024)
Mean Y	0.913	0.912	0.912	0.862	0.921
Observations	7372	5361	5361	2369	2986
<i>Attended Some College</i>					
Head Start	0.038 (0.024)	0.039 (0.029)	0.046 (0.033)	-0.016 (0.036)	0.120** (0.053)
Other preschool	0.068*** (0.019)	0.069*** (0.023)	0.034 (0.039)	-0.011 (0.046)	0.043 (0.047)
Mean Y	0.531	0.532	0.532	0.396	0.556
Observations	7372	5361	5361	2369	2986
<i>Fraction of last 5 yrs not on Food Stamps/SNAP, age 30</i>					
Head Start	-0.018 (0.015)	0.011 (0.017)	0.043 (0.033)	0.042 (0.037)	0.076 (0.055)
Other preschool	-0.003 (0.007)	0.007 (0.009)	-0.019 (0.018)	-0.019 (0.047)	-0.015 (0.019)
Mean Y	0.936	0.929	0.929	0.831	0.949
Observations	4186	2805	2175	887	1285
<i>Never on AFDC/TANF by age 30</i>					
Head Start	-0.028* (0.016)	-0.015 (0.018)	-0.009 (0.020)	-0.001 (0.023)	0.001 (0.034)
Other preschool	0.022*** (0.008)	0.026*** (0.009)	0.004 (0.011)	-0.010 (0.025)	0.005 (0.012)
Mean Y	0.938	0.940	0.940	0.819	0.962
Observations	7372	5361	5361	2369	2986
<i>Fraction of last 5 yrs with positive earnings, age 30</i>					
Head Start	0.041*** (0.015)	0.035** (0.017)	0.061 (0.038)	0.026 (0.034)	0.088 (0.072)
Other preschool	0.013 (0.011)	0.008 (0.013)	0.015 (0.019)	-0.047 (0.048)	0.027 (0.020)
Mean Y	0.895	0.898	0.898	0.845	0.907
Observations	4378	3295	2800	1054	1740
<i>Mean Inc. Rel. Pov. in last 5 years, age 30</i>					
Head Start	-29.579*** (10.548)	-27.953** (12.160)	-16.953 (14.369)	5.860 (12.890)	-24.477 (23.499)
Other preschool	42.704** (18.606)	46.790*** (17.411)	-1.326 (16.118)	-4.147 (17.769)	0.923 (18.924)
Mean Y	385.831	385.933	385.933	224.651	412.236
Observations	5293	4068	3694	1514	2175
<i>Fraction of last 5 yrs no unemployment, age 30</i>					
Head Start	-0.007 (0.015)	-0.001 (0.016)	0.005 (0.030)	-0.013 (0.031)	0.056 (0.049)
Other preschool	-0.017 (0.012)	-0.013 (0.014)	-0.029 (0.027)	0.022 (0.029)	-0.040 (0.032)
Mean Y	0.854	0.850	0.850	0.807	0.857
Observations	4259	3184	2670	981	1683

Notes: This table shows estimates from regressions of the inputs to the Economic Sufficiency Index at age 30 on an indicator for participation in Head Start together with control variables described in the text. Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.4: Effect of Head Start on Inputs to Economic Sufficiency Index at age 40 (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>Fraction of last 5 yrs not on Food Stamps/SNAP, age 40</i>					
Head Start	0.001 (0.019)	0.009 (0.020)	0.045 (0.033)	0.054 (0.044)	0.051 (0.049)
Other preschool	0.001 (0.010)	0.003 (0.013)	-0.010 (0.023)	-0.013 (0.062)	-0.008 (0.023)
Mean Y	0.957	0.957	0.957	0.866	0.971
Observations	1972	1423	1213	564	647
<i>Never on AFDC/TANF by age 40</i>					
Head Start	0.008 (0.020)	0.022 (0.023)	-0.009 (0.030)	-0.010 (0.039)	0.002 (0.048)
Other preschool	0.016 (0.010)	0.019 (0.012)	0.018 (0.021)	-0.034 (0.062)	0.025 (0.021)
Mean Y	0.932	0.933	0.933	0.778	0.959
Observations	4085	2845	2503	1065	1435
<i>Fraction of last 5 yrs with positive earnings, age 40</i>					
Head Start	0.026 (0.031)	0.022 (0.038)	0.021 (0.062)	0.073 (0.053)	-0.180 (0.130)
Other preschool	-0.004 (0.027)	-0.012 (0.033)	-0.026 (0.051)	-0.135*** (0.052)	0.003 (0.060)
Mean Y	0.850	0.849	0.849	0.856	0.847
Observations	1829	1369	1078	445	633
<i>Mean Inc. Rel. Pov. in last 5 years, age 40</i>					
Head Start	1.769 (21.347)	3.447 (26.202)	32.738 (30.410)	27.251 (24.095)	-11.620 (56.148)
Other preschool	97.953** (38.986)	101.861** (47.085)	24.513 (40.157)	17.035 (22.343)	26.140 (50.412)
Mean Y	436.769	434.280	434.280	234.965	466.741
Observations	2152	1613	1272	540	732
<i>Fraction of last 5 yrs no unemployment, age 40</i>					
Head Start	-0.003 (0.022)	-0.022 (0.027)	-0.028 (0.047)	-0.033 (0.056)	-0.046 (0.083)
Other preschool	-0.011 (0.017)	-0.011 (0.021)	-0.026 (0.037)	-0.053 (0.060)	-0.016 (0.044)
Mean Y	0.906	0.902	0.902	0.841	0.911
Observations	1825	1365	1073	440	633
<i>Fraction of last 5 yrs owned home, age 40</i>					
Head Start	-0.022 (0.049)	-0.024 (0.056)	0.045 (0.056)	-0.058 (0.054)	0.070 (0.121)
Other preschool	-0.041 (0.037)	-0.053 (0.044)	-0.057 (0.058)	-0.079 (0.079)	-0.058 (0.074)
Mean Y	0.500	0.522	0.522	0.324	0.554
Observations	2292	1625	1391	642	747

Notes: This table shows estimates from regressions of the inputs to the Economic Sufficiency Index at age 40 on an indicator for participation in Head Start together with control variables described in the text. Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.5: Effect of Head Start on Inputs to Good Health Index at age 30 (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>Fraction of last 5 yrs smoked less than 1 cigarette/day, age 30</i>					
Head Start	-0.064*	-0.031	0.021	-0.127*	0.049
	(0.035)	(0.039)	(0.080)	(0.072)	(0.110)
Other preschool	-0.017	0.017	-0.011	-0.181**	0.012
	(0.021)	(0.024)	(0.052)	(0.091)	(0.056)
Mean Y	0.745	0.755	0.755	0.785	0.750
Observations	2267	1742	1174	376	796
<i>Fraction of last 5 yrs reported good or better health, age 30</i>					
Head Start	-0.001	0.001	0.042	0.047	0.039
	(0.012)	(0.013)	(0.031)	(0.034)	(0.052)
Other preschool	0.008	0.004	0.005	-0.009	0.010
	(0.008)	(0.010)	(0.016)	(0.035)	(0.017)
Mean Y	0.948	0.950	0.950	0.890	0.959
Observations	3763	2806	2292	829	1459
<i>Negative Mean BMI in last 5 years, age 30</i>					
Head Start	-1.063**	-0.982*	-0.485	1.408	-1.514
	(0.436)	(0.506)	(0.765)	(0.984)	(1.128)
Other preschool	0.046	-0.096	-0.332	-0.357	-0.202
	(0.266)	(0.313)	(0.441)	(1.069)	(0.472)
Mean Y	-26.569	-26.615	-26.615	-28.826	-26.267
Observations	3248	2528	1978	689	1286

