

Understanding Program Complementarities: Estimating the Dynamic Effects of Head Start with Multiple Alternatives

Marc K. Chan* Antonio Dalla-Zuanna † Kai Liu‡

September 11, 2024

Abstract

We use experimental data from the Head Start Impact Study to examine the effects of childcare program pathways on cognitive outcomes. Using a sequential choice model with state dependence and dynamic selection, we identify and estimate joint returns and complementarities across skill investment programs. Our results indicate differential returns by the sequencing of program participation, and support engaging low-ability children with prolonged center-based care and high-ability children with some home care. We then use a dynamic structural model to evaluate a counterfactual reform which limits children to one year of Head Start, showcasing the importance of accounting for dynamic behavior.

*Faculty of Business and Economics, University of Melbourne. Email: marc.chan@unimelb.edu.au.

†Bank of Italy. Email: antonio.dalla.zuanna@gmail.com

‡Faculty of Economics, University of Cambridge and IZA. kai.liu@econ.cam.ac.uk

1 Introduction

Many social programs are evaluated by experiments which offer initial access to a focal program via lottery. In addition to the well-known issue of imperfect compliance, the existence of close alternatives complicates the experimental evaluation, especially if most control group individuals end up receiving similar services to the focal program (Heckman, Hohmann, Smith, and Khoo, 2000). Another challenge, which is important in evaluating programs that span multiple periods, is that people may choose different paths following the initial period of randomization. For example, those in the treatment group are not obliged to stay in the focal program and may switch to an alternative program; those in the control group who initially enrolled in an alternative program may switch to the focal program later.

This paper analyzes returns to the Head Start Program (HS), the largest and longest-running preschool education program for disadvantaged children aged 3–5 in the United States. Using experimental data from the Head Start Impact Study (HSIS), we confront the challenge arising when parents choose from multiple childcare options in every preschool year. The HSIS randomly assigned eligible HS applicants to treatment and control groups. The treatment group was allowed to enroll in HS; the control group was not to be admitted to HS during the first program year, but may select other types of care and re-apply for HS in the following year.

A close examination of childcare use for the age-3 HSIS cohort suggests that experimental estimates from comparing treatment- and control-group outcomes are difficult to interpret. First, almost 90% of control group children receive center-based care (HS or others) in at least one of two pre-school years. Second, treatment group children exhibit significant variations in the time spent in HS. Among children who enrolled in HS at age 3, less than three-quarters stay in HS at age 4, and about 20% switch to other center-based care at age 4. These observations suggest the need for a rich and well-defined set of causal returns beyond the experimental impacts.¹

We develop a potential outcome framework for understanding the causal returns to *program se-*

¹Puma et al. (2010) find positive but relatively small experimental impacts on children’s test scores during the pre-school period. Subsequent studies emphasize the heterogeneity in the impacts by children’s characteristics (e.g., Bitler, Hoynes, and Domina, 2014) and characteristics of the HS centers (e.g., Walters, 2015). These positive effects are generally found to be fading out over time, but several studies that use non-experimental methods (not relying on HSIS data) find positive effects of HS on later outcomes, such as education, health and income (e.g., Ludwig and Miller, 2007; Deming, 2009; Carneiro and Ginja, 2014; De Haan and Leuven, 2020). Griffen and Todd (2017) compare experimental estimates derived from the HSIS data to nonexperimental estimates obtained from comparable nonexperimental data.

quences with multiple program options in each period. We distinguish two types of returns: (1) the effect of a childcare program relative to home care at one age while fixing the rest of the program sequences to home care (“partial return”), and (2) the effect of a program at age 3 followed by another program (same or different) at age 4 relative to home care at both ages (“joint return”). The differential between a joint return and the sum of the relevant partial returns identifies the *cross return* between two programs, a concept closely related to dynamic complementarity in the child development literature (Cunha and Heckman, 2007). A positive cross return signals that the two programs are complementary, indicating that the return to one program is higher when combined with the other.

In a fully nonparametric model, the causal effects of program sequences are not identified with a single instrument from the HSIS experiment. We provide semiparametric identification by exploiting choice index restrictions and a factor structure. More specifically, building on the empirical approach in Heckman and Navarro (2007) and Heckman, Humphries, and Veramendi (2016), we specify a quasi-structural threshold model with approximate decision rules, combined with multiple exclusion restrictions and a system of measurements for the child’s unobserved skill endowment.² We assume that (i) the measurements (at pre-treatment) provide a valid proxy for the unobserved skill factor, and (ii) selection patterns at each period are driven by the same factor. Under these assumptions, we show that the causal effects of program sequences can be identified using the experimental variation, because the HSIS experiment modifies the selection mix of individuals into each and every program sequence.

Our estimates suggest that enrolling in HS at either age 3 or age 4 alone (while fixing the rest of the program sequences to home care) improves performance in reading and literacy tests at the end of age-4 year, by 0.30 and 0.34 standard deviations (sd) on average. After one year of HS at age 3, an additional year in HS (at age 4) only leads to a small improvement to mean test scores relative to home care, implying a negative cross return from two periods of HS participation. Other-center care produces positive returns that are comparable to HS (when home care is the benchmark), whereas the joint returns from mixed program experiences can vary a lot depending on their sequence. For instance, enrollment in HS at age 3 combined with other center care at age 4 leads to the largest gain in mean test scores. By contrast, reversing the sequence – other center care at age 3 followed by HS at age

²Recent applications of this estimation approach include, for example, Fruehwirth, Navarro, and Takahashi (2016) (the effect of grade retention), Heckman, Humphries, and Veramendi (2018) (the effect of school choices on earnings and health) and Rodríguez, Saltiel, and Urzúa (2018) (the effect of on-the-job training). In these papers, there are only two alternatives at each decision node.

4 – generates a very low average joint return, highlighting the importance of *sequencing* of program participation.

Exploring heterogeneity in program returns by children’s latent baseline ability (measured prior to the chosen sequence), we find that two years of center-based care (both HS and other-center care) improve the outcomes of low-ability children the most, suggesting a “compensatory” role of prolonged enrollment in center-based care. By comparison, high-ability children gain more from mixing center-based care at age 3 with home care at age 4 than low-ability children, which is consistent with the “skills-beget-skills” hypothesis. Overall, our estimates support both the “compensatory” hypothesis and the “skills-beget-skills” hypothesis that were previously examined by Bitler, Hoynes, and Domina (2014). We show that the specific mechanism depends on the sequence of childcare programs.

The estimates of the threshold model reveal state dependence in program choices, with age-3 enrollment in HS (other-center care) strongly predicting the use of the same program in age 4. State dependence implies that a one-off incentive for HS enrollment given at the start of age 3 (to the treatment group) continues to affect the demand for center-based care at age 4. Our estimated threshold model predicts that about 60% of the HSIS compliers change their entire path of program participation when assigned to the treatment group (“path compliers”). The average effects of HS experience differ widely between path compliers and other compliers who change their program choice temporarily in the first year. For instance, for children who would have otherwise been cared for at home at age 3, the average effects of HS is much larger among the path compliers than the temporary compliers.

The threshold model does not impose information sets and specify how people discount the future. These assumptions are not necessary to identify ex-post program returns, but they are essential when evaluating certain counterfactual reforms that involve a fundamental change in program features. One important yet unanswered policy question from the HSIS is the consequence of limiting HS to one year as opposed to two years (Puma et al., 2010). Such a change in program feature cannot be analyzed using the threshold model – the proposed policy reform changes the continuation value of choosing HS in age 3, implying that forward-looking individuals should adjust their program choice as soon as the information about the reform is revealed.

To this end, we build and estimate a dynamic structural model of program participation. The structural model’s key distinction is that it focuses on ex ante choice decisions via an explicit information

set and a model of expectations, which allows the utility flow and continuation value of each alternative to be disentangled. With additional exclusion assumptions, we identify preference parameters that measure the tradeoff between flow utilities and continuation values. Our estimates lend support for forward-looking behavior, showing that individuals incorporate the expected future outcomes as part of the information set and act upon it.

We use the estimated structural model to evaluate the hypothetical policy reform where the maximum duration of participation in HS is limited to one year, but individuals may choose whether to enrol at age 3 or 4. While the policy has limited impact on the average outcome of children, it exacerbates inequality. Individuals with low baseline ability perform substantially worse (relative to the status-quo of no limit in HS duration) while those with high baseline ability perform better. By removing the option of two consecutive years of HS participation, many individuals move away from HS at age 3 preemptively. This preemptive response is more common among low-ability individuals and the majority of them switch to home care at age 3, which leads to worse outcomes. In comparison, preemptive responses are less common among high-ability individuals; those who respond preemptively tend to switch to home care, which leads to better outcomes.

This paper is closely related to two papers analyzing the role of program substitution in program evaluation. Heckman, Hohmann, Smith, and Khoo (2000) show the importance of considering substitution and dropout biases in estimating the effects of job training. They point out that the experimental evaluation can only identify the effect of the program (the effect relative to the alternative training programs taken by individuals in the control group), but not the effect of training relative to no training. The dropout behavior in their paper refers to dropping out of the program prior to receiving training, whereas we consider the issue arising from individuals who drop out after having enrolled in a program. More recently, Kline and Walters (2016) use HSIS to estimate the returns and cost effectiveness of the HS program, taking into account that individuals may self-select into an alternative educational program that might be a close substitute to HS. They estimate a single-period selection model where both the choices and outcomes are measured within the first year of entry into HSIS (in Spring 2003). We not only distinguish multiple outside options like Kline and Walters (2016) but also study childcare choices and outcomes beyond the first year. This has non-trivial implications; for example, their framework implies that an expansion of HS can only crowd out enrollments in other com-

peting center care because the two programs are mutually exclusive within a single period (hence they are either substitutes or independent from each other), but across periods they can be complements. Our framework allows us to obtain estimates of cross returns between two different programs which would otherwise not be possible from a single-period model. In addition, we also consider the dynamic structural model described above, which can deliver important policy-relevant parameters beyond the scope of the threshold model.

By considering a multiple-period framework, this paper also relates to a strand of the program evaluation literature which incorporates dynamic selection into programs, including Ham and LaLonde (1996), Eberwein, Ham, and LaLonde (1997), Abbring and Van den Berg (2003), Heckman and Navarro (2007) and Heckman, Humphries, and Veramendi (2016).³ Studies in this literature only consider a binary program choice (one program versus no program). We contribute to this literature by combining a model of dynamic selection with multiple outside options in every decision period, which enables us to explain richer substitution patterns and understand how the returns to one program may be affected by participation in another program.

Our interest in the cross returns between two childcare programs coincides with a small but growing literature which tests for dynamic complementarities using experimental or quasi-experimental variations. To secure identification, papers in this literature use two independent sources of exogenous variation, one for each type of human capital investment (see, among others, Bhalotra and Venkataramani, 2015; Johnson and Jackson, 2019; Rossin-Slater and Wüst, 2020; Goff, Malamud, Pop-Eleches, and Urquiola, 2022).⁴ We complement this literature by presenting the partial, joint and cross returns to program participation that are easily comparable and interpretable. Most papers in this literature report the “reduced-form” impact without accounting for program take-up, and hence do not identify the cross returns to program participation.⁵ Among the few papers that attempt to use the two exoge-

³Ham and LaLonde (1996) and Eberwein, Ham, and LaLonde (1997) consider the effect of participating in training programs on the duration of subsequent employment and unemployment spells. They assume that the treatment effect is constant with respect to the time spent in the program, whereas we are also interested in understanding the effect of program duration. Abbring and Van den Berg (2003) allow for the effect to change with the duration of the spell in the program, but this is at the expense of ruling out the heterogeneous effect of the treatment with respect to individuals’ unobservable characteristics, which determine individuals’ choices in terms of duration in the program. In our framework, we explicitly allow for this type of heterogeneity.

⁴This identification strategy is also known as the “lightning strikes twice”: two identification strategies affecting the same cohort but at adjacent developmental stages (Almond and Mazumder, 2013).

⁵For instance, Rossin-Slater and Wüst (2020) find that the effect of *access* to preschool was smaller among those who had *access* to home visits during infancy. Goff, Malamud, Pop-Eleches, and Urquiola (2022) find weak evidence of negative effects of the interaction between exogenously driven improvements in family environment and being quasi-randomly *offered* a seat at a better school.

nous variations as instruments, the identified cross return is relevant only for a very specific subgroup of population which may not be policy relevant.⁶

The paper is organized as follows. Section 2 describes the HSIS and the sample and presents descriptive evidence on childcare choices and experimental impacts. Section 3 defines the joint, partial and cross returns from program sequences, laying the groundwork for our empirical analysis. In Section 4, we present the sequential threshold model and discuss identification and estimation, before reporting the estimates of program returns and their implications in Section 5. Section 6 discusses the identification and estimation of a dynamic structural model of program participation, as well as the estimated structural parameters and the counterfactual policy reform. Section 7 concludes.

2 Experimental Evaluation of the Head Start Program

2.1 The Head Start Program and the Head Start Impact Study

Head Start (HS) is a U.S. federal program that offers year-long care to children between three and five years of age, with the aim of fostering their early reading and maths skills to be ready for school. Launched in 1965, it enrolled almost one million children in 2003 (Puma et al., 2010). HS is administered by local agencies, which compete for funding and are required to adhere to national standards. It targets poor families with income below the federal poverty line, although local agencies are allowed to admit up to 10% of children from wealthier families.

The vast majority of HS providers offer center-based care. HS centers offer a large number of activities to support children’s literacy and maths skills. They also aim at supporting parents, with some centers providing transportation services for children and training/housing assistance for parents themselves; in addition, families receive at least two visits per year from HS staff. Children who attend HS are not charged any fee.

In the 1998 reauthorization, the Congress mandated an evaluation of Head Start’s impact on children’s school readiness and parental practices that support children’s development. The Head Start Impact Study (HSIS) followed this mandate: in fall 2002, around 5,000 newly entering children aged 3

⁶Johnson and Jackson (2019) examine the complementarity between Head Start spending and K-12 spending. Although not directly related to program participation, they use exposure to the rollout of Head Start across counties as a shock to Head Start spending and the implementation of court-ordered school finance reforms as a shock to K-12 spending, finding a positive effect of the interaction of increased spending in education at different ages.

and 4 were randomly assigned to either a treatment group or a control group. The randomization was implemented at the center level – at each HS center, eligible applicants were randomly assigned to the treatment and control groups.⁷ The treatment group was allowed to enroll in HS; the control group was not to be admitted to HS during the 2002-03 program year, but may select other types of care and re-apply next year. Therefore, the embargo period lasted one year only. For the age-3 cohort in the random assignment, this means that the control group could enroll in HS at age 4 in fall 2003 (and children in the treatment group could choose a different program at age 4). Despite efforts to maintain the integrity of the control group, some children assigned to the control group enrolled in HS in the 2002-03 program year (“crossovers”). These crossovers occur because local staff intentionally enrolled control group children into HS, and more commonly, parents applied to another nearby HS program as information on HSIS was not shared with programs not involved in the study (Puma et al., 2010). Overall, children who were aged 3 at the time of randomization could spend up to two years in HS before kindergarten at age 5, while those who were aged 4 could spend one year in HS.

2.2 Data

We use the first three waves of data from the HSIS. Given our interest in the returns to program sequences, we focus on children who joined HSIS at age 3 (i.e., the age-3 cohort). The baseline survey conducted at the time of randomization (Fall 2002) contains the experimental assignment, baseline characteristics of the children and their families, as well as characteristics of childcare centers. From the first follow-up survey (Spring 2003), we obtain information on the type of childcare use at age 3 ($t = 1$). From the second follow-up survey (Spring 2004), we use information on the type of childcare use at age 4 ($t = 2$) and children’s cognitive test scores.⁸

To measure children’s skills, we use the Peabody Picture Vocabulary Test (PPVT), the Woodcock-Johnson III (WJIII) Spelling test and WJIII Letter-Word Identification test. These tests measure children’s language and early literacy skills, and their levels and scales are comparable across years.⁹

⁷The sample of children in the HSIS consists of HS applicants to a nationally representative sample of 84 grantee/delegate agencies.

⁸To investigate the extent of within-year changes in program participation, we cross tabulate program participation using the Fall 2002 survey (at the start of the first HS year) and the Spring 2003 survey (at the end of the first HS year). We find a high degree of overlap in childcare use – over 90% of children report identical childcare use in the two surveys.

⁹PPVT and WJIII scores obtained from HSIS data have been transformed using item response theory to make them comparable across ages.

We use the Fall 2002 survey for measuring children’s baseline skills, and the Spring 2004 survey to construct their skill outcomes at the end of the age-4 year. The main outcome of interest is a composite score that averages PPVT, WJIII Spelling and WJIII Letter-Word Identification test scores, with scores from each test scaled down by 100. Section 5.1 examines the sensitivity of our results by including WJIII Applied Problems test in our analysis. The Applied Problems test is widely used to measure early numeracy skills.

As characteristics of childcare centers, we focus on center quality and whether the center provides transportation services. Center quality is a composite index derived (by the HSIS) based on responses to parent, center director, teacher and provider interviews. It combines information on center characteristics (teacher and center director education and qualifications, child/staff ratio) and practices (variety of literacy and math activities, home visiting, health and nutrition). For Head Start, the quality index and transportation availability refer to a child’s HSIS center of random assignment (i.e., the center at which the child was randomized into a experimental group). For other center care, we use the average characteristics of all other center care in the local area (proxied by the child’s assigned HS center identifier).¹⁰

Children who did not choose HS had other program options aside from home care, including private care or other federal or state-funded programs. The following key characteristics differentiate HS from other center-based care obtained by children in our sample (Puma et al., 2010): (a) HS offers more education-oriented activities and support for parents (see also Section 2.1), (b) health-related services are also included, (c) teachers are more likely to hold a bachelor’s degree, have taken college ECE courses and receive training on a monthly basis, (d) centers are filled to capacity, and (e) centers are generally of higher quality as certified by specific indicators (ECERS-R, Arnett, and quality composite scores).

We drop 573 children with missing information on childcare type and test scores. Our final sample consists of 1,876 children who were aged 3 in fall 2002, with 62.8% in the treatment group and 37.2% in the control group.¹¹ The treatment and control groups in our sample are very similar in demographic

¹⁰For each HSIS center of random assignment, there is an average of 2.5 non-HS other childcare centers. Characteristics of HS center are drawn from the Spring 2003 Child Experiences Data. To increase sample size, characteristics of other center-based care are drawn from the Spring 2003 and the Spring 2004 Child Experiences Data.

¹¹The entire HSIS data have 2,449 children in the age-3 cohort, 60% of whom are in the treatment group and 40% are in the control group.

characteristics, baseline test scores, and characteristics of the HS centers that they applied to (see Appendix Table A.1).

Table 1 shows the proportion of children in different types of care in our sample. The HS offer induced most children in the treatment group ($Z = 1$) to receive HS at $t = 1$ (85.1%), while a non-negligible proportion chose other center care (5.4%) or home care (9.5%). 14.4% of control-group children ($Z = 0$) crossed-over to receive HS at $t = 1$. These crossovers represent always-takers in the experiment. This implies an experimental compliance rate (the population share of compliers due to the random assignment) of 70.6%(=85.1%-14.4%).

Table 1 also highlights rich transition patterns in the type of care between $t = 1$ and $t = 2$. In both the control and treatment groups, among those who received HS at $t = 1$, slightly less than three quarters stay in HS at $t = 2$, about 20% transition to other center care, and about 6% transition to home care. By contrast, among those who were in other center care or home care at $t = 1$, the transition patterns are markedly different between the control and treatment groups. In the control group, among those who received other center care at $t = 1$, 37.6% transitioned to HS at $t=2$; this proportion is lower at 21.9% in the treatment group. The same pattern is observed in home care to HS transitions (49.8% in the control group versus 31.2% in the treatment group). These results support the notion of a rebound effect on HS receipt following the end of the embargo period. Yet even though the embargo period is over, control group children transition into HS at a considerably lower rate than the initial HS take-up rate of the treatment group (85.1%). Overall, the proportion receiving HS at $t = 2$ is 50.2% (10.7 + 9.3 + 30.2) in the control group, versus 66.2% (62.0 + 1.2 + 3.0) in the treatment group.

While the table shows that the majority of control group children chose home care at $t = 1$ (60.7%), this does not imply that they do not receive any formal care before starting kindergarten. The vast majority (49.8 + 30.7 = 80.5%) of these children made a transition to HS or other center care at $t = 2$. Therefore, most children in the control group enrolled in formal center-based care prior to the kindergarten age, making the experimental impact of HS difficult to interpret.

2.3 Experimental Impact of Head Start

Table 2 reports the experimental impact on language and early literacy skills at age 4. Column (1) shows that children in the treatment group gained 2.9 points relative to the control group, or 0.10 of a standard deviation (sd) of the test scores in the control group.¹² Given the proportion of experimental compliers (70.6%; see above), this implies that HS participation at age 3 increases the age-4 test score among compliers by 0.041 (p-val=0.008).

In column (2) we explore heterogeneity in the experimental impact by regressing test score on the treatment group indicator Z , and interactions of Z with indicators of period-2 choice of other center care (c) or HS (h). Because childcare choices at $t = 2$ are endogenous, these estimates do not have any causal interpretation. Nevertheless, column (2) shows large heterogeneity in the experimental impact across types of care selected at $t = 2$. The estimated coefficient on the interaction of Z and c indicates that the experimental impact is significantly larger among individuals who choose other center care at $t = 2$ than those who choose home care. In fact, this subgroup drives the positive experimental effect because the outcomes of the treatment group individuals who choose either home care or HS at $t = 2$ are not significantly different from the outcomes of the control group. These results are suggestive that subsequent program experience in $t = 2$ may be relevant in explaining the experimental impact.

We also investigate the heterogeneity in the proportion of compliers and in the experimental effect by children’s baseline skills. Consistent with Bitler, Hoynes, and Domina (2014), Appendix Figure A.1 shows that the experimental impact declines significantly with children’s baseline skills, whereas the compliance rate is only weakly increasing with children’s baseline skills.¹³ The underlying mechanisms are less clear. Children with different baseline skills may benefit differently from the same childcare services or may choose different childcare services over time.

¹²Similar to Kline and Walters (2016), we control for household size, number of siblings, dummies for whether the child is female, black, hispanic, uses English as home language, is living in urban area, is living with both parents, is in need of special education, is the child of a teen mother, is the child of a mother who never married or is separated, is the child of a mother with high school or more than high school. Our point estimates are similar if these baseline characteristics are excluded from the regression. Our estimates are close to Kline and Walters (2016), who estimate a significant increase in test score for the treated group at the end of the age-4 year by 0.09 of a s.d. in the control group. Their test score measure is derived from the average of PPVT scores and WJIII Pre-Academic Skills scores (a summary index of Letter-Word Identification, Spelling, and Applied Problems subtests).

¹³In Figure A.2 we show that the LATE is also declining with baseline skills.

3 Return to Program Sequence

In this section we define the joint, partial and cross returns from program sequences, laying the groundwork for our empirical analysis. Consider potential outcome $Y_i^{j,j'}$ for an individual i who is externally assigned to program j at $t = 1$ and program j' at $t = 2$. There are three program alternatives in every period: HS (h), other center care (c), and home care (n). We define the *joint return* from a sequence of programs as the difference between the potential outcome upon completing sequence (j, j') , where $j, j' \in \{h, c\}$, and the potential outcome of home care in both periods (i.e., (n, n)). The individual's joint return from program sequence (j, j') can be decomposed into

$$Y_i^{j,j'} - Y_i^{n,n} = \underbrace{(Y_i^{j,n} - Y_i^{n,n})}_{\text{Partial Return of } j \text{ at } t=1} + \underbrace{(Y_i^{n,j'} - Y_i^{n,n})}_{\text{Partial Return of } j' \text{ at } t=2} + \underbrace{[(Y_i^{j,j'} - Y_i^{n,j'}) - (Y_i^{j,n} - Y_i^{n,n})]}_{\text{Cross Return between } j \text{ and } j'}. \quad (1)$$

The first two terms are *partial returns*, which capture the effect of completing a program at a specific age as compared to home care at that age, while fixing the rest of the program sequences to home care. Note that the two partial returns may be different even when $j = j'$; for instance, attending HS at age 3 may have a different effect than attending HS at age 4.

Of particular interest are the cross returns, which indicate the extent of complementarity between HS and alternative childcare use. Cross return arises because the joint return from program j at $t = 1$ and program j' at $t = 2$ may be different from the sum of their partial returns. A positive cross return would in fact signal that the two programs are *complements*, implying that the interaction of the two programs gives an even higher return than the combined partial returns from each single program. A negative cross return, on the other hand, implies that the two programs are *substitutes*, with the impact of one being offset by having enrolled in the other. This may signal that some of the competences learnt in one program are replicated by the other, thus resulting in no additional knowledge.¹⁴

¹⁴Our definition of program cross returns is closely related to the concept of dynamic complementarity in the context of the child development literature. Formally, one can define dynamic complementarity between investments at different ages for individual i as

$$\frac{\partial^2 Y_{i,t}}{\partial x_{i,t} \partial x_{i,t-1}} = \frac{\partial^2 Y_{i,t}}{\partial x_{i,t} \partial Y_{i,t-1}} \frac{\partial Y_{i,t-1}}{\partial x_{i,t-1}},$$

where $x_{i,t}$ are the investment made at time t (Cunha and Heckman, 2007; Aizer and Cunha, 2012). In our setting we explore complementarity (or substitutability) not only between investments made at different points in time but also between different *types* of investments, i.e. between different programs *relative to* a benchmark program (home care). Importantly, the cross returns in our context may vary across persons because they are defined from potential outcomes allowing for flexible individual-specific program returns.

Table 3 lists the cross-returns in the HSIS context, consisting of two within-program cross returns ($j = j'$) and two between-program cross returns ($j \neq j'$). Note that the between-program cross returns can differ from each other as they are defined from different program sequences. For instance, return to h at $t = 1$ can be boosted by enrolling in c at $t = 2$, yet return to h at $t = 2$ may be reduced by enrolling in c at $t = 1$.¹⁵ More generally, equation (1) does not restrict the heterogeneity in returns: for a given individual, the return may vary depending on the sequences being compared; and for a given pair of sequences, the return may vary across individuals.

The main empirical obstacle to estimate the average returns defined by equation (1) for a population of interest is endogenous self-selection into program sequences. We only observe each child’s potential outcome given the chosen sequence; the rest of potential outcomes are unobserved. There are at least three approaches that correct for selection: (i) instrumental variables (IV), (ii) a quasi-structural approach where we specify reduced form or approximate decision rules for the program choices and estimate those equations jointly with the potential outcomes, or (iii) we could estimate a fully-specified dynamic structural model. A conventional approach would be to use an IV strategy. However, there are generally insufficient number of instruments relative to the large number of (endogenous) program paths, making the IV approach impractical. In a fully nonparametric model, the causal returns are not identified with a binary instrument (such as the HSIS experiment).

In Sections 4 and 5, we identify the causal returns using a quasi-structural threshold model, combined with multiple exclusion restrictions and a system of measurements for the child’s unobserved skill endowment. Although the quasi-structural model makes more assumptions than IV, returns to program sequences can be identified with a limited number of instruments and the identified returns are relevant to the broader population rather than a specific subpopulation of compliers defined by the set of instruments. Unlike a dynamic structural model, the quasi-structural model does not impose information sets and specify how people discount the future, which are attractive to researchers who may worry about the sensitivity of estimated program returns to these assumptions. More broadly, the approximate decision rules can nest behaviors from structural models involving different assumptions

¹⁵The cross-return parameters can be interpreted in two equivalent ways: the effect of enrolling in program j at $t = 1$ on the return to program j' at $t = 2$ or the effect of enrolling in program j' at $t = 2$ on the return to program j at $t = 1$. This can be easily seen by switching the terms on the right hand side of equation (1): $(Y_i^{j,j'} - Y_i^{n,j'}) - (Y_i^{j,n} - Y_i^{n,n}) = (Y_i^{j,j'} - Y_i^{j,n}) - (Y_i^{n,j'} - Y_i^{n,n})$.

of the decision process.¹⁶ Therefore, if the goal is to identify ex-post returns to program sequences, then the quasi-structural model, which tends to be more flexible than the dynamic structural model, is sufficient.¹⁷

The dynamic structural model can evaluate counterfactual policies that are not possible for the estimated threshold model. In Section 6, we estimate a dynamic structural model and use it to evaluate a hypothetical policy reform where the maximum participation duration allowed in HS is reduced from two years to one year. Importantly, the dynamic structural model allows individuals to re-optimize into different program paths as soon as information of the reform is revealed to the individuals. These types of dynamic changes in behavior cannot be captured by a quasi-structural model because we need to disentangle the continuation value from flow utility for each program alternative. In our view, this a key strength of the dynamic structural model in program evaluation.

