# COMPARED TO WHAT? VARIATION IN THE IMPACTS OF EARLY CHILDHOOD EDUCATION BY ALTERNATIVE CARE TYPE[1]

BY AVI FELLER[*], TODD GRINDAL[†], LUKE MIRATRIX[‡]
AND LINDSAY C. PAGE[§]

*University of California, Berkeley[*], Abt Associates[†], Harvard University[‡]
and University of Pittsburgh[§]*

Early childhood education research often compares a group of children who receive the intervention of interest to a group of children who receive care in a range of different care settings. In this paper, we estimate differential impacts of an early childhood intervention by alternative care type, using data from the Head Start Impact Study, a large-scale randomized evaluation. To do so, we utilize a Bayesian principal stratification framework to estimate separate impacts for two types of Compliers: those children who would otherwise be in other center-based care when assigned to control and those who would otherwise be in home-based care. We find strong, positive short-term effects of Head Start on receptive vocabulary for those Compliers who would otherwise be in home-based care. By contrast, we find no meaningful impact of Head Start on vocabulary for those Compliers who would otherwise be in other center-based care. Our findings suggest that alternative care type is a potentially important source of variation in early childhood education interventions.

**1. Introduction.** Access to publicly funded prekindergarten in the United States has expanded substantially in recent years. In the last decade, the percentage of U.S. four-year-old children enrolled in public preschool has increased by one-third—from 31 to 40 percent—with some states now serving nearly 90 percent of all four-year-old children through publicly funded preschool programs [Barnett et al. (2014)]. Many cities, such as Boston, Los Angeles, New York, and Washington, D.C., have added to this expansion through locally-funded prekindergarten programs. The Obama Administration has called for additional funds to support even greater access to high-quality early childhood education across the country.

Those who support the expansion of publicly funded preschool point to nearly 50 years of research indicating that participation in high-quality pre-school programs can yield individual and societal benefits in both the short and long term,

often highlighting historically important interventions such as the Perry Preschool Project [e.g., Barnett (1995), Heckman (2006)]. Opponents argue that current public preschool programs, especially Head Start, the largest and most prominent public preschool program in the United States, have failed to replicate these initial successes at scale [e.g., Coulson (2013), Whitehurst (2013a)]. This belief stems in part from the results of the Head Start Impact Study (HSIS), a randomized evaluation that found that the opportunity to enroll in Head Start improved children's performance on short-term measures of cognitive and social–emotional development but that, in general, these initial impacts were no longer apparent after children finished first grade [Puma et al. (2010a)].

Researchers and policymakers have posited a wide range of explanations for differences between the Head Start results and those of early model programs like Perry preschool, including differences in program features, program intensity, and program targeting [Barnett (2011), Bitler, Hoynes and Domina (2014), Elango et al. (2015)]. We focus on one prominent explanation: that the care settings of control group children attenuated the reported effects for Head Start [e.g., National Forum on Early Childhood Policy and Programs (2010)]. In the Perry Preschool Project, all control group children were cared for in their homes by a parent or other adult. By contrast, in the Head Start Impact Study, roughly one-third of children not in Head Start enrolled in other center-based care, with services similar to those provided by Head Start, instead of receiving care in a home-based setting.

In this paper, we conduct a comprehensive analysis of the differential impact of enrolling in Head Start by the setting in which children would otherwise receive care. Our main result is that enrollment in Head Start yields strong, positive short-term effects on a measure of receptive vocabulary among those children who would enroll in Head Start when offered the opportunity to do so but who would otherwise be cared for by a parent or other caregiver at home or in a home-based setting. For this group of children, we estimate that, after one year, enrollment in Head Start improved children's performance by over 0.2 standard deviations, more than 50 percent larger than the corresponding intent-to-treat estimates reported in Puma et al. (2010a). By contrast, we find no meaningful impact of Head Start for those children who would otherwise enroll in non-Head Start center-based care.[2]

Our analysis makes three main substantive contributions. First, we find meaningful impact variation by alternative care type that is masked by the HSIS topline results. This suggests that sweeping claims of Head Start's ineffectiveness [e.g., Whitehurst (2013b)] are misplaced, at least in terms of impact on receptive vocabulary. At the same time, we find no evidence that other center-based alternatives are more effective than Head Start on average, despite research arguing that this might be the case [Gormley et al. (2010)]. Second, this pattern of impact variation broadly holds across outcome quantiles [Bitler, Hoynes and Domina (2014)]

---

[2]These results are corroborated in independent work by Kline and Walters (2016), who find the same general pattern using a structural model. We compare our approaches in Section 7.

and within key subgroups [Bloom and Weiland (2014)], although these estimates are imprecise. We find especially large impacts among Dual-Language Learner children who would otherwise be in home-based care. Third, consistent with the HSIS results [Puma et al. (2010a)], we find that, while the impact of Head Start indeed declines over time, it is a gradual decline rather than the rapid attenuation identified by prior work [Gibbs, Ludwig and Miller (2013)]. We also find modest evidence of positive impacts of Head Start through first grade.

Our paper also makes several methodological contributions. First, we set up an approach for identifying and estimating impacts in the presence of multiple counterfactual treatment options, which is common in early childhood education studies and in program evaluation more generally [e.g., Duncan and Magnuson (2013), Heckman et al. (2000)]. To do so, we use the *principal stratification* framework of Frangakis and Rubin (2002), which is a generalization of the usual instrumental variables (IV) approach for noncompliance in randomized experiments [Angrist, Imbens and Rubin (1996)]. In the standard IV case, the goal is to estimate the impact of randomization for Compliers, known as the Local Average Treatment Effect (LATE). In HSIS, Compliers are children who would enroll in Head Start under treatment and would not enroll in Head Start under control. In our analysis, we are instead interested in two different types of Compliers: Center-based Compliers, children who would enroll in Head Start under treatment and would enroll in other center-based care under control, and Home-based Compliers, children who would enroll in Head Start under treatment and would otherwise enroll in home-based care. This approach yields two LATEs, rather than just one.

Identifying and estimating impacts for these subgroups is challenging. Extending results from the IV setting [Abadie (2003), Imbens and Rubin (1997a)], we first show that a range of quantities of interest can be immediately estimated using moment-based methods, including the relative sample shares of Center- and Home-based Compliers and the outcome distributions for these groups under control. The outcome distributions under treatment, however, are more difficult to estimate. To overcome these obstacles, we therefore utilize a hierarchical Bayesian modeling approach [e.g., Imbens and Rubin (1997b)]. In addition to providing a natural paradigm for causal inference with potential outcomes, this approach easily allows us to account for many of the real-world complications in the Head Start Impact Study, including missing data and a multilevel structure, with children nested within Head Start centers. We estimate this model via an implementation of Hamiltonian Monte Carlo called Stan [Stan Development Team (2014)], which builds on recent advances in Bayesian computation. To the best of our knowledge, this is the first implementation of a principal stratification model with site-level random effects.

We organize the paper as follows. Section 2 gives background on Head Start and the principal stratification approach. Section 3 describes the HSIS data. Sections 4 and 5 provide an overview of the analytic framework and give some descriptive

information about the principal strata. Section 6 gives an overview of our identification and estimation approaches. Section 7 presents our results. We close with a discussion of the substantive implications for this work for early childhood policy and reflect on the broader methodological implications. We defer all detailed technical discussions and proofs to the Supplementary Material [Feller et al. (2016)].

## 2. Background.

2.1. *Background on Head Start and the Head Start impact study.* Originally launched in the summer of 1965 as a two-month intervention to help low-income children prepare for kindergarten, Head Start programs across the United States currently provide early childhood education and family support services to more than 900,000 low-income children and their families each year. Head Start services are administered by nearly 1600 local grantee agencies that receive a total of $8 billion in annual state and federal funds [Administration for Children and Families (2014)]. Today, Head Start programs must adhere to a set of performance standards that specify requirements for program services, curricula, teacher preparation and professional development. For example, current Head Start classes serving four- or five-year-olds can have no more than 20 children, and those serving three-year-olds can have no more than 17 children. Programs must screen all enrolled children for developmental, sensory and behavioral disabilities and have a written curricula to support each child's cognitive and language development. Head Start programs are also required to engage in collaborative partnership-building with parents through processes that include structured home visits, parenting education classes and assistance in accessing food, housing, clothing and transportation.

Researchers and policy makers have debated the effectiveness of Head Start since the program's inception. In their summary of the initial research on Head Start from the 1960s, Zigler and Muenchow (1992) show that children enrolled in these early evaluations of Head Start exhibited large gains on measures of cognitive achievement between their initial enrollment and program completion. Excitement regarding these impressive findings was soon tempered, however, by additional research indicating that the effects of Head Start participation were no longer apparent once children reached elementary school [Westinghouse Learning Corporation (1969)]. Nevertheless, many of the quasi-experimental studies that followed over the next four decades indicated positive impacts of Head Start on a range of outcomes from short-term academic skill development to long-term outcomes measured in adulthood [e.g., Carneiro and Ginja (2014), Currie and Thomas (1995), Deming (2009), Garces, Thomas and Currie (2002), Ludwig and Miller (2007)].

The mixed results of the randomized Head Start Impact Study did little to settle this debate [National Forum on Early Childhood Policy and Programs (2010)]. Nonetheless, the rich HSIS data has led to a host of secondary analyses. Bloom and Weiland (2014) and Walters (2015), for example, examine impact variation across Head Start centers, finding substantial heterogeneity. Bitler, Hoynes and Domina

(2014) use quantile regression to examine impact variation across the entire outcome distribution, finding substantially larger effects for children with low scores. Bitler, Hoynes and Domina (2014) and Bloom and Weiland (2014) also examine heterogeneity across important subgroups, with both studies highlighting significantly larger effects among Dual-Language Learners than among native English speaking students. Finally, other studies, such as Gelber and Isen (2013) and Miller et al. (2014), find that parents play an important role in the effects of Head Start.

2.2. *Heterogeneity by alternative care type.* The goal of this paper is to explore a specific type of impact heterogeneity: whether or not the impact of Head Start varies by alternative care type. There is substantial evidence in the literature suggesting that this might be the case. First, a recent meta-analysis of 28 studies of Head Start conducted between the program's inception and 2007 found that much of the variation in the findings regarding Head Start's impact on child achievement and cognitive development could be explained by differences in the types of preschool services used by the control group [Shager et al. (2013)]. Although studies of Head Start programs yielded overall positive effects on short-term indicators of children's cognitive skills and achievement, with average effect sizes of 0.27, those studies in which the children in the control group experienced other forms of center-based care yielded significantly smaller effects as compared to those studies of Head Start in which control group children received no additional services [see also Duncan and Magnuson (2013) for a broader discussion of the counterfactual problem]. Zhai, Brooks-Gunn and Waldfogel (2011) find a similar result using longitudinal data from the Fragile Families and Child Wellbeing Study, concluding that impacts of Head Start were largest relative to noncenter-based care.

Second, a few authors have used HSIS data to address this question. Using a matching approach, Zhai, Brooks-Gunn and Waldfogel (2014) find significant effects of Head Start compared to parent care and relative/nonrelative care, but find no meaningful differences in outcomes between Head Start and other center-based care. Using variation across sites, Walters (2015) finds that impacts are smaller for Head Start centers that draw more children from other center-based programs rather than from home-based care. Finally, using a structural model, Kline and Walters (2016) find that the effects of Head Start are larger relative to home-based care than relative to other center-based care. We discuss the relationship between our results and those of Kline and Walters (2016) in Section 7.

At the same time, some authors have argued against alternative care type as an important source of impact variation. Bitler, Hoynes and Domina (2014), for example, find no relationship between observed impacts and the distribution of counterfactual care type across a range of subgroups in HSIS. Barnett (2011) points to the Abecedarian study, initially launched in 1972, which demonstrated large, sustained program impacts, even though roughly two-thirds of control group children attended high-quality center care.

2.3. *Principal stratification.*    There is a small but growing literature on the use of model-based principal stratification in social science applications. Page et al. (2015) provide a recent nontechnical review [see also Schochet, Puma and Deke (2014)]. Some previous education examples include Barnard et al. (2003) on the effect of a randomized lottery for private school voucher use in New York City with complex noncompliance patterns [see also Jin and Rubin (2009)]; Page (2012) on the relative importance of student exposure to the labor market in career academy high schools; and Schochet (2013) on student mobility in school-based randomized trials. Outside of education, several studies have used principal stratification to analyze the JobCorps evaluation [e.g., Frumento et al. (2012), Zhang, Rubin and Mealli (2009)] and JOBS II evaluation [Mattei, Li and Mealli (2013)]. Finally, a separate series of papers use a *principal score* approach, rather than model-based inference, to estimate similar quantities of interest. The key assumption with this approach is *principal ignorability*: conditional on covariates, stratum membership is ignorable. Examples include Hill, Waldfogel and Brooks-Gunn (2002), who analyze the Infant Home Development Program, Schochet and Burghardt (2007), who analyze the JobCorps data, Jo and Stuart (2009), who analyze the JOBS II data, and Scott-Clayton and Minaya (2014), who analyze student employment data.

## 3. Head Start impact study.

3.1. *Overview.*    Our primary source of data is the HSIS, which was conducted within oversubscribed Head Start centers throughout the U.S. In the HSIS, children randomized to treatment were offered enrollment in a Head Start program for the 2002–2003 school year, while children randomized to control were not offered enrollment. In total, 4440 children, aged either three or four years old, were randomized to treatment or control across 351 Head Start centers. We exclude all children from Puerto Rico because they are not available in the public use data set. The randomization itself was complex; treatment probabilities varied by the child's age, the date the child was first put on a Head Start center wait list, and the distribution of eligible children across neighboring Head Start centers.[3] While it is infeasible to recreate the true randomization procedure using currently available data, we can approximately account for the complex structure of the randomization by analyzing the data as if randomization were conducted separately within each center. After excluding children from centers that did not have at least one child in each experimental condition, we obtain a data set with 4385 children across 340 Head Start centers. We refer to the first year of the study as the Head Start year.

---

[3]The official HSIS report also uses a complex set of weights to extrapolate the experimental results to a "nationally representative" population of potentially eligible Head Start children [see Gibbs, Ludwig and Miller (2013) for a discussion]. We do not use those weights here, instead focusing on the results for the experimental sample.

3.2. *Outcomes.* The HSIS research team collected a wide array of outcomes on children in the sample. A key requirement of our analytic approach, however, is the ability to find a close parametric approximation to the underlying outcome distribution. Therefore, we currently cannot assess several important cognitive outcomes, such as the Woodcock–Johnson III Applied Problems test, and social–emotional outcomes, such as externalizing behavior, since they are poorly suited to typical parametric approximations, even conditional on covariates.