Notes: This table shows estimates from regressions of the inputs to the Good Health Index at age 30 on an indicator for participation in Head Start together with control variables described in the text. Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * $p < .10$, ** $p < .05$, *** $p < .01$. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.6: Effect of Head Start on Inputs to Good Health Index at age 40 (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>Fraction of last 5 yrs smoked less than 1 cigarette/day, age 40</i>					
Head Start	-0.022 (0.047)	0.013 (0.050)	0.002 (0.075)	0.074 (0.077)	0.099 (0.148)
Other preschool	0.003 (0.039)	0.041 (0.047)	-0.033 (0.126)	0.218** (0.097)	-0.104 (0.150)
Mean Y	0.738	0.728	0.728	0.713	0.731
Observations	1280	930	698	300	398
<i>Fraction of last 5 yrs reported good or better health, age 40</i>					
Head Start	0.010 (0.034)	0.008 (0.039)	0.013 (0.059)	0.021 (0.061)	0.002 (0.144)
Other preschool	0.016 (0.029)	0.010 (0.035)	0.026 (0.023)	0.026 (0.065)	0.025 (0.023)
Mean Y	0.919	0.922	0.922	0.871	0.930
Observations	1463	1116	884	398	486
<i>Negative Mean BMI in last 5 years, age 40</i>					
Head Start	-1.218** (0.613)	-1.297* (0.731)	-0.976 (0.867)	-0.475 (1.055)	0.501 (1.251)
Other preschool	-0.330 (0.424)	-0.741 (0.518)	-1.861*** (0.647)	1.271 (1.503)	-2.360*** (0.693)
Mean Y	-27.504	-27.433	-27.433	-29.491	-27.095
Observations	2037	1486	1116	413	703

Notes: This table shows estimates from regressions of the inputs to the Good Health Index at age 40 on an indicator for participation in Head Start together with control variables described in the text. Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * $p < .10$, ** $p < .05$, *** $p < .01$. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.7: Impact of Head Start on High School, College, Earnings, and Criminal Behavior (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>A. Completed High School</i>					
Head Start	0.007 (0.018)	-0.002 (0.021)	-0.011 (0.026)	-0.024 (0.031)	-0.015 (0.045)
Other preschool	-0.002 (0.011)	-0.008 (0.014)	0.036* (0.021)	-0.012 (0.048)	0.046* (0.024)
R-Squared	0.098	0.105	0.028	0.050	0.038
Observations	7372	5361	5361	2369	2986
<i>B. Completed Some College</i>					
Head Start	0.038 (0.024)	0.039 (0.029)	0.046 (0.033)	-0.016 (0.036)	0.120** (0.053)
Other preschool	0.068*** (0.019)	0.069*** (0.023)	0.034 (0.039)	-0.011 (0.046)	0.043 (0.047)
R-Squared	0.213	0.233	0.050	0.056	0.057
Observations	7372	5361	5361	2369	2986
<i>C. Ln Earnings 23-25</i>					
Head Start	0.040 (0.056)	0.032 (0.066)	0.064 (0.109)	0.057 (0.142)	0.113 (0.158)
Other preschool	0.064 (0.045)	0.035 (0.052)	0.084 (0.098)	0.174 (0.173)	0.070 (0.110)
R-Squared	0.151	0.161	0.131	0.095	0.152
Observations	4351	3309	2726	986	1736
<i>D. Not Booked/Charged with Crime</i>					
Head Start	-0.007 (0.025)	-0.012 (0.031)	-0.008 (0.033)	0.028 (0.028)	-0.068 (0.064)
Other preschool	-0.006 (0.014)	0.007 (0.017)	-0.002 (0.033)	-0.022 (0.036)	0.002 (0.039)
R-Squared	0.055	0.062	0.089	0.074	0.106
Observations	5005	3591	3206	1366	1836

Notes: This table shows estimates from regressions of high school graduation (panel A), some college attainment (panel B), ln earnings between ages 23 and 25 (panel C) and not being charged with a crime (panel D) on an indicator for participation in Head Start together with control variables described in the text. Among the 7,372 individuals in the sample, 1098 individuals are in families that have variation in the Head Start variable (347 families), among those for whom we observe completed education; 887 black (277 black families), and 211 white individuals (70 white families). Crime sample limited to individuals age ≥ 16 at the time of interview in 1995. Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * $p < .10$, ** $p < .05$, *** $p < .01$. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.8: Impact of Head Start on Main Outcomes by Sex

	Black		White	
	Female	Male	Female	Male
<i>High School</i>				
Head Start x Sex	0.008 (0.033)	-0.062 (0.042)	0.005 (0.059)	-0.043 (0.054)
P-value of Difference	0.092		0.497	
<i>Some College</i>				
Head Start x Sex	-0.012 (0.044)	-0.021 (0.045)	0.102 (0.074)	0.145*** (0.053)
P-value of Difference	0.873		0.582	
<i>Ln Earnings 23-25</i>				
Head Start x Sex	0.265 (0.171)	-0.238 (0.202)	0.133 (0.217)	0.078 (0.174)
P-value of Difference	0.037		0.834	
<i>No Crime</i>				
Head Start x Sex	0.038 (0.035)	0.016 (0.041)	-0.036 (0.073)	-0.112 (0.089)
P-value of Difference	0.661		0.448	
<i>Economic Sufficiency Index, age 30</i>				
Head Start x Sex	-0.052 (0.099)	-0.197** (0.090)	-0.099 (0.112)	0.078 (0.141)
P-value of Difference	0.148		0.252	
<i>Economic Sufficiency Index, age 40</i>				
Head Start x Sex	-0.021 (0.173)	-0.363** (0.164)	0.058 (0.140)	-0.271 (0.184)
P-value of Difference	0.098		0.099	
<i>Good Health Index, Age 30</i>				
Head Start x Sex	0.042 (0.159)	-0.004 (0.218)	-0.198 (0.278)	-0.361 (0.378)
P-value of Difference	0.838		0.690	
<i>Good Health Index, Age 40</i>				
Head Start x Sex	0.349 (0.273)	-0.672** (0.271)	0.605 (0.378)	-1.099** (0.480)
P-value of Difference	0.014		0.004	

This table shows estimates from regressions of our main outcomes on an indicator for participation in Head Start interacted with an indicator for being female or male. The estimated interactions between Head Start and female (male) are shown in columns 1 and 3 (2 and 4). Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * $p < .10$, ** $p < .05$, *** $p < .01$. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table C.9: Regression: Interaction with Cohort (Linear) (PSID)

	All	Sibs	Mom FE	Blk, FE	Wht, FE
<i>Economic Sufficiency Index, age 30</i>					
Head Start	-0.054 (0.066)	-0.033 (0.073)	-0.038 (0.086)	-0.081 (0.104)	0.094 (0.153)
Head Start x trend	-0.010** (0.005)	-0.010* (0.006)	-0.009 (0.007)	-0.007 (0.008)	-0.017 (0.013)
Mean Y	0.094	0.096	0.096	-0.552	0.213
Observations	7372	5361	5361	2369	2986
<i>Economic Sufficiency Index, age 40</i>					
Head Start	-0.042 (0.084)	-0.038 (0.093)	-0.030 (0.104)	-0.155 (0.136)	-0.026 (0.118)
Head Start x trend	-0.014 (0.012)	-0.015 (0.013)	-0.031* (0.017)	-0.050** (0.025)	-0.029 (0.019)
Mean Y	0.020	0.025	0.025	-0.670	0.142
Observations	4085	2845	2503	1065	1435
<i>Good Health Index, Age 30</i>					
Head Start	-0.318*** (0.064)	-0.291*** (0.065)	-0.113 (0.161)	-0.087 (0.167)	0.018 (0.293)
Head Start x trend	-0.004 (0.007)	-0.004 (0.007)	-0.007 (0.019)	0.034** (0.017)	-0.044 (0.034)
Mean Y	0.004	0.017	0.017	-0.357	0.074
Observations	4749	3600	3114	1150	1959
<i>Good Health Index, Age 40</i>					
Head Start	-0.135 (0.149)	-0.110 (0.167)	-0.129 (0.210)	0.066 (0.188)	0.422 (0.513)
Head Start x trend	-0.028 (0.024)	-0.034 (0.026)	-0.026 (0.037)	0.067 (0.044)	-0.186** (0.083)
Mean Y	0.011	0.015	0.015	-0.290	0.062
Observations	2228	1673	1306	511	795

Notes: This table shows estimates from regressions of the Economic Sufficiency and Good Health Indices on an indicator participation in Head Start interacted with a normed linear trend in year of birth (year of birth minus 1966, where 1966 represents the first year that Head Start was available). Estimates are weighted to be representative of 1995 population; see text for details. Standard errors are clustered at 1968 family id in column 1 and on mother id level otherwise. * p < .10, ** p < .05, *** p < .01. Source: Panel Study of Income Dynamics, 1968-2011 waves.