4 A Sequential Model of Program Participation

Choices. Individual i starts in an initial state, Z_i , that has been randomly assigned by an experiment. Let $Z_i \in \{0, 1\}$ indicate the individual's experimental status (= 1 if treatment group and 0 if control group). At the start of period $t = 1$ (age 3), the perceived value of each program alternative is given by:

$$U_{i1}^h = \psi_1^h + \mathbf{X}'_{ih}\boldsymbol{\beta}_1^h + f(Z_i, \mathbf{X}_{ih}, \theta_i) + \lambda_1^h\theta_i + \nu_{i1}^h \quad (2)$$

$$U_{i1}^c = \psi_1^c + \mathbf{X}'_{ic}\boldsymbol{\beta}_1^c + \lambda_1^c\theta_i + \nu_{i1}^c \quad (3)$$

$$U_{i1}^n = 0 \quad (4)$$

where U_{i1}^h is the value of HS, U_{i1}^c is the value of other center care, and the value of home care (U_{i1}^n) is normalized to zero. The value of each alternative is determined by characteristics of the respective program ($\mathbf{X}_{ih}, \mathbf{X}_{ic}$), the experimental status (Z_i), the unobserved factor (θ_i) which will be measured by a set of baseline test scores, as well as idiosyncratic shocks (ν_{i1}^h and ν_{i1}^c). The alternative-specific

¹⁶For example, the threshold decision rule can be consistent with a fully static multinomial choice model where individuals commit to program sequence chosen at the start of the initial period, or a dynamic structural model that involves a forward-looking agent under the restriction that the choice shocks are serially uncorrelated.

¹⁷For a recent application of the quasi-structural model on child skill production function, see Chan and Liu (forthcoming).

idiosyncratic shocks reflect unobserved shocks to demand for HS and other center care relative to home care. They are assumed independent from $\theta_i, Z_i, \mathbf{X}_{ih}, \mathbf{X}_{ic}$ and over time; within the same period, the unobserved demand shocks for HS and other center care may be correlated.

The perceived value of HS varies by $\mathbf{X}_{ih} = (q_{ih}, tp_{ih})$, including the quality index (q_{ih}) and transportation availability (tp_{ih}) of the child’s center of random assignment.¹⁸ Similarly, the perceived value of other-center care varies by $\mathbf{X}_{ic} = (q_{ic}, tp_{ic})$, the mean quality (q_{ic}) and transportation availability (tp_{ic}) of non-HS childcare centers in the local area. One concern is that these center characteristics may reflect neighborhood characteristics and hence ultimately correlate with unobserved children characteristics due to, for instance, endogenous sorting. To alleviate this concern, center characteristics are first residualized against a large set of socio-demographic characteristics (which reflect neighbourhood characteristics), and hence reflect residual variations in center quality net of these characteristics.¹⁹ Hereafter, “center characteristics” (\mathbf{X}_{ih} and \mathbf{X}_{ic}) refer to residualized center characteristics. After standardizing, they have means of zero and standard deviations of 1. We assume that \mathbf{X}_{ih} and \mathbf{X}_{ic} are exogenous to the underlying latent factor (θ_i). We observe a set of baseline test scores (before the HSIS experiment) which is a good proxy for the latent factor. Appendix Figure A.3 shows that there is little correlation between (residualized) center characteristics and children’s baseline test scores, providing empirical support to this assumption.

We approximate $f(Z_i, \mathbf{X}_{ih}, \theta_i)$ with a flexible function:

$$f(Z_i, \mathbf{X}_{ih}, \theta_i) = (\beta_Z + \mathbf{X}'_{ih}\beta_{Zx} + \beta_{Z\theta}\theta_i) Z_i \quad (5)$$

Equation (5) captures the effects of the experimental status on HS enrollment. Being in the HSIS treatment group improves ease of access to HS, which is captured by coefficient β_Z . Empirically we find that the take-up rates of HS in the treatment group relative to the control group varies by center

¹⁸We observe \mathbf{X}_{ih} for all children because the randomization of HSIS is implemented at the HS center level (see Section 2.2). At each HS center, applicants are randomly assigned to either the treatment group or the control group.

¹⁹More specifically, we obtain residuals from a regression of each center characteristic on household size, number of siblings, dummies for whether the child is female, black, hispanic, uses English as home language, is living in urban area, is living with both parents, is in need of special education, is the child of a teen mother, is the child of a mother who never married or is separated, is the child of a mother with high school or more than high school. An alternative approach is to include the relevant family characteristics as additional covariates. By comparison, our approach can reduce the dimensionality of the estimation problem, which is particularly useful when it comes to estimating the dynamic structural model in Section 6.

and individual characteristics.²⁰ To capture such heterogeneity, we interact Z_i with the vector of HS center characteristics (\mathbf{X}_{ih}) and the unobserved factor (θ_i). Then, the HS enrollment variation by the unobserved factor in the control group is captured by λ_1^h in equation (2). Similarly, β_1^h in the same equation captures heterogeneity in HS enrollment probability by center characteristics for the control group. The covariates in equation (5) serve as exclusion restrictions in the model (to discuss below).²¹

Denote the individual's program choice at t by $K_{it} = \underset{k \in \{h,c,n\}}{\operatorname{argmax}} U_{it}^k$. At the start of period $t = 2$ (age 4), individuals choose from the same menu of programs, while taking into account the program choice they made at $t = 1$:

$$U_{i2}^h = \psi_2^h + \psi_2^{hh} \mathbf{1}(K_{i1} = h) + \psi_2^{ch} \mathbf{1}(K_{i1} = c) + \mathbf{X}'_{ih} \beta_2^h + \lambda_2^h \theta_i + \nu_{i2}^h \quad (6)$$

$$U_{i2}^c = \psi_2^c + \psi_2^{hc} \mathbf{1}(K_{i1} = h) + \psi_2^{cc} \mathbf{1}(K_{i1} = c) + \mathbf{X}'_{ic} \beta_2^c + \lambda_2^c \theta_i + \nu_{i2}^c \quad (7)$$

$$U_{i2}^n = 0 \quad (8)$$

The value of HS in $t = 2$ no longer depends on Z_i directly because the experimental offer was assigned at the start of $t = 1$. We emphasize that Z_i can still affect the choices at $t = 2$, via changing the decision at $t = 1$. In fact, this is a key mechanism that our model examines – the perceived value of each childcare alternative in $t = 2$ can differ by program experience in $t = 1$. For example, if $\psi_2^{hh}, \psi_2^{ch} > 0$, then an individual values h more in $t = 2$ if she chose h or c last period as opposed to n . Note that all parameters in the choice equations, with one exception, are also time-specific. For instance, ψ_2^h, ψ_2^c may differ from ψ_1^h, ψ_1^c , capturing that there were more formal childcare facilities available for age-4 children than for age-3 children (Puma et al., 2010; Bitler, Hoynes, and Domina, 2014). The only exception is the alternative-specific idiosyncratic shocks, where ν_{i2}^h and ν_{i2}^c are drawn from the same distribution as the choice shocks in $t = 1$.

In Appendix Section B.1.1, we estimate a more flexible specification where the coefficients on the unobserved factor in $t = 2$ may differ by the program choice made in $t = 1$. This relaxes the constant effect assumption in our current model where program choice in $t = 1$ only makes an intercept shift to

²⁰For instance, we find that lower-quality HS centers tend to admit control group individuals.

²¹We assume that Z_i and \mathbf{X}_{ih} have no effect on the value of other center care (relative to home care) in period 1. Note that in a structural model with forward-looking agents, \mathbf{X}_{ih} may indirectly affect the values of care in period 1 via the expected value function (e.g., high-quality HS center raises the expected value of all modes of care now due to a possibility of choosing HS in the future). However, the symmetry of our decision tree and additional robustness checks suggest that any shifts in the expected value function have roughly the same magnitude in the “c” and “n” branches at $t = 1$. This motivates the exclusion restriction in the sequential choice model.

the perceived value of program alternatives in $t = 2$. Effectively, individuals with different θ are now subject to different degrees of “state dependence” in preferences and hence perturbing initial childcare choice can have different consequences for subsequent choice among different types of individuals. However, the data does not lend support to this richer specification – we cannot reject the hypothesis that the marginal values of the unobserved factor are homogeneous across lagged program choices. In light of this evidence and other sensitivity analysis reported in Appendix Section B, we focus on the current specification in the rest of the paper.²²

Potential outcomes. Denote individual i ’s sequence of choices from period 1 up to period t by $D_{it} \equiv (K_{i1}, \dots, K_{it})$. Let Y_{it}^d denote the potential outcome of child i at the end of period t if she is externally assigned to program sequence $D_{it} = d$. We only observe the realized outcome Y_{it} (a composite test score; see Section 2.2), which is equal to the potential outcome of the actual program sequence choice. We parameterize the potential outcomes as a function of program characteristics, the unobserved factor and idiosyncratic shocks:

$$Y_{it}^d = \alpha_t^d + \beta_{th}^d q_{ih} + \beta_{tc}^d q_{ic} + \gamma_t^d \theta_i + \varepsilon_{it}^d, \quad (9)$$

where q_{ih} and q_{ic} are HS and other center quality, respectively, and ε_{it}^d is the idiosyncratic shock that is i.i.d across d and independent from θ and all other components in the model. The factor loadings (γ_t^d) describe how potential outcomes vary with the unobservable factor, and they are informative about the nature of selection and program returns.²³ Note that the experimental status Z_i and its interaction with center characteristics are excluded from equation (9). This is consistent with the standard exclusion restriction assumption that the experimental offer is unrelated to potential outcomes once the agent is externally assigned to $D = d$ (e.g., Angrist, Imbens, and Rubin (1996)).²⁴

²²Although the setup is not most general, the sequential structure of the model naturally corresponds to the features of HSIS, and the perceived values of different choices are easily comparable with the values estimated in the dynamic structural model in Section 6 (which is also sequential as shocks are revealed over time). The model specification is also simpler, e.g., if we “flatten” it to a multinomial model instead we would allow $3 \times 3 - 1 = 8$ alternative-specific choice shocks to be mutually correlated without a priori restrictions.

²³For instance, a test of $\gamma_t^d = 0, \forall d$ indicates the presence of selection and testing the null hypothesis $\gamma_t^j = \gamma_t^k, \forall j \neq k$ indicates the presence of constant returns. See Section 5.1 for details. Also see Heckman, Urzua, and Vytlačil (2006); Fruehwirth, Navarro, and Takahashi (2016).

²⁴Given that the variables on transportation availability (tp_{ih}, tp_{ic}) do not enter into the potential outcomes, they are exclusion restrictions in the model. In Appendix Section B.1.3, we show that our estimates change very little if we do not use them as exclusion restrictions (also see Section 5).

All parameters in equation (9) vary by program sequence. They measure the technological returns to participation in childcare programs over ages 3 and 4, where the returns can vary by individual’s baseline skills (θ_i) and childcare center quality. HS center quality affects potential outcomes only if individual i enrolls in at least 1 period of HS, meaning that $\beta_{th}^{nn} = \beta_{th}^{nc} = \beta_{th}^{cn} = \beta_{th}^{cc} = 0$. We further restrict the HS center quality to have homogeneous effects across program sequences involving 1 period of HS, by constraining $\beta_{th}^{hn} = \beta_{th}^{nh} = \beta_{th}^{ch} = \beta_{th}^{hc}$. This restriction has minimal effects on the point estimates and cannot be rejected statistically. We make similar assumptions on the coefficients related to other-center quality, where $\beta_{tc}^{nn} = \beta_{tc}^{nh} = \beta_{tc}^{hn} = \beta_{tc}^{hh} = 0$ and $\beta_{tc}^{cn} = \beta_{tc}^{nc} = \beta_{tc}^{ch} = \beta_{tc}^{hc}$.

Measurement equations. Following the literature on dynamic treatment effects (e.g., Heckman, Humphries, and Veramendi (2016)), we supplement the outcomes with a measurement system to proxy the unobserved factor and correct for the effects of measurement error in the proxy:

$$M_i^M = \alpha^M + \gamma^M \theta_i + e_i^M, M = \{pp, ws, wl\} \tag{10}$$

where M_i^M is a baseline test score that is observed around the time of randomization (in Fall 2002). We utilize three types of baseline test scores which measure language and literacy skills, including PPVT ($M = pp$), WJIII Spelling test ($M = ws$) and WJIII Letter-Word Identification test ($M = wl$).²⁵ Given that these measurements are observed around the time of randomization, the parameters affecting M_i^M are invariant to program sequence. The measurement error terms e_i^M are independent of each other, the unobserved factor and all other components in the model.²⁶

4.1 Identification

Our identification argument consists of four blocks: (1) measurement equations, (2) choice equations at $t = 1$, (3) choice equations at $t = 2$, and (4) outcome equations. We focus on key intuitions below and more details are in Appendix C.

²⁵As a sensitivity check, we also include WJIII Applied Problems test score (measuring early numeracy) from the baseline survey as the fourth measurement. See Section 5.1 for details.

²⁶In their preferred model, Kline and Walters (2016) instead control for baseline test scores as covariates in both the selection equations and potential outcomes equations. Therefore, differently from our model, the unobservables driving selection in Kline and Walters (2016) are interpreted as unobserved tastes (conditional on baseline test scores of children).

Measurement equations. Location and scale of the factor are not identified, so normalizations are required. We normalize the mean of the factor (θ_i) and measurement errors (e^{pp} , e^{ws} , e^{wl}) to zero and the factor loading in the baseline PPVT equation to one ($\gamma^{pp}=1$). By linking the factor to the baseline PPVT score, it can be interpreted as the latent ability of the individual. The intercepts α^{pp} , α^{ws} , α^{wl} are identified from the first moments of measurements. With θ , e^{pp} , e^{ws} , e^{wl} mutually independent, covariances of the measurements identify the factor loadings γ^{ws} , γ^{wl} and factor variance σ_θ^2 :

$$Cov(M_i^{pp}, M_i^{ws}) = \gamma^{ws} \sigma_\theta^2 \quad (11)$$

$$Cov(M_i^{pp}, M_i^{wl}) = \gamma^{wl} \sigma_\theta^2 \quad (12)$$

$$Cov(M_i^{ws}, M_i^{wl}) = \gamma^{ws} \gamma^{wl} \sigma_\theta^2 \quad (13)$$

Then, utilizing the joint distribution of $(M_i^{pp} - \alpha^{pp}, \frac{M_i^{ws} - \alpha^{ws}}{\gamma^{ws}})$, the distribution of θ_i is identified via Kotlarski (1967).²⁷ The distributions of e_i^{pp} , e_i^{ws} , e_i^{wl} are then identified by deconvolution of θ_i from $M_i^{pp} - \alpha^{pp}$, $\frac{M_i^{ws} - \alpha^{ws}}{\gamma^{ws}}$ and $\frac{M_i^{wl} - \alpha^{wl}}{\gamma^{wl}}$, respectively.²⁸

Choice equations at $t = 1$. Identification follows the literature on semiparametric discrete choice models with index restrictions. For example, the probability of choosing n , which is also an expectation, is a function of two indices consisting of the exogenous covariates:

$$p(K_{i1} = n | \tilde{\mathbf{X}}_{ih}, \mathbf{X}_{ic}) = G_{(\lambda_1^h \theta_i + \nu_{i1}^h, \lambda_1^c \theta_i + \nu_{i1}^c)}(-\psi_1^h - \tilde{\mathbf{X}}_{ih}' \tilde{\boldsymbol{\beta}}_1^h, -\psi_1^c - \mathbf{X}_{ic}' \boldsymbol{\beta}_1^c) \quad (14)$$

where $G(\cdot)$ is the joint distribution function of composite choice errors, $\tilde{\mathbf{X}}_{ih} \equiv (\mathbf{X}_{ih}, Z_i, Z_i \mathbf{X}_{ih})$ is a vector of all covariates in the h-equation.²⁹ To identify $\tilde{\boldsymbol{\beta}}_1^h$, $\boldsymbol{\beta}_1^c$ up to scale, key conditions are: (1) each index contains at least one distinct continuous variable with nonzero coefficient, not contained in the other index (e.g., q_{ih}, q_{ic}), (2) the conditional density of the distinct continuous variable on other covariates (discrete/continuous) is absolutely continuous and positive (possibly bounded support), (3) full rank in covariates, and (4) $G(\cdot)$ is continuously differentiable. With a scale normalization (e.g.,

²⁷See also Theorem 1 of Cunha, Heckman, and Schennach (2010). Evdokimov and White (2012) show that Kotlarski's result holds even when the common latent variable θ_i has a vanishing characteristic function (e.g., compact support).

²⁸Appendix Section C shows the identification results when only two types of baseline test scores are available in the data. See also Fruehwirth, Navarro, and Takahashi (2016).

²⁹Assume away $Z_i \theta$ for simplicity, but this can be accommodated by repeating the analysis by $Z_i = 0, 1$.

coefficient on the distinct continuous variable set to one), identification of the index coefficients $\tilde{\beta}_1^h$, β_1^c follows Lee (1995) and Klein and Spady (1993) without specifying the error distribution, based on quasi-likelihood and nonparametric kernel regression for the conditional expectation in equation (14) (similar analysis for c-choice). With support of the indices covering that of the composite errors, $G(\cdot)$ can be identified by tracing out the indices along this conditional expectation. The by-product $G(\cdot|M_i^{PP})$ is identified (since we observe M_i^{PP}), hence the joint distribution of composite errors and M_i^{PP} is also identified. The covariances therein identify factor loadings λ_1^h , λ_1^c (σ_θ^2 is known from equations (11)-(13)):

$$Cov(M_i^{PP}, \lambda_1^h \theta_i + \nu_{i1}^h) = \lambda_1^h \sigma_\theta^2 \quad (15)$$

$$Cov(M_i^{PP}, \lambda_1^c \theta_i + \nu_{i1}^c) = \lambda_1^c \sigma_\theta^2 \quad (16)$$

The intercepts ψ_1^h , ψ_1^c cannot be identified without a location normalization. A natural restriction is to set the means of ν_{i1}^h and ν_{i1}^c to zero. Lewbel (1997) shows identification of intercepts with unknown $G(\cdot)$ and large index support.

Finally, the joint distribution of (ν_{i1}^h, ν_{i1}^c) is identified via deconvolution. The joint characteristic function of the composite errors is:

$$\varphi(\lambda_1^h \theta_i + \nu_{i1}^h, \lambda_1^c \theta_i + \nu_{i1}^c)(s, w) = E(e^{i[s(\lambda_1^h \theta_i + \nu_{i1}^h) + w(\lambda_1^c \theta_i + \nu_{i1}^c)]}) \quad (17)$$

$$= E(e^{i[(s\lambda_1^h + w\lambda_1^c)\theta_i + s\nu_{i1}^h + w\nu_{i1}^c]}) \quad (18)$$

$$= E(e^{i(s\lambda_1^h + w\lambda_1^c)\theta_i}) E(e^{i(s\nu_{i1}^h + w\nu_{i1}^c)}) \quad \text{by } \theta \perp \nu_{i1}^h, \nu_{i1}^c \quad (19)$$

$$= \varphi_\theta(s\lambda_1^h + w\lambda_1^c) \varphi_{(\nu_{i1}^h, \nu_{i1}^c)}(s, w) \quad (20)$$

where $i = \sqrt{-1}$, and the characteristic function $\varphi_\theta(\cdot)$ and loadings λ_1^h , λ_1^c are known. The joint density of (ν_{i1}^h, ν_{i1}^c) is obtained by inverse Fourier transform: $(2\pi)^{-2} \iint \frac{\varphi(\lambda_1^h \theta_i + \nu_{i1}^h, \lambda_1^c \theta_i + \nu_{i1}^c)(s, w)}{\varphi_\theta(s\lambda_1^h + w\lambda_1^c)} e^{-i(s\nu_{i1}^h + w\nu_{i1}^c)} ds dw$.³⁰

³⁰In cases where $\varphi_\theta(\cdot)$ vanishes (e.g., compact support), ridge deconvolution is available (Hall and Meister (2007)).

Choice equations at t=2. The conditional density of factor θ_i given period-1 choice K_{i1} , baseline measurement M_i^{pp} and covariates $\tilde{\mathbf{X}}_i(\equiv (\tilde{\mathbf{X}}_{ih}, \mathbf{X}_{ic}))$, which captures selection over time, is identified:

$$\begin{aligned} f_{\theta_i|K_{i1}, M_i^{pp}, \tilde{\mathbf{X}}_i}(\theta|k, m, x) &= \frac{p(K_{i1} = k|\theta_i = \theta, M_i^{pp} = m, \tilde{\mathbf{X}}_i = x) f_{M_i^{pp}|\theta_i}(m|\theta) f_{\theta_i}(\theta)}{p(K_{i1} = k|M_i^{pp} = m, \tilde{\mathbf{X}}_i = x) f_{M_i^{pp}}(m)} \\ &= \frac{p(K_{i1} = k|\theta_i = \theta, \tilde{\mathbf{X}}_i = x) f_{e_i^{pp}}(m - \theta) f_{\theta_i}(\theta)}{p(K_{i1} = k|M_i^{pp} = m, \tilde{\mathbf{X}}_i = x) f_{M_i^{pp}}(m)} \end{aligned} \quad (21)$$

where the conditional choice probability $p(K_{i1} = k|\theta_i = \theta, \tilde{\mathbf{X}}_i = x)$ is known from the identified joint distribution of (ν_{i1}^h, ν_{i1}^c) and parameters in the period-1 choice equations, $f_{e_i^{pp}}(\cdot)$ and $f_{\theta_i}(\cdot)$ are known, and the denominator terms are directly identified from data. Importantly, the choice equations (6)-(8) then represent a mixture model where both the mixture density $f_{\theta_i|K_{i1}, M_i^{pp}, \tilde{\mathbf{X}}_i}(\cdot)$ and the joint error density (ν_{i2}^h, ν_{i2}^c) are known (note that (ν_i^h, ν_i^c) is i.i.d. over time). Hence using period-2 choice probabilities conditional on $K_{i1}, M_i^{pp}, \tilde{\mathbf{X}}_i$, one can specify a likelihood to estimate all the period-2 parameters without scale normalization. While it is standard to assume unique optimum in the population criterion function, the intuition for identification stems from the following variations: (1) Given $K_{i1} = n$ and $\tilde{\mathbf{X}}_i = 0$, variation in M_i^{pp} shifts choices via the mixture, targeting loadings λ_2^h, λ_2^c and intercepts ψ_2^h, ψ_2^c ; (2) Given $\tilde{\mathbf{X}}_i = 0$ and $M_i^{pp} = m$, variation in K_{i1} shifts choices via the mixture, and since loadings λ_2^h, λ_2^c are known, this targets state dependence parameters $\psi_2^{hh}, \psi_2^{hc}, \psi_2^{ch}, \psi_2^{cc}$; (3) Given K_{i1} and M_i^{pp} , variation in $\tilde{\mathbf{X}}_i$ shifts choices via the indices and mixture, targeting β_2^h, β_2^c .

To further illustrate identification, the score-weighted moment equals:

$$E \left[\mathbf{1}\{K_2 = n\} \frac{-\partial \ln f(\theta, M^{pp}|K_1=n, \tilde{\mathbf{X}}=0)}{\partial M^{pp}} \right] = \lambda_2^h E \left[\frac{\partial G_{\nu^h, \nu^c}(s_1, s_2)}{\partial s_1} \right] + \lambda_2^c E \left[\frac{\partial G_{\nu^h, \nu^c}(s_1, s_2)}{\partial s_2} \right] \quad (22)$$

with expectation taken over θ and M^{pp} conditional on $K_1 = n$ and $\tilde{\mathbf{X}} = 0$. In a similar spirit to the average derivative estimator of Stoker (1986), this is a function of λ_2^h, λ_2^c weighted by the average derivative of the distribution function G_{ν^h, ν^c} on the index for h and c , respectively. If G_{ν^h, ν^c} were unknown, λ_2^h, λ_2^c would have been identified up to scale only. But since G_{ν^h, ν^c} is known, λ_2^h and λ_2^c can be identified using moments involving period-2 choices (h, c or n).

Outcome equations. Conditional on choice sequence D_{it} , treatment group indicator Z_i , and other covariates $\mathbf{X}_i(\equiv (\mathbf{X}_{ih}, \mathbf{X}_{ic}))$, the expected potential outcome is:

$$E(Y_{it}^d | D_{it} = d, Z_i = z, \mathbf{X}_i = x) = \alpha_t^d + \beta_{th}^d q_{ih} + \beta_{tc}^d q_{ic} + \gamma_t^d \lambda(d, z, x) + E(\varepsilon_{it}^d | D_{it}=d, Z_i=z, \mathbf{X}_i=x) \quad (23)$$

where the conditional expected factor $\lambda(d, z, x)$, which captures selection, is identified:

$$\begin{aligned} \lambda(d, z, x) \equiv E(\theta_i | D_{it} = d, Z_i = z, \mathbf{X}_i = x) &= \int \theta f(\theta | D_{it} = d, Z_i = z, \mathbf{X}_i = x) d\theta & (24) \\ &= \int \theta \frac{p(D_{it} = d | \theta_i = \theta, Z_i = z, \mathbf{X}_i = x)}{p(D_{it} = d | Z_i = z, \mathbf{X}_i = x)} f_{\theta_i}(\theta) d\theta & (25) \end{aligned}$$

since the conditional choice probability in the numerator is known from the identified choice equations at $t = 1, 2$, $f_{\theta_i}(\cdot)$ is known, and the denominator term is directly identified from data. The conditional density of ε_{it}^d can be simplified due to mutual independence of ε^d , (ν^h, ν^c) , θ , Z , \mathbf{X} :

$$f_{\varepsilon_{it}^d | D_{it}, \theta, Z_i, \mathbf{X}_i} = \frac{f_{\varepsilon_{it}^d, D_{it} | \theta, Z_i, \mathbf{X}_i}}{f_{D_{it} | \theta, Z_i, \mathbf{X}_i}} = \frac{f_{\varepsilon_{it}^d | \theta, Z_i, \mathbf{X}_i} f_{D_{it} | \theta, Z_i, \mathbf{X}_i}}{f_{D_{it} | \theta, Z_i, \mathbf{X}_i}} = f_{\varepsilon_{it}^d | \theta, Z_i, \mathbf{X}_i} = f_{\varepsilon_{it}^d} \quad (26)$$

Hence $E(\varepsilon_{it}^d | D_{it}, Z_i, \mathbf{X}_i) = 0$ by iterated expectations and a location normalization to zero.

In a similar spirit to Hansen, Heckman, and Mullen (2004), the factor loading γ_t^d is identified via exclusion restrictions in relation to $\lambda(d, z, x)$. For example, utilizing experimental variation Z_i (see also Kline and Walters (2016)), we obtain:

$$\gamma_t^d = \frac{E(Y_{it}^d | D_{it} = d, Z_i = 1, \mathbf{X}_i = x) - E(Y_{it}^d | D_{it} = d, Z_i = 0, \mathbf{X}_i = x)}{\lambda(d, 1, x) - \lambda(d, 0, x)} \quad (27)$$

where the numerator terms are identified from the observed outcome Y_{it} . The denominator is non-zero if Z_i modifies the selection mix of individuals into D_{it} by shifting choices. In the empirical model, we use Z_i , $Z_i \mathbf{X}_{ih}$, tp_{ih} , tp_{ic} as exclusion restrictions that affect choices and are excluded from the outcome equations (Table 4).³¹ We use the data and the estimated choice model to compute $\lambda(d, z, x)$, and show that $\lambda(d, 1, x) \neq \lambda(d, 0, x)$ for all d and all major values of x (Figure A.4).

³¹ Additionally, q_{ih} and q_{ic} are exclusion restrictions in the outcome equations where β_{th}^d and/or β_{tc}^d are set to zero.