We therefore restrict our analysis to the Peabody Picture Vocabulary Test (PPVT), a standardized measure of children's receptive vocabulary in which the evaluator shows the child a page containing three to four pictures and asks the child to identify the picture that best represents the meaning of a word presented orally by the assessor. See Puma et al. (2010b), Sections 3–10, for additional details on the exact form of the PPVT. The PPVT is our outcome of choice for two reasons. First, the PPVT, which is derived from an item response theory score, is unimodal and roughly bell-shaped. Second, the PPVT is a widely used assessment and is predictive of key skills later in life [Romano et al. (2010)]. Based on results from the pretest, the average child at the beginning of the HSIS performed at roughly the 30th percentile of national PPVT performance, reflecting this group's relative disadvantage in pre-academic skills.

An important complication in the HSIS is the high proportion of missing outcomes. Overall, around 18 percent of PPVT scores are missing in the Head Start year, increasing to around 22 percent two years later. Twenty five percent of PPVT pretest scores are missing. Furthermore, treatment group children are much more likely to have observed outcomes than control group children: in the Head Start year, 24 percent of control group children have missing PPVT scores, compared to just 13 percent of treatment group children. Around 40 percent of children are missing at least one PPVT score from the pretest, the Head Start year, the first follow-up year or the second follow-up year; around 10 percent do not have an observed PPVT score for any of these four tests.

3.3. *Covariates.* Covariates play a particularly important role in principal stratification models. Zhang, Rubin and Mealli (2009) point to two main functions. First, covariates can be predictive of the outcome and stratum membership. Second, parametric assumptions can often be more plausible conditional on covariates than marginally. For additional discussion, see Ding et al. (2011), Feller (2015), Hirano et al. (2000), Jo (2002), Jo and Stuart (2009).

Thankfully, the HSIS data set includes a rich set of covariates on child and family characteristics. As part of a broader research effort, we also appended center-level characteristics and neighborhood-level variables for the area around each child's Head Start center of random assignment. Neighborhood-level information includes geocoded data from the 2000 Census, the 2002 Business Census, the Department of Education and the FBI crime database [McCoy et al. (2015)].

Table 1 assesses balance across conditions for the HSIS covariates we use in our analysis. The left column shows the covariate mean for those children assigned to the control group. The middle column shows the difference between covariate means in the treatment and control groups. Finally, the right column shows the normalized differences, a standardized measure of covariate balance across treatment conditions [Imai, King and Stuart (2008), Imbens and Rubin (2015)]. There is excellent covariate balance between treatment and control groups, with all normalized differences below 0.1 in absolute value.

Overall, HSIS children had diverse background characteristics (reporting control group means for simplicity): around 30 percent identified as Black, 37 percent as Hispanic, 29 percent spoke a non-English language at home, roughly half lived with both biological parents, and one-fifth had a mother who was a recent immigrant. The children generally come from disadvantaged households: around 70 percent have a mother with at most a high school degree or GED, and around 80 percent have an assessed family risk that is moderate to high.[4] As would be expected, the children's households are generally situated in disadvantaged neighborhoods. Based on the census data for the Head Start centers, nearly one-quarter of neighborhood households were in poverty. Further, while the national unemployment rate in the US was roughly four percent in 2000, the unemployment rate in these communities was nearly eleven percent, although there is substantial heterogeneity across neighborhoods [McCoy et al. (2015)].

3.4. *Child care setting.* Standard practice in early childhood education research is to divide care settings into home-based versus center-based care [e.g., Gormley (2007)]. Given our main substantive question, we therefore categorize care settings into three main groups: Head Start, non-Head Start center care and home care. Home care encompasses a variety of home-based settings including being cared for by a parent at home (73 percent), being cared for in a nonrelative home-based child care setting (11 percent), being cared for by a relative in that relative's home (9 percent) and being cared for by a nonparent in the family's home (6 percent). Although it may be of some substantive interest to separate out these different home-based settings, it was not feasible given the small sample sizes.

Table 2 shows the distribution of observed child care settings in the Head Start year for children in the HSIS treatment and control groups. Among treatment group children, 77 percent took up the offered slot and enrolled in Head Start in the treatment year. Approximately eight percent of children assigned to treatment enrolled in a non-Head Start center, and nine percent were cared for by a

---

[4]Family risk in HSIS is based on the sum of five variables: "(1) whether the household received food stamps or TANF in Fall 2002; (2) if neither parent was a high school graduate; (3) if neither parent is working; (4) if the mother was a teen mother; (5) and if the mother is a single mother" [Puma et al. (2010b)].

TABLE 1
*Covariate balance at baseline*

| | Control Mean | T-C Diff. | Norm. Diff. |
|---|---|---|---|
| *Child characteristics* | | | |
| PPVT pretest (std.) | 0.03 | −0.05 | −0.04 |
| Bottom third by pretest | 0.32 | 0.02 | 0.03 |
| Three-year old | 0.55 | – | 0.01 |
| Male | 0.51 | – | −0.01 |
| Black | 0.30 | 0.01 | 0.02 |
| Hispanic | 0.37 | 0.01 | 0.01 |
| Dual-language learner | 0.29 | 0.01 | 0.03 |
| Special needs | 0.11 | 0.03 | 0.08 |
| *Caregiver and family characteristics* | | | |
| Caregiver age: <25 | 0.32 | −0.02 | −0.05 |
| Caregiver age: 25–29 | 0.31 | – | – |
| Caregiver age: 30–39 | 0.29 | 0.01 | 0.02 |
| Caregiver age: 40+ | 0.07 | 0.02 | 0.06 |
| Teen mother | 0.19 | −0.03 | −0.07 |
| High school dropout | 0.39 | −0.02 | −0.04 |
| Only high school diploma/GED | 0.33 | 0.01 | 0.02 |
| Married | 0.45 | −0.01 | −0.01 |
| Previously married | 0.16 | – | – |
| Urban | 0.84 | – | – |
| Family risk: medium/high | 0.22 | 0.03 | 0.06 |
| Lives with both biological parents | 0.49 | – | – |
| Recent immigrant | 0.19 | – | 0.01 |
| Any older sibling attended Head Start | 0.37 | 0.04 | 0.09 |
| Oldest child | 0.45 | −0.03 | −0.06 |
| *Head Start center of random assignment characteristics* | | | |
| Provides transportation | 0.63 | – | – |
| At least four home visits per year | 0.21 | – | −0.01 |
| Full day child care | 0.64 | – | 0.01 |
| Student–teacher ratio | 6.75 | −0.02 | −0.01 |
| All teachers certified in early childhood | 0.41 | – | – |
| All teachers have mentors | 0.46 | – | – |
| Center is always filled | 0.48 | – | – |
| Number of children randomized | 17 | – | – |
| *Neighborhood and state characteristics* | | | |
| Percent in poverty | 0.25 | – | – |
| Percent minority | 0.44 | – | – |
| Percent unemployed | 0.11 | – | – |
| Percent commute by car | 0.82 | – | – |
| Number of crimes per 1000 people | 44 | 0.1 | 0.01 |
| State has DOE Pre-K | 0.64 | – | 0.01 |
| State per-child spending ($'000) | 3.9 | – | 0.01 |
| State Head Start teacher salary ($'000) | 21.8 | – | 0.01 |

*Child care setting by treatment group, based on responses from the Spring* 2003 *parent reports. "Head Start* (*admin.*)*" refers to the administrative records collected as part of HSIS and is the compliance rate used in* Puma et al. (2010a)

| | Treatment | Control | Difference |
|---|---|---|---|
| Head Start | 0.77 | 0.11 | 0.66 |
| Other center-based care | 0.08 | 0.26 | −0.18 |
| Home-based care | 0.09 | 0.47 | −0.38 |
| Missing | 0.06 | 0.16 | −0.10 |
| Head Start (admin.) | 0.81 | 0.12 | 0.69 |

parent or other relative or enrolled in a home-based childcare program. In principle, children randomized to the control group were free to take up any available early childhood program except for that provided by the Head Start center to which they had applied and had not been offered enrollment. In practice, among control group children with an observed care setting, 13 percent enrolled in a Head Start center (most in the center in which they had lost the lottery), 31 percent enrolled in a non-Head Start center and 56 percent were cared for by a parent, a relative or within a home-based childcare program. Note that the HSIS sample consists entirely of families who actively sought to enroll a child in Head Start. Thus, there was at least some initial indication of a preference for Head Start.

**4. Analytic framework.** We next outline the technical aspects of our Bayesian principal stratification framework. We begin with a general setup for the problem, review the case with binary treatment compliance—that is, Head Start vs. not Head Start—and then extend this setup to the more general multi-valued treatment setting. Additional technical details are deferred to the Supplementary Materials.

4.1. *Overview of Bayesian principal stratification.* Following Splawa-Neyman (1990) and Rubin (1974), we set up our problem using the potential outcomes notation. Thus, the causal effects of interest are defined regardless of the mode of inference. With this setup, we explore two common inferential approaches: moment-based and model-based. In the moment-based approach, the idea is to equate the causal quantities of interest with population moments, and then introduce identifying assumptions to create valid moment estimators. In this setting, a parameter is said to be point-identified if the moment equations and identifying assumptions yield a single estimate [see Zhang and Rubin (2003) for relevant discussion]. In the Bayesian model-based approach, by contrast, unobserved potential outcomes are treated as unknown parameters to be estimated given the

model and the observed data. Importantly, identification issues are quite different from this perspective. In a Bayesian setting, proper prior distributions always yield proper posterior distributions. Thus, lack of identification results in regions of flatness of the posterior [Imbens and Rubin (1997b)], and identifying assumptions are not strictly necessary. Rather, introducing these assumptions sharpens the resulting inference.

Our primary approach in this paper is the parametric Bayesian paradigm, which has become widespread for principal stratification analysis [e.g., Hirano et al. (2000), Mattei, Li and Mealli (2013)]. First, the Bayesian approach is attractive for causal inference with potential outcomes, which is essentially a missing data problem. Second, as Imbens and Rubin (1997a) discuss, parsimoniously parameterized models can often lead to better practical performance (in the sense of lower root-mean-squared error) than corresponding moment-based approaches. Finally, we face a range of real-world complications in the HSIS example: missing data and study attrition; stratified randomization across many, small Head Start centers; and a mix of child- and center-level covariates. Addressing these issues is natural in a full Bayesian model, but would be quite difficult with moment-based approaches.

At the same time, we still find it useful to articulate the assumptions necessary for a moment-based analysis. First, while hierarchical Bayesian modeling is a powerful inferential tool, it is often difficult to determine what "drives" such models in practice. Indeed, Cox and Donnelly (2011), page 96, warn that "if an issue can be addressed nonparametrically, then it will often be better to tackle it parametrically; however, if it cannot be resolved nonparametrically, then it is usually dangerous to resolve it parametrically." We therefore believe it is useful to assess the level of danger we face. By thinking through the nonparametric approach, we highlight the importance of the Normality assumption.

4.2. *Setup and ITT.*   We observe $N$ children, $N_1$ of whom are randomized to receive the opportunity to enroll in Head Start, with treatment indicator $Z_i = 1$ for child $i$, and $N_0$ of whom are not, with $Z_i = 0$. We analyze the HSIS data as a stratified randomized experiment, with child-level randomization conducted separately within each Head Start center.

In order to use the potential outcomes notation, we first make the standard Stable Unit Treatment Value Assumption [SUTVA; Rubin (1980)], which states that the treatment assignment of one child does not affect the outcome of another child. Next, we define the relevant potential outcomes. First, let $D_i^{\text{obs}} \in \mathcal{D}$ denote the observed care setting for child $i$, where $\mathcal{D}$ is the set of possible care settings and $D_i(z)$ is the care setting child $i$ would have received if that child had been assigned to treatment condition $z$. Second, let $Y_i^{\text{obs}} \in \mathbb{R}$ denote the observed outcome of interest (e.g., PPVT), with corresponding potential outcomes, $Y_i(z)$. With this setup, $Y_i^{\text{obs}} = Z_i Y_i(1) + (1 - Z_i) Y_i(0)$ and $D_i^{\text{obs}} = Z_i D_i(1) + (1 - Z_i) D_i(0)$.

We now formalize the assumption that randomization is valid, which is sensible given that HSIS is indeed a randomized experiment [Imbens and Rubin (2015)].

TABLE 3
*Possible principal strata in the Head Start Impact Study with binary $D^*$: Head Start vs. No Head Start*

|  |  | $Z = 0$ | |
|---|---|---|---|
|  |  | **Head Start** | **Not Head Start** |
| $Z = 1$ | Head Start | Always Head Start | Complier |
|  | Not Head Start | (*Defier*) | Never Head Start |

ASSUMPTION R (Random assignment). Treatment assignment probabilities do not depend on the potential outcomes:

$$Z_i \perp\!\!\!\perp \big( Y_i(0), Y_i(1), D_i(0), D_i(1) \big).$$

Finally, we define the Intent-to-Treat (ITT) estimand as

$$\text{ITT} = \frac{1}{N} \sum Y_i(1) - Y_i(0).$$

Under assumption R, we can estimate the ITT with the usual difference-in-means estimator. We note that our estimands of interest are defined for the finite sample of $N$ children observed in HSIS, which is straightforward to estimate in the Bayesian paradigm [see Imbens and Rubin (2015) for further discussion]. However, since we also present moment-based results, we present all assumptions in terms of a superpopulation for convenience.

4.3. *IV*: $D_i^* \in \{$*Head Start*, *Not Head Start*$\}$. To introduce the overall approach, we briefly walk through the assumptions necessary to identify the Local Average Treatment Effect, following Angrist, Imbens and Rubin (1996). Let $D_i^*$ be a binary indicator for whether or not child $i$ participated in Head Start in the first year. Define child $i$'s compliance type, $S_i^*$, via the joint values $(D_i^*(0), D_i^*(1))$, as shown in Table 3. For continuity with the next section, we refer to these compliance types by the more general term, principal strata, taking values $S_i^* \in \{$Always Head Start, Never Head Start, Complier, Defier$\}$. As usual, we define the LATE as the impact of randomization on the Compliers:

$$\text{LATE} = \frac{1}{N_c} \sum_{i:S_i^*=c} Y_i(1) - Y_i(0).$$

The two standard assumptions for IV are as follows: (1) the "no defiers" assumption; and (2) the exclusion restrictions for Always Head Start and Never Head Start children.

ASSUMPTION IV-1 (IV monotonicity/No defiers). There are no individuals with $\{D_i^*(0) = 1, D_i^*(1) = 0\}$.

The monotonicity assumption states that there are no children who would enroll in Head Start when denied access to the program but who would not enroll when explicitly offered a position. While such behavior is possible in other settings, it is unlikely in the context of HSIS, where enrolling in Head Start without an available position is already quite difficult.

ASSUMPTION IV-2 (IV exclusion restrictions).    For $S_i^* \in$ {Always Head Start, Never Head Start}, $Y_i(0) = Y_i(1)$.