D Replication of GTC (2002)

D.1 Summary

In this appendix we describe the results of our replication of Garces, Thomas and Currie (2002) (GTC). We describe our replication methods in Section D.2.

Table D.1 below shows the summary statistics corresponding to Table 1 of GTC for our sample. We include GTC Table 1 for comparison as Table D.2. In general, the results across the two tables are similar, albeit not identical. The most notable difference is that we find a lower share of respondents participate in Head Start, although the difference is smaller for the sibling sample. The shares of respondents who graduate high school and college are higher in our sample than in GTC. We report average earnings from age 23-25 in nominal terms as well as adjusted to 1999 dollars. Our adjusted earnings are consistently higher than GTC's reported adjusted earnings, but our unadjusted earnings are quite close to their mean adjusted earnings. We suspect that GTC may have reported unadjusted earnings, although it is also possible that the discrepancy is due to a slightly larger sample of individuals with earnings in GTC's sample. Again, the number of observations we report in the final row of the table is based on the number of individuals responding to the Head Start participation question.

Our replication of the main regression results in GTC are shown in Table D.3. We include GTC's Table 2 as Table D.4 for comparison. Our regression results are qualitatively similar, especially for the larger samples (panels A, B, and C). GTC found few statistically significant results, one of which was a negative effect of Head Start on high school completion before including controls. We, too, find this negative and significant result, though ours is slightly smaller. The result in Column (6), which GTC find to be positive and significant, we do not find to be significant. Our results for the college outcomes are aligned with the findings in GTC. The magnitudes that we report are not statistically different from GTC and in particular we replicate the key finding that Head Start influences college going for white children and not for black children. Our replication of Panel C is qualitatively similar to GTC. We do not find a statistically significant decrease in black crime rates as GTC do, although our point estimates are consistently negative for blacks. Otherwise, our estimates are quite imprecise and not statistically different from GTC's.

Our earnings results (panel C replication) are quite different from GTC, but this may be due to differences in how we defined earnings rather than differences in our samples. This is apparent in the fact that we have many fewer observations than GTC beginning from column 2 onward, about 24% smaller in column 2 and 48% smaller in column 8.

D.2 Replication Methodology

This section documents the process of replicating Garces, Thomas and Currie (2002) (GTC) for future scholars wishing to repeat our steps. We describe three stages of the replication: construction

of the dataset, iterations to identify the likely variable definitions, and our final decisions based on these iterations. We also include information about the mechanics of downloading the data and the variables we use.

D.2.1 Construction of Dataset

We begin by assembling data from the Panel Study of Income Dynamics (PSID), a nationally representative longitudinal dataset that forms the basis for the analysis in GTC. The PSID consists of the survey responses of household heads and their wives, which compose the annual household-level datasets (“family files”), as well as a smaller database of responses of all individuals in the household to a small set of questions (“cross-year individual files”). We merge the family files to the cross-year individual files using the “case id” number, which is present both on the individual and family files. We also merge responses of an individual’s mother and father from the crossyear file for those individuals whose mother or father have been identified in the PSID crossyear file.

The result is a dataset with 71,285 individual observations, each of which contains the personal responses of an individual over time, the responses (usually given by the head of household) to the family interview questions for each year, and the responses of an individual’s parents to the cross-year survey. The base dataset includes the Survey of Economic Opportunity “poverty oversample” and the Latino oversample, two populations specifically targeted by the PSID in order to improve the representativeness of the survey. We proceed by excluding the Latino oversample in accordance with GTC’s footnote 4.

Next, we construct the variables needed to define our sample. GTC delineate the specifications for their sample throughout the paper, and in particular we rely on their descriptions in Section II and footnote 7. A key stratifying variable in GTC is race, which is also a limiting factor for the sample size since the GTC sample is restricted to only black and white individuals (see footnote 4 of GTC). Unfortunately, the PSID does not assign a race to each individual, so race must be imputed from the annual family responses about race. Specifically, the PSID surveys families about the race of the head and wife of the head of household, so an individual’s race can only be identified if that individual becomes a head of household or his wife. Otherwise we must infer the race of the individual through their relation to the head of household or his wife.

The process of identifying race from the responses of other family members can be done at any age and from a variety of different family members, so we have experimented with using more and less restrictive definitions. We establish five definitions of race based on the relations through which we allow inference and the survey years over which we make the inference. These definitions are summarized over those two dimensions below in Table D.5.

The second limiting criterion is the age of individuals. GTC include respondents aged 18 and over in 1995, which results in a sample of respondents born between 1965 and 1977. They exclude the 1964 and 1965 cohorts. Since this sample restriction can be defined and replicated in a few different ways with PSID variables, we develop three candidate limitations on age and year of birth

for individuals in our sample. We describe the criteria which define these alternative candidates in Table D.6.

The third criterion is to identify sets of siblings within the remaining sample that comprise the “siblings subsample.” Since the identification strategy relies on the inclusion of a mother fixed effect, we define siblings as any two individuals who satisfy the race and age criteria for the sample and have the same unique mother identification number. The mother identification number is a combination of a family identifier and a personal identifying number which is assigned by the PSID. Individuals that do not have a mother identification number are excluded from the sibling subsample.

Next, we flag observations from the SEO poverty oversample with the intention of excluding them as GTC do. We ultimately do not exclude these observations because comparisons of the sample statistics with and without the SEO sample make us speculate that the results in GTC were generated from a sample that included the SEO sample.

We construct sample weights using CPS weights to make the sample representative of the 1995 white and African-American populations. Specifically, we collapsed the 1995 CPS weights to age-race-sex cells (year of birth is not available) and merge the cell weight onto each observation of our sample. Then, we divide the cell weight by the number of individuals in that age-race-sex cell who are in our sample and the resulting individual weight is what we use for our analysis.

D.2.2 Search for identical dataset construction

As mentioned previously, the sample construction criteria are clearly documented in GTC. For some dimensions, we could think of a few ways to define variables and samples in accordance with their descriptions. Therefore, we conducted tests to determine the procedures that would yield a dataset consistent with GTC, as well as to assess the stability of the results.

Our search iterations hinge on four parameters: inclusion or exclusion of the SEO oversample; the algorithm for identifying an individual’s race; the criteria for age; and the order in which we dropped observations and weighted the sample. For this last parameter, we weighted the sample before dropping the Latino oversample as well as after. We do not present the results for the variations on this final parameter because the exercise clearly indicated that dropping the Latino oversample best matched GTC’s results regardless of how the first three parameters were defined.

Table D.7 below shows the results of our iteration of the summary statistics results for a select set of variables. Our goal was to match the results to Table 1 in GTC, reproduced on the first row of the table. The number of observations we report is for the variable for Head Start participation, although some variables have fewer observations. For example, over half the observations for the income variable are missing. GTC also report one N for each column, although they also likely had fewer observations for variables like income.

Our sample is weighted based on race, gender, and age variables from the CPS, so we expect that the mean values for the weighted PSID sample should be similar to the CPS means. We

include the CPS means for the three variables as a comparison. The definitions for age and race are as described in the previous section. There are a number of conclusions we draw from this table. First, we speculate that the 25.17 percent black reported in GTC is, in fact, 15.17 percent, which is much closer to the CPS means. Second, inclusion of the SEO oversample adds approximately 1,500 observations to our sample and brings us quite close to the size of the sample and sample means reported in GTC.