The intercept α_t^d is identified:

$$\alpha_t^d = E(Y_{it}^d | D_{it} = d, Z_i = z, \mathbf{X}_i = 0) - \gamma_t^d \lambda(d, z, 0) \quad (28)$$

as are β_{th}^d and β_{tc}^d by conditioning on two values of q_{ih} and q_{ic} , respectively. Using equation (26), express the conditional density of Y_{it}^d as a convolution:

$$f_{Y_{it}^d | D_{it}=d, Z_i=z, \mathbf{X}_i=x}(y|d, z, x) = \int f_{\varepsilon_{it}^d}(y - \alpha_t^d - \beta_{th}^d q_{ih} - \beta_{tc}^d q_{ic} - \gamma_t^d \theta) f_{\theta_i | D_{it}=d, Z_i=z, \mathbf{X}_i=x}(\theta|d, z, x) d\theta \quad (29)$$

The density $f_{\varepsilon_{it}^d}(\cdot)$ can then be identified by deconvolution since $f_{Y_{it}^d | D_{it}=d, Z_i=z, \mathbf{X}_i=x}$ is observed and $f_{\theta_i | D_{it}=d, Z_i=z, \mathbf{X}_i=x}$ is known.³² Finally, the joint distribution of potential outcomes is identified due to ε_{it}^d i.i.d. across d :

$$f_{Y_{it}^1, \dots, Y_{it}^{\bar{d}} | \theta_i=\theta, Z_i=z, \mathbf{X}_i=x}(y^1, \dots, y^{\bar{d}} | \theta, z, x) = \prod_{d=1}^{\bar{d}} f_{\varepsilon_{it}^d}(y^d - \alpha_t^d - \beta_{th}^d q_{ih} - \beta_{tc}^d q_{ic} - \gamma_t^d \theta) \quad (30)$$

Discussion. Our approach complements the existing literature that examines semiparametric identification of the joint system of measurement, choice and outcome equations (e.g., Carneiro, Hansen, and Heckman (2003), Heckman and Navarro (2007), Fruehwirth, Navarro, and Takahashi (2016)). While our model is more complicated in that it incorporates multiple alternatives in a sequential model, our identification approach is in line with this literature by exploiting exogenous measurements, choice index restrictions and a factor structure in errors for identification.

The general framework of Heckman and Navarro (2007), which focuses on stopping-time sequential models, considers identification with period-specific instruments via limit set arguments as in Heckman and Smith (1998).³³ Yet they note that this may be a demanding requirement, and they develop two strategies that impose additional structure: (1) Time-invariant instruments can be used instead if their coefficients differ across periods, which is motivated by the model's shrinking time horizon;

³²Few have examined deconvolution involving selection. An exception is Camirand Lemyre, Carroll, and Delaigle (2022), who consider $W = X + U$ where W is observed when $W > 0$, $p(W > 0) = H(\beta_0 + \beta_1 X)$ with known $H(\cdot)$, latent variable X and error U are unobserved and distributions unknown, and $U \perp X$. They show that if there are two observations per individual and f_U is symmetric and continuous, then f_X , f_U , β_0 , β_1 can be identified. They do not have other sources of identifying information like ours.

³³In the limit, all individuals are shifted to the same choice (or path) by an instrument, which eliminates the selection problem involving that choice, allowing the joint distribution of choice errors and the outcome error to be identified.

(2) A factor structure has the crucial benefit of identifying the joint distribution of counterfactual outcomes (treatment effects), which cannot be achieved using instruments alone. Our approach can be viewed as a mix of these strategies. While we impose stronger dependence assumptions on choice and outcome errors via a factor structure, we have period-1 specific instruments $(Z_i, Z_i \mathbf{X}_{ih})$, time-invariant instruments (tp_{ih}, tp_{ic}) , and alternative-specific covariates (e.g., q_{ih}, q_{ic}) for identification.³⁴ This analysis is useful for applications in more typical empirical settings, including randomization experiments.

4.2 Estimation

Although the identification does *not* require any distributional assumption on the unobserved factor, to facilitate estimation, we assume that the unobserved factor follows a normal distribution, and consider sensitivity of our estimates to alternative distributional assumptions of the unobserved factor (see Section 5.1). Alternative-specific choice shocks are assumed to follow a bivariate normal distribution with zero means, variance of 1 (scale normalization) for ν_{it}^h , variance of $\sigma_{\nu^c}^2$ for ν_{it}^c , and correlation coefficient ρ_{hc} . We also assume that the potential-outcome-specific idiosyncratic error terms (ε_{it}^d) and the measurement error terms (e_i^M) are normally distributed. Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n) , to avoid overfitting we restrict the potential outcome for (c, n) to have the same intercept and factor loading as (h, n) (i.e., $\alpha_t^{cn} = \alpha_t^{hn}$, $\gamma_t^{cn} = \gamma_t^{hn}$). The model is estimated using the method of maximum likelihood. Details of the estimation procedure are presented in Appendix Section F.2.

5 Results from the Sequential Threshold Model

This section reports the results from the sequential threshold model. Section 5.1 presents the model parameter estimates and sensitivity checks. Section 5.2 discusses the estimated joint, partial and cross returns from sequential program participation and how these parameters can be utilized to improve the average outcome in the population. Section 5.3 analyzes the characteristics and returns of HSIS experimental compliers.

³⁴While θ_i is unidimensional, our strategy applies to the multidimensional setting, with tedious derivations. A necessary condition for identification is that the number of measurements is at least $2 \times (\text{number of factors}) + 1$.

5.1 Estimated Model Parameters

Goodness of fit. Table 5 shows that program choices and mean outcomes in different program paths between model predictions and data and by experimental group line up closely. The model predictions are evaluated by simulating 30 paths of potential choices and associated outcomes for each individual in the sample with and without assigning an experimental offer.³⁵ Appendix Figure A.5 shows that the predicted and actual distribution of outcomes also overlap closely. In addition, the predicted experimental return from the estimated model is very close to the experimental return that is nonparametrically identified in the data.³⁶ Appendix Figure A.6 further shows that the close correspondence between model-predicted and actual experimental return and compliance share holds even across different groups stratified by baseline test scores, which we have documented large heterogeneity in Section 2.3. One avenue where the model can be improved is that it under-predicts the experimental returns for children in the bottom two deciles of the baseline test scores distribution (corresponding to the observations with the largest estimated experimental returns).

Measurement equations. Using the estimated measurement equation parameters (panels B and C in Table 7), we compute the signal-to-noise ratio for each measurement.³⁷ The baseline PPVT score has the largest signal to noise ratio (37%), followed by the baseline WJIII Letter-Word score (31%) and WJIII Spelling score (17%). Therefore, all three measurements include a substantial amount of information about the unobserved factor. At the same time, they also demonstrate the importance of allowing for measurement error.

Choice equations. The estimates and standard errors of the parameters of the choice equations are reported in Table 6. Given the normalization ($\gamma^{PP} = 1$), the unobservable factor θ_i can be interpreted as child baseline ability. The factor loading in the period-1 perceived value of c (λ_1^c) is positive and large, whereas the factor loading in the perceived value of h in the same period is much smaller and insignificantly different from zero. This indicates that children choosing c in the first period are

³⁵Each simulation consists of random draws of the unobserved factor and utility shocks. For each simulation path, we first simulate all choice shocks across periods, and then switch Z from 0 to 1 keeping all shocks fixed.

³⁶The experimental return is 0.027 in the data (without controlling for any covariates). The model-predicted return is 0.024.

³⁷The expression for the signal-to-noise ratio of measurement M is $\frac{(\gamma^M)^2 \sigma_\theta^2}{(\gamma^M)^2 \sigma_\theta^2 + \sigma_{eM}^2}$. It is the ratio of the variance of the latent factor to the variance of each measurement.

positively selected by ability. Although no significant pattern of selection emerges among those choosing h relative to n , there is negative selection into h relative to c because children are positively selected into c . Not surprisingly, being in the treatment group ($Z = 1$) increases the probability of enrolling in HS. The increments are larger among HS centers that are rated as high quality or provide transportation services. By contrast, children in the control group ($Z = 0$) are less likely to enroll in high-quality HS centers or HS centers with transportation services, presumably because these centers are more likely to be filled to capacity and generate fewer crossovers.

Our estimates reveal significant state dependence in program choices, as $\psi_2^{jj'} > 0, \forall j, j' \in \{h, c\}$. To quantify the magnitude of state dependence, we compute the average marginal response of program choices in period 2 when period 1 enrollment changes. When all children are assigned to h instead of n in period 1, the period-2 proportion of children choosing h increases by 28 percentage points (pp) while the proportion of choosing c decreases by 10 pp (and hence the probability of choosing home care also declines, by 18 pp). Similarly, if all children are assigned to c instead of n in period 1, the period-2 proportion of enrolling in c increases by 26 pp and the proportion of enrolling in h declines by 9 pp. Therefore, all else being equal, enrollment in HS or other center care in age 3 increases the likelihood of participating in the same type of center-based care and decreases the likelihood of using home care in age 4.

The correlation between HS and other-center care choice shocks (ρ_{hc}) is positive and precisely estimated, rejecting the IIA assumption underlying the logit-based choice models.

Potential outcomes. Panel A in Table 7 reports parameter estimates of the potential outcome equations. The benchmark sequence is two periods of home care (n, n). For all other program sequences, we report the parameters in terms of difference relative to the benchmark sequence. Given that the covariates are normalized to have unconditional mean zero, the pairwise difference in the intercepts $\alpha_t^d - \alpha_t^{nn}$ reflects the average return to program sequence d relative to two periods of home care.

We find that experience with any center-based care over the 2-year period is beneficial to children overall relative to home care alone, as evidenced by the positive average returns across program paths. The relative sizes of the average returns are informative about the overall effectiveness of different program paths. For instance, if all children are enrolled in the HS program at age 3, their average outcome would be higher if they are assigned to other center care at age 4 rather than staying in HS

for another year ($\alpha_t^{hc} > \alpha_t^{hh}$). However, this does not necessarily imply that age-4 enrollment in HS has zero return. In fact, if all children were assigned to home care at age 3, they would be better off by enrolling in HS instead of home care at age 4 ($\alpha_t^{nh} > \alpha_t^{nn}$). Section 5.2 provides an extensive analysis on the program returns and their implications for policy.

Although not all of the estimated coefficients on the factor are precise, we can reject the hypotheses of no selection ($H_0 : \gamma_t^d = 0, \forall d$; p-val < 0.01) and constant returns ($H_0 : \gamma_t^j = \gamma_t^k, \forall j \neq k$; p-val < 0.01). The evidence for selection on gains is mixed. We find evidence for Roy’s selection on gains into (h, h) relative to (h, n) , because low-ability children are more likely to select h relative to n in period 2 and their outcomes improve more when shifted from n to h than high-ability children ($\gamma_t^{hh} - \gamma_t^{hn} < 0$). On the contrary, high-ability children are more likely to self-select into c relative to n in both periods, yet they have smaller than average gains when switching from n to c (e.g., $\gamma_t^{cc} - \gamma_t^{nn} < 0$ and $\gamma_t^{hc} - \gamma_t^{hn} < 0$). This is consistent with reverse-Roy pattern of negative selection.

Sensitivity Analysis. In Appendix Section B, we report sensitivity analysis of our estimates against functional form and distributional assumptions. We show that our parameter estimates are robust to a richer specification of state dependence in the decision rules in $t = 2$ (Appendix Section B.1.1), excluding transportation services and their interaction with Z from choice equations in both periods (Appendix Section B.1.3), and relaxing the factor distribution used in estimation by modelling it as a mixture normal distribution (Appendix Section B.2). The choice equations in the sequential threshold model are approximations to the perceived value of each alternative. For forward-looking individuals, these perceived values are a combination of utility flow and continuation value, meaning that the perceived values may be nonlinear in parameters due to incorporation of uncertainty of future payoffs.³⁸ To account for such nonlinearity, we also estimate a model by specifying a quadratic function of factor θ in the period-1 choice equations for h and c (see Appendix Section B.1.2). Our estimates are robust across these alternative specifications.

Finally, Appendix Section B.3 reports the results when WJIII Applied Problems test are included in the measurement system and outcomes. The estimated parameters change little; if anything, incorporating early numeracy/math skills tends to strengthen the effects of certain program sequences for

³⁸For instance, in the dynamic structural model constructed in Section 6, we show that the expected value function in the first decision period is a nonlinear function of parameters (also see Appendix Section F).

low-ability children.

5.2 Joint, Partial and Cross Returns of Childcare Programs

Average program returns. Figure 1 shows the average partial, joint, and cross returns from program sequences. These parameters are population averages of the individual-specific causal returns defined in Table 3, abstracting from endogenous selection.³⁹ The first bar (HH) in panel (a) shows that enrolling in HS at either age 3 or age 4 has a positive effect on test scores compared to home care. The partial return is slightly smaller for HS at $t = 1$ (0.082, or 0.30sd), relative to HS at $t = 2$ (0.092, or 0.34sd).⁴⁰ This bar also shows that the joint return from enrolling in HS at both $t = 1$ and $t = 2$ is positive (0.098, or 0.36sd), yet it is smaller than the sum of the two partial returns of HS. In fact, this joint return is only slightly higher than the partial return from HS obtained at $t = 1$. This indicates a negative and significant within-program cross return from two periods of HS, as shown in the first bar of panel (b). Age-4 enrollment in HS reduces the effectiveness of HS received at age 3 on average, suggesting that one period of HS enrollment is a substitute for another.⁴¹

We find a large and positive joint return from HS at $t = 1$ combined with other center care at $t = 2$ (0.141, or 0.52sd; column (HC)), which is *higher* than the joint return from two periods of HS (p-val=0.01) and also (h, n) (p-val=0.06). Given that the partial return of HS at $t = 1$ is also positive, we conclude that enrollment in HS at age 3, regardless of the type of the follow-up program in age 4, leads to higher average returns than receiving home care at both ages. There is also evidence that the competences offered by HS and by other center care are less substitutable than two periods of HS, as the estimated cross return of (h, c) is larger than the cross return of (h, h), although the difference is not statistically significant (p-val=0.26).

The remaining bars in panel (a) report average program returns where all individuals are assigned to other center care in the first period. Children in both (h, c) and (c, h) experience one year of HS and one year of other center care by the end of age 4. However, the joint returns differ widely between these

³⁹Given the covariates in the potential outcome equations are normalized to have mean zero (see Section 4), they correspond to the intercepts in the outcome equations.

⁴⁰1sd is defined as the standard deviation of the age-4 test scores in the control group (1sd=0.269).

⁴¹Note that all program returns are defined relative to the benchmark sequence (n, n). The benchmark sequence (n, n) is also a type of investment which may change the child's cognitive outcomes. The average joint return from (h, h) being smaller than the average partial return of (h, n) does not imply human capital destruction because (h, n) is an alternative sequence of investments, rather than merely HS and no other investments.

two program sequences – the return from (h, c) is significantly higher than (c, h) (p-val=0.003), with the latter return small and insignificantly different from 0 (p-val=0.14). Such difference is largely driven by the differential cross returns; the cross return of (c, h) is much lower (significant at 5% level) than that of (h, c) . This result highlights the importance of considering the *sequence* of program participation.

Finally, for two periods of enrollment in other center care (c, c) , both the joint return (panel (a), fourth bar) and the cross return (panel (b), fourth bar) are quite similar to those from two periods of HS (h, h) . The partial return from enrolling in c in $t = 2$ is also not significantly different from the partial return from enrolling h in the same period (p-val=0.76).⁴² Therefore, our estimates indicate that HS provides similar returns to other-center care if children can only experience a single type of care in pre-school years. When children can switch between HS and other-center care in pre-school years (as in reality), our estimates suggest that the combined joint returns from different programs can vary a lot depending on their sequence.

Heterogeneity in program returns. The average returns hide relevant heterogeneity. To illustrate how program returns differ by child latent baseline ability (θ), Figure 2 shows the estimated returns from program sequences – (h, h) , (h, c) , (c, h) , or (c, c) – for individuals at different percentiles of the θ -distribution. Additionally, Appendix Figure A.7 presents the partial, joint, and cross returns and their 95% confidence intervals when θ is fixed at the 10th and 90th percentile of the factor distribution.⁴³

Figure 2 shows that the lowest-ability children gain the most from prolonged formal care participation – the joint returns across four program sequences are the highest for children with the lowest θ and declining with θ . For example, in the (h, h) sequence, children at the 10th percentile of θ gain 0.16 (0.59sd, p-val=0.004), while children at the 90th percentile of θ only gain 0.03 (0.13sd) which is insignificantly different from 0. This results in a narrowing of 0.13 in the age-4 test score, which constitutes close to one-fourth of the 90/10 gap when individuals are assigned to home care in both periods instead.⁴⁴ These results indicate compensatory effects of two periods of center-based care for

⁴²Recall that we restrict the potential outcomes to be common across program sequence (h, n) and (c, n) , due to the small sample size of parents making (c, n) choice in the data (see Section 4.2).

⁴³We have also explored how program returns differ by the quality of the childcare center. Appendix Figure A.8 shows that enrolling in high-quality centers in both $t = 1$ and $t = 2$ boost the cross returns for (h, h) and (c, c) .

⁴⁴Under the baseline scenario of home care at both age 3 and 4 (i.e., fixing the program sequence to (n, n)), moving from the 10th to the 90th percentile of θ implies an increase in age-4 test score by 0.55. This is computed from the (n, n) potential outcome equation: $3.268 + 0.974\theta$. At the 10th and 90th percentiles of θ , the scores are 2.98 and 3.54, respectively.

those at the bottom of baseline ability distribution. Notably, the partial returns to HS at $t = 1$ are positive and increasing with θ , suggesting that enrollment in HS in age 3 (while fixing age 4 to home care) can benefit children with high baseline ability. Therefore, for the program sequence (h, n) , there is evidence in favor of the skills-beget-skills hypothesis.

Overall, our estimates indicate that prolonged exposure to HS is important for low-ability children. One interpretation is that high-ability children may have adapted to the HS curriculum quickly and hence repeated exposure to the same curriculum at age 4 is as beneficial as engaging them with home care or other-center care.⁴⁵ By comparison, low-ability children at age 3 may not have developed the same level of school readiness as the high-ability children of the same age. In fact, for low-ability children, the partial returns to HS or other-center care at age 4 are larger than the corresponding partial returns at age 3 (although they are still smaller than either (h, h) or (h, c)).⁴⁶

Although the mean cross returns are negative, there is some weak evidence for positive cross returns among low-ability children. For instance, the estimated cross return of sequence (h, c) is positive (p-val=0.16) for children at the 10th percentiles of θ distribution (Panel (b) in Appendix Figure A.7). This indicates that receiving other center care at $t = 2$ boosts returns to HS at $t = 1$ for these children, an evidence that is consistent with dynamic complementarity.

Optimal pathways. Our results are informative of the “optimal” pathways for a policymaker who can assign different children to different program paths in order to maximize their outcomes. For instance, for children with below-median θ , our estimates imply that the policymaker should assign them to HS at $t = 1$ followed by other center care at $t = 2$ in order to maximize their expected outcomes. High- θ children should instead be assigned to any center-based care at $t = 1$ followed by home care at $t = 2$.

Appendix Section D formalizes this discussion by considering an optimal program assignment prob-

⁴⁵The Head Start Performance Standard mandates delegate agencies to implement a curriculum. The majority of children in HS in 2000 followed either High Scope or Creative Curriculum (Shaul, Ward-Zukerman, Edmondson, Moy, Moriarity, and Picyk, 2003).

⁴⁶Another interpretation is that the quality of home care differs by the baseline ability of the child. Relative to center care, existing evidence suggests that children from disadvantaged environments are exposed to a substantially less rich vocabulary than children from more advantaged families (Hart and Risley (1995), Fernald, Marchman, and Weisleder (2013)). Chan and Liu (2018) provide additional evidence from estimation of children’s cognitive ability production function. Therefore, home care could be less effective than center care for low-ability children than for high-ability children. This could also explain the high returns from prolonged participation in center-based care for low-ability children and from exposure to some home care for high-ability children.

lem facing a social planner whose objective is to maximize the average test score of the population at the end of age 4. If the planner observes each child’s baseline ability θ_i perfectly, then assigning individuals to their optimal program sequence will lead to an average age-4 test score of 3.424 (Table 8, second row). This is an increase of 0.053 (0.20sd) relative to the average outcome when all children are in the treatment group and choose their program sequences according to the estimated threshold model (Table 8, first row). The improvement is felt across ability levels, with high-ability children gaining more from the external assignment relative to making their individual choice (Table 8, columns (2)–(4)). This is consistent with our previous finding that high-ability children tend to self-select into center-based care in both periods, which does not maximize their test scores.

Yet θ_i cannot be perfectly observed by the social planner. The remaining rows in Table 8 show that the improvements in test scores are smaller when the social planner can only assign the optimal program path using baseline test scores instead of θ_i . The efficiency loss depends on the signal-to-noise ratio of the specific measurement being used. For instance, the improvement is the smallest (0.039, or 0.14sd) when the social planner relies on the baseline WJIII-Spelling test, which has the lowest signal-to-noise ratio among the three baseline measurements (as reported in Section 5.1).

5.3 Program Returns for HSIS Compliers

In this section we analyze the program returns for HSIS compliers, i.e., children who respond to the HS offer by enrolling in HS. Standard classification of experimental compliers focuses on program choice in period 1 only. For instance, Kline and Walters (2016) partition the HSIS compliers into two types, those who switch to HS from home care and other-center care when receiving the experimental offer. A novelty of our framework is that an increase in the availability of HS at $t = 1$ can change the demand for center-based care at $t = 2$, because childcare experience in $t = 1$ directly affects childcare choice in $t = 2$ (state dependence). This means that we can predict returns for multiple types of compliers, who differ in their program choices made with and without the experimental offer and over time.⁴⁷

⁴⁷As in the sequential threshold model, we assume that the experimental offer Z_i affects the valuation for h in $t = 1$ only; importantly, the experimental offer does not affect the perceived value of c and n in the same period, nor does it directly affect the perceived value of any childcare choices beyond $t = 1$. This assumption restricts period-1 behaviour in that anyone who responds to the HS offer does so to attend HS in $t = 1$. The same restriction is imposed by Kline and Walters (2016) and by Kirkeboen, Leuven, and Mogstad (2016) in a different context. The compliance shares and compliance-specific returns are computed from by simulating 30 paths of potential choices and associated outcomes for each individual in the sample, with ($Z = 1$) and without ($Z = 0$) the experimental offer (as in Section 5.1).

We classify the HSIS experimental compliers into two groups: (1) temporary switchers, where *only* the choice at $t = 1$ is changed due to the HS offer at $t = 1$, and (2) path switchers, where the entire *trajectory* of choice is changed due to the HS offer at $t = 1$. The incidence of path switchers signifies the importance of state dependence in program choice.

Further distinguishing the type of care from which compliers switch to HS in $t = 1$, we partition the experimental compliers into the following 4 types:

- temporary n-compliers: $D_i^0 = (n, j), D_i^1 = (h, j), \forall j \in \{h, c, n\}$
- temporary c-compliers: $D_i^0 = (c, j), D_i^1 = (h, j), \forall j \in \{h, c, n\}$
- path n-compliers: $D_i^0 = (n, j), D_i^1 = (h, j'), \forall j, j' \in \{h, c, n\}$ and $j \neq j'$
- path c-compliers: $D_i^0 = (c, j), D_i^1 = (h, j'), \forall j, j' \in \{h, c, n\}$ and $j \neq j'$

where $D_i^z = (K_{i1}^z, K_{i2}^z)$ is individual i 's potential sequence of choices when the experimental offer is $Z_i = z$ (with K_{it}^z being the potential choice at time t with $Z_i = z$). For each individual, we observe either D_i^0 if the experimental variation puts her in the control group, or D_i^1 if she belongs to the treatment group.

Column (1) of Table 9 reports the population shares of four compliance types. Considering the potential choices in $t = 1$ only, the HS offer shifts 50.84% of the population from n to h ($=29.39+21.45$; n-compliers), and 19.86% from c to h ($=7.92+11.94$; c-compliers).⁴⁸ Therefore, when considering only the first period, the HS program crowds out other center care and home care.

Considering potential choices in a multi-period framework enriches substitution patterns. We find that *path switchers* constitute about 60% ($\frac{11.94+29.39}{70.70}$) of all compliers – that is, more than half of the compliers also change the program choice at age 4 due to the HS offer at age 3. Appendix Table A.2 further examines the compliance behaviors based on the individual potential program sequences when in the treatment group and when in the control group. Among c-compliers, 43% (or 8.33% of the population) change their period-2 choice from c to h ($D^0 = (c, c), D^1 = (h, h)$). This implies that

⁴⁸These compliance types can also be non-parametrically identified from the data directly. More specifically, $P(K_{i1} = n|Z_i = 0) - P(K_{i1} = n|Z_i = 1)$ identifies the proportion of n-compliers and $P(K_{i1} = c|Z_i = 0) - P(K_{i1} = c|Z_i = 1)$ identifies the proportion of c-compliers. According to this calculation, n-CP compliers are 51.2% of the population and c-CP compliers are 19.3% (see Table 1, which shows that the HS offer reduces the share of children in other center care from 25% to 5%, and reduces the share of children in home care from 61% to 10%). These estimates are very close to the model predictions, lending further credibility to the model. They are also comparable with Kline and Walters (2016) (n -compliers is 45% and c -compliers is 23%), whose sample also includes the age-4 children cohort.

receiving the HS offer not only reduces the demand for c at $t = 1$ but also moves a substantial fraction of children away from c towards h at $t = 2$. Among n-compliers, 40% (or 20.19% of the population) change their period-2 choice to h ($D^0 = \{(n, c), (n, n)\}, D^1 = (h, h)$). Note that path switchers do not necessarily change their period-2 choice to HS. For example, some of them actually switch from h to c at $t = 2$, both among the c-compliers ($D^0 = (c, h), D^1 = (h, c)$, 6% of c-compliers) and the n-compliers ($D^0 = (n, h), D^1 = (h, c)$, 9% of the n-compliers).

The remaining columns in Table 9 report the estimated average return (column 3) and the average value of θ (column 2) for each complier type. The overall LATE (=0.034) masks large heterogeneity in the returns across different types of compliers. c-compliers are dominated by high- θ children, who are likely to have larger gain from switching to HS temporarily in $t = 1$ rather than changing their entire path of childcare choices (the latter may reduce home care use which is beneficial to high-ability children). The pattern is reversed for n-compliers, who are relatively less able. These children instead benefit from a permanent change in program path induced by the experimental variation, as the new program path is likely to involve prolonged exposure to center-based care.

6 A Dynamic Structural Model with Forward-looking Behavior

The sequential threshold model represents reduced-form decision rules that approximate choices from a structural model. We now build and estimate a dynamic structural model that involves a forward-looking agent. The structural model's key distinction is that it focuses on ex ante choice decisions via an explicit information set and a model of expectations, which allows the utility flow and continuation value of each alternative to be disentangled. We then use the estimated structural model to examine the effect of restricting HS enrolment to a maximum of one year. We show that specifying forward-looking behavior is important to understand the policy implications from such a fundamental change to program features.

To facilitate comparison, the structural model is set up to be closely connected to the sequential threshold model in Section 4. Flow utilities are linear functions of center characteristics and factor θ_i as in the threshold model (see Appendix E and discussions below). Importantly, the center characteristics covariates, θ_i , as well as the technology of child skill production are assumed to be known to the individual throughout; the choice and outcome shocks are revealed to the individual sequentially.

In each period, the individual chooses an alternative based on the realization of preference shocks, comparing the current utility flow and expected future value embedded in each alternative. To express this intertemporal optimization problem in recursive form, consider an individual having chosen $k_1 \in \{h, c, n\}$ at $t = 1$. The value function at $t = 2$ is:

$$V_{i2}(K_{i1} = k_1, \mathbf{X}_{ih}, \mathbf{X}_{ic}, \theta_i) := \max_{k_2 \in \{h, c, n\}} \left[\tilde{u}_{i2}^{k_2}(k_1, \mathbf{X}_{ih}, \mathbf{X}_{ic}, \theta_i) + \kappa E_2 Y_{iT}^{k_1, k_2}(q_{ih}, q_{ic}, \theta_i) \right] \quad (31)$$

The value function $V_{i2}(\cdot)$ has state variables K_{i1} , \mathbf{X}_{ih} , \mathbf{X}_{ic} , and θ_i . The *utility flow* from each alternative in $t = 2$ is denoted by $\tilde{u}_{i2}^{k_2}$. Besides the utility flow, the individual receives a terminal value as a function of the period- T outcome Y_{iT}^d given program sequence $d = (k_1, k_2)$ (where $k_1, k_2 \in \{h, c, n\}$). Potential outcome Y_{iT}^d , given by equation (9), represents the technology of child skill production and is subject to idiosyncratic outcome shocks ε_{iT}^d . The expected outcome $E_2 Y_{iT}^{k_1, k_2}(q_{ih}, q_{ic}, \theta_i) = \bar{Y}_{iT}^{k_1, k_2}(q_{ih}, q_{ic}, \theta_i)$, which is the outcome exclusive of the idiosyncratic shocks.