The exclusion restriction for Never Head Start children states that there is no effect of randomization on those children who would never enroll. While this is not always a plausible assumption [e.g., Jo and Stuart (2009)], in the context of Head Start, there is no reason to expect that turning down the offer of enrollment will have any effect on test scores. The exclusion restriction for Always Head Start children states that there is no effect of randomization on those children who would enroll in Head Start regardless of random assignment. As Gibbs, Ludwig and Miller (2013) argue, this exclusion restriction might not hold in practice. In particular, roughly half of Always Head Start children enroll in Head Start centers other than the center of random assignment. If these alternative centers systematically differ from centers of random assignment, then the exclusion restriction might not hold for this group.

From a moment-based perspective, Assumptions IV-1 and IV-2 are necessary to identify the LATE, as we discuss in Section 6.2 [Angrist, Imbens and Rubin (1996)]. From a Bayesian model-based perspective, these assumptions are not strictly necessary for inference. Therefore, it is possible to assess these assumptions by relaxing them in the model [e.g., Hirano et al. (2000), Imbens and Rubin (1997b), Mattei, Li and Mealli (2013)]. As these questions are not central to our main substantive point, however, we do not explore them further here.

4.4. *Principal stratification*: $D_i \in$ {*Head Start*, *Other Center*, *Home*}.    The IV approach allows us to estimate the impact of Head Start among Compliers. However, we wish to estimate differential impacts for children within this group. Our inferential goal is to divide the overall LATE into one LATE for those who would otherwise receive care in another non-Head Start center and a second LATE for those who would otherwise receive care in a home-based setting. To do so, we disaggregate the binary indicator, $D_i^*$, to three levels: $D_i \in$ {Head Start, Other Center, Home}. We also disaggregate the set of three standard compliance types, $\mathcal{S}^*$, into a more complete set of principal strata, $\mathcal{S}$. Table 4 shows the nine possible combinations of care types under both treatment and control. Column headings correspond to the type of care each child would experience if assigned to the control condition; row headings correspond to the type of care each would experience if assigned to the treatment condition.

TABLE 4
*Possible principal strata in the Head Start Impact Study with multi-valued $D$: Head Start, Other Center-based care, Home-based care*

| | | Z = 0 | | |
| | | **Head Start** | **Center Care** | **Home Care** |
| --- | --- | --- | --- | --- |
| $Z = 1$ | Head Start | Always Head Start | Center Complier | Home Complier |
| | Center Care | (A) | Always Center Care | (B) |
| | Home Care | (C) | (D) | Always Home Care |

As in the standard IV case, we make two key types of assumptions: monotonicity assumptions and exclusion restrictions. The standard monotonicity assumption from the IV setting becomes a statement about four strata rather than just one. We break this statement into two parts.

ASSUMPTION PS-1a (PS monotonicity/No defiers). There are no individuals with $\{D_i(0) = \text{HS}, D_i(1) = \text{Center}\}$ or $\{D_i(0) = \text{HS}, D_i(1) = \text{Home}\}$.

Assumption PS-1a states that there are no children who would take up Head Start under assignment to control but not under assignment to treatment. Therefore, strata A and C in Table 4 do not exist. This is a natural extension of Assumption IV-1 to multi-valued $D$. As with Defiers in the IV setup, these two types of Defiers are unlikely to exist in HSIS.

ASSUMPTION PS-1b (Irrelevant alternatives). There are no individuals with $\{D_i(0) = \text{Center}, D_i(1) = \text{Home}\}$ or $\{D_i(0) = \text{Home}, D_i(1) = \text{Center}\}$.

Assumption PS-1b states that the Head Start offer does not change the care setting for families choosing between non-Head Start options. Therefore, strata B and D in Table 4 do not exist. Walters (2015) motivates this assumption with a revealed preference argument: since the availability of non-Head Start preschool is not affected by a Head Start offer, preferences among non-Head Start care options should not be affected either. While this is an unverifiable assumption, it is likely that, if such families do exist, they make up only a very small fraction of the overall population.

This yields five possible principal strata: Always Head Start (ahs), Always Center (ac), Always Home (ah), Center Complier (cc), and Home Complier (hc). As in the IV case, we can naturally make exclusion restrictions for principal strata unaffected by randomization. In particular, we assume zero treatment effect for the Always Head Start, Always Center and Always Home strata.

ASSUMPTION PS-2 (PS exclusion restrictions). For $S_i \in \{$Always Head Start, Always Center, Always Home$\}$, $Y_i(0) = Y_i(1)$.

The exclusion restriction for Always Head Start children here is identical to the exclusion restriction for Always Head Start children in the IV case. The exclusion restriction for Never Head Start children in the IV case directly implies the exclusion restrictions for Always Center and Always Home children here.

The remaining strata are Center Compliers and Home Compliers. Our goal is to estimate the impacts of randomization for these groups, which are the effects of receiving Head Start versus receiving other center-based care and home-based care, respectively:

$$\text{LATE}_{\text{cc}} = \frac{1}{N_{\text{cc}}} \sum_{i:S_i=\text{cc}} Y_i(1) - Y_i(0),$$

$$\text{LATE}_{\text{hc}} = \frac{1}{N_{\text{hc}}} \sum_{i:S_i=\text{hc}} Y_i(1) - Y_i(0).$$

As with the overall LATE, these are local effects since they are only defined for specific subgroups. In other words, we cannot interpret the difference between $\text{LATE}_{\text{cc}}$ and $\text{LATE}_{\text{hc}}$, as the causal effect of other center-based care versus home care—these two subgroups are not the same children. They differ across a range of unobserved and observed characteristics, such as child pretest scores and family characteristics.

Finally, the overall LATE is a weighted average of these two estimands:

$$\text{LATE} = \frac{\pi_{\text{cc}}}{\pi_{\text{cc}} + \pi_{\text{hc}}} \text{LATE}_{\text{cc}} + \frac{\pi_{\text{hc}}}{\pi_{\text{cc}} + \pi_{\text{hc}}} \text{LATE}_{\text{hc}},$$

where $\pi_s$ denotes the proportion of children in stratum $s$.

**5. Describing principal strata.** In most subgroup analyses, the groups themselves are known and fixed. For example, we can easily estimate the differential impact of Head Start for boys and girls: after collecting baseline data, each child's gender is known. While principal strata are well-defined subgroups, just like three- and four-year-olds, we cannot directly observe subgroup membership for all children.

Fortunately, we can extend some results from the IV case to provide useful descriptions of the principal strata themselves. In particular, we can nonparametrically identify the overall distribution of principal strata as well as the distribution of covariates within each stratum. We could therefore use moment-based methods to estimate these distributions. However, as discussed in Section 4.1, we instead use a Bayesian model-based approach which allows us to address important study complications. Unsurprisingly, we find that these principal strata indeed differ across observed characteristics and that this variation is consistent with intuition and results in the early childhood literature.

The Supplementary Materials give further details for the results we present below, along with proofs of all the lemmas.

TABLE 5
*Distribution of principal strata (posterior medians, with missing care type imputed)*

| Noncompliers | | | Compliers | |
|---|---|---|---|---|
| Always HS | Always Center | Always Home | Center Complier | Home Complier |
| 0.11 | 0.11 | 0.12 | 0.20 | 0.45 |

5.1. *Overall distribution of principal strata.* Extending the standard results from the IV case [Angrist, Imbens and Rubin (1996)], we can estimate the overall size of each principal stratum.

LEMMA 1 (Distribution of principal strata). *Under Assumptions R, PS-1a and PS-1b, the distribution of principal strata, $\pi_s \equiv \mathbb{P}\{S_i = s\}$, is nonparametrically identified for all $s$.*

For intuition on Lemma 1, it is useful to see the analogue in the IV setting: we first estimate the proportion of Always Head Start children in the control group and Never Head Start children in the treatment group, and then subtract to estimate the proportion of Compliers. Table 5 shows point estimates for the distribution of principal strata in the sample. Roughly one-third of all children are noncompliers of various types; each noncomplier stratum is around 10 percent of the overall sample. The remaining two-thirds are split between the two Complier groups; Home Compliers total around 70 percent of all Compliers.

5.2. *Using covariates to predict stratum membership.* Since HSIS is a randomized experiment, we can examine the distribution of principal strata for specific subgroups, such as for all boys in the sample. Following Hill, Waldfogel and Brooks-Gunn (2002), we define the *principal score* as $\pi_{s|\mathbf{x}} \equiv \mathbb{P}(S_i = s | \mathbf{X}_i = \mathbf{x})$, the probability that a child belongs to principal stratum $s$ given that child's observed covariates [see also Abadie (2003), Jo and Stuart (2009)]. Note that this is a simple generalization of modeling the "first stage" in the standard IV setting as a function of the covariates [e.g., Angrist (2004)].[5]

For HSIS, we estimate the principal score using multinomial logistic regression and a simple data augmentation procedure.[6] Figure 1 shows the resulting logistic

---

[5]Unlike the usual first stage model, $\mathbb{P}\{D^{*,\mathrm{obs}} | X_i = \mathbf{x}\}$, the principal score is vector-valued, since $S_i$ is discrete rather than binary.

[6]This approach improves on simpler versions of this model fit by Walters (2015) and Zhai, Brooks-Gunn and Waldfogel (2014). Walters (2015) effectively estimates the share of Center-based Compliers and Home-based Compliers for each Head Start center, doing so via two separate logistic regressions, rather than via multinomial logistic regression. Zhai, Brooks-Gunn and Waldfogel (2014) estimate a multinomial logistic regression using covariates to predict $D(0)$ rather than stratum membership, therefore conflating Always Center-based children and Center Compliers under control and conflating Always Home-based children and Home Compliers under control.
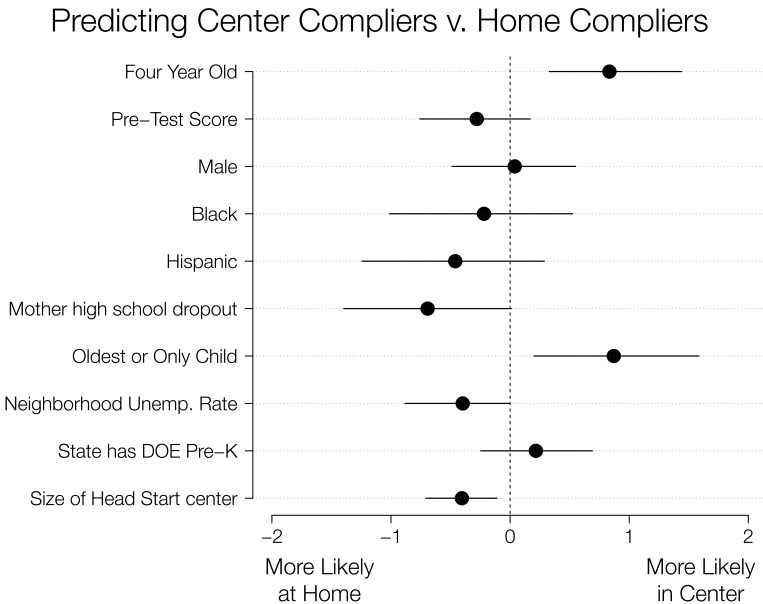
## Predicting Center Compliers v. Home Compliers



FIG. 1.  *Logistic regression coefficients predicting Center vs. Home Compliers, generated from a multinomial logistic regression predicting all types. All continuous covariates are standardized. Point estimates and error bars show posterior medians and 95% credible intervals.*

regression coefficients for select covariates that are predictive of being a Center-based vs. Home-based Complier. We discuss these results below.

5.3. *Distribution of covariates by principal stratum.*    We can also estimate the distribution of covariates for each principal stratum.

LEMMA 2 (Distribution of covariates by principal stratum).    *Under Assumptions* R, PS-1a *and* PS-1b, $\mathbb{P}\{\mathbf{X}_i = \mathbf{x}|S_i = s\}$ *is nonparametrically identified for all* $s$.

This lemma is a simple extension of the comparable IV result in Abadie (2003) and allows us to make concrete observations about otherwise unobservable groups [see also Angrist and Pischke (2008), Frumento et al. (2012)]. Table 6 shows the means for select covariates for each stratum; Figure 2 separately shows the means for pretest score by principal stratum. There are key differences in observable characteristics across the latent groups. Columns 1–3 on Table 6 show variation in pre-treatment covariates across the different types of noncompliers. Overall, these results suggest that children who always enroll in a non-Head Start center-based setting outperform their counterparts who would always be in Head Start or in a home-based setting. For example, as shown in Figure 2, Always Center-based children strongly outperform Always Head Start and Always Home-based children on