As we had hoped, moving from iteration to iteration substantially changes the number of observations, which suggest which decisions produced the sample of GTC. For example, holding SEO and age definitions constant, moving from our conservative definition of race (2) to the liberal definition (4) adds approximately 30 to 50 observations, an approximately 1.5 percent increase in sample size. The specification of age is also important for defining the sample size. For example, the movement from row 1, 1, 2 (N=3,286) to 1, 2, 2 (N=3,548) is an eight percent increase, and the subsequent movement to row 1, 3, 2 (N=4,187) is an 18 percent increase.

Despite the variability in sample size, our sample characteristics are not sensitive to the decisions along each of these dimensions. Additionally, while our results for these select variables are at times statistically different from those of GTC, we remain close to the magnitudes that they report. The race, gender, and age means are very similar across the specifications, likely on account of the weighting. The preschool participation and high school graduation rates are nearly identical throughout, especially when we include the SEO oversample. The exception to this pattern is Head Start participation. The SEO oversample increases the share of respondents who were in Head Start to close to nine percent, which is still lower than the 10.57 percent reported in GTC. We were unable to replicate this high incidence of Head Start participation throughout the iteration process, including in iterations not reported here.

We also performed iterations on the regression models from GTC's Table 2. GTC conduct a similar regression for each of four outcome variables: high school graduation, college graduation, crime, and later earnings. The first of these three are fairly similar: they are defined by one variable in the PSID. In this comparison table we only show results for high school graduation. On the other hand, compiling a consistent variable for earnings is trickier. Here we present results for one of our regressions, but in general we were not able to replicate the findings for this outcome variable.

There are eight different models in GTC. The first three are on the full sample, the sibling sample, and the sibling sample with controls. The next five models use mother fixed effects: first on the full sample, then the full sample split by whether the mother was white or black, and finally for the subset of mothers with less than a high school education, also split by race.

Table D.8 shows a comparison of the results. We show iterations on the same three age restrictions as above, as well as race definitions for definitions 4 and 5 as defined in the previous section. For each regression the corresponding result from GTC is shown on the first row.

Our regression results are qualitatively similar, especially for the larger samples (panel A). GTC found few statistically significant results, one of which was a negative effect of Head Start on high

school completion (result A.1). We, too, replicate this negative and significant result, though ours are smaller. As can be noted in result A.4, our models using later earnings were similar to those in the paper. The result in B.4, which GTC find to be positive and significant, we do not find significant. However, all of our replications of this result fall within the confidence interval they use.

Among our various iterations, the results are stable. Only the result in A.1 has a difference of one standard error between estimates, with the rest of these results never straying more than half a standard error from each other.

D.3 Final dataset restrictions

Given our iteration exercises, our preferred sample definition includes the SEO poverty over-sample, uses age definition 1 and uses race definition 5 as explained in the first section of this appendix. Our choice of age and race definitions is appropriate for three reasons. First, they replicate the GTC adequately. Second, they are a reasonable method for a researcher not attempting to replicate findings. Third, they result in large samples, which is important for additional analyses.

D.4 More on the data

We downloaded the data files from <http://simba.isr.umich.edu/Zips/ZipMain.aspx>. Table D.9 shows the variables we downloaded.

Table D.1: Replication of Garces, Thomas, Currie (2002) Summary Statistics

	All	Head Start	No Head Start	Sibling Sample
Head Start	0.0873 (.282)	1 (0)	0 (0)	0.103 (.304)
Other preschool	0.266 (.442)	0 (0)	0.291 (.454)	0.281 (.45)
Fraction completed hs	0.851 (.356)	0.752 (.432)	0.860 (.347)	0.854 (.353)
Fraction attended some college	0.468 (.499)	0.339 (.474)	0.481 (.5)	0.482 (.5)
Avg. Earnings age 23-25	18543.5 (14929)	13361.3 (12057)	18962.7 (15062)	20116.3 (17141)
Avg. Earnings age 23-25 (CPI adjusted)	20367.9 (15646)	14730.7 (12950)	20823.9 (15758)	21734.8 (17521)
Fraction booked/charged with crime	0.0998 (.3)	0.124 (.33)	0.0975 (.297)	0.106 (.308)
Fraction African-American	0.150 (.357)	0.619 (.486)	0.105 (.307)	0.162 (.369)
Fraction female	0.502 (.5)	0.533 (.499)	0.499 (.5)	0.475 (.5)
Age in 1995	23.67 (3.44)	23.14 (3.5)	23.72 (3.43)	23.14 (3.28)
Fraction eldest child in family	0.345 (.475)	0.335 (.472)	0.346 (.476)	0.364 (.481)
Fraction low birth weight	0.0608 (.239)	0.110 (.314)	0.0553 (.229)	0.0560 (.23)
Mother's yrs education	11.36 (2.58)	10.00 (2.44)	11.49 (2.56)	11.17 (2.54)
Fraction whose mother completed hs	0.772 (.419)	0.585 (.493)	0.790 (.407)	0.770 (.421)
Father's yrs education	11.46 (3.01)	9.806 (2.78)	11.60 (2.98)	11.37 (3)
Fraction whose father completed hs	0.725 (.446)	0.475 (.5)	0.747 (.435)	0.717 (.451)
Family income (age 3-6) (CPI adjusted)	48040.3 (27470)	30253.9 (15498)	49699.4 (27756)	48580.8 (29193)
Had a single mother at age 4	0.119 (.324)	0.320 (.467)	0.0998 (.3)	0.108 (.31)
Household size at age 4	4.659 (1.81)	5.109 (2.18)	4.616 (1.76)	4.831 (1.71)
Observations	3399	552	2847	1541

Notes: Weighted to be representative of 1995 population; see text for details.

Table D.2: GTC Table 1: Summary Statistics

	All sample	Head Start	Not in Head Start	Sibling Sample
Head Start	0.1057 (.0053)	1 (0)	0 (0)	0.1089 (.0073)
Other preschool	0.2834 (.0077)	0.1333 (.0151)	0.3011 (.0085)	0.2771 (.0105)
Pct. completed hs	0.7660 (.0074)	0.6465 (.0216)	0.7803 (.0079)	0.7721 (.0101)
Pct. attended some college	0.3714 (.0085)	0.2508 (.0196)	0.3859 (.0093)	0.3880 (.0117)
Average earnings between age 23-25	-	-	-	-
Average earnings between age 23-25 - CPI adjusted	17290 (690)	12100 (670)	17810 (760)	17310 (1000)
Pct. booked/charged with crime	0.0969 (.0051)	0.1104 (.00139)	0.0953 (.0054)	0.1004 (.0070)
Pct. African-American	0.2517 (.0074)	0.7532 (.00192)	0.1924 (.0078)	0.2285 (.0098)
Pct. female	0.5149 (.0085)	0.5641 (.0220)	0.5091 (.0093)	0.5075 (.0117)
Age in 1995	23.66 (.06)	23.35 (.15)	23.70 (.06)	23.65 (.08)
Pct. eldest child in family	0.5311 (.0056)	0.5089 (.0141)	0.5337 (.0061)	0.5057 (.0076)
Pct. low birth weight	0.0699 (.0037)	0.1040 (.0124)	0.0659 (.0038)	0.0669 (.0056)
Mother's yrs education	12.14 (.04)	11.33 (.09)	12.24 (.04)	12.30 (.05)
Pct. whose mother completed hs	0.7037 (.0078)	0.5552 (.0221)	0.7212 (.0083)	0.7815 (.0097)
Father's yrs education	11.60 (.06)	10.19 (.14)	11.76 (.06)	12.23 (.07)
Pct. whose father completed hs	0.5612 (.0085)	0.2638 (.0196)	0.5964 (.0091)	0.6330 (.0113)
Family income (age 3-6) - CPI adjusted	46230 (460)	26620 (580)	48540 (500)	47330 (670)
Had a single mother at age 4	0.1642 (.0061)	0.4035 (.0216)	0.1359 (.0061)	0.1306 (.0079)
Household size at age 4	4.59 (.03)	4.97 (.09)	4.55 (.03)	4.84 (.04)
Observations	3255	489	2766	1742