The terminal value scaling factor, κ , quantifies how sensitive the choices are to expected test scores. A test of $\kappa = 0$ versus $\kappa \neq 0$ informs whether the individual incorporates the expected future outcomes as part of her information set and acts upon it (Abbring and Heckman, 2007). Hence the decision problem at $t = 2$ is simply an intertemporal tradeoff between the current utility flow and the expected future outcome.

At $t = 1$, the value function is:

$$V_{i1}(Z_i, \mathbf{X}_{ih}, \mathbf{X}_{ic}, \theta_i) := \max_{k_1 \in \{h, c, n\}} \left[\tilde{u}_{i1}^{k_1}(Z_i, \mathbf{X}_{ih}, \mathbf{X}_{ic}, \theta_i) + \delta E_1 V_{i2}(k_1, \mathbf{X}_{ih}, \mathbf{X}_{ic}, \theta_i) \right] \quad (32)$$

with state variables Z_i , \mathbf{X}_{ih} , \mathbf{X}_{ic} , and θ_i ; δ is the discount factor that determines the importance of continuation values in ex ante choice decisions; $E_1(\cdot)$ is an expectation that integrates out the preference shocks at $t = 2$. Appendix Section F derives the semi-closed form solution to $E_1 V_{i2}(\cdot)$, which is a nonlinear function of parameters and incorporates uncertainty of preference shocks at $t = 2$ and expected returns in outcomes.⁴⁹ $\tilde{u}_{i1}^{k_1}$ captures utility flow specific to each alternative in $t = 1$. We emphasize that, as long as $\delta \neq 0$ and $\kappa \neq 0$, the flow utilities are different from the perceived values

⁴⁹ Across all periods, we assume that the preference shocks have the same statistical properties as in the sequential threshold model. More specifically, the preference shocks, $(\tilde{v}_{ij}^h, \tilde{v}_{ij}^c)$, follow a bivariate normal distribution with zero mean and covariance matrix $[1 \quad \tilde{\rho}_{hc} \tilde{\sigma}_{\nu^c}; \quad \tilde{\rho}_{hc} \tilde{\sigma}_{\nu^c} \quad \tilde{\sigma}_{\nu^c}^2]$, and are independent across time.

defined in the sequential threshold model.

Identification of the Structural Model. If the flow utility is as flexibly specified as the perceived values in the threshold model, then the structural model is underidentified. For instance, equations (31) and (32) suggest that the intercepts and state dependence parameters in the flow utility cannot be disentangled from the terminal and expected values. Also, both values are affected by the factor and center characteristics, yet these same variables also affect the individual’s flow utility.⁵⁰

Therefore, the identification of the structural model requires additional exclusion restrictions – the excluded variable should affect decisions via expected future outcomes only. Specifically, compared to the threshold model (equations (2)-(8)), center quality variables (q_{ih}, q_{ic}) are excluded from the utility flow in $t = 2$ (see Appendix E for details). That is, we assume that center quality drives childcare choice in period 2 only via expectations of improved future outcomes. These exclusions identify the terminal value scaling factor κ .

To identify the discount factor δ , we impose exclusions in the period-1 flow utility. The excluded variables include the quality of other-center care (q_{ic}) and the HS center quality of the treated group ($q_{ih} \times Z_i$). The rationale for the latter is driven by the institutional background. In period 1, control group individuals do not have an experimental HS offer, but they may “crossover” and enrol primarily in low-quality HS centers (see Sections 2.1 and 5). Due to this supply-side barrier cost to crossover, we keep $q_{ih} \times (1 - Z_i)$ in the period-1 utility flow for h among control group individuals (hence it is not an exclusion variable). By contrast, treatment group individuals do not face this barrier cost because they all have an experimental HS offer.⁵¹

Empirically, we find that δ is harder to identify than κ . This may be due to the high HS takeup rate in the treatment group at $t = 1$. In the baseline structural model, we set $\delta = 1$. In a sensitivity analysis, we estimate δ jointly with other parameters in the structural model. The point estimate

⁵⁰Without loss of generality, the intercept of the terminal value function is normalized to zero. The intercept affects the terminal values of all states by an equal amount so it does not affect program choices. See also Keane and Wolpin (2001), Fang and Silverman (2009) and Chan and Liu (2018).

⁵¹An alternative approach is to restrict the form of state dependence and unobserved heterogeneity by reducing the number of intercept and loading parameters in the utility function. This approach is popular in conventional structural models for theoretical or computational considerations, and not necessarily (or explicitly) for the purpose of identification. The intuition for identification is also less transparent when we impose these restrictions. Therefore, we keep the utility function as closely comparable to the sequential threshold model as possible, and rely solely on the exclusion variables to identify δ and κ . More ideally, one would find experimental variation in a policy that affects expected future outcomes but not current utility; see, e.g., Chan (2017) for an example from welfare reform.

of δ is 0.81 (se=0.628), which provides reassurance that the exclusion restriction contains reasonable identifying information. Our results also change little if we set the discount factor δ to 0.95 (Appendix Table A.4).

Limitations. We acknowledge that this structural model is a simplified version of many richer models that incorporate other important behavioral mechanisms. We view this as a middle ground where certain key behavioral features such as continuation values are explicitly modeled, but others are kept implicit and subsumed into the preference parameters. For example, consider the tradeoff between child outcome and the parent’s labor supply.⁵² A high-education mother (whose child may have higher θ_i) may be aware that putting her child in center care may worsen the child’s outcome relative to taking care of the child on her own. However, she may still put her child in center care after weighing the benefits of working (e.g., better market opportunities or higher work preference) against the deterioration of child outcome. This tradeoff is subsumed into the factor loadings in the choice and outcome equations of our models.

6.1 Structural Model Estimation Results

Table 10 reports the estimates of the flow utility and terminal value scaling factor κ in the structural model. As discussed in Section 4, the technological parameters in the baseline measurement equations, outcome equations and the factor distribution have been identified from the sequential threshold model. To facilitate comparison of the results in the choice equations, when we estimate the structural model we fix all technological coefficients to be the same as in the sequential threshold model (see Table 7). Details of the estimation procedure are in Appendix Section F.⁵³

Comparing the estimated flow utility with the perceived value (estimated from the threshold model) reveals the relative importance of contemporaneous and continuation value in determining individuals’ program choices. For instance, in $t = 1$, the alternative-specific intercepts for h and c are considerably

⁵²See, for example, Griffen (2019) who estimates a full structural model of Head Start participation that captures this tradeoff explicitly. Like many other structural models, to maintain tractability the extra complexity comes at the expense of modeling certain parts of the model (outcome equation) in a highly stylized manner. Although our approach is less complex, it recognizes the sequential nature of participation in different programs and it models outcome equations flexibly.

⁵³Appendix Table A.3 shows that the predicted program choices align well with data, across experimental groups. The simulated experimental return also aligns closely with the experimental return in the data (0.025 predicted by the estimated structural model vs. 0.027 in the data).

more negative (-1.481 and -2.681, respectively) than those in the sequential threshold model, which suggests that average control-group individuals: (1) prefer home care (n) if they consider period-1 utility only, and (2) have higher continuation values for h and c than n . By contrast, the factor loadings in the h equation ($\tilde{\lambda}_1^h$ and also $\tilde{\beta}_{Z\theta}$) are considerably larger than that in the sequential threshold model, which suggests that high- θ individuals: (1) prefer HS (h) relative to n if they consider period-1 utility only, and (2) have lower continuation values for h than the other options. Indeed, as shown in Appendix Figure A.9, we find that the flow utility of h (relative to that of n and excluding choice shocks) increases with θ whereas its continuation value decreases with θ , leading to the total perceived value at $t = 1$ relatively flat with respect to θ . This provides a structural interpretation to the relatively small factor loading in the choice equation for h in the sequential choice model.⁵⁴

Similar to $t = 1$, the alternative-specific intercepts for h and c in $t = 2$ are smaller than those in the sequential threshold model. Yet there is stronger evidence of state dependence in the choice of center-based care, as period-1 program choice of h or c increases the period-2 flow utilities of h and c considerably more than the perceived values ($\tilde{\psi}_2^{hh}, \tilde{\psi}_2^{hc}, \tilde{\psi}_2^{ch}, \tilde{\psi}_2^{cc}$ larger than their counterparts in the threshold model). There is also a stronger gradient of flow utility on children's baseline ability than the threshold model, as indicated by the larger factor loadings for c and h . Note also that the period-2 utilities form part of the continuation values in period 1, which drives period-1 choice. For example, for high- θ individuals, choosing h at $t = 1$ may be attractive because it opens up the option to choose their strongly preferred c in period 2 (accounting for uncertainty in preference shocks).

The terminal value scaling factor κ has a point estimate of 3.512 and is significant at the 1-percent level. This suggests that individuals incorporate the expected future outcomes as part of the information set and act upon it. Specifically, they are more likely to choose an option that yields better expected outcomes, keeping preferences fixed.⁵⁵

⁵⁴Appendix Figure A.9 can be easily modified to show the flow utility and continuation values of h for individuals in the control group. The experimental offer Z only affects the flow utility of HS at $t = 1$ (and does not affect the continuation value).

⁵⁵To interpret its magnitude, consider an example of an individual in period 2 who chose home care in period 1. The difference in expected outcome is $0.082 + 0.227\theta$ if she chooses h relative to n ($Y^{nh} - Y^{nn}$ in Table 7). This difference in expected outcomes is equivalent to a difference in period-2 utility of $3.512 \times [0.082 + 0.227\theta]$ between h and n , implying that high- θ individuals are incentivized to choose h over n due to higher expected outcome of (n, h) relative to (n, n) .

6.2 Counterfactual policy reform

An important policy debate about Head Start is the desirability of two-year versus one-year enrolment. In this section, we use the structural model to perform a hypothetical policy reform, which restricts all treatment-group individuals to *at most* one period of HS only, while still allowing them to decide when to enrol in HS. Importantly, this fundamental change to program feature is revealed to all individuals at the beginning of period 1 (age 3). The key insight from the dynamic structural model is that individuals will adjust their program choice in $t = 1$ as soon as the information about the reform is revealed, because the policy reform changes the continuation value of choosing HS in $t = 1$.

We examine the effects of this policy on program choices, outcomes, and welfare. Table 11 reports the results, which are based on 30 simulated paths per treatment-group individual under the baseline and counterfactual scenarios and keeping the realization of shocks identical across scenarios. Relative to the no-restriction baseline scenario (Column 1), the policy (Column 2) reduces period-1 HS enrolment by 25.5pp (down from 85.2%), the majority of which switch to home care (+18.8pp) and the rest to other center care (+6.7pp). This suggests that many individuals move away from HS in period-1 preemptively, even though they are not impacted by the restriction immediately. Such preemptive behavior is driven by the decline in the continuation value of HS in $t = 1$, which reduces the overall perceived value of HS relative to c and n .⁵⁶ In period 2, the policy reduces HS enrolment by 51pp (down from 68.1%), half of which switch to home care (+26.5pp) and the other half to other center care (+24.4pp). The effect is larger in period 2 because some individuals have become ineligible for HS.

Overall, in terms of program sequence, the elimination of the (h, h) path (from 61.9% to 0%) results in the largest shift to (h, n) (+21.4pp) followed by (h, c) (+15pp). The shifts to (n, h) (+8.7pp) and (c, h) (+2.3pp) are considerable and reflect the tendency to “bank up” HS eligibility for possible enrolment in period 2 (note: the period-2 utility shocks are not yet realized in period 1). There are also shifts to (n, c) (+5.4pp), (n, n) (+4.7pp) and (c, c) (+4pp), which involve no HS enrolment at all. This is driven by forward-looking behavior under uncertainty – while the individual “banks up” HS eligibility in period 1, it ends up not being used due to particular realizations of period-2 utility shocks. The policy reduces the average outcome slightly (from 3.371 to 3.366), but increases the standard deviation

⁵⁶Appendix Figure A.10 shows the continuation values and the overall perceived values of h and c relative to n . The flow utilities are not affected by the policy reform.

of outcome by 2.7 points (from 0.245 to 0.272). The changes in outcomes by program sequence (Column 5) reveal the complexity of selection effects; some paths such as (h, c) experience minimal changes in average outcomes, while others such as (h, n) experience large reductions in average outcomes.⁵⁷

To gauge the importance of forward-looking behavior, we also consider the case where individuals are assumed to *not* know about the enrolment limit policy at $t = 1$, and can only reoptimize at $t = 2$ when knowledge of the policy becomes part of the information set (Column 3 in Table 11). There is no pre-emptive response at $t = 1$: individuals switch program choice at $t = 2$ only, and the elimination of the (h, h) path results in shifts to (h, c) (+29pp) and (h, n) (+32.9pp) exclusively. Relative to the full-response scenario, the policy yields *better* outcomes in that the average outcome increases by 0.5 points (from 3.371 to 3.376, Column 6). This positive effect is largely driven by improvements in the outcomes of individuals who shifted from (h, h) to (h, c) .

Distributional impacts of the reform. While the policy has limited impact on the average outcome of children, it exacerbates inequality. Figure 3 reports the distributional effect of the HS enrolment limit policy on outcomes (the “full response” line). Individuals with low factor θ perform substantially worse while those with high θ perform better. To further examine the source of increased inequality in outcomes, in Appendix Table A.6 we report the effects on program choice separately for individuals with values of θ at the 10th and at the 90th percentiles, respectively. The preemptive response is much stronger among low- θ individuals; in period 1 the policy reduces HS enrolment by 39.1pp (down from 87.4%), the majority of which switch to home care (+31.5pp). Overall, low- θ individuals tend to switch from the (h, h) path to (n, \cdot) paths, all of which make their outcomes worse. High- θ individuals exhibit weaker preemptive responses and they tend to switch from the (h, h) path to (h, c) or (h, n) paths, which make their outcomes better. In addition, some high- θ individuals switch from the (h, h) path to (n, \cdot) paths, all of which make their outcomes better. As a comparison, we also report the case where individuals can re-optimize at $t = 2$ only (the “Reoptimize at $t = 2$ only” line). The effect is less adverse among low- θ individuals. In the full-response scenario, the preemptive response (switching

⁵⁷In Appendix Table A.5, we show that while the policy reduces the sum of present discounted ex-post utility (-0.54), it increases the period-1 utility slightly (+0.04) because the (h, h) path, which yields more utility in period 2 than 1 due to state dependence in the preference for h (see Table 10), is no longer feasible. In present-discounted terms, the policy reduces period-2 utility the most (-0.56) and reduces the terminal value marginally (-0.02)(includes scaling factor κ). We also show that, when individuals can only reoptimize at $t = 2$ (see discussions below), the sum of present discounted ex-post utility reduces by a larger amount (-0.63) due to the inability to smooth out the impact of the policy across $t = 1, 2$.

away from HS at $t = 1$) particularly worsens the outcomes for low- θ individuals (see also heterogeneity in technological returns in Figure 2).

If the policymaker can target the HS enrolment limit policy to subpopulations, Figure 3 shows that it should be imposed on children with θ above the 54th percentile, whose outcomes increase due to the policy. This makes the average outcome to increase by 2.3 points relative to the baseline. When the policymaker can only target the policy based on the baseline PPVT scores (a proxy for θ ; see Appendix Figure A.11), the overall improvement in outcome is smaller (+1.4 points) but remains substantial relative to the baseline. These results resonate with our discussion in Section 5 that policymakers should engage low- θ individuals with multiple periods of center-based care and high- θ individuals with at least one period of home care.

Appendix Figure A.12 shows the distributional effect separately for individuals whose HS center quality q_{ih} is at the highest quartile and lowest quartile, respectively. The policy tends to benefit the majority of children enrolled in a low-quality HS center; all children with θ above the 29th percentile are expected to gain from the policy. On the contrary, for high-quality HS centers, only children with high θ (above the 69th percentile) are expected to gain from the policy. These findings highlight the importance of considering the quality of HS centers in reforming HS program features. The enrolment limit policy should target low-quality HS centers as prolonged participation in these centers do not yield large returns for most children.

Predictions from the threshold model. What would happen if the sequential threshold model is used to evaluate the policy reform? Appendix Table A.7 reports the results, treating the threshold model as a sequentially static behavioral model and the perceived values as representing utilities. The predicted effects of the reform (columns 2, 4) are substantially different from the predictions from the dynamic structural model (columns 1, 3). In particular, the threshold model fails to predict the preemptive changes in program choices in period 1; in fact the predicted effects of the threshold model are almost identical to the prediction from the dynamic structural model for the partial response scenario when individuals can only reoptimize at $t = 2$ (Table 11, columns 3 and 6). The threshold model, while useful for accounting for program selection and estimating causal returns, can deliver wrong counterfactual results for policymakers when used as a sequential static behavioral model.

7 Conclusion

This paper provides a comprehensive evaluation of the Head Start program, the largest preschool education program for disadvantaged children in the United States. Using the HSIS experimental data, we analyzed the causal effects of sequential program participation when individuals face multiple childcare options in multiple periods. Our empirically-motivated threshold model accounts for selection into different modes of childcare use over time based on approximate decision rules. We estimated the model by exploiting the randomization of HS offer as a key source of identification and imposing a factor structure on the unobservables in characterizing the joint distribution of choices and outcomes.

We computed easily interpretable program returns. Relative to the estimated returns to HS that are available in the literature, our returns are defined on age-specific HS participation and use a specific alternative sequence of childcare use as benchmark comparison. For instance, our estimates of partial returns of HS indicate that, relative to using home care at both age 3 and age 4, one year of HS at either age 3 or age 4 (while fixing the other age to home care) raises the average test score at the end of the pre-school period by 0.30 sd and 0.34 sd, respectively. Two years of HS participation also improves children’s average cognitive outcomes relative to home care at both ages, but the improvement is not as high as the combined partial returns to HS at age 3 and 4. This is because a second year of HS participation lowers the return to HS in the first year, suggesting that returns to a single program spanning multiple periods are not necessarily linearly additive and separable. In addition, we showed important interactions between programs in the context of sequential participation. For instance, compared with two years of HS, we found that children can achieve an even better average outcome if they are assigned to HS at age 3 followed by other center care at age 4. However, this does not necessarily imply that a combination of HS and other center care always yields better outcomes. Reversing the sequence – other center care at age 3 followed by HS at age 4 – instead generates a very low average return.

Our estimates imply that policy makers should engage low-ability individuals with prolonged exposure to center-based care and high-ability individuals with some home care. One central question then is how to redesign HS to induce more individuals of heterogenous ability to choose their “socially desirable” program sequence. We analyzed a counterfactual policy reform which had not been considered before but deemed important in the HSIS evaluation report: limiting the maximum duration of

participation in HS to one year only (but keeping open the option to enrol at age 3 or 4). By specifying and estimating a dynamic structural model that involves a forward-looking agent who knows the imposed time limit from the beginning, we find that the policy reduces the average outcome slightly but worsens inequality. This is driven by pre-emptive responses at age 3 to “bank up” HS eligibility for possible enrolment at age 4, which is more common among low-ability children and worsens their outcomes. By way of comparison, had the policy not been known to individuals until age 4, which rules out pre-emptive responses, the policy becomes more attractive in that the average outcome increases slightly. These results highlight the importance of the specification of information set and modeling forward-looking behavior, which is a strength of the dynamic structural model over the sequential threshold model considered in the first part of the paper.

Our modeling framework can be applied to other contexts, especially in the analysis of RCTs that span multiple periods and that rely on a “encouragement design” (e.g., Duflo, 2001), where the randomization is based on providing initial access to a particular program. Similar to HS, in this type of experiments the duration of program participation may be endogenous, and participation in one program may affect decisions to participate in other programs in subsequent periods. Therefore, the experimental-oriented parameters may be hard to interpret and become less relevant for policy in the presence of multiple outside options and dynamic selection.

References

- ABBRING, J. H., AND J. J. HECKMAN (2007): “Econometric evaluation of social programs, part III: Distributional treatment effects, dynamic treatment effects, dynamic discrete choice, and general equilibrium policy evaluation,” *Handbook of econometrics*, 6, 5145–5303.
- ABBRING, J. H., AND G. J. VAN DEN BERG (2003): “The nonparametric identification of treatment effects in duration models,” *Econometrica*, 71(5), 1491–1517.
- AIZER, A., AND F. CUNHA (2012): “The production of human capital: Endowments, investments and fertility,” Discussion paper, National Bureau of Economic Research.
- ALMOND, D., AND B. MAZUMDER (2013): “Fetal origins and parental responses,” *Annual Review of Economics*, 5(1), 37–56.
- ANGRIST, J. D., G. W. IMBENS, AND D. B. RUBIN (1996): “Identification of causal effects using instrumental variables,” *Journal of the American statistical Association*, 91(434), 444–455.
- BERNDT, E. K., B. H. HALL, R. E. HALL, AND J. HAUSMAN (1974): “Estimation and inference in nonlinear structural models,” *Annals of Economic and Social Measurement*, 3(4), 653–665.
- BHALOTRA, S. R., AND A. VENKATARAMANI (2015): “Shadows of the captain of the men of death: Early life health interventions, human capital investments, and institutions,” *Human Capital Investments, and Institutions (August 8, 2015)*.
- BITLER, M. P., H. W. HOYNES, AND T. DOMINA (2014): “Experimental evidence on distributional effects of Head Start,” Discussion paper, National Bureau of Economic Research.
- BONHOMME, S., AND J.-M. ROBIN (2009): “Consistent noisy independent component analysis,” *Journal of Econometrics*, 149(1), 12–25.
- CAMIRAND LEMYRE, F., R. J. CARROLL, AND A. DELAIGLE (2022): “Semiparametric Estimation of the Distribution of Episodically Consumed Foods Measured With Error,” *Journal of the American Statistical Association*, 117, 469–481.
- CARNEIRO, P., AND R. GINJA (2014): “Long-term impacts of compensatory preschool on health and behavior: Evidence from Head Start,” *American Economic Journal: Economic Policy*, 6(4), 135–73.
- CARNEIRO, P., K. HANSEN, AND J. HECKMAN (2003): “Estimating distributions of treatment effects with an application to the returns to schooling and measurement of the effects of uncertainty on college choice,” *International Economic Review*, 44(2), 361–422.
- CHAN, M. K. (2017): “Welfare dependence and self-control: An empirical analysis,” *The Review of Economic Studies*, 84(4), 1379–1423.

- CHAN, M. K., AND K. LIU (2018): “Life-cycle and intergenerational effects of child care reforms,” *Quantitative Economics*, 9(2), 659–706.
- (forthcoming): “Changing families: family relationships, parental decisions and child development,” *Journal of Labor Economics*.
- CUNHA, F., AND J. HECKMAN (2007): “The technology of skill formation,” *American Economic Review*, 97(2), 31–47.
- CUNHA, F., J. HECKMAN, AND S. SCHENNACH (2010): “Estimating the Technology of Cognitive and Noncognitive Skill Formation,” *Econometrica*, 78(3), 883–931.
- DE HAAN, M., AND E. LEUVEN (2020): “Head start and the distribution of long-term education and labor market outcomes,” *Journal of Labor Economics*, 38(3), 727–765.
- DEMING, D. (2009): “Early childhood intervention and life-cycle skill development: Evidence from Head Start,” *American Economic Journal: Applied Economics*, 1(3), 111–34.
- DUFLO, E. (2001): “Schooling and Labor Market Consequences of School Construction in Indonesia: Evidence from an Unusual Policy Experiment,” *American Economic Review*, 91, 795–813.
- EBERWEIN, C., J. C. HAM, AND R. J. LALONDE (1997): “The impact of being offered and receiving classroom training on the employment histories of disadvantaged women: Evidence from experimental data,” *The Review of Economic Studies*, 64(4), 655–682.
- EVDOKIMOV, K., AND H. WHITE (2012): “Some Extensions of A Lemma of Kotlarski,” *Econometric Theory*, 28, 925–932.
- FANG, H., AND D. SILVERMAN (2009): “Time-Inconsistency and Welfare Program Participation: Evidence from the NLSY,” *International Economic Review*, 50, 1043–1077.
- FERNALD, A., V. A. MARCHMAN, AND A. WEISLEDER (2013): “SES differences in language processing skill and vocabulary are evident at 18 months,” *Developmental science*, 16(2), 234–248.
- FRUEHWIRTH, J. C., S. NAVARRO, AND Y. TAKAHASHI (2016): “How the timing of grade retention affects outcomes: Identification and estimation of time-varying treatment effects,” *Journal of Labor Economics*, 34(4), 979–1021.
- GOFF, L., O. MALAMUD, C. POP-ELECHES, AND M. URQUIOLA (2022): “Interactions between Family and School Environments: Access to Abortion and Selective Schools,” *mimeo*.
- GRIFFEN, A. S. (2019): “Evaluating the effects of childcare policies on children’s cognitive development and maternal labor supply,” *Journal of Human Resources*, 54(3), 604–655.
- GRIFFEN, A. S., AND P. E. TODD (2017): “Assessing the performance of nonexperimental estimators for evaluating Head Start,” *Journal of Labor Economics*, 35(S1), S7–S63.

- HALL, P., AND A. MEISTER (2007): “A Ridge-Parameter Approach to Deconvolution,” *Annals of Statistics*, 35(4), 1535–1558.
- HAM, J. C., AND R. J. LALONDE (1996): “The effect of sample selection and initial conditions in duration models: Evidence from experimental data on training,” *Econometrica*, 64, 175–206.
- HANSEN, K., J. HECKMAN, AND K. MULLEN (2004): “The Effect of Schooling and Ability on Achievement Test Scores,” *Journal of Econometrics*, 121, 39–98.
- HART, B., AND T. R. RISLEY (1995): *Meaningful differences in the everyday experience of young American children*. Paul H Brookes Publishing.
- HECKMAN, J., N. HOHMANN, J. SMITH, AND M. KHOO (2000): “Substitution and dropout bias in social experiments: A study of an influential social experiment,” *The Quarterly Journal of Economics*, 115(2), 651–694.
- HECKMAN, J. J., J. E. HUMPHRIES, AND G. VERAMENDI (2016): “Dynamic treatment effects,” *Journal of econometrics*, 191(2), 276–292.
- (2018): “Returns to education: The causal effects of education on earnings, health, and smoking,” *Journal of Political Economy*, 126(S1), S197–S246.
- HECKMAN, J. J., AND S. NAVARRO (2007): “Dynamic discrete choice and dynamic treatment effects,” *Journal of Econometrics*, 136(2), 341–396.
- HECKMAN, J. J., AND J. A. SMITH (1998): “Evaluating the Welfare State,” *Econometrics and Economic Theory in the Twentieth Century: The Ragnar Frisch Centennial Symposium*, edited by S. Strom.
- HECKMAN, J. J., S. URZUA, AND E. VYTLACIL (2006): “Understanding instrumental variables in models with essential heterogeneity,” *The Review of Economics and Statistics*, 88(3), 389–432.
- JOHNSON, R. C., AND C. K. JACKSON (2019): “Reducing inequality through dynamic complementarity: Evidence from Head Start and public school spending,” *American Economic Journal: Economic Policy*, 11(4), 310–49.
- KEANE, M. (1994): “A computationally practical simulation estimator for panel data,” *Econometrica: Journal of the Econometric Society*, pp. 95–116.
- KEANE, M. P., AND K. I. WOLPIN (2001): “The Effect of Parental Transfers and Borrowing Constraints on Educational Attainment,” *International Economic Review*, 42(4), 1051–1103.
- KIRKEBOEN, L. J., E. LEUVEN, AND M. MOGSTAD (2016): “Field of study, earnings, and self-selection,” *The Quarterly Journal of Economics*, 131(3), 1057–1111.
- KLEIN, R., AND R. SPADY (1993): “An Efficient Semiparametric Estimator for Binary Response Models,” *Econometrica*, 61(2), 387–421.