TABLE 6
*Covariate means by principal stratum*

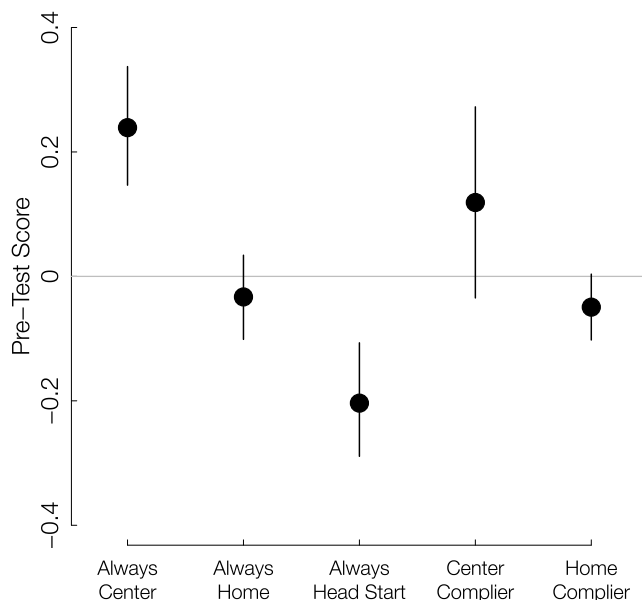| | Always Head Start | Always Center | Always Home | Center Complier | Home Complier |
|---|---|---|---|---|---|
| *Child characteristics* | | | | | |
| PPVT pretest (std.) | −0.20 | 0.24 | −0.03 | 0.12 | −0.05 |
| Bottom third by pretest | 0.36 | 0.28 | 0.35 | 0.30 | 0.35 |
| Three-year-old | 0.63 | 0.43 | 0.53 | 0.47 | 0.59 |
| Male | 0.54 | 0.53 | 0.57 | 0.48 | 0.48 |
| Black | 0.35 | 0.40 | 0.23 | 0.31 | 0.30 |
| Hispanic | 0.42 | 0.33 | 0.39 | 0.35 | 0.37 |
| Dual-language learner | 0.37 | 0.28 | 0.27 | 0.33 | 0.29 |
| Special needs | 0.15 | 0.17 | 0.12 | 0.13 | 0.11 |
| *Caregiver and family characteristics* | | | | | |
| Caregiver age: <25 | 0.29 | 0.31 | 0.38 | 0.24 | 0.32 |
| Caregiver age: 25–29 | 0.29 | 0.31 | 0.30 | 0.35 | 0.31 |
| Caregiver age: 30–39 | 0.34 | 0.28 | 0.26 | 0.30 | 0.29 |
| Caregiver age: 40+ | 0.08 | 0.09 | 0.07 | 0.11 | 0.08 |
| Teen mother | 0.16 | 0.19 | 0.20 | 0.15 | 0.17 |
| High school dropout | 0.44 | 0.27 | 0.47 | 0.35 | 0.38 |
| Only high school diploma/GED | 0.28 | 0.35 | 0.31 | 0.30 | 0.36 |
| Married | 0.46 | 0.42 | 0.47 | 0.45 | 0.44 |
| Previously married | 0.14 | 0.18 | 0.16 | 0.17 | 0.16 |
| Urban | 0.90 | 0.87 | 0.86 | 0.86 | 0.81 |
| Family risk: medium/high | 0.27 | 0.16 | 0.22 | 0.24 | 0.25 |
| Lives with both biological parents | 0.51 | 0.47 | 0.48 | 0.48 | 0.51 |
| Recent immigrant | 0.23 | 0.20 | 0.18 | 0.22 | 0.17 |
| Any older sibling attended Head Start | 0.40 | 0.34 | 0.38 | 0.34 | 0.43 |
| Oldest child | 0.43 | 0.47 | 0.45 | 0.50 | 0.39 |
| *Head Start center of random assignment characteristics* | | | | | |
| Provides transportation | 0.44 | 0.61 | 0.60 | 0.62 | 0.68 |
| At least four home visits per year | 0.15 | 0.17 | 0.21 | 0.18 | 0.25 |
| Full day child care | 0.69 | 0.75 | 0.59 | 0.67 | 0.61 |
| Student–teacher ratio | 6.66 | 6.89 | 6.58 | 7.07 | 6.64 |
| All teachers certified in early childhood | 0.50 | 0.43 | 0.41 | 0.44 | 0.38 |
| All teachers have mentors | 0.38 | 0.49 | 0.43 | 0.46 | 0.48 |
| Center is always filled | 0.50 | 0.43 | 0.44 | 0.48 | 0.49 |
| Number of children randomized | 14 | 18 | 15 | 16 | 18 |
| *Neighborhood and state characteristics* | | | | | |
| Percent in poverty | 0.27 | 0.25 | 0.23 | 0.27 | 0.24 |
| Percent minority | 0.55 | 0.49 | 0.40 | 0.45 | 0.40 |
| Percent unemployed | 0.12 | 0.11 | 0.10 | 0.11 | 0.10 |
| Percent commute by car | 0.72 | 0.77 | 0.82 | 0.81 | 0.85 |
| Number of crimes per 1000 people | 49 | 45 | 42 | 47 | 43 |
| State has DOE Pre-K | 0.72 | 0.69 | 0.59 | 0.68 | 0.62 |
| State per-child spending ($'000) | 3.4 | 3.8 | 3.4 | 4.1 | 4.2 |
| State Head Start teacher salary ($'000) | 21.1 | 21.7 | 21.3 | 21.9 | 22.1 |

FIG. 2. *PPVT pretest score by principal stratum. Point estimates and error bars show posterior medians and* 95% *credible intervals.*

the PPVT pretest. Other covariates also sensibly predict differences among the noncomplier types. For example, Always Center children are much more likely to live in a state that has state-funded preschool than Always Home children. In general, this ordering is consistent with the selection results from Deming (2009), who finds that families of children in non-Head Start preschools have higher income and maternal education than families of children in Head Start or in no preschool.

We can compare our two Complier groups by examining columns 4 and 5 of Table 6, which are the complement to the logistic regression coefficients in Figure 1. Consistent with research that has found that parents typically prefer center-based care for four-year-olds [e.g., Huston, Chang and Gennetian (2002), Rose and Elicker (2010)], roughly 60 percent of Home Compliers are three years old, compared to only 45 percent of Center Compliers. We also find that Home Compliers enter the study with lower pre-academic skills. Home Compliers exhibit lower PPVT performance at the beginning of the study and are more likely to be in the bottom third of PPVT performance compared to Center Compliers. Home Compliers are additionally more likely to have a mother with less than a high school education. As above, Center Compliers are more likely to live in states that, during the time of the HSIS, provided state-funded prekindergarten. Note that we do not find meaningful differences between these two groups based on race or ethnicity or based on Dual Language Learner status.

Overall, these differences in covariate means by principal stratum underscore that children in different principal strata do, indeed, differ in terms of their baseline

characteristics. Therefore, while estimates of causal effects *within* each principal stratum are valid, comparisons *between* principal strata are descriptive rather than causal, in the same way that comparing treatment effects for males and females is descriptive rather than causal. In other words, differential impacts across strata could also be due to differences in observed or unobserved characteristics other than care type. See Gallop et al. (2009) for a discussion of using principal stratification for mediation analysis, which generally requires much stronger assumptions than those presented here.

**6. Overview of identification and estimation.** This section provides an overview of the identification and estimation strategies used in this paper. Interested readers can find greater detail in the Supplementary Materials, which give an in-depth discussion of possible identification approaches, our hierarchical Bayesian estimation procedure, robustness to different parametric assumptions and other technical information. Conversely, readers can skip to Section 7 for a discussion of the results.

6.1. *Identification*.    The identification strategy rests on the idea that we can identify the outcome distributions for each principal stratum. This builds on earlier work in the IV case from Imbens and Rubin (1997a) and Abadie (2003). We provide a brief sketch of the idea here. The Supplementary Materials provide additional discussion of identification in principal stratification models [see also Zhang, Rubin and Mealli (2009)].

To illustrate the identification approach, first consider a standard subgroup analysis, for example, estimating the impact of Head Start for the subgroup of boys. Formally, we can achieve this in two distinct steps. The first step is to identify the distribution of outcomes for boys in the treatment group, which we denote $g_{\text{boys } 1}(y)$, and the corresponding distribution of outcomes for boys in the control group, which we denote $g_{\text{boys } 0}(y)$. Since HSIS is a randomized experiment, and since we directly observe which children are boys, we can nonparametrically identify both $g_{\text{boys } 0}(y)$ and $g_{\text{boys } 1}(y)$ from the corresponding sample (e.g., via kernel density estimation). We can then obtain the average impact of Head Start on boys by comparing the means of the two distributions. While not necessarily practical, this is nonetheless a valid procedure for identifying an average treatment effect for a subgroup.

6.1.1. *Instrumental variables*.    While we directly observe gender, we do not directly observe compliance type for all children. We therefore must adopt a different approach for estimating the outcome distributions by compliance type. For illustration, we again begin with the standard IV set up for noncompliance, where we compare Head Start versus not Head Start:

TABLE 7
*Relationship between Observed Care Type and Principal
Strata for binary $D^*$: Head Start vs. No Head Start*

| Z | D* | Possible Principal Strata |
|---|-----|---------------------------|
| 1 | HS | Always Head Start, Complier (treat) |
| 1 | Not HS | Never Head Start |
| 0 | HS | Always Head Start |
| 0 | Not HS | Never Head Start, Complier (control) |

- *Always Head Start and Never Head Start.* Under monotonicity, we know that any children in the control group who enroll in Head Start must be Always Head Start children. As a result, we can directly estimate the outcome distribution for the Always Head Start subgroup under control, $g_{\text{ahs } 0}(y)$. Since we assume that there is no treatment effect for this group (i.e., that the exclusion restriction holds for Always Head Start children), then $g_{\text{ahs } 1}(y) = g_{\text{ahs } 0}(y) = g_{\text{ahs}}(y)$. We can repeat this approach for Never Head Start children in the treatment group, which yields $g_{\text{nhs } 1}(y) = g_{\text{nhs } 0}(y) = g_{\text{nhs}}(y)$.
- *Compliers.* We must take a different approach for the Compliers. First, we cannot directly observe which children are Compliers. Second, since we are interested in the LATE, we can no longer assume that Compliers have the same outcome distribution under treatment and control. The key insight is to focus on the relationship between the observed treatment and the unobserved compliance type; Table 7 shows these relationships for the IV case. For example, children in the control group who do not enroll in Head Start are either Compliers or Never Head Start children. In other words, the observed outcome distribution for these children is a mixture of $g_{\text{nhs}}(y)$ and $g_{\text{co } 0}(y)$. Formally,

$$(6.1) \qquad f_{00}(y) = \frac{\pi_{\text{nhs}}}{\pi_{\text{nhs}} + \pi_{\text{co}}} g_{\text{nhs}}(y) + \frac{\pi_{\text{co}}}{\pi_{\text{nhs}} + \pi_{\text{co}}} g_{\text{co } 0}(y),$$

where $f_{zd}(y)$ is the observed outcome distribution for children with treatment assignment $Z_i = z$ and treatment received $D_i^* = d$. For example, $f_{00}(y)$ is the observed outcome distribution for children assigned to the control condition who do not experience Head Start. Since we can directly observe $f_{00}(y)$, $\pi_{\text{nhs}}$, $\pi_{\text{co}}$ and $g_{\text{nhs}}(y)$, we can re-arrange terms to identify $g_{\text{co } 0}(y)$, the outcome distribution for Complier children in the control group. We can repeat this with the mixture of Always Head Start and Compliers under treatment to obtain $g_{\text{co } 1}(y)$. Therefore, we can nonparametrically identify both $g_{\text{co } 0}(y)$ and $g_{\text{co } 1}(y)$, even though we cannot observe these distributions directly; see Imbens and Rubin (1997a) for additional discussion.

Once we have all the outcome distributions, we can immediately obtain the average outcomes by principal stratum, $\mu_{sz}$, and finally obtain LATE $= \mu_{\text{co } 1} - \mu_{\text{co } 0}$.

TABLE 8
*Relationship between Observed Care Type and Principal Strata for multi-valued D: Head Start, Other Center-based care, Home-based care*

| Z | D | Possible Principal Strata |
|---|---|---|
| 1 | HS | Always Head Start, Center Complier (treat), Home Complier (treat) |
| 1 | Center | Always Center |
| 1 | Home | Always Home |
| 0 | HS | Always Head Start |
| 0 | Center | Always Center, Center Complier (control) |
| 0 | Home | Always Home, Home Complier (control) |

Also see Kling, Liebman and Katz (2007) for an example in which the Complier means are substantively meaningful in their own right. More generally, Abadie (2003) shows that we can use this approach to identify a broad range of features by compliance type, including covariate distributions.

6.1.2. *Principal stratification.* We now extend the argument from the IV case to identify the outcome distributions for our principal strata of interest. We again have observed mixtures, as shown in Table 8.

- *Always Head Start, Always Center-based and Always Home-based.* Just as with the Always Head Start and Never Head Start groups, we directly observe the outcome distributions for the Always Head Start, Always Center-based and Always Home-based strata. For example, we directly observe the Always Home-based children under treatment and can therefore nonparametrically identify $g_{ah\ 1}(y)$. Since we assume that there is no impact of randomization on this group, $g_{ah\ 1}(y) = g_{ah\ 0}(y) = g_{ah}(y)$. We repeat this for the Always Head Start and Always Center-based strata, yielding nonparametric identification for $g_{ahs}(y)$, $g_{ac}(y)$ and $g_{ah}(y)$.
- *Center-based Compliers (control) and Home-based Compliers (control).* As in the IV case, we cannot directly observe the outcome distributions for Center-based Compliers and Home-based Compliers and must instead identify these distributions indirectly. We begin with the outcome distribution for Home-based Compliers under control, $g_{hc\ 0}(y)$. Analogous to equation (6.1), the outcome distribution for control group children in home-based care is a mixture of $g_{ah}(y)$ and $g_{hc\ 0}(y)$:

$$(6.2) \qquad f_{0\ \text{Home}}(y) = \frac{\pi_{ah}}{\pi_{ah} + \pi_{hc}} g_{ah}(y) + \frac{\pi_{hc}}{\pi_{ah} + \pi_{hc}} g_{hc\ 0}(y),$$

where we have previously identified $g_{ah}(y)$. Similarly, we rearrange terms to nonparametrically identify $g_{hc\ 0}(y)$ and repeat this procedure for the Center-based Compliers under control, $g_{cc\ 0}(y)$.

- *Center-based Compliers (treated) and Home-based Compliers (treated).* Identifying the corresponding Complier distributions under treatment requires additional steps. As in the IV case, we can reduce the problem to estimating a mixture of two types:

$$(6.3) \qquad f^*_{1\,\mathrm{HS}}(y) = \frac{\pi_{\mathrm{cc}}}{\pi_{\mathrm{cc}} + \pi_{\mathrm{hc}}} g_{\mathrm{cc}\,1}(y) + \frac{\pi_{\mathrm{hc}}}{\pi_{\mathrm{cc}} + \pi_{\mathrm{hc}}} g_{\mathrm{hc}\,1}(y),$$

where $f^*_{1\,\mathrm{HS}}(y)$ is the observed outcome distribution after "backing out" the Always Head Start outcome distribution. Unlike the IV case, however, neither mixture component is known, which leads to a two-component finite mixture. Without additional assumptions, the component densities, $g_{\mathrm{cc}\,1}(y)$ and $g_{\mathrm{hc}\,1}(y)$, are not identified.

Therefore, the key inferential challenge, at least implicitly, is estimating the parameters of a two-component finite mixture. Once we obtain the relevant component means, $\mu_{\mathrm{cc}\,1}$ and $\mu_{\mathrm{hc}\,1}$, we can then estimate $\mathrm{LATE}_{\mathrm{cc}} = \mu_{\mathrm{cc}\,1} - \mu_{\mathrm{cc}\,0}$ and $\mathrm{LATE}_{\mathrm{hc}} = \mu_{\mathrm{hc}\,1} - \mu_{\mathrm{hc}\,0}$.

There are many possible approaches to disentangle the finite mixture model. Since we adopt a Bayesian parametric framework here, it is natural to assume that the component densities, $g_{\mathrm{cc}\,1}(y)$ and $g_{\mathrm{hc}\,1}(y)$, follow a parametric distribution, namely, Normality. In a classic result, Pearson (1894) showed that the component parameters are all identified under this assumption. Similar results hold for a broad class of parametric models [Frühwirth-Schnatter (2006)] and for distributions with shape restrictions, such as symmetry [Bordes, Mottelet and Vandekerkhove (2006), Hunter, Wang and Hettmansperger (2007)]. Note that, as we discuss in the next section, our model imposes the Normality assumption on the outcome residuals (i.e., conditional on covariates) rather than on the marginal outcome distributions.