Table D.3: Replication of Garces, Thomas, Currie (2002) Regressions

	All	Sibs	Controls	Mom FE	Blk, FE	Wht, FE
<i>Panel A. High School</i>						
Head Start	-0.064*	-0.017	0.009	0.031	-0.017	0.093
	(0.034)	(0.043)	(0.040)	(0.057)	(0.063)	(0.092)
Other Preschool	0.082***	0.076***	0.014	0.028	0.068	0.021
	(0.013)	(0.022)	(0.021)	(0.035)	(0.072)	(0.038)
Observations	3399	1541	1541	1541	615	923
<i>Panel B. College</i>						
Head Start	-0.027	-0.021	0.033	0.100*	-0.039	0.232**
	(0.035)	(0.053)	(0.045)	(0.059)	(0.059)	(0.094)
Other Preschool	0.200***	0.219***	0.098***	0.047	-0.062	0.059
	(0.025)	(0.034)	(0.033)	(0.044)	(0.101)	(0.049)
Observations	3399	1541	1541	1541	615	923
<i>Panel C. Earnings</i>						
Head Start	-0.139*	-0.142	-0.056	-0.041	0.427*	-0.322
	(0.074)	(0.108)	(0.113)	(0.191)	(0.245)	(0.261)
Other Preschool	0.067	-0.023	-0.125*	-0.013	0.286	-0.017
	(0.062)	(0.072)	(0.074)	(0.116)	(0.448)	(0.118)
Observations	2118	972	972	779	236	541
<i>Panel D. No Crime</i>						
Head Start	-0.028	0.069	-0.055	-0.086	0.065	-0.222*
	(0.028)	(0.050)	(0.049)	(0.070)	(0.044)	(0.125)
Other Preschool	-0.000	-0.020	0.004	-0.046	0.059	-0.059
	(0.015)	(0.019)	(0.020)	(0.038)	(0.052)	(0.043)
Observations	3387	1537	1537	1535	614	918

Notes: * $p < .10$, ** $p < .05$, *** $p < .01$. Weighted to be representative of 1995 population; see text for details. SE clustered at 1968 family id in column 1 and at mother id level otherwise.

Table D.4: GTC Table 2: Regressions

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	All	Sibs	Controls	Mom FE	Blk, FE	Wht, FE	Blk, l.e. HS	Wht, l.e. HS
<i>Completed high School</i>								
Head Start	-0.089 (0.026)	-0.075 (0.035)	0.006 (0.034)	0.037 (0.053)	-0.025 (.065)	0.203 (0.098)	0.000 (0.071)	0.283 (0.119)
Other Preschool	0.085 (0.016)	0.073 (0.022)	0.003 (0.021)	-0.032 (0.038)	-0.056 (0.064)	-0.014 (0.048)	-0.080 (0.077)	-0.019 (0.067)
Difference	-0.174	-0.148	0.003	0.069	0.031	0.217	0.081	0.302
S.E Difference	0.028	0.037	0.039	0.062	0.085	0.105	0.097	0.126
N	3255	1742	1742	1742	706	1036	554	677
<i>Attended Some College</i>								
Head Start	-0.038 (0.023)	-0.016 (0.033)	0.075 (0.033)	0.092 (0.056)	0.023 (.066)	0.281 (0.108)	0.031 (0.067)	0.276 (0.120)
Other Preschool	0.142 (0.019)	0.149 (0.027)	0.023 (0.026)	0.050 (0.040)	-0.007 (0.064)	0.095 (0.052)	0.022 (0.072)	0.0103 (0.068)
Difference	-0.180	-0.165	0.052	0.042	0.030	0.186	0.009	0.173
S.E Difference	0.028	0.040	0.041	0.065	0.085	0.115	0.092	0.127
N	3255	1742	1742	1742	706	1036	554	677
<i>ln(earnings 23-25)</i>								
Head Start	-0.034 (0.090)	0.053 (0.116)	0.170 (0.117)	0.194 (0.257)	0.073 (0.321)	0.566 (0.459)	0.051 (0.357)	1.004 (0.516)
Other Preschool	0.173 (0.063)	0.174 (0.086)	0.002 (0.082)	0.079 (0.171)	-0.087 (0.287)	0.146 (0.219)	0.124 (0.341)	0.136 (0.306)
Difference	-0.207	-0.122	0.167	0.115	0.160	0.420	-0.073	0.868
S.E Difference	0.104	0.138	0.144	0.302	0.420	0.504	0.482	0.548
N	1383	728	728	728	272	456	216	320
<i>Booked or charged with crime</i>								
Head Start	0.023 (0.018)	0.041 (0.026)	0.012 (0.026)	-0.053 (0.039)	-0.116 (0.045)	0.122 (0.077)	-0.126 (0.050)	0.058 (0.095)
Other Preschool	-0.017 (0.011)	-0.022 (0.016)	-0.001 (0.017)	0.032 (0.028)	0.000 (0.045)	0.063 (0.036)	-0.023 (0.056)	0.147 (0.054)
Difference	0.040	0.063	0.013	-0.085	-0.117	0.059	-0.103	-0.089
S.E Difference	0.020	0.028	0.030	0.045	0.059	0.082	0.070	0.100
N	3255	1742	1742	1742	706	1036	554	677

SE in parentheses.

Table D.5: Alternative Definitions of Race

Defn.	Survey Years		Relation to Head (or Wife)				
	1995	1985-1996	Head	Wife	Child	Parent	Sibling
1	X		X	X	X		
2	X		X	X	X	X	X
3		X	X	X	X		
4		X	X	X	X	X	X
5		X	X	X	X	X	X

Table D.6: Candidate limitations on birth year and age

Defn.	BirthYears			Age in 1995		
	1966-1977	Not 1965, 1978	No Restriction	>18	17-29	17-30
1	X			X		
2		X			X	
3			X			X

Table D.7: Iterations for Summary Statistics Table

	Black	Female	Age	Head Start	Preschool	High School	N		
GTC(2002)	0.252	0.515	23.660	0.106	0.283	0.766	3255		
CPS 1995	0.150	0.505	23.686						
<i>Sample Iterations</i>									
SEO	Age	Race							
0	1	2	0.149	0.497	22.952	0.078	0.302	0.822	1708
0	1	4	0.149	0.497	22.950	0.079	0.299	0.820	1735
0	2	2	0.154	0.499	22.859	0.079	0.309	0.811	1855
0	2	4	0.154	0.499	22.857	0.080	0.306	0.809	1883
0	3	2	0.150	0.503	23.713	0.076	0.286	0.820	2173
0	3	4	0.150	0.503	23.712	0.076	0.284	0.818	2204
1	1	2	0.153	0.498	22.959	0.089	0.290	0.788	3286
1	1	4	0.153	0.498	22.958	0.089	0.288	0.787	3333
1	2	2	0.157	0.500	22.926	0.087	0.292	0.782	3548
1	2	4	0.157	0.500	22.925	0.087	0.290	0.781	3597
1	3	2	0.150	0.503	23.710	0.082	0.276	0.788	4187
1	3	4	0.120	0.503	23.710	0.082	0.274	0.787	4244

Notes: First row corresponds to selections from Garces, Thomas and Currie (2002) table 1. Second row corresponds to 1995 CPS means, as described in the text of the appendix. The next 12 columns correspond to sample iterations on three criteria. The first is the inclusion (SEO=1) or exclusion (SEO=0) of the Survey of Economic Opportunity sample. The three age criteria and two race criteria are explained in detail in the previous table. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table D.8: Iterations for Regressions Table