- KLINE, P., AND C. R. WALTERS (2016): “Evaluating public programs with close substitutes: The case of Head Start,” *The Quarterly Journal of Economics*, 131(4), 1795–1848.
- KOTLARSKI, I. (1967): “On characterizing the gamma and the normal distribution,” *Pacific Journal of Mathematics*, 20(1), 69–76.
- LEE, L.-F. (1995): “Semiparametric Maximum Likelihood Estimation of Polychotomous and Sequential Choice Models,” *Journal of Econometrics*, 65, 381–428.
- LEWBEL, A. (1997): “Semiparametric Estimation of Location and Other Discrete Choice Moments,” *Econometric Theory*, 13, 32–51.
- LUDWIG, J., AND D. L. MILLER (2007): “Does Head Start improve children’s life chances? Evidence from a regression discontinuity design,” *The Quarterly journal of economics*, 122(1), 159–208.
- PUMA, M., ET AL. (2010): “Head Start Impact Study. Final Report.,” *Administration for Children & Families, U.S. Department of Health and Human Services*.
- REIERSØL, O. (1950): “Identifiability of a linear relation between variables which are subject to error,” *Econometrica: Journal of the Econometric Society*, pp. 375–389.
- RODRÍGUEZ, J., F. SALTIEL, AND S. S. URZÚA (2018): “Dynamic Treatment Effects of Job Training,” Discussion paper, National Bureau of Economic Research.
- ROSENBAUM, S. (1961): “Moments of a truncated bivariate normal distribution,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 23(2), 405–408.
- ROSSIN-SLATER, M., AND M. WÜST (2020): “What is the added value of preschool for poor children? Long-term and intergenerational impacts and interactions with an infant health intervention,” *American Economic Journal: Applied Economics*, 12(3), 255–86.
- SHAUL, M. S., B. WARD-ZUKERMAN, S. EDMONDSON, L. MOY, C. MORIARITY, AND E. PICYK (2003): “Head Start: Curriculum Use and Individual Child Assessment in Cognitive and Language Development. Report to Congressional Requesters.,” .
- STOKER, T. (1986): “Consistent Estimation of Scaled Coefficients,” *Econometrica*, 54(6), 1461–1481.
- TALLIS, G. M. (1961): “The moment generating function of the truncated multi-normal distribution,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 23(1), 223–229.
- WALTERS, C. R. (2015): “Inputs in the production of early childhood human capital: Evidence from Head Start,” *American Economic Journal: Applied Economics*, 7(4), 76–102.

Table 1: Proportion of Children in Different Childcare Programs

Program Sequence		Proportion (%)			
$t = 1$	$t = 2$	$Z = 0$		$Z = 1$	
Head Start (h)	Head Start (h)	10.7	<i>74.3</i>	62.0	<i>72.9</i>
	Other Center Care (c)	2.7	<i>18.8</i>	18.1	<i>21.3</i>
	Home Care (n)	1.0	<i>6.9</i>	5.0	<i>5.9</i>
Total (h)		14.4%		85.1%	
Other Center Care (c)	Head Start (h)	9.3	<i>37.6</i>	1.2	<i>21.9</i>
	Other Center Care (c)	13.9	<i>56.1</i>	3.7	<i>68.8</i>
	Home Care (n)	1.6	<i>6.4</i>	0.5	<i>9.4</i>
Total (c)		24.8%		5.4%	
Home Care (n)	Head Start (h)	30.2	<i>49.8</i>	3.0	<i>31.2</i>
	Other Center Care (c)	18.6	<i>30.7</i>	2.8	<i>29.5</i>
	Home Care (n)	11.9	<i>19.6</i>	3.7	<i>39.3</i>
Total (n)		60.7%		9.5%	

Notes : This table reports the percentages of children selecting different combinations of childcare programs, separately for the control ($Z = 0$) and for the treatment group ($Z = 1$). Numbers in italic show the proportions selecting the program at $t = 2$ (age 4) *conditional* on the choice made at $t = 1$ (age 3). “Total (•)” rows show the overall proportion making each choice at $t = 1$ (i.e., the sum of the three rows above).

Table 2: Experimental Impact of Head Start: Reduced Form Estimates

	(1)	(2)
Z	0.029***	-0.020
	(0.011)	(0.023)
$Z \times$ other-center care in $t = 2$		0.109***
		(0.026)
$Z \times$ HS in $t = 2$		0.033
		(0.023)
Constant	3.281***	3.287***
	(0.034)	(0.034)
Observations	1,876	1,876
R-squared	0.197	0.209

Notes : $Z = 1$ if the child is assigned to the treatment group. Each column is from a separate regression of the outcome variable (average of the three test scores – PPVT, WJIII Letter-Word and WJIII Spelling– at the end of the age-4 year) on a dummy variable for having been randomly assigned to HS (column (1)) and its interaction with indicator variables for the type of childcare use at $t = 2$ (column (2)). We include controls for household size, number of siblings, dummies for whether the child is female, black, hispanic, use English as home language, living in urban area, living with both parents, in need of special education, child of a teen mother, child of a mother who never married or is separated, child of a mother with high school or more than high school. Standard errors are in parentheses. We adjust the PPVT and WJIII test scores by dividing the raw score by 100. Significance level (t-test for the parameter being 0): *** 1%, ** 5%, * 10%

Table 3: Within- and Between-Program Cross Returns

Program at $t = 1$	Program at $t = 2$	
	Head Start (h)	Other Center Care (c)
Head Start (h)	$(Y_i^{hh} - Y_i^{nh}) - (Y_i^{hn} - Y_i^{nn})$	$(Y_i^{hc} - Y_i^{nc}) - (Y_i^{hn} - Y_i^{nn})$
Other Center Care (c)	$(Y_i^{ch} - Y_i^{nh}) - (Y_i^{cn} - Y_i^{nn})$	$(Y_i^{cc} - Y_i^{nc}) - (Y_i^{cn} - Y_i^{nn})$

Table 4: List of Covariates and Exclusion Variables

	Period 1 choice		Period 2 choice		Potential outcome
	HS	Other center	HS	Other center	
<i>A. Exclusion restrictions</i>					
Experimental offer (Z_i)	x				
$Z_i \times$ HS center quality	x				
$Z_i \times$ HS center transport	x				
HS center transport	x		x		
Other center transport		x		x	
<i>B. Other covariates</i>					
HS center quality	x		x		x
Other center quality		x		x	x
<i>C. Unobserved heterogeneity</i>					
Child skill endowment (θ_i)	x	x	x	x	x

Note: Boldface **x** indicates that the variable serves as an exclusion variable.

Table 5: Goodness of Fit of the Sequential Threshold Model

		Proportions			Outcome		
		Program $t = 2$					
		h	c	n	h	c	n
Program $t = 1$		Control Group			Control Group		
Simulation	h	10.4	3.0	0.9	3.32	3.41	3.34
Data	h	<i>10.7</i>	<i>2.7</i>	<i>1.0</i>	<i>3.29</i>	<i>3.46</i>	<i>3.56</i>
Simulation	c	8.5	15.0	1.8	3.33	3.43	3.48
Data	c	<i>9.3</i>	<i>13.9</i>	<i>1.6</i>	<i>3.32</i>	<i>3.43</i>	<i>3.51</i>
Simulation	n	27.5	18.7	14.1	3.32	3.38	3.26
Data	n	<i>30.2</i>	<i>18.6</i>	<i>11.9</i>	<i>3.34</i>	<i>3.36</i>	<i>3.25</i>
		Treatment Group			Treatment Group		
Simulation	h	62.0	17.8	5.2	3.36	3.45	3.39
Data	h	<i>62.0</i>	<i>18.1</i>	<i>5.0</i>	<i>3.35</i>	<i>3.45</i>	<i>3.35</i>
Simulation	c	1.7	3.3	0.4	3.36	3.46	3.47
Data	c	<i>1.2</i>	<i>3.7</i>	<i>0.5</i>	<i>3.41</i>	<i>3.44</i>	<i>3.74</i>
Simulation	n	4.4	2.9	2.4	3.29	3.36	3.23
Data	n	<i>3.0</i>	<i>2.8</i>	<i>3.7</i>	<i>3.29</i>	<i>3.36</i>	<i>3.23</i>

Notes : Numbers in italics are the population equivalents (computed directly from the data). Each cell shows the proportion or the average outcome (average between PPVT score, WJIII Letter-Word score and WJIII Spelling score at the end of age-4 year) for the group of children who at $t = 1$ enrol in the program reported in the second column and at $t = 2$ enrol in the program reported in the second row, separately for the treatment and for the control group. h is HS, c is other center care, n is home care. We perform tests of equality between model-predicted and empirical moments, separately for choice proportions and outcomes and by treatment and control group. The p-values for a test of the null hypothesis of equality between the data and the choice proportions implied by the estimated model (Chi-square test, $df=8$) are 0.66 ($Z = 0$) and 0.04 ($Z = 1$); the p-values for the null hypothesis of equality in the 9 mean outcomes between data and model predictions (F test) are 0.28 ($Z = 0$) and 0.42 ($Z = 1$).

Table 6: Estimates of the Parameters of the Choice Equations from the Sequential Threshold Model

Covariate	Perceived Value					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept (ψ_1^h, ψ_1^c)	-0.575	(0.148)	***	-1.153	(0.243)	***
Factor θ (λ_1^h, λ_1^c)	0.218	(0.322)		1.965	(0.506)	***
Treatment group Z	1.895	(0.152)	***			
HS center quality $\times Z$	0.315	(0.068)	***			
HS center transport $\times Z$	0.148	(0.068)	**			
$\theta \times Z$	0.343	(0.394)				
HS center quality	-0.199	(0.055)	***			
HS center transport	-0.145	(0.049)	***			
Other center quality				0.126	(0.067)	*
Other center transport				-0.092	(0.066)	
<i>Period 2 choice parameters:</i>						
Intercept (ψ_2^h, ψ_2^c)	0.613	(0.104)	***	-0.063	(0.147)	
Period 1 in h (ψ_2^{hh}, ψ_2^{hc})	0.930	(0.103)	***	0.375	(0.159)	**
Period 1 in c (ψ_2^{ch}, ψ_2^{cc})	0.731	(0.172)	***	1.723	(0.307)	***
Factor θ (λ_2^h, λ_2^c)	-0.317	(0.252)		0.736	(0.424)	*
HS center quality	0.076	(0.031)	**			
HS center transport	0.044	(0.033)				
Other center quality				0.189	(0.056)	***
Other center transport				0.214	(0.055)	***
<i>Other parameters:</i>						
σ_{ν^h}	1.000	-				
σ_{ν^c}	2.079	(0.334)	***			
ρ_{hc}	0.748	(0.135)	***			

Notes : N=1876. Log-likelihood = -3035.72. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of the age-4 year as the outcome. Alternative-specific choice shocks assumed to follow a bivariate normal distribution with zero means, variance of 1 (normalization) for ν_{it}^h , variance of $\sigma_{\nu^c}^2$ for ν_{it}^c , and correlation coefficient ρ_{hc} . Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table 7: Estimates of the Parameters of the Potential Outcome and Measurement Equations from the Sequential Threshold Model

<i>A. Potential outcome equations</i>						
	Intercept			Factor		
Potential outcomes:						
Y^{nn}	3.268	(0.018)	***	0.974	(0.089)	***
$Y^{hn} - Y^{nn}$	0.082	(0.032)	***	0.227	(0.134)	*
$Y^{nh} - Y^{nn}$	0.092	(0.024)	***	-0.108	(0.108)	
$Y^{nc} - Y^{nn}$	0.099	(0.025)	***	0.041	(0.106)	
$Y^{hh} - Y^{nn}$	0.098	(0.02)	***	-0.219	(0.087)	**
$Y^{hc} - Y^{nn}$	0.141	(0.024)	***	-0.175	(0.104)	*
$Y^{ch} - Y^{nn}$	0.049	(0.033)		-0.104	(0.159)	
$Y^{cc} - Y^{nn}$	0.098	(0.026)	***	-0.227	(0.104)	**
Other covariates						
1 period quality HS	0.006	(0.008)				
2 periods quality HS	0.022	(0.007)	***			
1 period quality other center	0.004	(0.011)				
2 periods quality other center	0.023	(0.017)				
<i>B. Measurement equations</i>						
	Intercept			Factor		
PPVT	2.299	(0.009)	***	1.000	-	
WJIII Spelling	3.348	(0.007)	***	0.500	(0.034)	***
WJIII Letter-Word	2.939	(0.006)	***	0.555	(0.033)	***
<i>C. Factor and error SD</i>						
Factor	0.223	(0.009)	***			
Outcome	0.158	(0.005)	***			
Baseline measurements:						
PPVT	0.293	(0.006)	***			
WJIII Spelling	0.243	(0.004)	***			
WJIII Letter-Word	0.185	(0.004)	***			

Notes : N=1876. Log-likelihood = -3035.72. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of age-4 year as the outcome. The intercept and the factor loading of the outcome equation differ depending on the program sequence (see equation 9), and are all relative to the parameters of the outcome equation for the nn sequence. h is for HS, c is for other center care, n is for home care. Baseline measurements are test scores measured around the time of randomization (in Fall 2002). The factor loading in the measurement equation for the baseline PPVT test score is normalized to one ($\gamma^{pp}=1$). Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n) , to avoid overfitting we restrict $\alpha_t^{cn} = \alpha_t^{hn}$ and $\gamma_t^{cn} = \gamma_t^{hn}$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): ***1%, ** 5%, *10%.

Table 8: Optimal Program Assignment: Predicted Outcomes

	Population	θ Distribution		
	Average	p(10)	p(50)	p(90)
	(1)	(2)	(3)	(4)
Individual choice, $Z = 1$	3.371	3.141	3.371	3.597
Optimal path when observing:				
Factor (θ_i)	3.424	3.181	3.410	3.689
PPVT	3.413	3.179	3.405	3.650
WJIII Spelling	3.410	3.179	3.407	3.639
WJIII Letter-Word	3.412	3.178	3.406	3.650

Notes : Given the estimated model, we simulate 9 potential outcomes (corresponding to 9 possible program combinations) for 30 random draws of the unobserved factor and outcome shocks for each individual in the sample (following equation (9)). Thus the distributions of center characteristics, factor and outcome shocks are identical between the simulated data and the sample. Outcome is the average of PPVT score, WJIII Letter-Word score and WJIII Spelling score at the end of age-4 year. Individual choice refers to the predicted outcome when all children are assigned to the treatment group ($Z = 1$) and choose program sequence according to the estimated sequential threshold model. Optimal path refers to the assignment of children to their optimal program sequence which maximizes the expected outcome (see Appendix D for details). Column (1) is the average predicted outcome in the population. Columns (2), (3) and (4) refer to the outcomes for children with baseline ability corresponding to the 10th, the 50th and the 90th percentile of the distribution of the unobserved factor.

Table 9: Proportion, Mean of the Unobserved Factor and Returns for Compliance Groups

Compliance type	Proportion	θ	Return
	(1)	(2)	(3)
temporary c-compliers	7.92	0.044	0.049
path c-compliers	11.94	0.061	0.013
temporary n-compliers	21.45	-0.033	0.021
path n-compliers	29.39	0.001	0.049
All compliers	70.7	0.006	0.034
c-compliers	19.86	0.054	0.028
n-compliers	50.84	-0.013	0.037

Notes : Columns (1) and (2) report proportions and the average baseline ability (θ) of each compliance group. Column (3) reports the average return (in terms of the average of PPVT score, WJIII Letter-Word score and WJIII Spelling score at the end of age-4 year) for the specific compliance type. See Section 5.3 for definitions of the compliance groups. The row “*All compliers*” refers to all the compliers from the HSIS experiment (who select h at $t = 1$ if in the treatment group and either c or n at $t = 1$ if in the control group).

Table 10: Estimates of the Parameters of the Utility Function from the Structural Model

Covariate	Flow Utility					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept ($\tilde{\psi}_1^h, \tilde{\psi}_1^c$)	-1.481	(0.204)	***	-2.681	(0.389)	***
Factor θ ($\tilde{\lambda}_1^h, \tilde{\lambda}_1^c$)	0.628	(0.298)	**	1.984	(0.458)	***
Treatment group Z	1.859	(0.150)	***			
HS center transport $\times Z$	0.153	(0.064)	**			
$\theta \times Z$	0.437	(0.372)				
HS center quality $\times (1-Z)$	-0.239	(0.057)	***			
HS center transport	-0.152	(0.046)	***			
Other center transport				-0.174	(0.071)	**
<i>Period 2 choice parameters:</i>						
Intercept ($\tilde{\psi}_2^h, \tilde{\psi}_2^c$)	0.326	(0.123)	***	-0.445	(0.177)	**
Period 1 in h ($\tilde{\psi}_2^{hh}, \tilde{\psi}_2^{hc}$)	1.183	(0.121)	***	0.530	(0.162)	***
Period 1 in c ($\tilde{\psi}_2^{ch}, \tilde{\psi}_2^{cc}$)	1.199	(0.208)	***	2.195	(0.316)	***
Factor θ ($\tilde{\lambda}_2^h, \tilde{\lambda}_2^c$)	0.718	(0.361)	**	1.561	(0.453)	***
HS center transport	0.046	(0.030)				
Other center transport				0.223	(0.054)	***
<i>Other parameters:</i>						
$\tilde{\sigma}_{\nu^h}$	1.000	-				
$\tilde{\sigma}_{\nu^c}$	2.203	(0.314)	***			
$\tilde{\rho}_{hc}$	0.783	(0.123)	***			
κ	3.512	(0.987)	***			

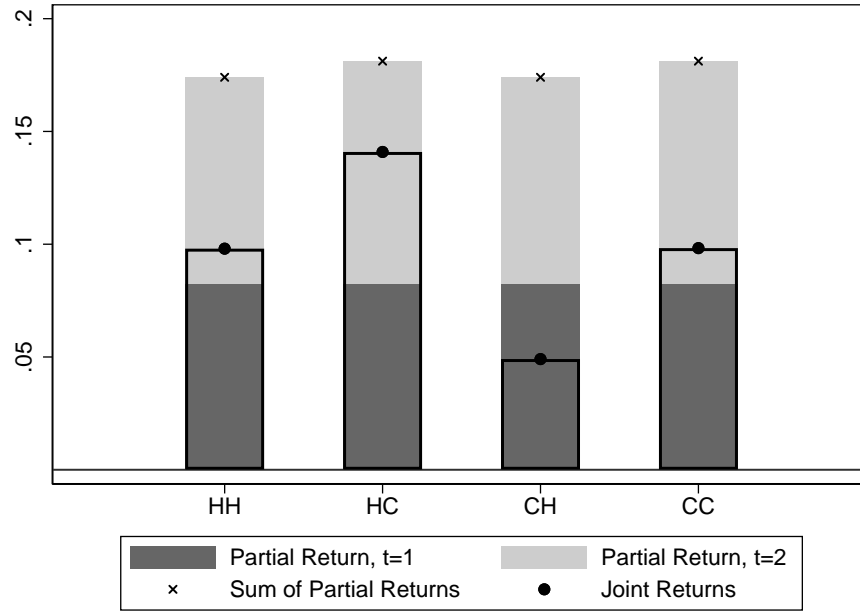
Notes : N=1876. Log-likelihood = -3048.47. Flow utilities are specified in Appendix Section E. The technological parameters in the baseline measurement equations, outcome equations and the factor distribution are fixed at the estimated values from the sequential threshold model (see Table 7). κ is the terminal value scaling factor. The discount factor δ is set to 1. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table 11: Effects of Limiting Head Start Enrolment to One Period Only

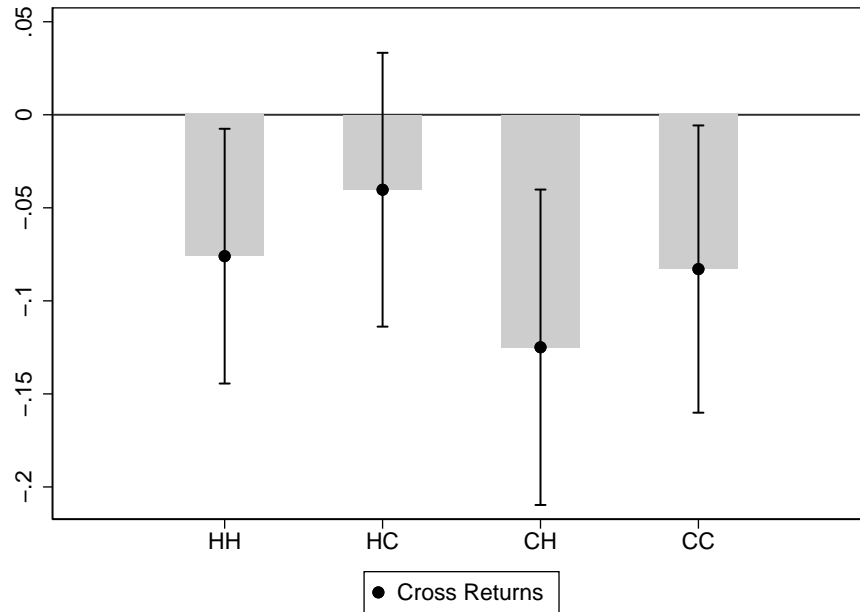
	Program choice (%)			Average Outcome		
	Baseline	Counterfactual		Baseline	Counterfactual	
	(1)	Full Resp.	Partial Resp.	(4)	Full Resp.	Partial Resp.
	(1)	(2)	(3)	(4)	(5)	(6)
Choices at t=1						
<i>h</i>	85.2	-25.5	0.0			
<i>c</i>	5.3	+6.7	0.0			
<i>n</i>	9.4	+18.8	0.0			
Choices at t=2						
<i>h</i>	68.1	-51.0	-61.9			
<i>c</i>	23.9	+24.4	+29.0			
<i>n</i>	8.0	+26.5	+32.9			
Program sequence						
<i>h, h</i>	61.9	-61.9	-61.9	3.35	-	-
<i>h, c</i>	18.0	+15.0	+29.0	3.44	-0.01	-0.03
<i>h, n</i>	5.3	+21.4	+32.9	3.45	-0.07	-0.10
<i>c, h</i>	1.8	+2.3	0.0	3.37	-0.04	0.00
<i>c, c</i>	3.2	+4.0	0.0	3.43	-0.03	0.00
<i>c, n</i>	0.4	+0.4	0.0	3.51	-0.06	0.00
<i>n, h</i>	4.4	+8.7	0.0	3.30	-0.01	0.00
<i>n, c</i>	2.7	+5.4	0.0	3.36	-0.02	0.00
<i>n, n</i>	2.3	+4.7	0.0	3.19	-0.02	0.00
Overall mean				3.371	3.366	3.376
Overall s.d.				0.245	0.272	0.273

Notes : All scenarios are simulated from the structural model with 30 paths per treatment-group individual. The “Baseline” columns show the proportion of the population making each choice and their average outcome under the baseline scenario, i.e. when individuals are allowed to select two periods of HS. The “Counterfactual” columns report the difference in the proportion of population selecting each childcare option and in the outcomes between the baseline scenario and the counterfactuals: “Full resp.” (full response) allows individuals to select at most one period in HS and assumes that they can reoptimize both at $t = 1$ and $t = 2$; “Partial resp.” (partial response) assumes individuals at $t = 1$ do not know about the HS enrolment limit, so they can reoptimize at $t = 2$ only.

Figure 1: Average Partial, Joint and Cross Returns from Sequential Program Participation



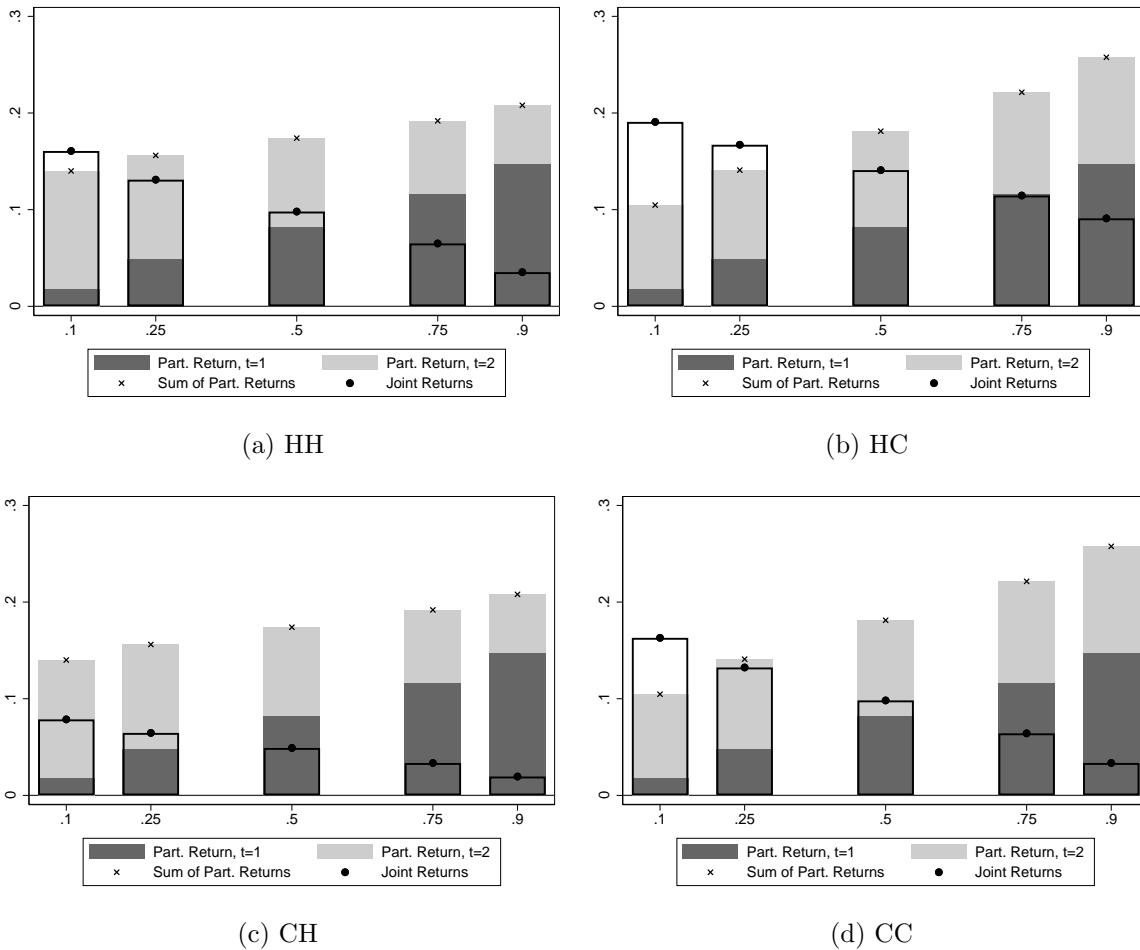
(a) Partial and Joint Returns



(b) Cross-Returns (with 95% C.I.)

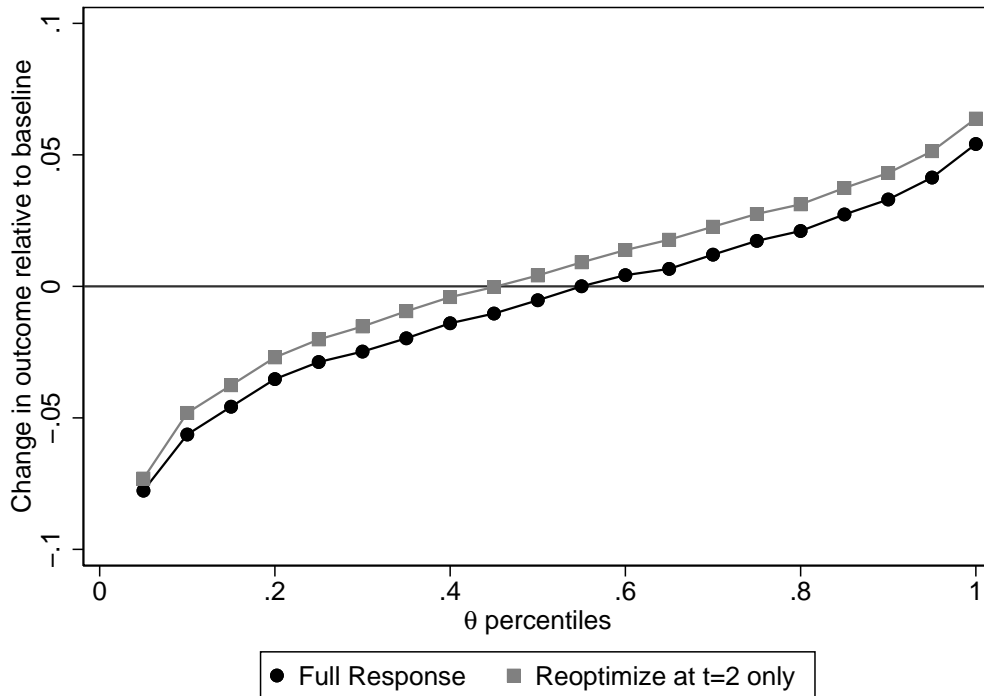
Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. In Panel (a) each bar represents the joint and partial return of different program sequences. In Panel (b) each bar represents the cross-returns, i.e. the difference between the joint return and the sum of the partial returns at $t=1$ and at $t=2$, with the respective 95% confidence intervals. The first bar (HH) shows the return of receiving HS both at $t=1$ and at $t=2$. The second bar (HC) shows the return of receiving HS at $t=1$ and other center care at $t=2$. The third bar (CH) shows the return of receiving other center care at $t=1$ and HS at $t=2$. The fourth bar (CC) shows the return of receiving other center care in both periods. All these returns use two periods in home care as the baseline program sequence for comparison (see Section 3 for the exact definition of partial, joint and cross-returns).