Finally, it is useful to briefly review some alternative strategies that leverage auxiliary covariates to disentangle the finite mixture model [Joffe, Small and Hsu (2007)]. First, researchers cam assume that, conditional on covariates, stratum membership is independent of potential outcomes, an assumption known as principal ignorability. This can be a sensible assumption in some settings [e.g., Hill, Waldfogel and Brooks-Gunn (2002), Schochet and Burghardt (2007), Scott-Clayton and Minaya (2014)], but seems somewhat implausible here, as we do not observe critical variables like parental preference for care type prior to randomization. Second, researchers can restrict the relationship between a special covariate and the outcome, for example, assuming that the treatment effect does not vary across site [Raudenbush, Reardon and Nomi (2012)]. While many such restrictions are possible [e.g., Ding et al. (2011), Jo (2002), Mealli and Pacini (2013)], there is no clear candidate for such a special covariate in HSIS, nor is it plausible to assume that the treatment effect is constant across Head Start centers. Finally, see Hall and Zhou (2003) and Mealli and Pacini (2013) for assumptions when there are multiple, independent outcomes.

6.2. *Estimation.* We now turn to model-based estimation. In practice, we could estimate the full parametric model from either a likelihood or Bayesian perspective. Indeed, some prominent applications of model-based principal stratification utilize a direct likelihood approach [e.g., Frumento et al. (2012), Zhang, Rubin and Mealli (2009)]. This approach is quite flexible and allows for straightforward comparisons between different models. It is especially attractive when specifying prior distributions is not desirable. An important feature of the Head Start data, however, is the multilevel structure of children nested within Head Start centers. Incorporating this structure is immediate with a Bayesian approach but can prove quite complex in a likelihood setting. In addition, accounting for uncertainty in the parameter estimates is natural with a Bayesian approach but can be more involved with a direct likelihood approach [see, for example, Frumento et al. (2016)]. While we use a Bayesian estimation approach, we would expect quite similar results using either method.

6.2.1. *Sketch of data augmentation.* To develop intuition, we first give a high-level sketch of a data augmentation procedure for estimating the parameters of interest. To focus on the core estimation problem, we initially ignore important complications, returning to them below. The key idea is to alternate between (1) estimating the vector of model parameters, $\theta$, given stratum membership, $S$, and (2) imputing each child's principal stratum membership, $S$, given $\theta$, beginning with an initial guess of principal stratum membership for each child:

- *Step* 1: *Given stratum membership, estimate model parameters.* We estimate model parameters via two submodels.
  – *Step* 1A: *Outcome submodel,* $g_{sz|\mathbf{x}}(y)$. First, we estimate the regression of $Y^{\mathrm{obs}}$ on $\mathbf{X}$ and $Z$ within each principal stratum, $S$. The critical assumption is that the residuals follow a Normal distribution.
  – *Step* 1B: *Principal score submodel,* $\pi_{s|\mathbf{x}}$. Second, we estimate a multinomial logistic regression predicting $S$ given $\mathbf{X}$.
- *Step* 2: *Given model parameters, predict stratum membership.* Given the outcome submodel, $g_{sz|\mathbf{x}}(y)$, and principal score submodel, $\pi_{s|\mathbf{x}}$, we can estimate the probability of stratum membership via Bayes' Rule. For example, if we observe a child in the control group who is in home-based care, the child's probability of being a Home Complier is

$$\mathbb{P}\{S_i = \mathrm{hc}|\mathrm{data}, \theta\} = \frac{\pi_{\mathrm{hc}|\mathbf{x}} \cdot g_{\mathrm{hc}\ 0|\mathbf{x}}(y)}{\pi_{\mathrm{hc}|\mathbf{x}} \cdot g_{\mathrm{hc}\ 0|\mathbf{x}}(y) + \pi_{\mathrm{ah}|\mathbf{x}} \cdot g_{\mathrm{ah}|\mathbf{x}}(y)}.$$

We then flip a weighted coin to predict $S_i$ for that child. By contrast, if we observe a child in the treatment group who is in home-based care, then the child must be in the Always Home-based stratum, and so $\mathbb{P}\{S_i = \mathrm{ah}|\mathrm{data}, \theta\} = 1$.

6.2.2. *Model details*. The actual model is considerably more complex. We highlight key issues here and defer additional technical details to the Supplementary Materials. First, the outcome models by principal stratum are as follows:

$$y_i^{\text{obs}}|(S_i = \text{ahs}, \theta, \mathbf{x}_i, z_i) \sim \mathcal{N}(\alpha_{\text{ahs}} + \beta_{\text{ahs}}\mathbf{x}_i + \psi_{j[i]}, \sigma_{\text{ahs}}^2),$$

$$y_i^{\text{obs}}|(S_i = \text{ac}, \theta, \mathbf{x}_i, z_i) \sim \mathcal{N}(\alpha_{\text{ac}} + \beta_{\text{ac}}\mathbf{x}_i + \psi_{j[i]}, \sigma_{\text{ac}}^2),$$

$$y_i^{\text{obs}}|(S_i = \text{ah}, \theta, \mathbf{x}_i, z_i) \sim \mathcal{N}(\alpha_{\text{ah}} + \beta_{\text{ah}}\mathbf{x}_i + \psi_{j[i]}, \sigma_{\text{ah}}^2),$$

$$y_i^{\text{obs}}|(S_i = \text{cc}, \theta, \mathbf{x}_i, z_i) \sim \mathcal{N}(\alpha_{\text{cc}} + \beta_{\text{cc}}\mathbf{x}_i + \psi_{j[i]} + \tau_{\text{cc}}z_i + \omega_{j[i],\text{cc}}z_i, \sigma_{\text{cc},z}^2),$$

$$y_i^{\text{obs}}|(S_i = \text{hc}, \theta, \mathbf{x}_i, z_i) \sim \mathcal{N}(\alpha_{\text{hc}} + \beta_{\text{hc}}\mathbf{x}_i + \psi_{j[i]} + \tau_{\text{hc}}z_i + \omega_{j[i],\text{hc}}z_i, \sigma_{\text{hc},z}^2),$$

where $j[i]$ denotes the site $j$ corresponding to child $i$. Within each stratum, this is essentially a varying intercept/varying slope model. To improve the stability of the model, the variance terms for the two complier groups under treatment are constrained to be equal, $\sigma_{\text{cc } 1}^2 = \sigma_{\text{hc } 1}^2$.[7] Given small sample sizes within each site, the random effects for site, $\{\psi_j\}$, are constrained to be equal across principal strata, although the treatment effects are allowed to differ. The site-level estimates follow a multivariate Normal distribution:

$$\begin{pmatrix} \psi_j \\ \omega_{j,\text{cc}} \\ \omega_{j,\text{hc}} \end{pmatrix} \sim \mathcal{N}\left( \begin{pmatrix} \boldsymbol{\gamma}^{\text{ctr}}\mathbf{w}_j \\ 0 \\ 0 \end{pmatrix}, \Sigma_y \right),$$

where $\mathbf{w}_j$ is a vector of site-level covariates and $\Sigma_y$ is an unconstrained covariance matrix. We include the proportion assigned to treatment, $\bar{z}_j$, as a site-level predictor in order to account for differing proportions randomized to treatment by site [see Bafumi and Gelman (2006), Raudenbush (2015)].

We also introduce a multilevel structure in the multinomial logistic regression model:

$$\mathbb{P}(S_i = s|\theta, \mathbf{x}_i) = \frac{\exp(\gamma_{s,j[i]} + \delta_s'\mathbf{x}_i)}{\sum_{s=1}^K \exp(\gamma_{s,j[i]} + \delta_s'\mathbf{x}_i)},$$

$$\gamma_{s,j} \sim \mathcal{N}(\mu_{\gamma,s} + \delta_s^{\text{ctr}}\mathbf{w}_j, \eta_{\gamma,s}^2),$$

where the site-level random effects are independent across strata; see the Supplementary Materials for additional details.

---

[7]Relaxing the constraint that $\sigma_{\text{cc } 1}^2 = \sigma_{\text{hc } 1}^2$ gives comparable results but leads to worse model fit since identification for these variance terms is rather weak. Alternatively, Imbens and Rubin (1997b) suggest modeling the variance based on treatment received rather than treatment assigned, which would lead to $\sigma_{\text{cc } 1}^2 = \sigma_{\text{hc } 1}^2 = \sigma_{\text{ahs}}^2$ in this context. While this is a stronger assumption than the equal variance case above, invoking this assumption reduces the number of unknown parameters in the mixture model by one; see also Griffin, McCaffrey and Morral (2008).

Three additional points are worth noting. First, as discussed in Section 3.2, there is considerable missingness in HSIS, especially in the outcomes. We address this by assuming that outcomes are Missing at Random (MAR) [Rubin (1976)],

$$\mathbb{P}\{M_i|Y_i, \mathbf{X}_i, Z_i, D_i^{\mathrm{obs}}\} = \mathbb{P}\{M_i|\mathbf{X}_i, Z_i, D_i^{\mathrm{obs}}\},$$

where $M_i$ is an indicator for missing outcome. In other words, given covariates, treatment assignment and observed child care setting, missing outcomes are just as likely to be low test scores as high test scores. While we address alternative assumptions in the Supplementary Materials, MAR is at least plausible for HSIS since the data collection procedures depended heavily on the child's actual care setting. Although implicit, this is also the assumption behind the nonresponse adjustment in the official HSIS report [Puma et al. (2010a)].

Second, as we discuss in Section 7.3, the treatment effect varies across observed covariates. Given the complexity of the base model, however, we report these treatment-by-covariate interactions one at a time. Since including multiple treatment-by-covariate interactions unsurprisingly yields poor model convergence, the main results are from a model that excludes such interactions. Finally, we use standard reference priors throughout; see the Supplementary Materials for additional details.

6.2.3. *Computational details.* While this data augmentation procedure helps to build intuition for estimation, convergence of the algorithm can be slow in practice. Instead, we estimate this model via Stan, a Bayesian programming language that implements a variant of Hamiltonian Monte Carlo [HMC; Hoffman and Gelman (2014), Stan Development Team (2014)]. Unlike, say, a classic Gibbs sampler, HMC-based samplers explore the space of the (log) posterior far more efficiently than more standard Markov chain Monte Carlo approaches, dramatically increasing the effective sample size of the same number of draws [Hoffman and Gelman (2014)]. One drawback of the HMC approach is that the log-posterior must have globally smooth gradients. As a result, Stan/HMC cannot incorporate discrete latent parameters, such as indicators for principal stratum membership that would be standard in a data augmentation scheme. Stan sidesteps this issue by maximizing the observed data log-posterior rather than the complete data log-posterior. While it is possible to couple a data augmentation Gibbs step with a bespoke HMC sampler, doing so would lose many of Stan's key advantages, including optimized C++ code and a powerful, flexible programming language. In the end, it is unlikely that this project would have been feasible without the development of Stan.

Each model was run with five separate chains with 500 "warm up" draws and 500 posterior draws. We assess model convergence in the usual way via traceplots, via Gelman–Rubin $\widehat{R}$ statistics at or near 1, and via measures of the effective sample size from each chain. All models reported here showed excellent convergence for parameters of interest. As with all hierarchical models, some hyperparameters were poorly estimated; we do not report those.

TABLE 9
*Impacts in the Head Start Year. Point estimates are
posterior medians, with 2.5 and 97.5 quantiles of
posterior distribution in parentheses. 95% posterior
intervals that exclude zero are printed in bold*

| **Panel A. ITT Model** | |
| --- | --- |
| ITT | **0.14** |
| | **(0.11, 0.16)** |
| **Panel B. IV Model** | |
| Overall LATE | **0.18** |
| | **(0.14, 0.23)** |
| **Panel C. Principal Stratification Model** | |
| LATE for Center Compliers | 0.00 |
| | (−0.13, 0.14) |
| LATE for Home Compliers | **0.23** |
| | **(0.15, 0.30)** |
| $\mathbb{P}\{\text{LATE}_{hc} > \text{LATE}_{cc}\}$ | 0.99 |

**7. Results.** We now summarize results for the Intent-to-Treat, Instrumental Variable and Principal Stratification models, beginning with impacts in the Head Start year. We then briefly explore impacts after the first year as well as additional impact heterogeneity, including distributional treatment effects. Finally, we report sensitivity and robustness checks.

7.1. *Impacts in the Head Start Year.* The first row of Table 9 shows the ITT estimate, the impact of opportunity to enroll in Head Start, on PPVT in effect size units (i.e., effects scaled by the SD of the control group). Consistent with the original Head Start results [Puma et al. (2010a)], we find that the overall impact of randomization to treatment is 0.14 in the Head Start year (posterior median). There is strong evidence that this impact is greater than zero.

In general, the results we present here give greater statistical evidence that the impacts are positive than the evidence presented in Puma et al. (2010a). Multiple factors contribute to these differences. First, unlike Puma et al. (2010a), we pool the three- and four-year-old cohorts, which roughly doubles the sample size. Second, unlike Puma et al. (2010a), we control for Head Start center of random assignment in the outcome model, which improves precision. Finally, we do not use the HSIS weights, which were created to generalize the experimental results to a particular population of Head Start children. As we estimate impacts for the finite sample of children in HSIS, the corresponding standard errors are smaller; see Bloom and Weiland (2014) for additional discussion.

The second row of Table 9 shows the corresponding LATE estimate from the IV model. Among Compliers, the impact of enrolling in Head Start on PPVT is

0.18. This estimate is nearly identical to that of Bloom and Weiland (2014), who conduct a similar analysis. As with the ITT, there is strong evidence that this impact is positive.[8] This effect is comparable to the average effects of early childhood education programs reported in a recent meta-analysis [effect size of 0.21; Duncan and Magnuson (2013)] and represents approximately one-quarter of the Black-White test score gap at the end of kindergarten [Fryer and Levitt (2004)].

The last three rows of Table 9 show the principal stratification results from the full model. For Home Compliers, we find a treatment effect of 0.23 on PPVT, with strong evidence that these impacts are greater than zero. This is much larger than the ITT effect. For Center Compliers, however, we find an effect of zero. Because we jointly estimate $\text{LATE}_{hc}$ and $\text{LATE}_{cc}$, we can calculate that $\mathbb{P}\{\text{LATE}_{hc} > \text{LATE}_{cc}\} = 0.99$. As we discussed above, this is a descriptive comparison—like claiming that the treatment effect is larger for boys than girls—but it nonetheless shows that impacts for these two latent groups are meaningfully different.

A useful check is to compare the implied LATE and ITT estimates from the principal stratification model with the corresponding estimates from the IV and ITT models, respectively. In particular, the implied LATE is 0.16, which is quite close to the IV model estimate of 0.18; the implied ITT is 0.11, again close to the ITT model estimate of 0.14. This similarity is reassuring given the additional flexibility and complexity of the principal stratification model.