<i>Panel A.</i>													
		HS, All			HS, Sib			HS, Mom FE			Log Earnings, All		
		b	se	N	b	se	N	b	se	N	b	se	N
GTC (2002)		-0.089	(0.026)	3255	-0.075	(0.035)	1742	0.037	(0.053)	1742	-0.034	(0.090)	1383
<i>Sample Iterations</i>													
Age	Race												
1	4	-0.075	(0.030)	3315	-0.035	(0.043)	1543	0.047	(0.075)	1543	-0.064	(0.106)	894
1	5	-0.071	(0.030)	3344	-0.025	(0.042)	1565	0.047	(0.075)	1565	-0.067	(0.105)	898
2	4	-0.073	(0.030)	3585	-0.034	(0.039)	1731	0.072	(0.077)	1731	-0.064	(0.104)	894
2	5	-0.067	(0.031)	3616	-0.024	(0.039)	1753	0.072	(0.076)	1753	-0.067	(0.104)	898
3	4	-0.052	(0.026)	4233	-0.046	(0.035)	2125	0.037	(0.063)	2125	-0.043	(0.092)	1132
3	5	-0.046	(0.027)	4264	-0.036	(0.035)	2147	0.036	(0.062)	2147	-0.046	(0.092)	1136

<i>Panel B.</i>													
		HS, Mom FE, Black			HS, Mom FE, White			HS, Mom<HS, Black			HS, Mom<HS, White		
		b	se	N	b	se	N	b	se	N	b	se	N
GTC (2002)		-0.025	(0.065)	706	0.203	(0.098)	1036	0	(0.071)	554	0.283	(0.119)	677
<i>Sample Iterations</i>													
Age	Race												
1	4	-0.030	(0.058)	625	0.133	(0.089)	898	-0.026	(0.058)	586	0.152	(0.099)	672
1	5	-0.030	(0.058)	625	0.133	(0.088)	920	-0.026	(0.058)	586	0.152	(0.098)	692
2	4	-0.028	(0.056)	702	0.181	(0.094)	1008	-0.025	(0.056)	649	0.203	(0.105)	759
2	5	-0.028	(0.056)	702	0.181	(0.092)	1030	-0.025	(0.056)	649	0.202	(0.104)	779
3	4	-0.043	(0.044)	858	0.120	(0.081)	1241	-0.045	(0.044)	797	0.136	(0.092)	961
3	5	-0.043	(0.044)	858	0.114	(0.079)	1263	-0.045	(0.044)	797	0.130	(0.088)	981

Notes: First row of each panel corresponds to selections from Garces, Thomas and Currie (2002) table 2. The three age criteria and two race criteria are explained in detail in the previous table. Source: Panel Study of Income Dynamics, 1968-2011 waves.

Table D.9: PSID Variables used in the analysis

Our Variable	PSID Original Variable	Description (derived variable)	Source
id1968	ER30001	Family identifier	Indiv. Cross year
pernum	ER30002	Personal identifier	Indiv. Cross year
relation1968- relation2001	ER30003, ER30022, ER30045, ER30069, ER30093, ER30119, ER30140, ER30162 ER30190, ER30219, ER30248 ER30285, ER30315, ER30345, ER30375, ER30401, ER30431, ER30465, ER30500, ER30537, ER30572, ER30608, ER30644 ER30691, ER30735, ER30808, ER33103, ER33203, ER33303, ER33403, ER33503, ER33603	Relation to head	Indiv. Cross year
caseid1968- caseid2011	ER30020 ER30043 ER30067 ER30091 ER30117 ER30138 ER30160 ER30188 ER30217 ER30246 ER30283 ER30313 ER30343 ER30373 ER30399 ER30429 ER30463 ER30498 ER30535 ER30570 ER30606 ER30642 ER30689 ER30733 ER30806 ER33101 ER33201 ER33301 ER33401 ER33501 ER33601 ER33701 ER33801 ER33901 ER34001 ER34101 ER33601	Fam. Interview Num- ber	Indiv. Cross year
edu1968- edu2011	ER30010 ER30052 ER30076 ER30100 ER30126 ER30147 ER30169 ER30197 ER30226 ER30255 ER30296 ER30326 ER30356 ER30384 ER30413 ER30443 ER30478 ER30513 ER30549 ER30584 ER30620 ER30657 ER30703 ER30748 ER30820 ER33115 ER33215 ER33315 ER33415 ER33516 ER33616 ER33716 ER33817 ER33917 ER34020 ER34119	Yrs. Education	Indiv. Cross year
age1995	ER33204	Age in 1995	Indiv. Cross year
birthyr1995	ER33206	Birthyear in 1995	Indiv. Cross year
headstart1995	ER33261	Head Start Response in 1995	Indiv. Cross year
preschool1995	ER33264	Preschool Response in 1995	Indiv. Cross year
preschool1995	ER33266	Crime Response in 1995	Indiv. Cross year
sex	ER32000	Sex	Indiv. Cross year
momid1968	ER32009	Mother's Family ID	Indiv. Cross year
mompernum	ER32010	Mother's Personal ID	Indiv. Cross year
dadid1968	ER32016	Father's Family ID	Indiv. Cross year
dadpernum	ER32017	Father's Personal ID	Indiv. Cross year
birthweight	ER32014	Birth weight	Indiv. Cross year

Our Variable	PSID Original Variable	Description (derived variable)	Source
crime1995	ER33266	Committed/Charged with Crime	Indiv. Cross year
parityofmom	ER32013	Parity of mom (Elders)	Indiv. Cross year
h_edu1968- h_edu2011	V313 V794 V1485 V2197 V2823 V3241 V3663 V4198 V5074 V5647 V6194 V6787 V7433 V8085 V8709 V9395 V11042 V12400 V13640 V14687 V16161 V17545 V18898 V20198 V21504 V23333 ER4158 ER6998 ER9249 ER12222 ER16516 ER20457 ER24148 ER28047 ER41037 ER46981 ER52405	Education of Head (Mom, Dad Education)	Family Interviews
w_edu1968, w_edu1972- w_edu2011	V246 V2687 V3216 V3638 V4199 V5075 V5648 V6195 V6788 V7434 V8086 V8710 V9396 V11043 V12401 V13641 V14688 V16162 V17546 V18899 V20199 V21505 V23334 ER4159 ER6999 ER9250 ER12223 ER16517 ER20458 ER24149 ER28048 ER41038 ER46982 ER52406	Education of Wife of Head (Mom Education)	Family Interviews
h_sex1968- h_sex2011	V119 V1010 V1240 V1943 V2543 V3096 V3509 V3922 V4437 V5351 V5851 V6463 V7068 V7659 V8353 V8962 V10420 V11607 V13012 V14115 V15131 V16632 V18050 V19350 V20652 V22407 ER2008 ER5007 ER7007 ER10010 ER13011 ER17014 ER21018 ER25018 ER36018 ER42018 ER47318	Sex of Head (Single mom)	Family Interviews
f_tanf1994- f_tanf2011	ER3262 ER6262 ER8379 ER11272 ER14538 ER18697 ER22069 ER26050 ER37068 ER43059 ER48381	Family Received AFDC/TANF last year	Family Interviews
f_fs1994- f_fs2011	ER3059 ER6058 ER8155 ER11049 ER14255 ER18386 ER21652 ER25654 ER36672 ER42691 ER48007	Family Received Food Stamps last year	Family Interviews
h_cigs1986, h_cigs1999- h_cigs2011	V13442 ER15544 ER19709 ER23124 ER27099 ER38310 ER44283 ER49621	Cigarettes Per Day of Head	Family Interviews
w_cigs1986, w_cigs1999- w_cigs2011	V13477 ER15652 ER19817 ER23251 ER27222 ER39407 ER45380 ER50739	Cigarettes Per Day of Wife of Head	Family Interviews
h_wlbs1999- h_wlbs2011	ER15552 ER19717 ER23132 ER38320 ER44293 ER49631	Weight of Head (BMI)	Family Interviews
w_wlbs1999- w_wlbs2011	ER15660 ER19825 ER23259 ER27232 ER39417 ER45390 ER50749	Weight of Wife of Head (BMI)	Family Interviews