Figure 2: Partial, Joint and Cross Returns from Sequential Program Participation, Different Values of θ



Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. Within each panel, the different bars represent the joint and partial returns of the same program sequence, but for different levels of baseline ability θ (“Part. Return” in label stands for “Partial Returns”). In particular, we report the returns for the 10th, the 25th, the 50th, the 75th and the 90th percentile of the θ distribution. Panel (a) is the return for two periods in HS (HH), panel (b) for one period in HS followed by other center care (HC), panel (c) for one period in other center care followed by HS (CH) and panel (d) for two periods in other center care (CC).

Figure 3: Effect of Limiting Head Start Enrolment to One Period Only for Different Values of θ

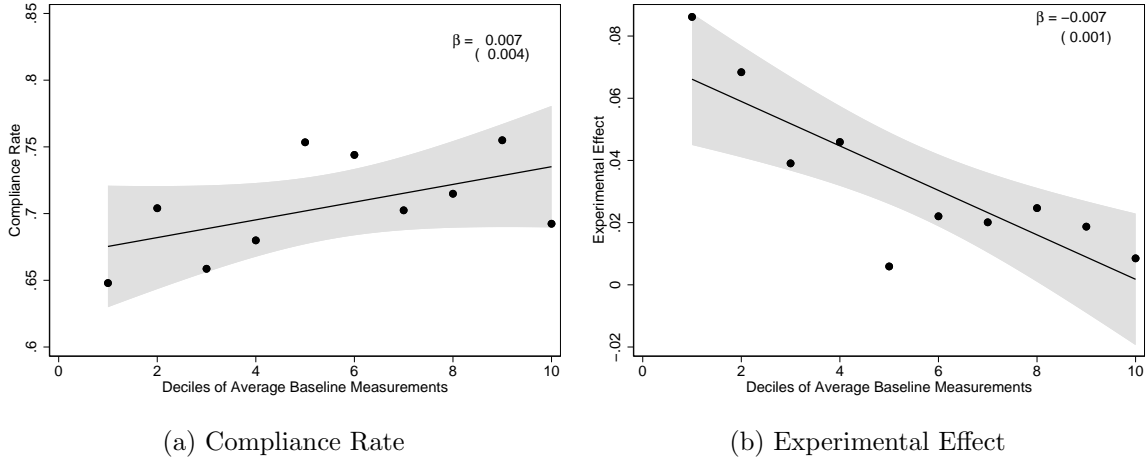


Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. Each curve represents deviations in outcomes from the baseline scenario in the treatment group, e.g., zero means no effect of the HS enrolment limit on the outcome. On each curve, each marker represents one ventile of the distribution of θ (i.e., baseline ability. E.g. the leftmost marker represents individuals with the lowest 5 percentiles of θ , the second leftmost marker represents individuals between the 5-10 percentiles, and so on). All scenarios are simulated from the structural model with 30 paths per treatment-group individual. “Full response”: allow individuals to reoptimize at t=1 and t=2. “Reoptimize at t=2 only”: assume individuals at t=1 do not know about the HS enrolment limit, so they can reoptimize at t=2 only.

ONLINE APPENDIX

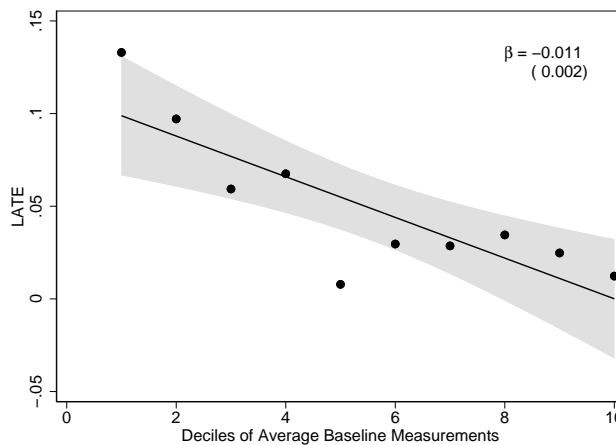
A Additional Results

Figure A.1: Compliance Rate and Experimental Effect by Decile of the Average Baseline Test Scores



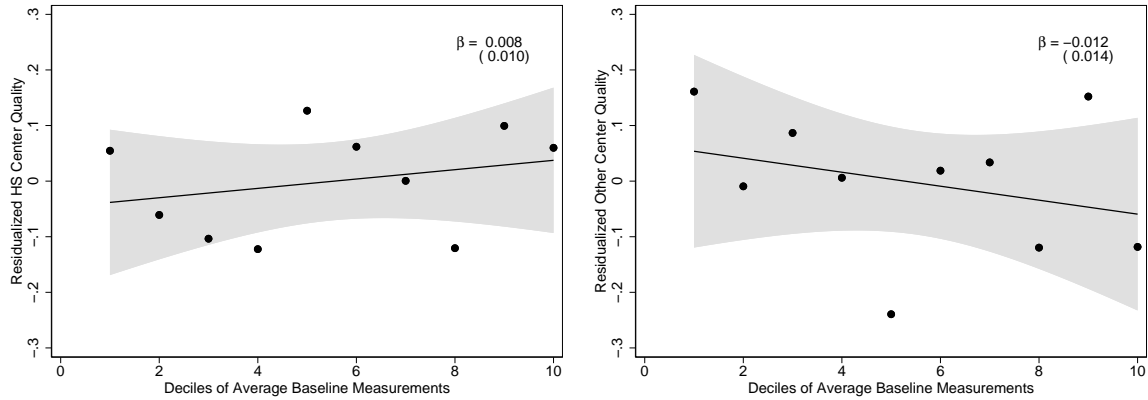
Notes : The average baseline measurements is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the time of randomization (in Fall 2002). We split the sample into 10 equal-sized bins based on the average baseline test score. Panel (a) shows the compliance rate in each test score bin, while panel (b) shows the experimental impact of the HS offer in each bin. Both figures report the line fitting the 10 points with its confidence interval while parameter β refers to the estimated slope of this line (s.e. is reported in parenthesis).

Figure A.2: Local Average Treatment Effect by Decile of the Average Baseline Test Scores



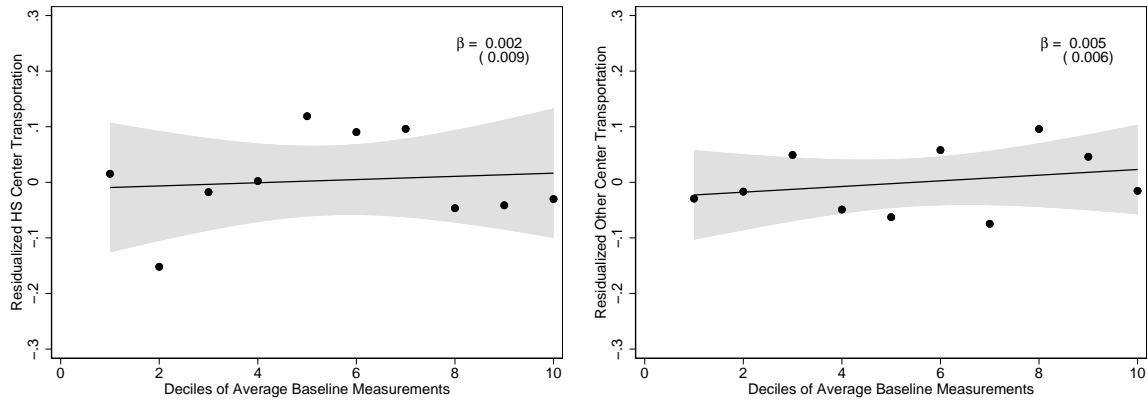
Notes : The average baseline test score is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the time of randomization (in Fall 2002). This average is divided in 10 equally sized bins. We also report the line fitting the 10 points with its confidence interval while parameter β refers to the estimated slope of this line (s.e. is reported in parenthesis).

Figure A.3: Residualized Center Quality and Transportation Availability by Decile of the Average Baseline Test Scores



(a) HS Center Quality

(b) Other Center Quality

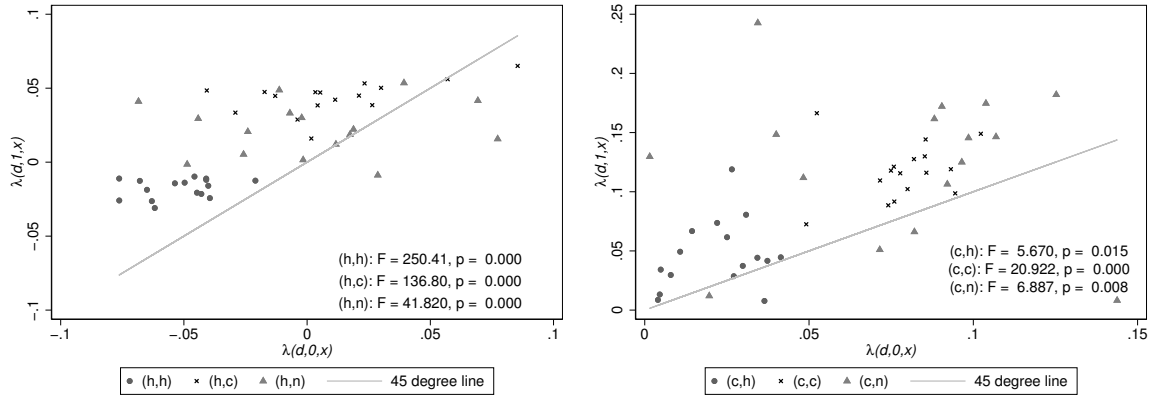


(c) HS Center Transportation

(d) Other Center Transportation

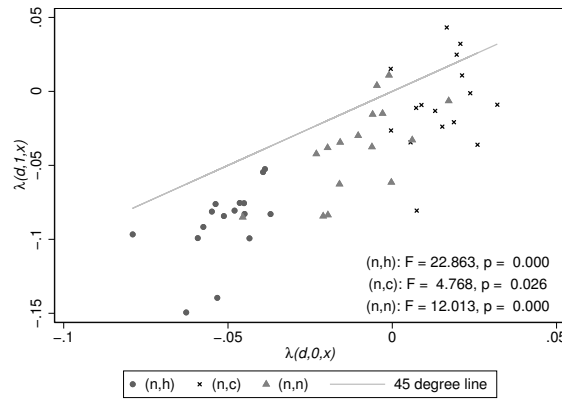
Notes : The average baseline test score is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the time of randomization (in Fall 2002). We split the sample into 10 equal-sized bins based on the average baseline test score. Panel (a) shows the residualized quality of the assigned HS center in each test score bin, panel (b) shows the residualized quality of the other center in each bin, panel (c) shows the residualized transportation availability of the assigned HS center in each bin and panel (d) shows the residualized transportation availability of the other center in each bin. Each variable is defined by taking the residuals of its regression on household size, number of siblings, dummies for whether the child is female, black, hispanic, use English as home language, living in urban area, living with both parents, in need of special education, child of a teen mother, child of a mother who never married or is separated, child of a mother with high school or more than high school. Each figure report the line fitting the 10 points with its confidence interval while parameter β refers to the estimated slope of this line (s.e. is reported in parenthesis)

Figure A.4: Difference between $\lambda(d, 1, x)$ and $\lambda(d, 0, x)$



(a) h at $t = 1$

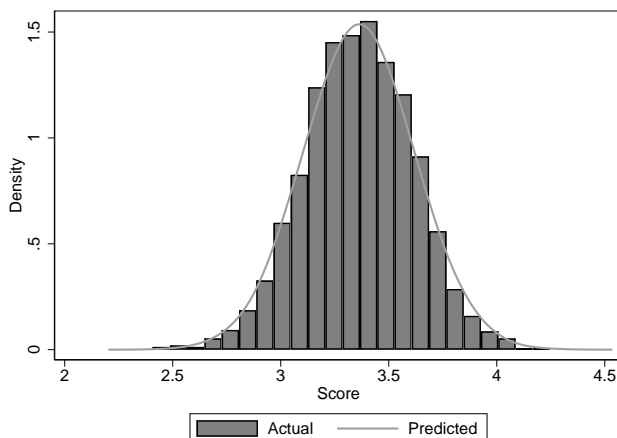
(b) c at $t = 1$



(c) n at $t = 1$

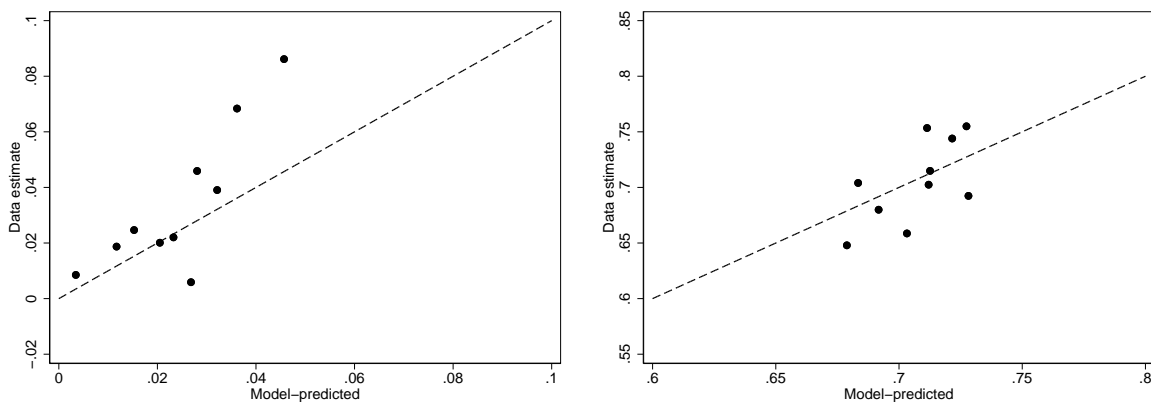
Notes : This figure plots the expected value of the unobserved factor θ conditional on the selected path d and covariates x . The horizontal axis shows the expected value without the HS offer, the vertical axis with the HS offer. Identification of the selection model requires the expected θ not to be the same with and without the offer, thus these values should not lie on the 45 degree line of the graph. Points show the means of θ for different combination of covariates value within each selected path. In particular, we define 16 mutually exclusive cells in terms of covariate values where cells are defined by interactions of below/above median bins of four covariates: HS center quality, HS center transportation, other center quality and other center transportation. Exploiting the simulated data, for each combination of covariate bin and selected path we compute the average value of the unobserved factor, separately when the experimental offer is on and when it is off. F-statistic and p-value come from the Wald tests of the hypothesis that a regression of $\lambda(d, 1, x)$ on $\lambda(d, 0, x)$ has constant=0 and slope=1.

Figure A.5: Goodness of Fit, Distribution of Test Score at Age 4



Notes : Comparison between the actual distribution of the test scores obtained at the end of the age-4 year (the test score measure is the average between PPVT score, WJIII Letter-Word score and WJIII Spelling score) as in the population and the distribution of the same test score as predicted by the sequential threshold model.

Figure A.6: Model-predicted Experimental Effects vs. Estimates from the Data

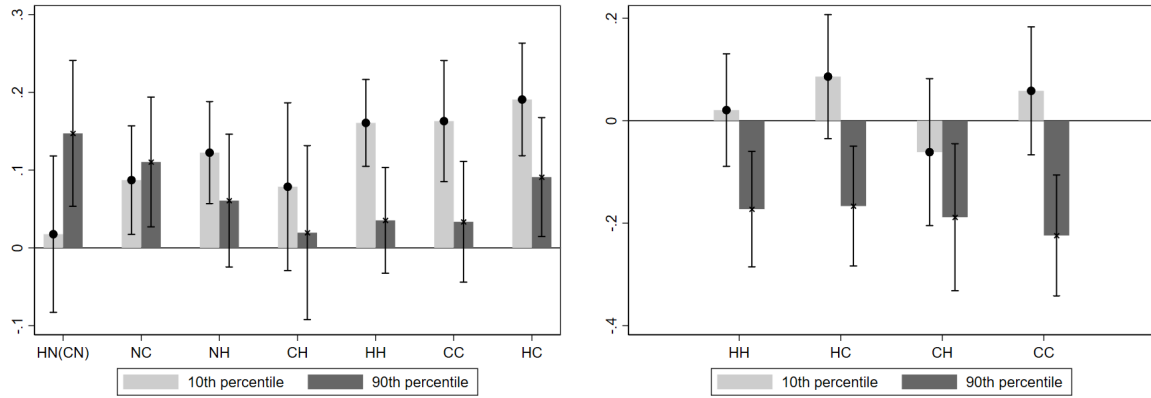


(a) Experimental Returns

(b) Experimental Compliance Shares

Notes : Comparison between the model-predicted and the actual sample estimated experimental returns (Panel (a)) and experimental compliance share (i.e. participation in HS at $t = 1$ when $Z = 1$, Panel(b)). The comparison is done within cells defined by the deciles of the baseline test score distribution (the average baseline test score is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the time of randomization). The dashed line corresponds to the 45-degree line. We test that the empirical and model-predicted returns and compliance shares are equal, treating the model predictions as fixed. For the experimental returns (Panel (a)) we obtain $F = 0.38$ with $p = 0.96$, for the compliance shares (Panel (b)) we obtain $F = 0.27$ and $p = 0.99$.

Figure A.7: Partial, Joint and Cross Returns from Sequential Program Participation, 10th and 90th Percentile of the distribution of θ

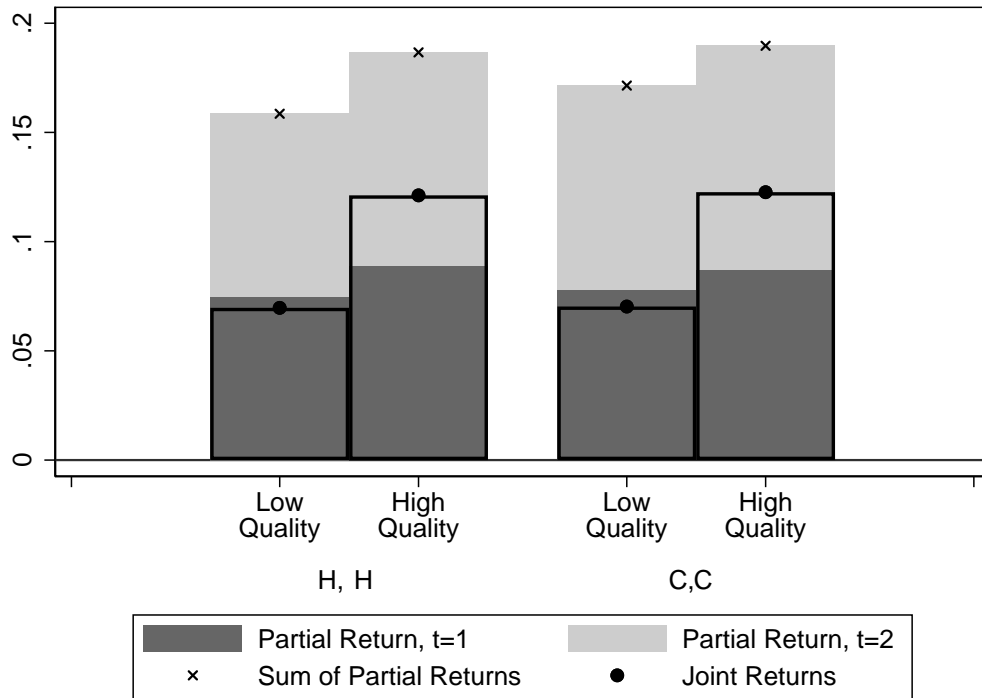


(a) Partial and Joint Returns (with 95% C.I.)

(b) Cross Returns (with 95% C.I.)

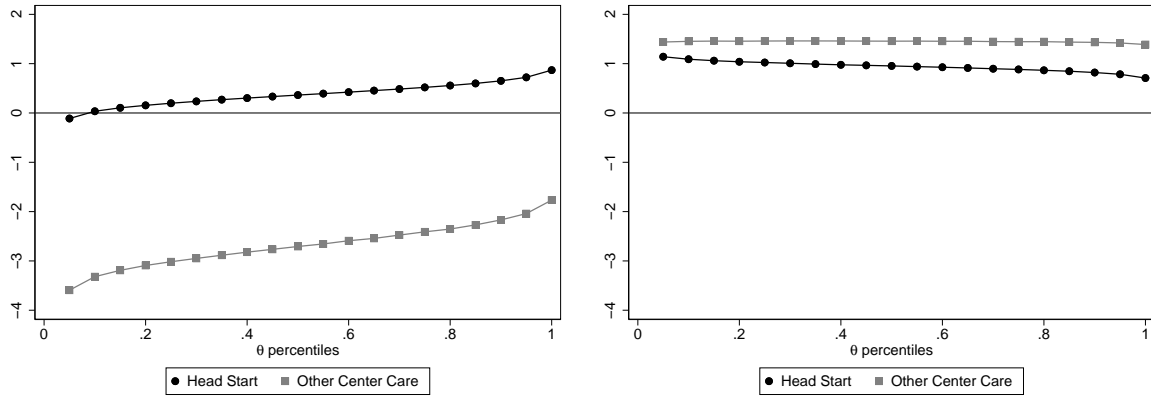
Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. In Panel (a) each bar represents the partial (HN(CN), NC, NH) and joint (CH, HH, CC, HC) returns of different program sequences, while in Panel (b) each bar represents the cross-returns, i.e. the difference between the joint return and the sum of the partial returns at $t=1$ and $t=2$, separately for children with a value of θ corresponding to the 10th and the 90th percentile of the distribution of θ . All these returns use two periods in home care as the baseline program sequence for comparison (see Section 5.2 for the exact definition of partial, joint and cross-returns).

Figure A.8: Partial, Joint and Cross Returns from Sequential Program Participation, by Different Center Quality



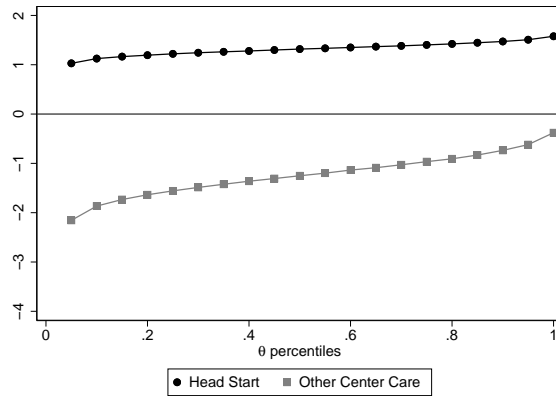
Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. The first two bars show the returns from two periods in HS (H,H) and the last two bars show the returns from two periods in other center care (C,C), by different center quality. For the (H,H) path, we consider the quality of the Head Start center, while for the (C,C) path we consider the quality of the other center care. As explained in Section 4, we define the quality of a center as the residualized quality due to the concern that the characteristics of the center may be correlated with family characteristics. We take the residuals of a regression of the center quality on household size, number of siblings, dummies for whether the child is female, black, hispanic, use English as home language, living in urban area, living with both parents, in need of special education, child of a teen mother, child of a mother who never married or is separated, child of a mother with high school or more than high school. In this graph, we define a low quality center as one with quality corresponding to the 10th percentile of the respective residualized quality distribution, while high quality is a center with quality corresponding to the 90th percentile.

Figure A.9: Flow Utility, Continuation Value and Overall Perceived Value of Attending Head Start and Other Center Care at $t = 1$ Relative to Home Care for Different Values of θ



(a) Flow Utility Relative to Home Care

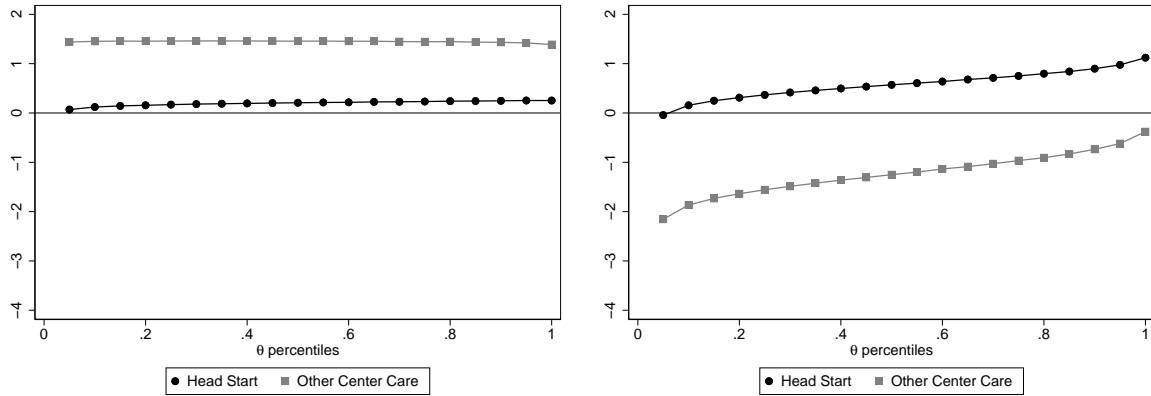
(b) Continuation Value Relative to Home Care



(c) Overall Perceived Value Relative to Home Care

Notes : On each curve, each marker represents the average within one ventile of the distribution of θ (i.e. baseline ability). Panel (a) shows the per-ventile average flow utility (excluding choice shocks) of attending HS or other center care at $t = 1$ as compared to the flow utility from home care (note that because these are means, the utility shock is averaged out). Panel (b) shows per-ventile average continuation value of attending HS or other center care at $t = 1$ as compared to the continuation value of home care. Panel (c) is the sum of the values in Panels (a) and (b) and shows the difference between the total perceived value of attending HS or other center care and the total perceived value of home care. All results are derived for treatment group individuals from the structural model.