Another useful check is to compare our results to those of Kline and Walters (2016), who use a structural model to estimate a range of different treatment effects for the HSIS data, including $\text{LATE}_{hc}$ and $\text{LATE}_{cc}$. Identification in the Kline and Walters (2016) paper comes from two main sources: (1) assuming that the choice of a child's care setting follows a multinomial Probit discrete choice model (i.e., that the latent choice utilities follow a multivariate Normal distribution), and (2) assuming that there is no interaction between covariates and $Z$. First, our multinomial logistic regression model is analogous to their multinomial Probit model, although our modeling choice is not critical for identification. Second, our assumption of Normality on the residuals broadly takes the place of their assumption of no interaction between covariates and $Z$: both place restrictions on the heterogeneity of the outcome distributions. Thus, while our approaches are quite different in formulation [see Mealli and Pacini (2008) for a comparison of selection models and principal stratification], the underlying assumptions are similar in spirit. It is therefore reassuring that Kline and Walters (2016) also find the same overall pattern of effects, with positive and significant impacts for Home Compliers and negligible impacts for Center Compliers. While their point estimate for $\text{LATE}_{hc}$ is somewhat larger than ours (0.35 vs. 0.23), it appears as though this discrepancy is largely due

---

[8]Note that this estimate differs from the usual Wald estimator for IV, $\frac{\text{ITT}}{\pi_c} = \frac{0.14}{0.7} = 0.20$. This is primarily due to the multi-site randomization and differences in compliance rates across Head Start centers. See Raudenbush, Reardon and Nomi (2012) and Reardon and Raudenbush (2013) for further discussion of this issue.

to a different choice of outcome; Kline and Walters (2016) estimate impacts on an index of outcomes, while we focus on PPVT alone.[9]

7.2. *Impacts after the Head Start Year.* A key feature of the HSIS design is that children in the three-year-old cohort control group were given access to the Head Start program in the second year of the study. In practice, nearly half of the control group took up the opportunity to enroll, with another 34 percent enrolling in other, non-Head Start center care during that year. Enrollment was similarly high for treatment group children: 64 percent enrolled in Head Start, with another 24 percent enrolling in other center care.[10] Therefore, by the second year of HSIS, the randomization only increased the probability of enrolling in Head Start by 16 percentage points and only increased the probability of enrolling in any center-based care setting by 6 percentage points.

There are several possible approaches to address this complication. First, we could expand the number of principal strata to allow for two years of enrollment in Head Start. However, this is impractical given the complexity of just modeling care setting in the first year. Another possibility is to redefine care setting to be $\mathcal{D} \in \{$Ever in Head Start, Home-based care, Center-based care$\}$; see, for example, Kline and Walters (2016). Consistent with the official report [Puma et al. (2010a)], we focus on the setting in which the child was cared for in the first year of the intervention, even for outcomes collected in subsequent years. We believe that this is a sensible definition, as the randomization encourages participation in Head Start in the first year only. Regardless, simply pooling cohorts after the Head Start year does not yield easily interpretable results.

Following Puma et al. (2010a), we therefore analyze the results separately by cohort to assess impacts after the Head Start year. Unfortunately, further dividing Center and Home Compliers into separate three- and four-year-old subgroups makes estimation more challenging. Sample sizes are relatively small. In addition, outcome missingness increases substantially over the course of the study, with roughly a quarter of all outcomes missing by the third year. Therefore, the cohort and subgroup results presented below should all be considered exploratory.

With this caveat in mind, Figure 3 shows the treatment effect on PPVT by cohort by assessment year for all Compliers, for Center Compliers and for Home Compliers.[11] Consistent with the official HSIS results, we find a decline in the

---

[9]We can assess the influence of the outcome choice with a simple back-of-the-envelope calculation. Their estimate of the overall LATE, which is nonparametrically identified, is roughly 40 percent larger than ours (0.25 vs. 0.18); their estimate of $LATE_{hc}$ is roughly 50 percent larger than ours (0.35 vs. 0.23).

[10]For the three-year-old cohort, 22 percent of control group children and 14 percent of treatment group children do not have an observed care setting in the second year of the study. Reported percentages are among children with observed care type.

[11]These are the normative grades for a given cohort. Children who began the study as three-year-olds were able to gain access to Head Start in year 2 and then enrolled in kindergarten in year 3.
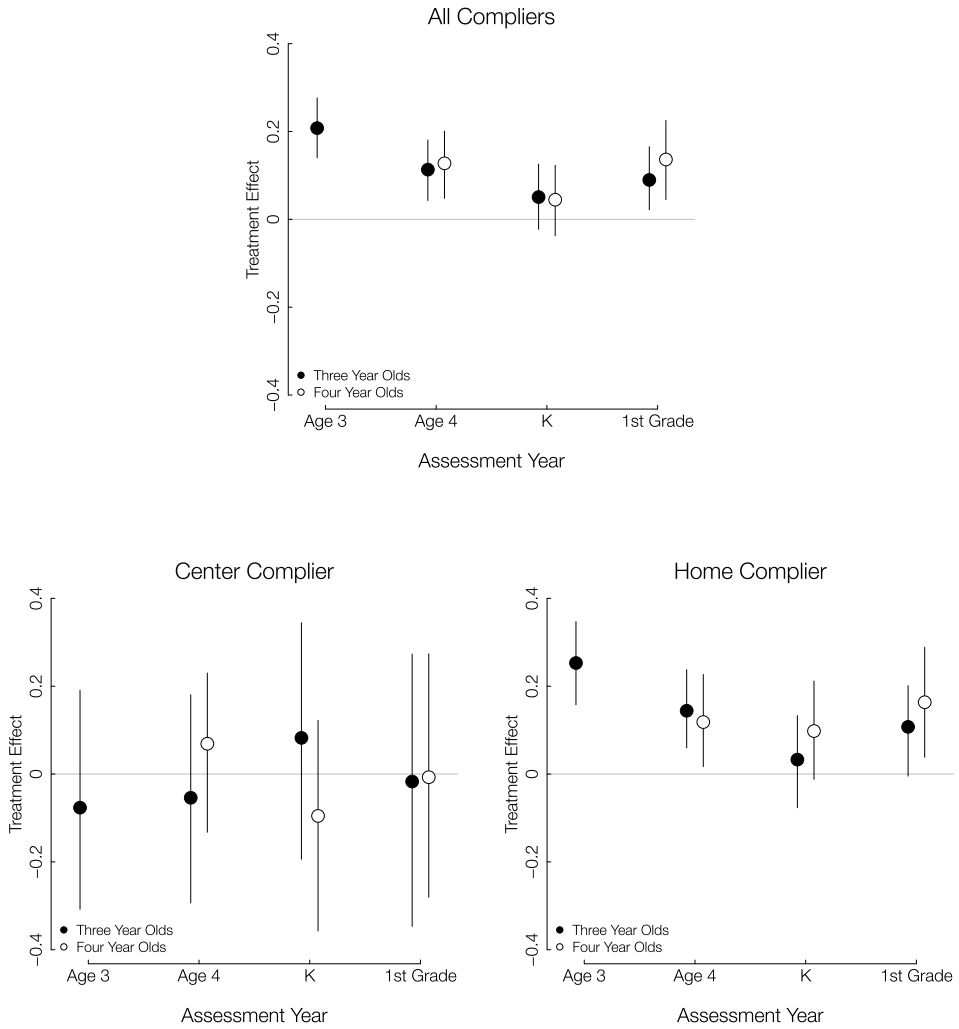
FIG. 3. *Impact estimates on PPVT by principal stratum and by three- and four-year-old cohort for each assessment year. Point estimates and error bars show posterior medians and* 95% *credible intervals. Effect sizes are calculated separately for each cohort in each assessment year.*

treatment effect as children age. Nonetheless, unlike in the official HSIS results, we find impacts that are positive and meaningfully different from zero by the time children are in 1st grade, with LATE estimates of 0.09 and 0.14 for the three- and four-year-old cohorts, respectively. The effects for Home Compliers follow

---

The four-year-olds transitioned to kindergarten and then first grade in the second and third years of the study. Therefore, by year 3, all children, if following a standard educational trajectory, were in elementary school.

the same decline as for the overall Compliers, albeit with slightly larger point estimates and with less precision. By contrast, the impacts among Center Compliers are best described as noise around zero, though this null result could be due to the limited sample size. Note that the pooled main effects in Table 9 are the (weighted) average of the impacts on three-year-olds at age 3 and four-year-olds at age 4.

While we regard these results as exploratory, they nonetheless suggest that the impact of Head Start might indeed persist into early elementary school, even if the magnitudes are modest. In particular, Gibbs, Ludwig and Miller (2013) argue that a key puzzle of the HSIS results is not that they decrease over time, but that they attenuate to zero as soon as children leave the program, much more rapidly than estimates based on quasi-experimental methods [e.g., Currie and Thomas (1995), Deming (2009)]. The results in Figure 3 show that the decline in treatment effects may not be nearly as rapid as in the reported topline results.

7.3. *Subgroup and quantile treatment effects.* Several recent papers have explored variation in Head Start's impact across observed subgroups and across quantiles of the outcome distribution. Since Center and Home Compliers differ across a range of observed and unobserved characteristics, an important question is therefore the extent to which these differences explain the different impacts for the two Complier groups. Again, these estimates should be considered exploratory.

First, we turn to variation across subgroups defined by pretreatment characteristics. Following Bloom and Weiland (2014) and Bitler, Hoynes and Domina (2014), we focus on variation by (1) whether a child is in the bottom third of pretest score by cohort, and (2) whether a child is a Dual-Language Learner (DLL). Table 10 shows the corresponding principal stratification estimates during the Head Start year. First, across all four subgroups, we observe the same pattern of positive, significant effects for Home Compliers and negligible effects for Center Compliers. While the smaller sample sizes limit statistical power, this consistency nonetheless bolsters the overall findings. Second, as in Bloom and Weiland (2014), we find larger Home Complier effects for children in the bottom third by pretest score and also for DLL students. The effect for DLL students is especially striking, with an effect size of around 0.35 SD in the Head Start year, more than double the point estimate for non-DLL students. This suggests that, at least in terms of vocabulary development, there is substantial impact of Head Start relative to a home-based setting in which English is likely not spoken; see Bloom and Weiland (2014) for additional discussion.

Another likely source of impact variation is heterogeneity across the outcome distribution [Bitler, Gelbach and Hoynes (2003)]. In a recent paper, Bitler, Hoynes and Domina (2014) estimate distributional effects for Head Start via quantile treatment effects, $G_{\mathrm{co}\,1}^{-1}(q) - G_{\mathrm{co}\,0}^{-1}(q)$, the difference between the $q$th quantiles of the outcome distributions for Compliers under treatment and control, respectively. The authors find that the impacts of Head Start on PPVT and other measures are largest at the bottom of the outcome distribution, both overall and among Compliers. As

TABLE 10
*Impacts in the Head Start Year for select subgroups. Point estimates are posterior medians, with 2.5 and 97.5 quantiles of posterior distribution in parentheses. 95% posterior intervals that exclude zero are printed in bold. Estimates are shown in effect size units, so point estimates might not average to the pooled estimate due to different outcome standard deviations*

|  | Center Compliers | Home Compliers |
|---|---|---|
| *Panel A. Bottom Third on Pretest* | | |
| Bottom Third | 0.19 | **0.30** |
|  | (−0.09, 0.47) | **(0.16, 0.45)** |
| Not Bottom Third | −0.06 | **0.21** |
|  | (−0.24, 0.16) | **(0.08, 0.31)** |
| *Panel B. DLL Status* | | |
| DLL Students | 0.06 | **0.36** |
|  | (−0.33, 0.42) | **(0.23, 0.49)** |
| Non-DLL Students | −0.04 | **0.15** |
|  | (−0.20, 0.12) | **(0.08, 0.23)** |

we discuss in the Supplementary Materials, we can leverage our framework both to replicate and to extend their results. Figure 4 shows the quantile treatment effect estimates for all Compliers, Center Compliers and Home Compliers during the Head Start year. As expected, our estimates for all Compliers are very close to those of Bitler, Hoynes and Domina (2014), showing large, positive effects at the bottom of the distribution of between 0.4 and 0.6 SD. The effects for Home Compliers are also positive and significant throughout, with larger effects at the bottom of the distribution. By contrast, the quantile treatment effects for Center Compliers are essentially zero across the entire distribution.

7.4. *Sensitivity checks.* We conducted robustness checks for our main results of impacts in the Head Start year, which we briefly discuss here. First, we assess sensitivity to our handling of missing data and refit the principal stratification model using only observed outcomes, which is approximately 80 percent of the overall sample. Table 11 shows the resulting complete case estimates, which are essentially unchanged from the full version. Second, as we discuss in Section 6, the Normality assumption plays a critical role in both identification and estimation. Table 11 shows the same model using a heavy-tailed Student $t_7$ distribution rather than a Gaussian. Again, the results are consistent.