Our Variable	PSID Original Variable	Description (derived variable)	Source
h_srhealth1984- h_srhealth2011	V10877 V11991 V13417 V14513 V15993 V17390 V18721 V20021 V21321 V23180 ER3853 ER6723 ER8969 ER11723 ER15447 ER19612 ER23009 ER26990 ER38202 ER44175 ER49494	Self-Reported Health of Head	Family Interviews
w_srhealth1984- w_srhealth2011	V10884 V12344 V13452 V14524 V15999 V17396 V18727 V20027 V21328 V23187 ER3858 ER6728 ER8974 ER11727 ER15555 ER19720 ER23136 ER27113 ER39299 ER45272 ER50612	Self Reported Health of Head of Wife	Family Interviews
f_rentown1968- f_rentown2011	V103 V593 V1264 V1967 V2566 V3108 V3522 V3939 V4450 V5364 V5864 V6479 V7084 V7675 V8364 V8974 V10437 V11618 V13023 V14126 V15140 V16641 V18072 V19372 V20672 V22427 ER2032 ER5031 ER7031 ER10035 ER13040 ER17043 ER21042 ER25028 ER36028 ER42029 ER47329	Family Rents/Owns Home	Family Interviews
h_wages1968- h_wages2011	V251 V699 V1191 V1892 V2493 V3046 V3458 V3858 V4373 V5283 V5782 V6391 V6981 V7573 V8265 V8873 V10256 V11397 V12796 V13898 V14913 V16413 V17829 V20178 V21484 V23323 ER4140 ER6980 ER9231 ER12080 ER16463 ER20443 ER24116 ER27931 ER40921 ER46829 ER52237	Earnings of Head	Family Interviews
w_wages1968- w_wages2011	V76 V516 V1198 V1899 V2500 V3053 V3465 V3865 V4379 V5289 V5788 V6398 V6988 V7580 V8273 V8881 V10263 V11404 V12803 V13905 V14920 V16420 V17836 V19136 V20436 V23324 ER4144 ER6984 ER9235 ER12082 ER16465 ER20447 ER24135 ER27943 ER40933 ER46841 ER52249	Earnings of Wife of Head	Family Interviews

E Functional form choices with Binary Treatment and Binary Outcome

We now consider potential sensitivity to functional form modeling assumptions. For binary outcomes the usual choice of specifications include linear probability model (LPM), logit, and probit. In the cross-sectional setting, the conventional wisdom is that the choice among these options is fairly innocuous, especially when the objective is to recover the ATE.⁶¹ We are not aware of any previous systematic exploration of these properties in extremely short-panel settings such as found in the FFE design. We demonstrate some complications that arise in such settings, and compare the performance of these estimators.

E.1 Specification choices

Empiricists commonly use LPM specification to estimate FE models. In our sample of papers, this is almost universally used as the primary, if not only, specification. We speculate that this is motivated by (1) the intuition carried over from the cross-sectional case that LPM models usually recover the ATE; (2) the benefit that the incidental parameters problem does not pollute the main parameters of interest (Chamberlain, 1980);⁶² (3) computational ease, especially when paired with other complications to the research design such as many fixed effects, instrumental variables, etc.); and (4) the fact that the estimated coefficient β_{LPM} directly gives the estimate of the ATE.

Obtaining ATE from a nonlinear specification is not only less common, but also sometimes less straightforward. The conditional logit model, sometimes referred to as logit FE, consistently estimates β_{Logit} by conditioning on the number of successes in a family, but does not have a paired method for obtaining treatment effects. To obtain ATE, Wooldridge (2010, section 15.8) recommends employing a regular logit model and including family-level-means of control variables, i.e. “Chamberlain-Mundlak controls,” (hereafter, Mundlak controls) rather than directly controlling for fixed effects (Mundlak, 1978; Chamberlain, 1980).⁶³

Fernandez-Val (2009) examines the probit FE model. He proposes a bias-correction approach, which is based on the large-T asymptotic bias resulting from the incidental parameters problem. He also derives a “small bias” property for uncorrected/naive estimates of marginal effects for the probit FE model, and demonstrates this for panels of length as short as $T=4$. However it is not clear that the results in Fernandez-Val (2009) should apply in the family FE setting. This is because: (1)

⁶¹See Angrist and Pischke (2009, pg. 107) and Wooldridge (2010, section 15.6). In contrast, Cameron and Trivedi (2005, pg. 471) recommend limiting LPM’s to exploratory analysis, and note that it does not do a good job making predicted probabilities for individual observations. In panel contexts, textbook treatments generally state that estimates should be fine using LPM (Wooldridge, 2010, pg. 608).

⁶²Because this inconsistency is based on the panel length being fixed, the problem may be especially acute for short panels.

⁶³The traditional implementation is to model the residual variance as having an i-level random effect, hence the terminology Correlated Random Effects given to this method. However, it is also possible to include family means of control variables and then estimate regular pooled logit or probit, as we will do.

we face extremely short panel lengths due to families commonly having only 2 children⁶⁴; (2) his results apply only to the leading (order 1/T) bias term, but with very short panels the subsequent bias terms could still be relevant; and (3) there is an unresolved challenge of how to address the extrapolation from the estimation sample to singletons (when they are in the target population).⁶⁵

Mundlak controls and naive-fixed-effects methods have the attractive properties of: (1) being easy to implement; (2) respecting the binary functional form of the left-hand-side (LHS) variable; and (3) straightforwardly obtaining ATEs. Nonetheless, empiricists’ use of either of these options is uncommon; in our sample of 58 papers discussed in section 2 these methods are not used.

An additional complication with conditional or fixed effect logit and probit models is that they use less variation relative to LPM. With these models, for any families that have no variation in outcomes, i.e “all successes” or “all failures”, the fixed effect parameters will be driven to +/-infinity, and these families will be dropped from estimation. This leaves only “double switchers”: families with variation in both the outcome variable and the treatment variable. This means that moving from LPM to nonlinear specification is automatically tied to a change in estimation sample, which can reduce the effective sample size and may exacerbate the issues discussed in Section 3. In our application for example we see a reduction from 2986 individuals in the overall white "siblings sample" to 211 individuals in the “RHS switchers" sample to 98 individuals (from only 27 families) in the double switchers sample. A related issue is that the LPM results will depend on the fraction of observations in families that are not LHS switchers, whereas the logit model estimates will be invariant to the number of these non-switchers.

E.2 Obtaining Marginal Effects from Conditional Logit

In order to address the challenge of translating the conditional logit coefficient, β_{Logit} , into ATE units that can be compared with LPM results, we introduce a “two-step logit” model. The first step is the usual conditional logit estimator, used to obtain a consistent coefficient $\hat{\beta}$ for variables that change within-family. The second step estimates a random effects logit model (over the full sample, including non-switchers), while imposing the coefficient on the treatment variable (and on other individual-level variables) from the first step model. The purposes of the second step are (1) to estimate coefficients on family-level variables, so as (2) to assign an estimated “logit index” value to each observation, and (3) to estimate the variance of the family-level random effect σ_u^2 . After the second step model is estimated, we then estimate the ATE using:

$$ATE_{2StepLogit} = \frac{1}{N} \sum_{i=1}^N \int_u (\hat{\beta}_{HeadStart} \cdot \Lambda(\hat{\beta}X_{if} + \hat{\gamma}Z_f + u) \cdot 1 - \Lambda(\hat{\beta}X_{if} + \hat{\gamma}Z_f + u) \phi(u) du \quad (29)$$

⁶⁴We have reproduced the results for mean bias from his Table 4 for Probit and LPM-FS. We then reduced the panel size to T=2, and we find a detectible bias of -6.4% of the true ATE for the Probit, and no bias for the LPM-FS

⁶⁵For singletons there is no ability to separately identify the value of the fixed effect from the idiosyncratic error term. This is not a problem when the target population is either RHS switcher or all siblings. For these target populations, the naive logit FE or probit FE model could be used following the reweighting ideas presented above.