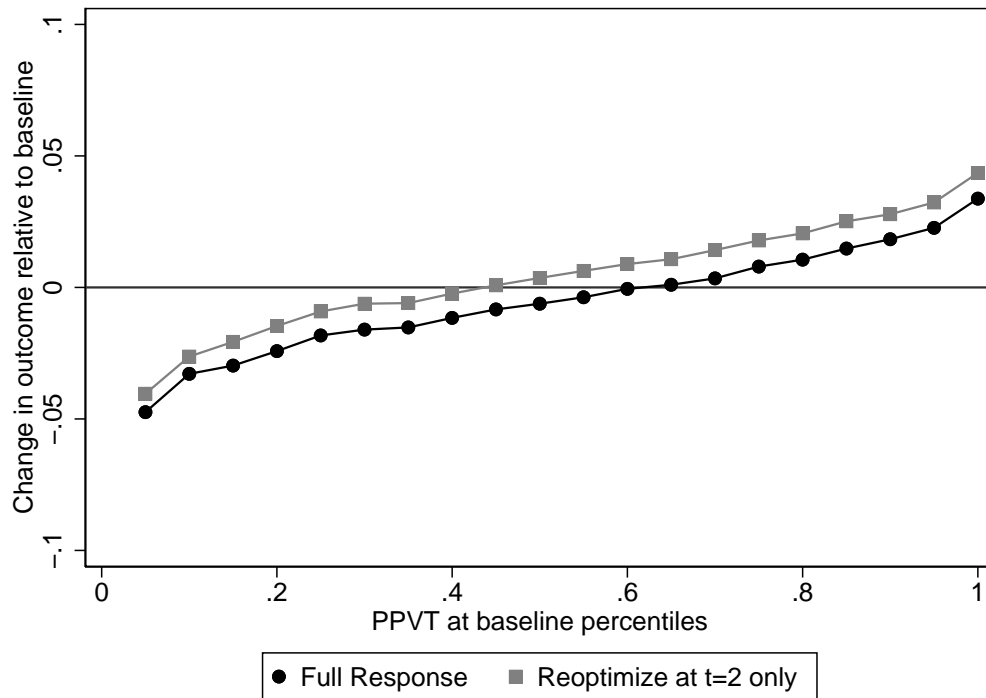
Figure A.10: Continuation Value and Overall Perceived Value of Attending Head Start and Other Center Care at $t = 1$ Relative to Home Care, Counterfactual Full Response Scenario



(a) Continuation Value Relative to Home Care (b) Overall Perceived Value Relative to Home Care

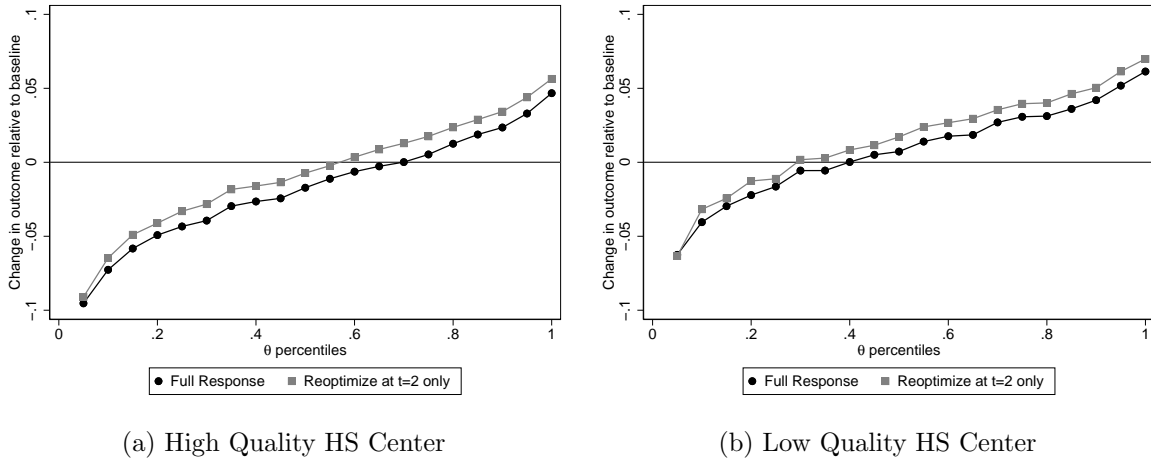
Notes : The figure shows the continuation values and the overall perceived values for the counterfactual scenario where individuals select at most one period in HS and assumes that they can reoptimize both at $t = 1$ and $t = 2$ (“Full response”). On each curve, each marker represents the average within one ventile of the distribution of θ (i.e. baseline ability). Panel (a) shows per-ventile average continuation value of attending HS or other center care at $t = 1$ as compared to the continuation value of home care. Panel (b) is the sum of the flow utility (same as the baseline scenario with no Head Start limits – see Figure A.9) and panel (a) and shows the difference between the total perceived value of attending HS or other center care and the total perceived value of home care. All results are derived from the structural model.

Figure A.11: Effect of Limiting Head Start Enrolment to One Period Only for Different Values of the Baseline PPVT Scores



Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. Each curve represents deviations in outcomes from the baseline scenario in the treatment group, e.g., zero means no effect of the HS enrolment limit on the outcome. On each curve, each marker represents one ventile of the distribution of the PPVT score obtained at baseline (e.g. the leftmost marker represents individuals with the lowest 5 percentiles of baseline PPVT score, the second leftmost marker represents individuals between the 5-10 percentiles, and so on). All scenarios are simulated from the structural model with 30 paths per treatment-group individual. “Full response”: allow individuals to reoptimize at t=1 and t=2. “Reoptimize at t=2 only”: assume individuals at t=1 do not know about the HS enrolment limit, so they can reoptimize at t=2 only.

Figure A.12: Effect of Limiting Head Start Enrolment to One Period Only for Different Values of θ , High and Low Quality of the HS Center



Notes : The outcome analyzed is the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year. Each curve represents deviations in outcomes from the baseline scenario in the treatment group, e.g., zero means no effect of the HS enrolment limit on the outcome. On each curve, each marker represents one ventile of the distribution of θ (i.e. baseline ability). Panel (a) shows the result for high quality HS centers, while panel (b) is for low quality (where high quality is defined as the top quartile of the quality distribution and low quality as the bottom quartile). As explained in Section 4, we define the quality of a center as the residualized quality due to the concern that the characteristics of the center may be correlated with family characteristics. We take the residuals of a regression of the center quality on household size, number of siblings, dummies for whether the child is female, black, hispanic, use English as home language, living in urban area, living with both parents, in need of special education, child of a teen mother, child of a mother who never married or is separated, child of a mother with high school or more than high school. All scenarios are simulated from the structural model with 30 paths per treatment-group individual. “Full response”: allow individuals to reoptimize at $t=1$ and $t=2$. “Reoptimize at $t=2$ only”: assume individuals at $t=1$ do not know about the HS enrolment limit, so they can reoptimize at $t=2$ only.

Table A.1: Descriptive Statistics in the Baseline Survey

	Control Group	Treatment Group	Difference
Household Size	4.53	4.50	0.02
Female=1	0.50	0.51	-0.00
Black=1	0.34	0.37	-0.03
Hispanic=1	0.32	0.33	-0.00
English as Home Language=1	0.75	0.75	0.00
In Need of Special Education=1	0.09	0.13	-0.03**
PPVT Score at Baseline	2.31	2.29	0.02
WJIII Score - Spelling at Baseline	3.35	3.34	0.01
WJIII Score - Letter-word at Baseline	2.93	2.94	-0.00
WJIII Score - Applied-problem at Baseline	3.69	3.70	-0.01
Number of Siblings	1.33	1.39	-0.06
Both Parents in the Household=1	0.47	0.48	-0.00
Teen Mother=1	0.17	0.13	0.03**
Mother not Married=1	0.42	0.43	-0.01
Mother Separated=1	0.14	0.13	0.00
Mother Education: High School Diploma=1	0.33	0.35	-0.02
Mother Education: More than High School=1	0.29	0.31	-0.01
Urban Area=1	0.80	0.82	-0.01
HS Center Quality Index	0.69	0.69	0.00
HS Center Transportation Available	0.68	0.67	0.01
Other Center Quality Index	0.59	0.59	-0.00
Other Center Transportation Available	0.52	0.54	-0.01
N	698	1,178	

Notes : This table reports the mean value of different individual characteristics comparing the control to the treatment group. The HS center quality index (obtained directly from the HSIS data) and the availability of transportation refer to the center of random assignment. The other centers quality index and the availability of transportation refer to the average of all other center care by the child's assigned HS center. Note that the center quality and the transportation variables used in actual estimation are residualized and standardized so to have standard deviation 1 (see Section 4 for details). We adjust the PPVT and WJIII test scores by dividing the raw score by 100. Significance level (t-test for the difference between the average for the treatment group and the one for the control group=0): *** 1%, ** 5%, * 10%

Table A.2: Proportion, Mean of the Unobserved Factor and Returns for each Compliance Group

D_i^0	D_i^1	Compliance Type		Proportion (1)	θ (2)	Return (3)	Population Return (4)
(h,h)	(h,h)	AT		10.45	-0.04	0	0
(h,c)	(h,c)	AT		2.98	0.01	0	0
(h,n)	(h,n)	AT		0.97	0.01	0	0
(c,h)	(h,h)	c-compliers	Temporary	5.03	0	0.052	0.048
(c,h)	(h,c)	c-compliers	Path	1.28	0.05	0.088	0.091
(c,h)	(h,n)	c-compliers	Path	0.37	0.06	0.054	0.033
(c,h)	(c,h)	NT		1.73	0.04	0	0
(c,c)	(h,h)	c-compliers	Path	8.33	0.05	0.001	0
(c,c)	(h,c)	c-compliers	Temporary	2.79	0.11	0.043	0.042
(c,c)	(h,n)	c-compliers	Path	0.71	0.1	0.028	-0.015
(c,c)	(c,c)	NT		3.23	0.11	0	0
(c,n)	(h,h)	c-compliers	Path	0.92	0.05	-0.009	0.015
(c,n)	(h,c)	c-compliers	Path	0.29	0.1	0.017	0.058
(c,n)	(h,n)	c-compliers	Temporary	0.09	0.09	0.001	0
(c,n)	(c,n)	NT		0.38	0.11	0	0
(n,h)	(h,h)	n-compliers	Temporary	17.13	-0.05	0.016	0.006
(n,h)	(h,c)	n-compliers	Path	4.49	0	0.049	0.049
(n,h)	(h,n)	n-compliers	Path	1.32	-0.01	-0.015	-0.009
(n,h)	(n,h)	NT		4.4	-0.08	0	0.007
(n,c)	(h,h)	n-compliers	Path	11.34	0	0	0
(n,c)	(h,c)	n-compliers	Temporary	3.57	0.06	0.028	0.042
(n,c)	(h,n)	n-compliers	Path	0.95	0.04	-0.008	-0.016
(n,c)	(n,c)	NT		2.83	0	0	0
(n,n)	(h,h)	n-compliers	Path	8.85	-0.01	0.104	0.098
(n,n)	(h,c)	n-compliers	Path	2.41	0.04	0.133	0.14
(n,n)	(h,n)	n-compliers	Temporary	0.74	0.02	0.086	0.082
(n,n)	(n,n)	NT		2.28	-0.03	0	0

Notes : D_i^0 and D_i^1 refer to the potential program sequence without and with the HS offer, respectively. Columns (1) and (2) report proportions and the average baseline ability (θ) of each compliance group. Column (3) reports the average return (in terms of the average between the score obtained by each child in the PPVT, WJIII Letter-Word and WJIII Spelling tests at the end of age-4 year) to program sequence D_i^1 relative to D_i^0 within each type of complier. Column (4) reports the difference in average potential outcome between program sequence D_i^1 and D_i^0 in the population. See Section 5.3 for definitions of the compliance groups.

Table A.3: Goodness of Fit of the Dynamic Structural Model

		Program $t = 2$		
		h	c	n
		Control Group		
Program $t = 1$				
Simulation	h	10.6	3.0	0.9
Data	h	<i>10.7</i>	<i>2.7</i>	<i>1.0</i>
Simulation	c	8.2	15.0	1.8
Data	c	<i>9.3</i>	<i>13.9</i>	<i>1.6</i>
Simulation	n	27.6	18.6	14.4
Data	n	<i>30.2</i>	<i>18.6</i>	<i>11.9</i>
		Treatment Group		
Simulation	h	61.9	18.0	5.2
Data	h	<i>62.0</i>	<i>18.1</i>	<i>5.0</i>
Simulation	c	1.7	3.2	0.4
Data	c	<i>1.2</i>	<i>3.7</i>	<i>0.5</i>
Simulation	n	4.4	2.7	2.4
Data	n	<i>3.0</i>	<i>2.8</i>	<i>3.7</i>

Notes : Numbers in italics are the population equivalents (computed directly from the data). Each cell shows the proportion for the group of children who at $t = 1$ enrol in the program reported in the second column and at $t = 2$ enrol in the program reported in the second row, separately for the treatment and for the control group. We perform tests of equality between model-predicted and empirical choice proportions by treatment and control group. The p-values (Chi-square test, df=8) are 0.56 ($Z = 0$) and 0.03 ($Z = 1$).

Table A.4: Estimates of the Parameters of the Utility Function from the Structural Model imposing $\delta=0.95$

Covariate	Flow Utility					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept ($\tilde{\psi}_1^h, \tilde{\psi}_1^c$)	-1.396	(0.194)	***	-2.718	(0.377)	***
Factor θ ($\tilde{\lambda}_1^h, \tilde{\lambda}_1^c$)	0.630	(0.291)	**	2.012	(0.469)	***
Treatment group Z	1.830	(0.145)	***			
HS center transport $\times Z$	0.150	(0.064)	**			
$\theta \times Z$	0.402	(0.368)				
HS center quality $\times (1-Z)$	-0.232	(0.056)	***			
HS center transport	-0.146	(0.046)	***			
Other center transport				-0.179	(0.072)	**
<i>Period 2 choice parameters:</i>						
Intercept ($\tilde{\psi}_2^h, \tilde{\psi}_2^c$)	0.333	(0.119)	***	-0.484	(0.184)	***
Period 1 in h ($\tilde{\psi}_2^{hh}, \tilde{\psi}_2^{hc}$)	1.176	(0.120)	***	0.515	(0.167)	***
Period 1 in c ($\tilde{\psi}_2^{ch}, \tilde{\psi}_2^{cc}$)	1.214	(0.204)	***	2.250	(0.315)	***
Factor θ ($\tilde{\lambda}_2^h, \tilde{\lambda}_2^c$)	0.712	(0.361)	**	1.569	(0.463)	***
HS center transport	0.045	(0.030)				
Other center transport				0.230	(0.057)	***
<i>Other parameters:</i>						
$\tilde{\sigma}_{\nu^h}$	1.000	-				
$\tilde{\sigma}_{\nu^c}$	2.317	(0.296)	***			
$\tilde{\rho}_{hc}$	0.817	(0.112)	***			
κ	3.558	(1.004)	***			

Notes : N=1876. Log-likelihood = -3048.76. Flow utilities are specified in Appendix Section E. The technological parameters in the baseline measurement equations, outcome equations and the factor distribution are fixed at the estimated values from the sequential threshold model (see Table 7). κ is the terminal value scaling factor. The discount factor δ is set to 0.95. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table A.5: Effects of Limiting Head Start Enrolment to One Period Only on the Present Value of ex-post Discounted Utility

	Baseline	Counterfactual	
		Full resp.	Partial resp.
First period	0.52	+0.04	0.00
Second period	1.64	-0.56	-0.65
Terminal period	11.84	-0.02	+0.02
Total	14.00	-0.54	-0.63

Notes : All scenarios are simulated from the structural model with 30 paths per treatment-group individual. The “Baseline” column shows the average present discounted utility of the population in each period under the baseline scenario, i.e. when individuals are allowed to select two periods of HS. The “Counterfactual” columns report the difference in the present discounted utility between the baseline scenario and the counterfactuals: “Full resp.” (full response) allows individuals to select at most one period in HS and assumes that they can reoptimize both at $t = 1$ and $t = 2$; “Partial resp.” (partial response) assumes individuals at $t = 1$ do not know about the HS enrolment limit, so they can reoptimize at $t = 2$ only.

Table A.6: Effects of Limiting Head Start Enrolment to One Period Only, by Children’s Ability

Program choice (%)	Low-ability children			High-ability children		
	Baseline	Counterfactual		Baseline	Counterfactual	
	(1)	Full resp.	Partial resp.	(4)	Full resp.	Partial resp.
	(1)	(2)	(3)	(4)	(5)	(6)
Choices at t=1						
<i>h</i>	87.4	-39.1	0.0	87.4	-18.5	0.0
<i>c</i>	3.7	+7.7	0.0	7.7	+6.6	0.0
<i>n</i>	8.9	+31.5	0.0	4.9	+11.9	0.0
Choices at t=2						
<i>h</i>	76.2	-48.9	-70.1	62.7	-50.5	-57.8
<i>c</i>	18.5	+20.3	+31.5	26.4	+26.2	+29.0
<i>n</i>	5.3	+28.6	+38.6	10.9	+24.3	+28.8
Program sequence						
<i>h, h</i>	70.1	-70.1	-70.1	57.8	-57.8	-57.8
<i>h, c</i>	14.6	+10.5	+31.5	20.3	+18.1	+29.0
<i>h, n</i>	2.7	+20.5	+38.6	9.2	+21.3	+28.8
<i>c, h</i>	1.2	+2.8	0.0	2.6	+1.9	0.0
<i>c, c</i>	2.1	+4.4	0.0	4.3	+4.1	0.0
<i>c, n</i>	0.4	+0.4	0.0	0.8	+0.6	0.0
<i>n, h</i>	4.8	+18.3	0.0	2.3	+5.5	0.0
<i>n, c</i>	1.8	+5.3	0.0	1.7	+4.0	0.0
<i>n, n</i>	2.3	+7.8	0.0	0.9	+2.4	0.0

Notes : All scenarios are simulated from the structural model with 30 paths per treatment-group individual. The “Baseline” columns show the proportion of the population making each choice under the baseline scenario, i.e. when individuals are allowed to select two periods of HS. The “Counterfactual” columns report the difference in the proportion of population selecting each childcare option between the baseline scenario and the counterfactuals: “Full resp.” (full response) allows individuals to select at most one period in HS and assumes that they can reoptimize both at $t = 1$ and $t = 2$; “Partial resp.” (partial response) assumes individuals at $t = 1$ do not know about the HS enrolment limit, so they can reoptimize at $t = 2$ only. Columns (1) – (3) (“Low-ability children”) report the results for children in the 10th percentile of the θ distribution, while Columns (4)–(6) (“High-ability children”) report the results for the children in the 90th percentile of the θ distribution.

Table A.7: Effects of Limiting Head Start Enrolment to One Period Only – Differences with Respective Baseline Scenarios

	Program choice (%)		Average Outcome	
	Dynamic model (Full resp.)	Threshold model	Dynamic model (Full resp.)	Threshold model
	(1)	(2)	(3)	(4)
Choices at $t = 1$				
h	-25.5	0.0		
c	+6.7	0.0		
n	+18.8	0.0		
Choices at $t = 2$				
h	-51.0	-62.1		
c	+24.4	+29.5		
n	+26.5	+32.5		
Program sequence				
h, h	-61.9	-62.1	-	-
h, c	+15.0	+29.5	-0.01	-0.03
h, n	+21.4	+32.5	-0.07	-0.05
c, h	+2.3	0.0	-0.04	0.00
c, c	+4.0	0.0	-0.03	0.00
c, n	+0.4	0.0	-0.06	0.00
n, h	+8.7	0.0	-0.01	0.00
n, c	+5.4	0.0	-0.02	0.00
n, n	+4.7	0.0	-0.02	0.00
Overall mean			3.366	3.374
Overall s.d.			0.272	0.274

Notes : Each column reports the predicted differences in the proportion of the population making each choice and their average outcome between the counterfactual scenarios and the baseline scenario where individuals are allowed to select two periods of HS. In the “Dynamic model” columns, the counterfactual is the “full response” scenario where individuals can reoptimize both at $t = 1$ and $t = 2$. In particular, column (1) corresponds to column (2) in Table 11, while column (3) corresponds to column (5) in Table 11. The predictions from the counterfactual and the baseline are computed from the estimated dynamic structural model from Section 6 (by simulating 30 paths per treatment-group individual). In the “Threshold model” columns, the counterfactual is computed by treating the threshold model as a sequentially static behavior model and the perceived values as representing utilities. The predictions from the counterfactual and the baseline are computed from the estimated threshold model from Section 4 (by simulating 30 paths per treatment-group individual).

B Robustness Checks of the Threshold Model

We conduct a range of robustness checks to gauge the sensitivity of our conclusions to various modeling and empirical specifications classified into three broad categories: (1) choice equations, (2) factor distribution and (3) measurement and outcome.

B.1 Choice equation specifications

B.1.1 Decision rules in $t = 2$

Our main model allows the perceived value of a given program alternative in $t = 2$ to differ by whether an individual experienced h, c or n in $t = 1$. However, we have assumed constant marginal value of lagged program choice; program choice in $t = 1$ only makes an intercept shift to the perceived value of program alternatives in $t = 2$.

We explore the sensitivity of our results by allowing for heterogeneity in the effects of lagged program choice. A convenient way to do so is to characterize separate decision rules in $t = 2$ by the program choice made in $t = 1$. Denote the set of decision nodes by $\mathcal{J} = \{\mathbf{o}, \mathbf{h}, \mathbf{c}, \mathbf{n}\}$, where \mathbf{o} is the decision node at $t = 1$, and \mathbf{h} , \mathbf{c} and \mathbf{n} are the decision nodes that the individual gets to at $t = 2$ after she chose h , c and n at $t = 1$. Figure B.1 illustrates the decision tree.

We specify a reduced-form threshold model for the choice process at each decision node. Upon arrival at node $j \in \mathcal{J}$, the perceived value from each alternative is given by:

$$U_{ij}^h = \psi_j^h + \mathbf{X}_{ih}'\boldsymbol{\beta}_j^h + f_j(Z_i, \mathbf{X}_{ih}, \theta_i) + \lambda_j^h\theta_i + \nu_{ij}^h \quad (\text{B.1})$$

$$U_{ij}^c = \psi_j^c + \mathbf{X}_{ic}'\boldsymbol{\beta}_j^c + \lambda_j^c\theta_i + \nu_{ij}^c \quad (\text{B.2})$$

$$U_{ij}^n = 0 \quad (\text{B.3})$$

where U_{ij}^h is the value of HS, U_{ij}^c is the value of other center care, and the value of home care (U_{ij}^n) is normalized to zero. As in our main specification, the random experimental offer Z_i affects the perceived value for h at decision node \mathbf{o} only but has no direct impacts on the perceived values of other options within the same decision node and of all options in other decision nodes:

$$f_j(Z_i, \mathbf{X}_{ih}, \theta_i) = \begin{cases} (\beta_Z + \mathbf{X}_{ih}'\boldsymbol{\beta}_{Zx} + \beta_{Z\theta}\theta_i) Z_i, & \text{for } t = 1(j = \mathbf{o}) \\ 0, & \text{for } t = 2(j \in \{\mathbf{h}, \mathbf{c}, \mathbf{n}\}) \end{cases} \quad (\text{B.4})$$

This specification relaxes the constant effect assumption – parameters $(\boldsymbol{\beta}_j^h, \lambda_j^h, \boldsymbol{\beta}_j^c, \lambda_j^c)$ can vary flexibly across the three decision nodes $j \in \{\mathbf{h}, \mathbf{c}, \mathbf{n}\}$ in $t = 2$. If these parameters are constant across the period-2 decision nodes, then it is easy to see that this alternative specification is equivalent to our

main model via a simple normalization of the intercepts:

$$\psi_n^h = \psi_2^h, \tag{B.5}$$

$$\psi_n^c = \psi_2^c, \tag{B.6}$$

$$\psi_h^h = \psi_2^h + \psi_2^{hh}, \tag{B.7}$$

$$\psi_h^c = \psi_2^c + \psi_2^{hc}, \tag{B.8}$$

$$\psi_c^h = \psi_2^h + \psi_2^{ch}, \tag{B.9}$$

$$\psi_c^c = \psi_2^c + \psi_2^{cc}. \tag{B.10}$$

In equations (B.5)-(B.10), the LHS corresponds to the node-specific intercepts in the alternative specification and the RHS corresponds to the parameters in the main model (equations (6)-(8)).

To gauge the sensitivity of the estimated program returns, we extend the main model by adding five parameters to $t = 2$ choice equations: $\lambda_h^h, \lambda_h^c, \lambda_c^h, \lambda_n^h$, and λ_n^c . In this extended model, the factor loadings in $t = 2$ equations differ by lagged program choices in $t = 1$.⁵⁸ Effectively, individuals with different θ are now subject to different degrees of “state dependence” in preferences and hence perturbing initial childcare choice can have different consequences for subsequent choice among different types of individuals. Appendix Tables B.1 and B.2 presents the estimates from this extended model. We cannot reject the null hypothesis that λ_h^h and λ_n^h are equal to the estimated coefficient λ_2^h in the main model (p-val=0.216). Similarly, we cannot reject the null hypothesis that λ_h^c, λ_c^c and λ_n^c are equal to the estimated coefficient λ_2^c in the main model (p-val=0.138). More importantly, it is reassuring that the estimated parameters in the potential outcome equations are similar to our main model.

B.1.2 Nonlinearity in θ

The choice equations in the sequential threshold model are approximations to the perceived value of each alternative. For forward-looking individuals, these perceived values are a combination of utility flow and perceived continuation value; even if the flow is linear in parameters, the perceived continuation value may be nonlinear in parameters due to incorporation of uncertainty of future payoffs (see Section 6). To assess this threat of functional form misspecification, we specify a quadratic function of factor θ in the period-1 choice equations for h and c to allow for nonlinear selection by θ .⁵⁹ The estimation results are reported in Tables B.3 and B.4. The quadratic coefficients are 1.618 (se=0.961) and 1.905 (se=1.922) in the choice equations for h and c , respectively. The other estimated parameters in the model are very similar to those in the baseline model.

⁵⁸Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n) , to avoid overfitting, we only estimate one of the two factor loadings associated with the decision node c . We choose to estimate λ_c^c and normalize $\lambda_c^h = \lambda_n^c = 0$. Our results are similar if we estimate λ_c^h but normalize $\lambda_c^c = \lambda_n^c = 0$.

⁵⁹As discussed in Section 6, the period-2 value function $V_{i2}(\cdot)$ can be expressed as a linear function of θ , therefore we focus on the threat of misspecification at $t = 1$.

B.1.3 Excluding transportation availability

Appendix Tables B.5 and B.6 reports the estimates when we exclude transportation options ($tp_{ih}, tp_{ic}, tp_{ih} \times Z_i$) from the choice equations in both periods. We find that our estimates change little if we do not use transportation options as exclusion restrictions. Although in our baseline model the availability of transportation services do not affect potential outcomes (see equation (9)) and hence become exclusion restrictions, in practice the transportation variables play little role in identifying the causal program returns.

B.2 Robustness to factor distribution

We relax the factor distribution used in estimation by modelling it as a mixture normal distribution instead. The estimation results are reported in Tables B.7 and B.8. The second mixture has a probability weight of 0.625 (se=0.10) with a mean of -0.078 (se=0.027) and standard deviation of 0.196 (se=0.014). The first mixture has a probability weight of 0.375 with a mean of 0.13 (implied by the zero-mean normalization of the factor distribution) and standard deviation of 0.218 (se=0.02). The other estimated parameters in the model remain very similar to those in the baseline model, with one notable exception: the average return for program sequence (c, h) becomes marginally significant, but still below the average returns from other program sequences.

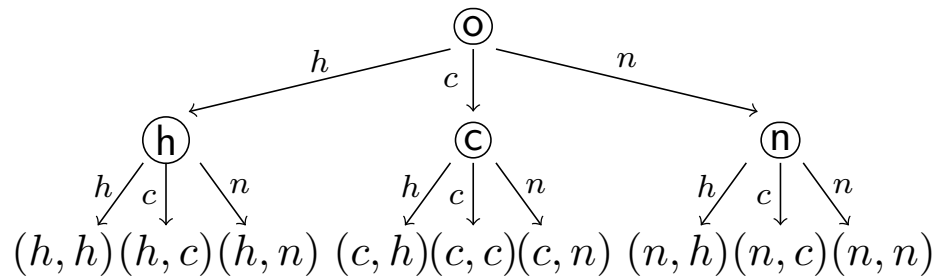
B.3 Robustness to measurement and outcome definition

We estimate the model by incorporating WJIII Applied Problems test into the production of child skills. WJIII Applied Problems test is the only assessment measuring early numeracy/math skills in pre-school years of the HSIS. We use the baseline WJIII Applied Problems test score as the 4th measurement of the latent factor. As outcome, we use the average of end-of-age-4 scores of PPVT and three WJIII subtests (Letter-Word, Spelling and Applied Problems). Given that scores from these three WJIII subtests form the WJIII Pre-Academic Skills scores, this outcome can be seen as a weighted average of PPVT and WJIII Pre-Academic Skills test scores.⁶⁰ It reflects a broader set of skills, including early mathematical skills in addition to literacy and reading skills.

The estimation results are reported in Tables B.9 and B.10, with the estimated program returns shown in Appendix Figure B.2 (average returns) and Figure B.3 (returns for different levels of θ). The estimated choice equations are similar to those in the baseline model. In the outcome equations, the intercept estimates are similar to those in the baseline model and the standard errors are about 10 percent smaller. The relative magnitudes of the factor loadings remain similar to those in the baseline model, although the levels of loadings associated with certain program sequence are reduced. The most notable difference is the loading for sequence (h, n) (reduced from 0.227 in the baseline to 0.099 and insignificant), and the loading for sequence (c, c) (reduced from -0.227 to -0.308). These changes

⁶⁰Kline and Walters (2016) use the simple average of PPVT and WJIII pre-academic skills test scores as the outcome variable. Note that we do not use standardized scores like they do.

Figure B.1: Decision Tree with Node-specific Perceived Values



Notes : \odot refers to a decision node. (j, j') refers to the program sequence where program $j \in (h, c, n)$ is obtained at period $t = 1$ and $j' \in (h, c, n)$ is obtained at period $t = 2$.

indicate that including early math skills tends to strengthen the gradient of certain program returns w.r.t baseline ability.