Finally, following Rubin (1984) and Gelman et al. (2013), we use posterior predictive checks to assess the fit of our full model to the observed data. Formally, let $y$ be the observed data and $\theta$ be the parameter vector. Define $y^{\text{rep}}$ as the replicated data that could have been observed if the study were replicated with the same

## All Compliers
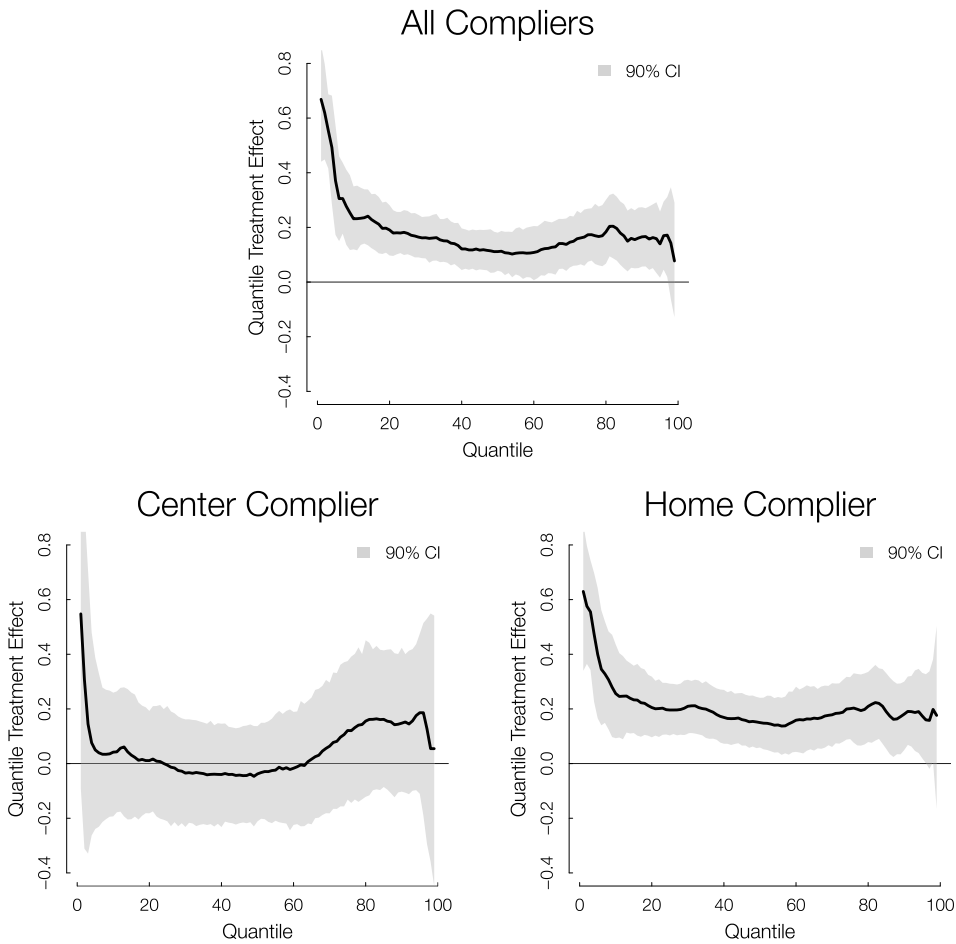


## Center Complier



## Home Complier



FIG. 4. *Quantile treatment estimates on PPVT by principal stratum for the Head Start year, with approximate 90 percent credible intervals.*

model and the same value of $\theta$ that produced $y$. We can estimate the distribution of $y^{\text{rep}}$ via the posterior predictive distribution,

$$p(y^{\text{rep}}|y) = \int p(y^{\text{rep}}|\theta) p(\theta|y) \, d\theta.$$

The intuition is to assess whether the replicated data produced from the model are similar to the observed data. In the Supplementary Materials, we assess this similarity in two ways. First, we visually inspect the observed and replicated data sets. Second, we compute posterior predictive $p$-values following a similar approach in Barnard et al. (2003) and Mattei, Li and Mealli (2013). Neither approach yields evidence that the model is a poor fit to the data.

*Sensitivity analysis for impacts in the Head Start Year. Point estimates are posterior medians, with 2.5 and 97.5 quantiles of posterior distribution in parentheses. 95% posterior intervals that exclude zero are printed in bold*

|                                        | Normal; Complete case | Student $t_7$; All observations |
|----------------------------------------|:---------------------:|:-------------------------------:|
| LATE for Center Compliers              | 0.03                  | 0.04                            |
|                                        | $(-0.07, 0.15)$       | (0.08)                          |
| LATE for Home Compliers                | **0.21**              | **0.21**                        |
|                                        | **(0.15, 0.27)**      | **(0.15, 0.26)**                |
| $\mathbb{P}\{\text{LATE}_{hc} > \text{LATE}_{cc}\}$ | 0.98     | 0.96                            |

**8. Discussion.** Our primary contribution is to develop a framework for estimating impact variation by alternative care setting and to apply this framework to the Head Start Impact Study. In particular, we find positive and meaningful impacts on key outcomes among Home-based Compliers, those children who would enroll in Head Start under treatment and who would otherwise be in home-based care. By contrast, we find no meaningful effects among Center-based Compliers, those children who would otherwise receive non-Head Start center care.

In doing so, we present a much more nuanced view of Head Start's impact than the topline experimental results indicate. We also refute sweeping generalizations made about Head Start, such as "Head Start does not improve the school readiness of children from low-income families" [Whitehurst (2013b)]. In addition, we do not find any evidence that available center-based alternatives are more effective than Head Start on average [e.g., Barnett and Haskins (2010), Gormley et al. (2010)]. In the HSIS sample, around half of the control group children who enrolled in some other form of center-based care did so in either a state-funded prekindergarten program or a prekindergarten program based in the public schools.[12] While statistical power is limited, the null finding for Center Compliers suggests that concerns over Head Start's comparative effectiveness may be misplaced.

In addition to showing larger impacts in the Head Start year, we also find that the fade out in treatment effects over time is gradual, not rapid [Gibbs, Ludwig and Miller (2013)]. This pattern closely resembles the observed fade out in other early childhood education studies [Leak et al. (2010), Magnuson, Ruhm and Waldfogel

---

[12]Like Head Start, these publicly funded programs typically feature minimum standards for important structural aspects of program quality such as teacher preparation, teacher-child ratio and curricula. This result is also consistent with a recent study in Tulsa that found that Head Start and a publicly funded prekindergarten program led to comparable school readiness [Jenkins et al. (2014)] and with the larger literature comparing quality for publicly funded versus private preschool programs [Kagan (1991), Morris and Helburn (2000)].

(2007)]. Further, while our estimates are imprecise, we find impacts between 0.10 to 0.15 for Home Complier children in first grade. These point estimates are very close to those in Deming (2009), who estimates Head Start impacts of 0.15 for children aged 5 to 6 and 0.13 for children aged 7 to 10.[13] Importantly, Deming (2009) observes outcomes for these same children in young adulthood, showing large long-term impacts. It is therefore possible that future follow up from the Head Start Impact Study will also find meaningful long-term impacts despite treatment effect fade out on short-term outcomes.

More generally, our analysis highlights the critical role that variation in counterfactual care type plays in early childhood education evaluations. Duncan and Magnuson (2013) argue that improving counterfactual conditions are a primary reason for a sharp decline in reported impacts of early childhood education interventions over the last half-century. We not only provide evidence consistent with this claim, but also outline a framework for reanalyzing other early childhood education studies to create comparable estimates. Of course, the issue of variation in counterfactual treatments is common in program evaluation settings, including for alternative schools [Bloom and Unterman (2014)] and job training programs [Heckman et al. (2000), Schochet, Burghardt and McConnell (2008)]. Our approach could easily be extended to these settings as well.

There are several promising avenues for future research. First, at present, we only analyze a single outcome of HSIS and analyze each follow-up year separately rather than jointly. Recent work from Mattei, Li and Mealli (2013) suggests that looking at multiple outcomes—either across different test scores or over time—could greatly improve inference for principal causal effects [see also Jo and Muthén (2001)]. In addition, repeated measures of the same outcome would likely make different assumptions about missingness more plausible [see, for example, Frumento et al. (2012)]. Second, while we conduct extensive sensitivity and robustness checks, inference with finite mixture models is notoriously difficult. In investigations for very simple mixture models, we have found that standard estimators can behave poorly when mixture components are not well separated [Day (1969), Feller et al. (2016)]. More work is needed to assess whether these same concerns apply to the much richer models we consider here, although Griffin, McCaffrey and Morral (2008) have taken an important step in this direction. That being said, the stability of the results to sensitivity checks, consistent patterns across subgroups and alignment with Kline and Walters (2016) are all encouraging.

In the end, our results support the argument that further efforts to improve the early skill development of US children through the expansion of publicly funded preschool programs should be targeted toward those who are currently not enrolling their children in center-based programs [for discussion, see Bassok, Fitzpatrick and Loeb (2013), Cascio and Schanzenbach (2013), Ludwig and Phillips

---

[13]The outcome in Deming (2009) combines PPVT with the Peabody Individual Achievement Tests (PIAT) for math and reading.

(2010)]. Nationwide, over 40 percent of eligible children are served by Head Start programs [Schmit et al. (2013)]. Although the availability of state and local prekindergarten has grown in recent years, many low-income children still spend their preschool years in home-based settings. In 2011, approximately 42 percent of three- and four-year-old children from low-income families enrolled in center-based prekindergarten compared to 59 percent of their nonlow-income peers [Burgess et al. (2014)]. Based on our results, shifting children from home-based care into formal care will likely lead to much larger effects than shifting children between preschool programs.

## SUPPLEMENTARY MATERIAL

**Supplement to "Compared to what? Variation in the impacts of early childhood education by alternative care type"** (DOI: 10.1214/16-AOAS910SUPP; .pdf). This files contains supporting material, additional results and proofs.

## REFERENCES

ABADIE, A. (2003). Semiparametric instrumental variable estimation of treatment response models. *J. Econometrics* **113** 231–263. MR1960380

ADMINISTRATION FOR CHILDREN AND FAMILIES (2014). Head Start Program Facts, Fiscal Year 2013. Available at https://eclkc.ohs.acf.hhs.gov/hslc/data/factsheets/docs/hs-program-fact-sheet-2013.pdf.

ANGRIST, J. D. (2004). Treatment effect heterogeneity in theory and practice. *Econ. J.* **114** C52–C83.

ANGRIST, J. D., IMBENS, G. W. and RUBIN, D. B. (1996). Identification of causal effects using instrumental variables. *J. Amer. Statist. Assoc.* **91** 444–455.

ANGRIST, J. D. and PISCHKE, J.-S. (2008). *Mostly Harmless Econometrics*: *An Empiricist's Companion*. Princeton University Press, Princeton, NJ.

BAFUMI, J. and GELMAN, A. E. (2006). Fitting multilevel models when predictors and group effects correlate. Unpublished manuscript.

BARNARD, J., FRANGAKIS, C. E., HILL, J. L. and RUBIN, D. B. (2003). Principal stratification approach to broken randomized experiments: A case study of school choice vouchers in New York City. *J. Amer. Statist. Assoc.* **98** 299–323. MR1995712

BARNETT, W. S. (1995). Long-term effects of early childhood programs on cognitive and school outcomes. *Future Child.* **5** 25.

BARNETT, W. S. (2011). Effectiveness of early educational intervention. *Science* **333** 975–978.

BARNETT, W. S. and HASKINS, R. (2010). *Investing in Young Children*: *New Directions in Federal Preschool and Early Childhood Policy*. The Brookings Institute, Washington, DC.

BARNETT, W. S., CAROLAN, M. E., SQUIRES, J. H. and BROWN, K. C. (2014). State of Preschool 2013: First Look. U.S. Dept. Education, National Center for Education Statistics, Washington, DC.

BASSOK, D., FITZPATRICK, M. and LOEB, S. (2013). Does state preschool crowd-out private provision? The impact of universal preschool on the childcare sector in Oklahoma and Georgia. NBER Working Paper 18605.

BITLER, M., GELBACH, J. and HOYNES, H. (2003). What mean impacts miss: Distributional effects of welfare reform experiments. *Am. Econ. Rev.* **96** 988–1012.

BITLER, M., HOYNES, H. and DOMINA, T. (2014). Experimental evidence on distributional effects of Head Start. Working paper.

BLOOM, H. S. and UNTERMAN, R. (2014). Can small high schools of choice improve educational prospects for disadvantaged students? *J. Policy Anal. Manage.* **33** 290–319.

BLOOM, H. S. and WEILAND, C. (2014). To what extent do the effects of Head Start on enrolled children vary across sites? Working paper.

BORDES, L., MOTTELET, S. and VANDEKERKHOVE, P. (2006). Semiparametric estimation of a two-component mixture model. *Ann. Statist.* **34** 1204–1232. MR2278356

BURGESS, K., CHIEN, N., MORRISSEY, T. and SWENSON, K. (2014). Trends in the use of early care and education, 1995–2011: Descriptive analysis of child care arrangements from national survey data. Report from the Office of the Assistant Secretary for Plannng and Evaluation, US Department of Health and Human Services.

CARNEIRO, P. and GINJA, R. (2014). Long term impacts of compensatory preschool on health and behavior: Evidence from Head Start. *Am. Econ. J. Appl. Econ.* **6** 135–173.

CASCIO, E. U. and SCHANZENBACH, D. W. (2013). The impacts of expanding access to high-quality preschool education. In *Brookings Papers on Economic Activity* 127–192. Brookings Institution, Washington, DC.

WESTINGHOUSE LEARNING CORPORATION (1969). *The Impact of Head Start: An Evaluation of the Effects of Head Start on Children's Cognitive and Affective Development, Vol. 1: Report to the Office of Economic Opportunity*. Westinghouse Learning Corporation and Ohio Univ., Athens, Ohio.

COULSON, A. J. (2013). Preschool's Anvil Chorus. Cato Institute, Washington, DC.

COX, D. R. and DONNELLY, C. A. (2011). *Principles of Applied Statistics*. Cambridge Univ. Press, Cambridge. MR2817147

CURRIE, J. and THOMAS, D. (1995). Does Head Start make a difference? *Am. Econ. Rev.* **85** 341–364.

DAY, N. E. (1969). Estimating the components of a mixture of normal distributions. *Biometrika* **56** 463–474. MR0254956

DEMING, D. (2009). Early childhood intervention and life-cycle skill development: Evidence from Head Start. *Am. Econ. J. Appl. Econ.* **1** 111–134.

DING, P., GENG, Z., YAN, W. and ZHOU, X.-H. (2011). Identifiability and estimation of causal effects by principal stratification with outcomes truncated by death. *J. Amer. Statist. Assoc.* **106** 1578–1591. MR2896858

DUNCAN, G. J. and MAGNUSON, K. (2013). Investing in preschool programs. *J. Econ. Perspect.* **27** 109–132.

ELANGO, S., GARCÍA, J. L., HECKMAN, J. J. and HOJMAN, A. (2015). Early childhood education. Technical report, National Bureau of Economic Research Working Paper No. 21766.

FELLER, A. (2015). Essays in public policy and causal inference. Ph.D. thesis, Harvard Univ., Cambridge, MA.

FELLER, A., GREIF, E., MIRATRIX, L. and PILLAI, N. (2016). Principal stratification in the Twilight Zone: Weakly separated components in finite mixture models. Available at arXiv:1602.06595.

FELLER, A., GRINDAL, T., MIRATRIX, L. and PAGE, L. C. (2016). Supplement to "Compared to what? Variation in the impacts of early childhood education by alternative care type." DOI:10.1214/16-AOAS910SUPP.

FRANGAKIS, C. E. and RUBIN, D. B. (2002). Principal stratification in causal inference. *Biometrics* **58** 21–29. MR1891039

FRÜHWIRTH-SCHNATTER, S. (2006). *Finite Mixture and Markov Switching Models*. Springer, New York. MR2265601

FRUMENTO, P., MEALLI, F., PACINI, B. and RUBIN, D. B. (2012). Evaluating the effect of training on wages in the presence of noncompliance, nonemployment, and missing outcome data. *J. Amer. Statist. Assoc.* **107** 450–466. MR2980057

FRUMENTO, P., MEALLI, F., PACINI, B. and RUBIN, D. B. (2016). The fragility of standard inferential approaches in principal stratification models relative to direct likelihood approaches. *Stat. Anal. Data Min.* **9** 58–70. MR3465093

FRYER, R. G. and LEVITT, S. D. (2004). Understanding the black–white test score gap in the first two years of school. *Rev. Econ. Stat.* **86** 447–464.