With $\hat{\beta}_{HeadStart}$ the coefficient on Head Start from the conditional logit first step; $\hat{\beta}$ the coefficient on i-level variables X_{if} from the conditional logit first step; $\hat{\gamma}$ the coefficients on family-level variables from the second step; and $\phi(u)$ the PDF from a normal distribution, with variance σ_u^2 estimated from the second step family-level random effects model. We have not yet found a prior implementation of this estimator in the literature; but it is similar in spirit to the two step fixed-effects logit proposed by Beck (2015).⁶⁶

E.3 Selection in Nonlinear Models

A desirable feature of the two-step logit and the Mundlak controls models is that both allow the marginal effect of treatment, $\frac{\partial Pr(Y=1)}{\partial HeadStart}$, to vary across individuals. In both cases, the treatment effect depends on the “index value” for each individual. However, the models maintain an assumption of constant treatment effects in logit units (β_{Logit}). If the model is misspecified, and instead there are variable treatment effects for different individuals, and that a reweighted estimation sample might produce more reliable results, especially when trying to measure the ATE for a pre-specified target group. This consideration is analogous to the treatment effect heterogeneity discussed in section 4.2.

We propose employing the in-regression weights $\widetilde{s}_f^{sw \rightarrow tg} \cdot v_f$ as discussed above in section 4.2. That is, the weights are a combination of (1) propensity score weights derived from a multinomial logit model predicting “RHS switcher” status and “in target population” status, and (2) inverse within-family conditional variance of the treatment variable of interest. For expediency, we continue to estimate this conditional variance from a linear model, and to apply it directly to the second stage logit estimation step.

We explore some of these models in the context of our empirical example, and find some differences in the point estimates and precision across linear and nonlinear specifications. Compared with LPM, we find somewhat smaller and less precise impacts of Head Start on some college when we use the 2-step approach (point estimate: 0.086 (se: 0.059)). We note that the slight decrease in precision here accompanies many fewer observations, which has fallen to 1200 for estimation of the logit beta instead of 2987 in the LPM.⁶⁷ The point estimate for the Mundlak controls is very similar to LPM, 0.126, but the standard errors are 20% larger (se: 0.053), so that the estimate is significant only at the 10% level.

E.4 Monte Carlo for Nonlinear Specifications

We next consider the bias of the different specifications in the context of a specific data generating process (DGP).

⁶⁶Beck’s second step is a logit FE (with dummies) estimator, with the β imposed from the conditional logit first stage. Then the estimated fixed effects are used to obtain the ATE.

⁶⁷Note that in the second step, the ATE is calculated over the full population. Another difference is that we weight the conditional logit regressions using family averages of individual weights, since conditional logit does not accommodate individual weights.

For our simulations, we continue with the PSID data setup presented above in Section 4.4. We take the original data, and estimate a logit model predicting some college attainment, using as regressors family level variables and family-level averages of individual variables. From this model we construct a family level logit index variable, x_f . For each simulation, the underlying logit index for each individual is equal to x_f , plus the Head Start dummy multiplied by the Head Start (logit) treatment effect. We then turn this index into $Pr(y = 1)$ using the logistic CDF, and then randomly draw outcomes y . We consider three DGPs. The first of these has a constant treatment effect (in logit units). The second has a treatment effect that is zero for small families (with 2 or 3 children), and a larger treatment effect for families with 4 or more children. The third DGP has a variable treatment effect which is decreasing linearly in x_f . For all of the DGPs, the treatment effect in terms of $Pr(y = 1)$ will vary across target populations because different children have different logit indices. For DGPs 2 and 3, there is additional variability stemming from family characteristics.

We run 2,500 Monte Carlo replications. In each replication we estimate a basic LPM, and LPM reweighted for the target population. We also estimate our two-step logit model and a logit model with Mundlak controls. For each of these we estimate both an unweighted version and a version that is reweighted for the target population. We consider the same four target populations, and present the results in Table E.1. The first panel shows results for DGP 1, with constant (in logit units) treatment effects. For this DGP, all models perform well for target groups of switchers, with biases that are small and usually not distinguishable from zero. When we target siblings, all children, or Head Start participants, the LPM model exhibits a detectable bias, which is slightly reduced by reweighting. The proposed 2-step logit model and Mundlak model do better, with small bias. However when they are reweighted with an aim to be representative of the target population, they too have a detectable bias.

In DGP 2 we now have treatment effects that vary with family size. Here all of the basic models perform poorly, both LPM and our two logit variations. Reweighting helps dramatically here, for all three models.

For DGP 3 all models give biased results when we target all children or all siblings. The three reweighted models perform roughly equally well. Each of the specifications does well for estimating treatment effects for switchers, Head Start participants, and Head Start siblings, with small biases.

In results not reported, we also explored a naive logit fixed effects specification for target groups of RHS switchers and sbilings. For these groups, this method performs similarly to the LPM, 2-step Logit, and Mundlak logit discussed above.

E.5 Discussion of Specification Choices

In our literature sample, use of OLS/LPM methods is ubiquitous. Based on the results of this section, we recommend continued use of this method. For researchers who want to pursue a logit type specification, we believe that either the two-step logit model (based off of a conditional logit estimation first step) or a logit with Mundlak controls can perform well.

Table E.1: Monte Carlo Experiments: Bias of Linear and Nonlinear Models Relative to Target ATE, and Effectiveness of Reweighting

	True ATE	LPM		Logit		Logit Reweight	
		FE Baseline	Reweight	2-Step	Mundlak	2-Step	Mundlak
<i>1. Constant TE</i>							
Switchers	86.4	-0.5	-0.4	-0.3	-1.6*	-1.5	-1.8*
Siblings	78.8	7.0*	5.6*	2.0*	0.6	5.5*	4.3*
All	78.8	7.1*	5.8*	2.0*	0.7	5.7*	4.3*
HS Participants	88.1	-2.3*	-2.1*	1.0	-0.2	-2.8*	-1.0
<i>2. Large Family TE</i>							
Switchers	79.6	-10.2*	0.0	-11.5*	-10.8*	-2.5*	-0.9
Siblings	44.5	24.9*	2.6*	-9.1*	20.0*	1.1	2.0*
All	36.1	33.2*	1.6	0.5	28.3*	0.5	1.1
HS Participants	40.1	29.2*	-0.6	40.7*	30.9*	-1.7*	-0.1
<i>3. TE linear in X_f</i>							
Switchers	102.2	0.1	0.8	-1.5	-1.3	-2.3*	-1.1
Siblings	84.3	18.1*	9.3*	3.8*	10.4*	7.3*	7.9*
All	84.2	18.2*	9.5*	9.6*	10.5*	7.6*	8.0*
HS Participants	101.9	0.4	-0.2	2.6*	2.5*	-2.8*	0.8

Notes: This table shows the results from 2,500 Monte Carlo simulations for three different DGPs of some college attainment, presented separately in each panel of the table, and four different target populations, shown in each row of the panel. The true DGP is a logit model, and is discussed in Section E.4. The first panel shows results where Head Start has a constant treatment effect (TE) (on the logit index) for all individuals; the second shows results where Head Start (HS) has no effect on individuals from small families (3 or fewer children) and a large effect for families with many children (4 or more children); and the third panel shows results where effects are linear in X_f . Column 1, “True Beta,” presents the true average increase in the probability of completing some college for participants in Head Start in the sample, which is a function of the DGP and sample composition. The remaining columns present the bias of various estimation strategies, defined as the difference between the estimated effects of Head Start and the true beta. Columns 2 and 3, LPM and LPM reweight, are defined as in Table 4. Columns 4 to 7 show the results from using the two step random effects estimator and Mundlak logit without and with propensity score weights, respectively. Reweighted estimates obtained using in-regression weighting, which accounts for the representativeness of switchers and the conditional variance of Head Start within families. All betas are multiplied by 1,000. * $p < .01$.