GALLOP, R., SMALL, D. S., LIN, J. Y., ELLIOTT, M. R., JOFFE, M. and TEN HAVE, T. R. (2009). Mediation analysis with principal stratification. *Stat. Med.* **28** 1108–1130. MR2662200

GARCES, E., THOMAS, D. and CURRIE, J. (2002). Longer-term effects of Head Start. *Am. Econ. Rev.* **92** 999–1012.

GELBER, A. and ISEN, A. (2013). Children's schooling and parents' behavior: Evidence from the Head Start Impact Study. *J. Public Econ.* **101** 25–38.

GELMAN, A., CARLIN, J. B., STERN, H. S., DUNSON, D. B., VEHTARI, A. and RUBIN, D. B. (2013). *Bayesian Data Analysis*. CRC press, Boca Raton, FL.

GIBBS, C., LUDWIG, J. and MILLER, D. L. (2013). Does Head Start do any lasting good? In *The War on Poverty*: *A* 50-*Year Retrospective* (M. J. Bailey and Sheldon Danziger, eds.). Russell Sage Foundation, New York.

GORMLEY, W. T. (2007). Early childhood care and education: Lessons and puzzles. *J. Policy Anal. Manage.* **26** 633–671.

GORMLEY, W. T., PHILLIPS, D., ADELSTEIN, S. and SHAW, C. (2010). Head Start's comparative advantage: Myth or reality? *Policy Stud. J.* **38** 397–418.

GRIFFIN, B. A., McCAFFREY, D. F. and MORRAL, A. R. (2008). An application of principal stratification to control for institutionalization at follow-up in studies of substance abuse treatment programs. *Ann. Appl. Stat.* **2** 1034–1055. MR2516803

HALL, P. and ZHOU, X.-H. (2003). Nonparametric estimation of component distributions in a multivariate mixture. *Ann. Statist.* **31** 201–224. MR1962504

HECKMAN, J. J. (2006). Skill formation and the economics of investing in disadvantaged children. *Science* **312** 1900–1902.

HECKMAN, J., HOHMANN, N., SMITH, J. and KHOO, M. (2000). Substitution and dropout bias in social experiments: A study of an influential social experiment. *Q. J. Econ.* **115** 651–694.

HILL, J., WALDFOGEL, J. and BROOKS-GUNN, J. (2002). Differential effects of high-quality child care. *J. Policy Anal. Manage.* **21** 601–627.

HIRANO, K., IMBENS, G. W., RUBIN, D. B. and ZHOU, X. H. (2000). Assessing the effect of an influenza vaccine in an encouragement design. *Biostatistics* **1** 69–88.

HOFFMAN, M. D. and GELMAN, A. (2014). The no-U-turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *J. Mach. Learn. Res.* **15** 1593–1623. MR3214779

HUNTER, D. R., WANG, S. and HETTMANSPERGER, T. P. (2007). Inference for mixtures of symmetric distributions. *Ann. Statist.* **35** 224–251. MR2332275

HUSTON, A. C., CHANG, Y. E. and GENNETIAN, L. (2002). Family and individual predictors of child care use by low-income families in different policy contexts. *Early Child. Res. Q.* **17** 441–469.

IMAI, K., KING, G. and STUART, E. A. (2008). Misunderstanding between experimentalists and observationalists about causal inference. *J. Roy. Statist. Soc. Ser. A* **171** 481–502. MR2427345

IMBENS, G. W. and RUBIN, D. B. (1997a). Estimating outcome distributions for compliers in instrumental variables models. *Rev. Econ. Stud.* **64** 555–574. MR1485828

IMBENS, G. W. and RUBIN, D. B. (1997b). Bayesian inference for causal effects in randomized experiments with noncompliance. *Ann. Statist.* **25** 305–327. MR1429927

IMBENS, G. W. and RUBIN, D. B. (2015). *Causal Inference for Statistics*, *Social*, *and Biomedical Sciences*: *An Introduction*. Cambridge Univ. Press, New York. MR3309951

JENKINS, J. M., FARKAS, G., DUNCAN, G. J., BURCHINAL, M. and VANDELL, D. L. (2014). Head Start at ages 3 and 4 versus Head Start followed by state pre-k: Which is more effective? Working paper.

JIN, H. and RUBIN, D. B. (2009). Public schools versus private schools: Causal inference with partial compliance. *J. Educ. Behav. Stat.* **34** 24–45.

JO, B. (2002). Estimation of intervention effects with noncompliance: Alternative model specifications. *Journal of Educational and Behavioral Statistics* **27** 385–409.

JO, B. and MUTHÉN, B. O. (2001). Modeling of intervention effects with noncompliance: A latent variable approach for randomized trials. In *New Developments and Techniques in Structural Equation Modeling* (G. A. Marcoulides and R. E. Schumacker, eds.) 57–87. Erlbaum Associates, Mahwah, NJ.

JO, B. and STUART, E. A. (2009). On the use of propensity scores in principal causal effect estimation. *Stat. Med.* **28** 2857–2875. MR2750169

JOFFE, M. M., SMALL, D. and HSU, C.-Y. (2007). Defining and estimating intervention effects for groups that will develop an auxiliary outcome. *Statist. Sci.* **22** 74–97. MR2408662

KAGAN, S. L. (1991). Examining profit and nonprofit child care: An odyssey of quality and auspices. *J. Soc. Issues* **47** 87–104.

KLINE, P. and WALTERS, C. (2016). Evaluating public programs with close substitutes: The case of Head Start. *Q. J. Econ.* To appear. DOI:10.1093/qje/qjw027.

KLING, J. R., LIEBMAN, J. B. and KATZ, L. F. (2007). Experimental analysis of neighborhood effects. *Econometrica* **75** 83–119.

LEAK, J., DUNCAN, G. J., LI, W., MAGNUSON, K. A., SCHINDLER, H. and YOSHIKAWA, H. (2010). Is timing everything? How early childhood education program impacts vary by starting age, program duration and time since the end of the program. Working paper.

LUDWIG, J. and MILLER, D. L. (2007). Does Head Start improve children's life chances? Evidence from a regression discontinuity design. *Q. J. Econ.* **122** 159–208.

LUDWIG, J. and PHILLIPS, D. A. (2010). Leave no (young) child behind: Prioritizing access in early childhood education. In *Investing in Young Children*: *New Directions in Federal Preschool and Early Childhood Policy* (R. Haskin and W. S. Barnett, eds.). Brookings and NIEER.

MAGNUSON, K. A., RUHM, C. and WALDFOGEL, J. (2007). The persistence of preschool effects: Do subsequent classroom experiences matter?. *Early Child. Res. Q.* **22** 18–38.

MATTEI, A., LI, F. and MEALLI, F. (2013). Exploiting multiple outcomes in Bayesian principal stratification analysis with application to the evaluation of a job training program. *Ann. Appl. Stat.* **7** 2336–2360. MR3161725

MCCOY, D. C., CONNORS, M. C., MORRIS, P. A., YOSHIKAWA, H. and FRIEDMAN-KRAUSS, A. H. (2015). Neighborhood economic disadvantage and children's cognitive and social-emotional development: Exploring Head Start classroom quality as a mediating mechanism. *Early Childhood Research Quarterly* **32** 150–159.

MEALLI, F. and PACINI, B. (2008). Comparing principal stratification and selection models in parametric causal inference with nonignorable missingness. *Comput. Statist. Data Anal.* **53** 507–516. MR2649105

MEALLI, F. and PACINI, B. (2013). Using secondary outcomes to sharpen inference in randomized experiments with noncompliance. *J. Amer. Statist. Assoc.* **108** 1120–1131. MR3174688

MILLER, E. B., FARKAS, G., VANDELL, D. L. and DUNCAN, G. J. (2014). Do the effects of Head Start vary by parental preacademic stimulation? *Child Dev.* **85** 1385–1400.

MORRIS, J. R. and HELBURN, S. W. (2000). Child care center quality differences: The role of profit status, client preferences, and trust. *Nonprofit Volunt. Sect. Q.* **29** 377–399.

NATIONAL FORUM ON EARLY CHILDHOOD POLICY AND PROGRAMS (2010). Understanding the Head Start impact study. Available at http://developingchild.harvard.edu.

PAGE, L. C. (2012). Principal stratification as a framework for investigating mediational processes in experimental settings. *Journal of Research on Educational Effectiveness* **5** 215–244.

PAGE, L. C., FELLER, A., GRINDAL, T., MIRATRIX, L. and SOMERS, M. A. (2015). Principal stratification: A tool for understanding variation in program effects across endogenous subgroups. *Am J. Eval*. **36** 514–531.

PEARSON, K. (1894). Contributions to the mathematical theory of evolution. *Philos*. *Trans*. *R*. *Soc*. *Lond*., *A* **185** 71–110.

PUMA, M., BELL, S. H., COOK, R., HEID, C. and SHAPIRO, G. (2010a). Head Start impact study. Final report, HHS, Administration for Children and Families.

PUMA, M., BELL, S. H., COOK, R., HEID, C. and SHAPIRO, G. (2010b). Head Start impact study. Technical report, HHS, Administration for Children and Families.

RAUDENBUSH, S. W. (2015). Estimation of means and covariance components in multi-site randomized trials. Unpublished manuscript.

RAUDENBUSH, S. W., REARDON, S. F. and NOMI, T. (2012). Statistical analysis for multisite trials using instrumental variables with random coefficients. *Journal of Research on Educational Effectiveness* **5** 303–332.

REARDON, S. F. and RAUDENBUSH, S. W. (2013). Under what assumptions do site-by-treatment instruments identify average causal effects? *Sociol*. *Methods Res*. **42** 143–163. MR3190727

ROMANO, E., BABCHISHIN, L., PAGANI, L. S. and KOHEN, D. (2010). School readiness and later achievement: Replication and extension using a nationwide Canadian survey. *Dev*. *Psychol*. **46** 995–1007.

ROSE, K. K. and ELICKER, J. (2010). Maternal child care preferences for infants, toddlers, and preschoolers: The disconnect between policy and preference in the USA. *Community Work Fam*. **13** 205–229.

RUBIN, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *J*. *Educ*. *Psychol*. **66** 688.

RUBIN, D. B. (1976). Inference and missing data. *Biometrika* **63** 581–592. MR0455196

RUBIN, D. B. (1980). Comment on "Randomization analysis of experimental data: The Fisher randomization test". *J*. *Amer*. *Statist*. *Assoc*. **75** 591–593.

RUBIN, D. B. (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician. *Ann*. *Statist*. **12** 1151–1172. MR0760681

SCHMIT, S., MATTHEWS, H., SMITH, S. and ROBBINS, T. (2013). Investing in young children: A fact sheet on early care and education participation, access, and quality. Fact sheet, New York, NY: National Center for Children in Poverty, Washington, DC: Center for Law and Social Policy.

SCHOCHET, P. Z. (2013). Student mobility, dosage, and principal stratification in school-based RCTs. *J*. *Educ*. *Behav*. *Stat*. **38** 323–354.

SCHOCHET, P. Z. and BURGHARDT, J. (2007). Using propensity scoring to estimate program-related subgroup impacts in experimental program evaluations. *Eval*. *Rev*. **31** 95–120.

SCHOCHET, P. Z., BURGHARDT, J. and MCCONNELL, S. (2008). Does job corps work? Impact findings from the national job corps study. *Am*. *Econ*. *Rev*. **98** 1864–1886.

SCHOCHET, P., PUMA, M. and DEKE, J. (2014). Understanding variation in treatment effects in education impact evaluations: An overview of quantitative methods (NCEE 2014–4017), Washington, DC: U.S. Department of Education, Institute of Education Sciences, National Center for Education Evaluation and Regional Assistance, Analytic Technical Assistance and Development.

SCOTT-CLAYTON, J. and MINAYA, V. (2014). Should student employment be subsidized? Conditional counterfactuals and the outcomes of work-study participation. Working Paper w20329, National Bureau of Economic Research.

SHAGER, H. M., SCHINDLER, H. S., MAGNUSON, K. A., DUNCAN, G. J., YOSHIKAWA, H. and HART, C. M. D. (2013). Can research design explain variation in Head Start research results? A meta-analysis of cognitive and achievement outcomes. *Educ*. *Eval*. *Policy Anal*. **35** 76–95.

SPLAWA-NEYMAN, J. (1990). On the application of probability theory to agricultural experiments. Essay on principles. Section 9. *Statist*. *Sci*. **5** 465–472. MR1092986

STAN DEVELOPMENT TEAM (2014). Stan: A C++ library for probability and sampling, Version 2.3.

WALTERS, C. R. (2015). Inputs in the production of early childhood human capital: Evidence from Head Start. *Am. Econ. J. Appl. Econ.* **7** 76–102.

WHITEHURST, G. J. (2013a). Obama's preschool plan. Brookings Institution.

WHITEHURST, G. J. (2013b). Can we be hard-headed about preschool? A look at Head Start. Brookings Institution.

ZHAI, F., BROOKS-GUNN, J. and WALDFOGEL, J. (2011). Head Start and urban children's school readiness: A birth cohort study in 18 cities. *Dev. Psychol.* **47** 134–152.

ZHAI, F., BROOKS-GUNN, J. and WALDFOGEL, J. (2014). Head Start's impact is contingent on alternative type of care in comparison group. *Dev. Psychol.* **50** 2572–2586.

ZHANG, J. L. and RUBIN, D. B. (2003). Estimation of causal effects via principal stratification when some outcomes are truncated by "death". *J. Educ. Behav. Stat.* **28** 353–368.

ZHANG, J. L., RUBIN, D. B. and MEALLI, F. (2009). Likelihood-based analysis of causal effects of job-training programs using principal stratification. *J. Amer. Statist. Assoc.* **104** 166–176. MR2663040

ZIGLER, E. and MUENCHOW, S. (1992). *Head Start*: *The Inside Story of America's Most Successful Educational Experiment*. Basic Books.

A. FELLER
GOLDMAN SCHOOL OF PUBLIC POLICY
UNIVERSITY OF CALIFORNIA, BERKELEY
2607 HEARST AVENUE
BERKELEY, CALIFORNIA 94702
USA
E-MAIL: afeller@berkeley.edu

L. MIRATRIX
HARVARD GRADUATE SCHOOL OF EDUCATION
14 APPIAN WAY
CAMBRIDGE, MASSACHUSETTS 02138
USA
E-MAIL: lmiratrix@g.harvard.edu

T. GRINDAL
ABT ASSOCIATES
55 WHEELER STREET
CAMBRIDGE, MASSACHUSETTS 02138
USA
E-MAIL: Todd_Grindal@abtassoc.com

L. C. PAGE
UNIVERSITY OF PITTSBURGH
SCHOOL OF EDUCATION
5918 WWPH
230 SOUTH BOUQUET STREET
PITTSBURGH, PENNSYLVANIA 15260
USA
E-MAIL: lpage@pitt.edu