Overall, the stability of these results is noteworthy given the non-trivial change in the measurement system and the outcome variable. In the measurement system, the loading on the baseline WJIII Applied Problems is sizable (0.596; a signal-to-noise ratio of 34%). This suggests that the latent ability captured by the factor θ remains most relevant to the skills embedded in PPVT (which has a signal-to-noise ratio of 37%, see Section 5.1).

Table B.1: Estimates of the Parameters of the Choice Equations, Different Factor Loadings at Different Nodes in $t = 2$

Covariate	Perceived Value					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept (ψ_1^h, ψ_1^c)	-0.476	(0.134)	***	-1.263	(0.235)	***
Factor θ (λ_1^h, λ_1^c)	0.293	(0.299)		2.096	(0.52)	***
Treatment group Z	1.802	(0.144)	***			
HS center quality $\times Z$	0.296	(0.066)	***			
HS center transport $\times Z$	0.134	(0.066)	**			
$\theta \times Z$	0.324	(0.378)				
HS center quality	-0.180	(0.052)	***			
HS center transport	-0.130	(0.045)	***			
Other center quality				0.129	(0.069)	*
Other center transport				-0.098	(0.067)	
<i>Period 2 choice parameters:</i>						
Intercept (ψ_2^h, ψ_2^c)	0.673	(0.081)	***	-0.097	(0.157)	
Period 1 in h (ψ_2^{hh}, ψ_2^{hc})	0.885	(0.097)	***	0.343	(0.168)	**
Period 1 in c (ψ_2^{ch}, ψ_2^{cc})	0.743	(0.154)	***	1.829	(0.294)	***
Factor θ after h in $t = 1$ (λ_h^h, λ_h^c)	-0.277	(0.332)		1.505	(0.615)	**
Factor θ after c in $t = 1$ (λ_c^c)	-			1.040	(0.808)	
Factor θ after n in $t = 1$ (λ_n^h, λ_n^c)	0.133	(0.365)		0.234	(0.695)	
HS center quality	0.070	(0.03)	**			
HS center transport	0.041	(0.032)				
Other center quality				0.198	(0.057)	***
Other center transport				0.220	(0.056)	***
<i>Other parameters:</i>						
σ_{ν^h}	1.000	-				
σ_{ν^c}	2.296	(0.278)	***			
ρ_{hc}	0.840	(0.105)	***			

Notes : N=1876. Log-likelihood = -3032.42. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of the age-4 year as the outcome. Alternative-specific choice shocks assumed to follow a bivariate normal distribution with zero means, variance of 1 (normalization) for ν_{ij}^h , variance of $\sigma_{\nu^c}^2$ for ν_{ij}^c , and correlation coefficient ρ_{hc} . In this extended model, the factor loadings in $t = 2$ equations differ by lagged program choices in $t = 1$. Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n) , to avoid overfitting, we only estimate one of the two factor loadings associated with the decision node c. We choose to estimate λ_c^c and normalize $\lambda_c^h = \lambda_c^n = 0$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table B.2: Estimates of the Parameters of the Potential Outcome and Measurement Equations, Different Factor Loadings at Different Nodes in $t = 2$ in the choice Equation

	<i>A. Potential outcome equations</i>					
	Intercept			Factor		
Potential outcomes:						
Y^{nn}	3.276	(0.019)	***	0.970	(0.089)	***
$Y^{hn} - Y^{nn}$	0.078	(0.034)	**	0.236	(0.134)	*
$Y^{nh} - Y^{nn}$	0.071	(0.025)	***	-0.114	(0.107)	
$Y^{nc} - Y^{nn}$	0.104	(0.027)	***	0.033	(0.106)	
$Y^{hh} - Y^{nn}$	0.092	(0.021)	***	-0.210	(0.087)	**
$Y^{hc} - Y^{nn}$	0.121	(0.025)	***	-0.161	(0.105)	
$Y^{ch} - Y^{nn}$	0.039	(0.035)		-0.100	(0.159)	
$Y^{cc} - Y^{nn}$	0.090	(0.027)	***	-0.224	(0.104)	**
	Other covariates					
1 period quality HS	0.006	(0.008)				
2 periods quality HS	0.021	(0.007)	***			
1 period quality other center	0.004	(0.011)				
2 periods quality other center	0.023	(0.017)				
	<i>B. Measurement equations</i>					
	Intercept			Factor		
PPVT	2.299	(0.009)	***	1.000	-	
WJIII Spelling	3.348	(0.007)	***	0.498	(0.034)	***
WJIII Letter-Word	2.939	(0.006)	***	0.551	(0.032)	***
	<i>C. Factor and error SD</i>					
Factor	0.223	(0.009)	***			
Outcome	0.157	(0.005)	***			
Baseline measurements:						
PPVT	0.293	(0.006)	***			
WJIII Spelling	0.243	(0.004)	***			
WJIII Letter-Word	0.185	(0.004)	***			

Notes : N=1876. Log-likelihood = -3032.42. In this extended model, the factor loadings in $t = 2$ equations differ by lagged program choices in $t = 1$. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of age-4 year as the outcome. The intercept and the factor loading of the outcome equation differ depending on the program sequence (see equation 9), and are all relative to the parameters of the outcome equation for the nn sequence. h is for HS, c is for other center care, n is for home care. Baseline measurements are test scores measured around the time of randomization (in Fall 2002). The factor loading in the measurement equation for the baseline PPVT test score is normalized to one ($\gamma^{pp}=1$). Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n), to avoid overfitting we restrict $\alpha_t^{cn} = \alpha_t^{hn}$ and $\gamma_t^{cn} = \gamma_t^{hn}$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): ***1%, ** 5%, *10%.

Table B.3: Estimates of the Parameters of the Choice Equations from the Sequential Threshold Model, Non-Linear Function of θ in the Choice Equation

Covariate	Perceived Value					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept (ψ_1^h, ψ_1^c)	-0.601	(0.149)	***	-1.328	(0.268)	***
Factor θ (λ_1^h, λ_1^c)	0.321	(0.303)		2.078	(0.542)	***
Factor θ squared	1.618	(0.961)	*	1.905	(1.922)	
Treatment group Z	1.856	(0.150)	***			
HS center quality $\times Z$	0.305	(0.067)	***			
HS center transport $\times Z$	0.140	(0.067)	**			
$\theta \times Z$	0.321	(0.400)				
HS center quality	-0.189	(0.052)	***			
HS center transport	-0.135	(0.046)	***			
Other center quality				0.129	(0.070)	*
Other center transport				-0.096	(0.067)	
<i>Period 2 choice parameters:</i>						
Intercept (ψ_2^h, ψ_2^c)	0.643	(0.090)	***	-0.088	(0.153)	
Period 1 in h (ψ_2^{hh}, ψ_2^{hc})	0.912	(0.101)	***	0.359	(0.166)	**
Period 1 in c (ψ_2^{ch}, ψ_2^{cc})	0.766	(0.164)	***	1.812	(0.300)	***
Factor θ (λ_2^h, λ_2^c)	-0.344	(0.252)		0.780	(0.444)	*
HS center quality	0.074	(0.029)	**			
HS center transport	0.045	(0.032)				
Other center quality				0.195	(0.057)	***
Other center transport				0.219	(0.056)	***
<i>Other parameters:</i>						
σ_{ν^h}	1.000	-				
σ_{ν^c}	2.225	(0.303)	***			
ρ_{hc}	0.803	(0.118)	***			

Notes : N=1876, Log-likelihood = -3033.93. A quadratic function of factor θ ("Factor θ squared") is included in the period-1 choice equations for h and c to allow for nonlinear selection by θ . The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of the age-4 year as the outcome. Standard errors are in parentheses. Alternative-specific choice shocks assumed to follow a bivariate normal distribution with zero means, variance of 1 (normalization) for ν_{ij}^h , variance of $\sigma_{\nu^c}^2$ for ν_{ij}^c , and correlation coefficient ρ_{hc} . Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table B.4: Estimates of the Parameters of the Potential Outcome and Measurement Equations from the Sequential Threshold Model, Non-Linear Function of θ in the Choice Equation

	<i>A. Potential outcome equations</i>					
	Intercept			Factor		
Potential outcomes:						
Y^{nn}	3.267	(0.018)	***	1.024	(0.093)	***
$Y^{hn} - Y^{nn}$	0.076	(0.032)	**	0.146	(0.136)	
$Y^{nh} - Y^{nn}$	0.094	(0.024)	***	-0.126	(0.112)	
$Y^{nc} - Y^{nn}$	0.098	(0.025)	***	0.032	(0.109)	
$Y^{hh} - Y^{nn}$	0.099	(0.019)	***	-0.272	(0.090)	***
$Y^{hc} - Y^{nn}$	0.140	(0.024)	***	-0.225	(0.109)	**
$Y^{ch} - Y^{nn}$	0.052	(0.032)		-0.161	(0.164)	
$Y^{cc} - Y^{nn}$	0.099	(0.026)	***	-0.290	(0.115)	**
	Other covariates					
1 period quality HS	0.005	(0.008)				
2 periods quality HS	0.022	(0.008)	***			
1 period quality other center	0.004	(0.010)				
2 periods quality other center	0.023	(0.017)				
	<i>B. Measurement equations</i>					
	Intercept			Factor		
PPVT	2.298	(0.008)	***	1.000	-	
WJIII Spelling	3.348	(0.007)	***	0.505	(0.034)	***
WJIII Letter-Word	2.939	(0.006)	***	0.559	(0.032)	***
	<i>C. Factor and error SD</i>					
Factor	0.221	(0.008)	***			
Outcome	0.158	(0.004)	***			
Baseline measurements:						
PPVT	0.294	(0.006)	***			
WJIII Spelling	0.242	(0.004)	***			
WJIII Letter-Word	0.185	(0.004)	***			

Notes : N=1876. Log-likelihood = -3033.93. A quadratic function of factor θ is included in the period-1 choice equations for h and c to allow for nonlinear selection by θ . The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of age-4 year as the outcome. The intercept and the factor loading of the outcome equation differ depending on the program sequence (see equation 9), and are all relative to the parameters of the outcome equation for the nn sequence. h is for HS, c is for other center care, n is for home care. Baseline measurements are test scores measured around the time of randomization (in Fall 2002). The factor loading in the measurement equation for the baseline PPVT test score is normalized to one ($\gamma^{pp}=1$). Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n), to avoid overfitting we restrict $\alpha_t^{cn} = \alpha_t^{hn}$ and $\gamma_t^{cn} = \gamma_t^{hn}$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): ***1%, ** 5%, *10%.

Table B.5: Estimates of the Parameters of the Choice Equations from the Sequential Threshold Model, Excluding Transportation

Covariate	Perceived Value					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept (ψ_1^h, ψ_1^c)	-0.522	(0.148)	***	-1.203	(0.243)	***
Factor θ (λ_1^h, λ_1^c)	0.248	(0.312)		2.007	(0.512)	***
Treatment group Z	1.845	(0.155)	***			
HS center quality $\times Z$	0.327	(0.068)	***			
$\theta \times Z$	0.324	(0.386)				
HS center quality	-0.208	(0.056)	***			
Other center quality				0.122	(0.067)	*
<i>Period 2 choice parameters:</i>						
Intercept (ψ_2^h, ψ_2^c)	0.644	(0.097)	***	-0.072	(0.150)	
Period 1 in h (ψ_2^{hh}, ψ_2^{hc})	0.906	(0.103)	***	0.383	(0.163)	**
Period 1 in c (ψ_2^{ch}, ψ_2^{cc})	0.743	(0.168)	***	1.759	(0.303)	***
Factor θ (λ_2^h, λ_2^c)	-0.272	(0.246)		0.754	(0.425)	*
HS center quality	0.072	(0.029)	**			
Other center quality				0.216	(0.057)	***
<i>Other parameters:</i>						
σ_{ν^h}	1.000	-				
σ_{ν^c}	2.178	(0.324)	***			
ρ_{hc}	0.790	(0.128)	***			

Notes : N=1876, Log-likelihood = -3050.69. Covariates including the transportation options ($tp_{ih}, tp_{ic}, tp_{ih} \times Z_i$) are excluded from the choice equations in both periods. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of the age-4 year as the outcome. Alternative-specific choice shocks assumed to follow a bivariate normal distribution with zero means, variance of 1 (normalization) for ν_{ij}^h , variance of $\sigma_{\nu^c}^2$ for ν_{ij}^c , and correlation coefficient ρ_{hc} . Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table B.6: Estimates of the Parameters of the Potential Outcome and Measurement Equations from the Sequential Threshold Model, Excluding Transportation from the Choice Equation

<i>A. Potential outcome equations</i>						
	Intercept			Factor		
Potential outcomes:						
Y^{nn}	3.268	(0.017)	***	0.971	(0.089)	***
$Y^{hn} - Y^{nn}$	0.084	(0.030)	***	0.228	(0.134)	*
$Y^{nh} - Y^{nn}$	0.091	(0.024)	***	-0.105	(0.108)	
$Y^{nc} - Y^{nn}$	0.099	(0.025)	***	0.044	(0.105)	
$Y^{hh} - Y^{nn}$	0.098	(0.019)	***	-0.217	(0.086)	**
$Y^{hc} - Y^{nn}$	0.142	(0.024)	***	-0.172	(0.104)	*
$Y^{ch} - Y^{nn}$	0.049	(0.032)		-0.103	(0.158)	
$Y^{cc} - Y^{nn}$	0.099	(0.026)	***	-0.224	(0.104)	**
Other covariates						
1 period quality HS	0.006	(0.008)				
2 periods quality HS	0.022	(0.007)	***			
1 period quality other center	0.004	(0.010)				
2 periods quality other center	0.023	(0.017)				
<i>B. Measurement equations</i>						
	Intercept			Factor		
PPVT	2.299	(0.008)	***	1.000	-	
WJIII Spelling	3.348	(0.007)	***	0.500	(0.034)	***
WJIII Letter-Word	2.939	(0.017)		0.554	(0.032)	***
<i>C. Factor and error SD</i>						
Factor	0.223	(0.008)	***			
Outcome	0.158	(0.004)	***			
Baseline measurements:						
PPVT	0.293	(0.006)	***			
WJIII Spelling	0.243	(0.004)	***			
WJIII Letter-Word	0.185	(0.004)	***			

Notes : N=1876. Log-likelihood = -3050.69. Covariates including the transportation options ($tp_{ih}, tp_{ic}, tp_{ih} \times Z_i$) are excluded from the choice equations in both periods. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of age-4 year as the outcome. The intercept and the factor loading of the outcome equation differ depending on the program sequence (see equation 9), and are all relative to the parameters of the outcome equation for the nn sequence. h is for HS, c is for other center care, n is for home care. Baseline measurements are test scores measured around the time of randomization (in Fall 2002). The factor loading in the measurement equation for the baseline PPVT test score is normalized to one ($\gamma^{pp}=1$). Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n), to avoid overfitting we restrict $\alpha_t^{cn} = \alpha_t^{hn}$ and $\gamma_t^{cn} = \gamma_t^{hn}$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): ***1%, ** 5%, *10%.

Table B.7: Estimates of the Parameters of the Choice Equations from the Sequential Threshold Model, with Factor Distributed as Mixture Normal

Covariate	Perceived Value					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept (ψ_1^h, ψ_1^c)	-0.624	(0.172)	***	-1.070	(0.268)	***
Factor θ (λ_1^h, λ_1^c)	0.173	(0.333)		1.702	(0.497)	***
Treatment group Z	1.935	(0.168)	***			
HS center quality $\times Z$	0.323	(0.070)	***			
HS center transport $\times Z$	0.152	(0.071)	**			
$\theta \times Z$	0.287	(0.395)				
HS center quality	-0.207	(0.057)	***			
HS center transport	-0.150	(0.052)	***			
Other center quality				0.120	(0.065)	*
Other center transport				-0.092	(0.064)	
<i>Period 2 choice parameters:</i>						
Intercept (ψ_2^h, ψ_2^c)	0.577	(0.131)	***	-0.044	(0.141)	
Period 1 in h (ψ_2^{hh}, ψ_2^{hc})	0.950	(0.110)	***	0.387	(0.154)	**
Period 1 in c (ψ_2^{ch}, ψ_2^{cc})	0.698	(0.190)	***	1.642	(0.342)	***
Factor θ (λ_2^h, λ_2^c)	-0.342	(0.259)		0.699	(0.421)	*
HS center quality	0.080	(0.032)	**			
HS center transport	0.044	(0.032)				
Other center quality				0.182	(0.054)	***
Other center transport				0.208	(0.054)	***
<i>Other parameters:</i>						
σ_{ν^h}	1.000	-				
σ_{ν^c}	1.941	(0.419)	***			
ρ_{hc}	0.697	(0.172)	***			

Notes : N=1876, Log-likelihood = -3024.11. We relax the factor distribution used in estimation by modelling it as a mixture normal distribution instead. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of the age-4 year as the outcome. Alternative-specific choice shocks assumed to follow a bivariate normal distribution with zero means, variance of 1 (normalization) for ν_{ij}^h , variance of $\sigma_{\nu^c}^2$ for ν_{ij}^c , and correlation coefficient ρ_{hc} . Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table B.8: Estimates of the Parameters of the Potential Outcome and Measurement Equations from the Sequential Threshold Model, with Factor Distributed as Mixture Normal

<i>A. Potential outcome equations</i>						
	Intercept			Factor		
Potential outcomes:						
Y^{nn}	3.266	(0.019)	***	0.998	(0.090)	***
$Y^{hn} - Y^{nn}$	0.086	(0.032)	**	0.261	(0.142)	*
$Y^{nh} - Y^{nn}$	0.094	(0.024)	***	-0.131	(0.107)	
$Y^{nc} - Y^{nn}$	0.095	(0.025)	***	-0.013	(0.101)	
$Y^{hh} - Y^{nn}$	0.099	(0.020)	***	-0.228	(0.087)	***
$Y^{hc} - Y^{nn}$	0.144	(0.024)	***	-0.197	(0.104)	*
$Y^{ch} - Y^{nn}$	0.060	(0.034)	*	-0.126	(0.163)	
$Y^{cc} - Y^{nn}$	0.099	(0.026)	***	-0.261	(0.101)	**
Other covariates						
1 period quality HS	0.003	(0.008)				
2 periods quality HS	0.021	(0.007)	***			
1 period quality other center	-0.001	(0.009)				
2 periods quality other center	0.025	(0.016)				
<i>B. Measurement equations</i>						
	Intercept			Factor		
PPVT	2.298	(0.008)	***	1.000	-	
WJIII Spelling	3.348	(0.007)	***	0.497	(0.034)	***
WJIII Letter-Word	2.939	(0.006)	***	0.556	(0.030)	***
<i>C. Factor paramters</i>						
Mixture 1 - SD	0.218	(0.020)	***			
Mixture 2 - Mean	-0.078	(0.027)	***			
Mixture 2 - SD	0.196	(0.014)	***			
Mixture 2 - Weight	0.625	(0.100)	***			
<i>D. Error SD</i>						
Outcome	0.154	(0.004)	***			
Baseline measurements:						
PPVT	0.292	(0.006)	***			
WJIII Spelling	0.243	(0.004)	***			
WJIII Letter-Word	0.184	(0.004)	***			

Notes : N=1876. Log-likelihood = -3024.11. We relax the factor distribution used in estimation by modelling it as a mixture normal distribution instead. The model uses the average of PPVT, WJIII Letter-Word test and WJIII Spelling test scores at the end of age-4 year as the outcome. The intercept and the factor loading of the outcome equation differ depending on the program sequence (see equation 9), and are all relative to the parameters of the outcome equation for the nn sequence. h is for HS, c is for other center care, n is for home care. Baseline measurements are test scores measured around the time of randomization (in Fall 2002). The factor loading in the measurement equation for the baseline PPVT test score is normalized to one ($\gamma^{pp}=1$). Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n), to avoid overfitting we restrict $\alpha_t^{cn} = \alpha_t^{hn}$ and $\gamma_t^{cn} = \gamma_t^{hn}$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): ***1%, ** 5%, *10%.

Table B.9: Estimates of the Parameters of the Choice Equations from the Sequential Threshold Model, Expanding the Measurement System and the Outcome to Include WJIII Applied Problems Test

Covariate	Perceived Value					
	Head Start (h)			Other center care (c)		
<i>Period 1 choice parameters:</i>						
Intercept (ψ_1^h, ψ_1^c)	-0.596	(0.150)	***	-1.122	(0.245)	***
Factor θ (λ_1^h, λ_1^c)	0.301	(0.326)		1.973	(0.501)	***
Treatment group Z	1.915	(0.151)	***			
HS center quality $\times Z$	0.320	(0.068)	***			
HS center transport $\times Z$	0.150	(0.068)	**			
$\theta \times Z$	0.321	(0.404)				
HS center quality	-0.204	(0.056)	***			
HS center transport	-0.147	(0.048)	***			
Other center quality				0.128	(0.065)	*
Other center transport				-0.096	(0.064)	
<i>Period 2 choice parameters:</i>						
Intercept (ψ_2^h, ψ_2^c)	0.596	(0.111)	***	-0.057	(0.143)	
Period 1 in h (ψ_2^{hh}, ψ_2^{hc})	0.943	(0.104)	***	0.375	(0.157)	**
Period 1 in c (ψ_2^{ch}, ψ_2^{cc})	0.724	(0.176)	***	1.697	(0.312)	***
Factor θ (λ_2^h, λ_2^c)	-0.404	(0.25)		0.730	(0.418)	*
HS center quality	0.078	(0.030)	**			
HS center transport	0.044	(0.032)				
Other center quality				0.189	(0.056)	***
Other center transport				0.211	(0.054)	***
<i>Other parameters:</i>						
σ_{ν^h}	1.000	-				
σ_{ν^c}	2.025	(0.347)	***			
ρ_{hc}	0.724	(0.142)	***			

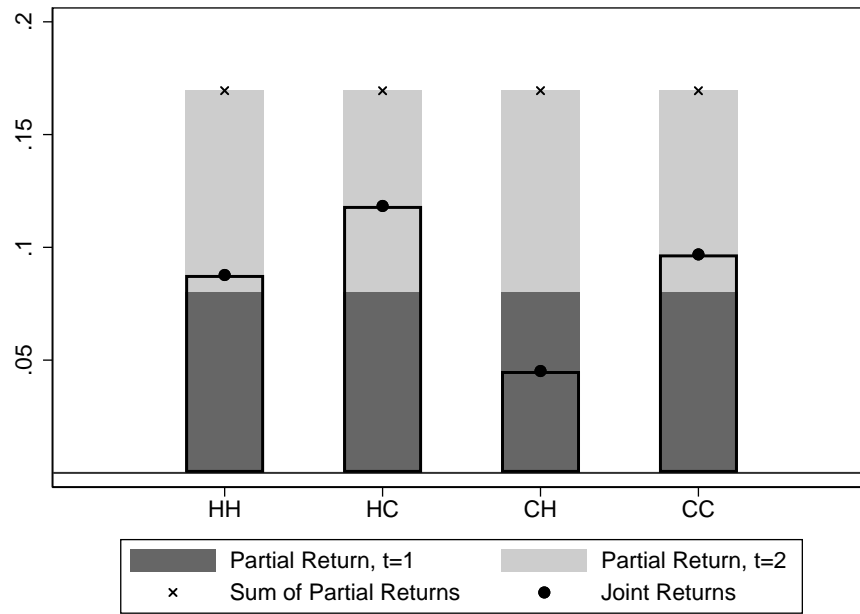
Notes : N=1876, Log-likelihood = -2928.92. The model uses the average of end-of-age-4 scores of PPVT and three WJIII subtests (Letter-Word, Spelling and Applied Problems) as the outcome, and the baseline WJIII Applied Problems test score as an additional measurement of the latent factor. Alternative-specific choice shocks assumed to follow a bivariate normal distribution with zero means, variance of 1 (normalization) for ν_{ij}^h , variance of $\sigma_{\nu^c}^2$ for ν_{ij}^c , and correlation coefficient ρ_{hc} . Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): *** 1%, ** 5%, *10%.

Table B.10: Estimates of the Parameters of the Potential Outcome and Measurement Equations from the Sequential Threshold Model, Expanding the Measurement System and the Outcome to Include WJIII Applied Problems Test

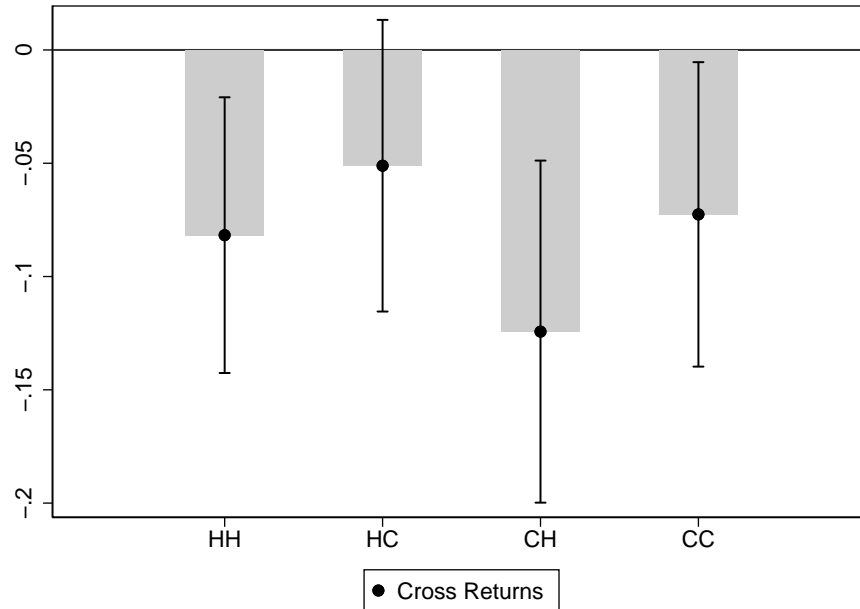
<i>A. Potential outcome equations</i>						
	Intercept			Factor		
Potential outcomes:						
Y^{nn}	3.437	(0.017)	***	0.942	(0.083)	***
$Y^{hn} - Y^{nn}$	0.080	(0.027)	***	0.099	(0.120)	
$Y^{nh} - Y^{nn}$	0.090	(0.021)	***	-0.101	(0.105)	
$Y^{nc} - Y^{nn}$	0.089	(0.021)	***	-0.013	(0.100)	
$Y^{hh} - Y^{nn}$	0.088	(0.017)	***	-0.245	(0.082)	***
$Y^{hc} - Y^{nn}$	0.118	(0.020)	***	-0.197	(0.098)	**
$Y^{ch} - Y^{nn}$	0.045	(0.029)		-0.117	(0.151)	
$Y^{cc} - Y^{nn}$	0.097	(0.023)	***	-0.308	(0.101)	***
Other covariates						
1 period quality HS	0.004	(0.008)				
2 periods quality HS	0.018	(0.007)	***			
1 period quality other center	0.005	(0.009)				
2 periods quality other center	0.023	(0.016)				
<i>B. Measurement equations</i>						
	Intercept			Factor		
PPVT	2.298	(0.008)	***	1.000	-	
WJIII Spelling	3.702	(0.007)	***	0.512	(0.039)	***
WJIII Letter-Word	3.348	(0.007)	***	0.593	(0.037)	***
WJIII Applied Problem	2.939	(0.006)	***	0.596	(0.035)	***
<i>C. Factor and error SD</i>						
Factor	0.212	(0.008)	***			
Outcome	0.151	(0.004)	***			
Baseline measurements:						
PPVT	0.300	(0.006)	***			
WJIII Spelling	0.242	(0.004)	***			
WJIII Letter-Word	0.235	(0.004)	***			
WJIII Applied Problem	0.183	(0.004)	***			

Notes : N=1876. Log-likelihood = -2928.92. The model uses the average of end-of-age-4 scores of PPVT and three WJIII subtests (Letter-Word, Spelling and Applied Problems) as the outcome, and the baseline WJIII Applied Problems test score as an additional measurement of the latent factor. The intercept and the factor loading of the outcome equation differ depending on the program sequence (see equation 9), and are all relative to the parameters of the outcome equation for the nn sequence. h is for HS, c is for other center care, n is for home care. Baseline measurements are test scores measured around the time of randomization (in Fall 2002). The factor loading in the measurement equation for the baseline PPVT test score is normalized to one ($\gamma^{pp}=1$). Because very few children ($\approx 1\%$) in our sample choose program sequence (c, n), to avoid overfitting we restrict $\alpha_t^{cn} = \alpha_t^{hn}$ and $\gamma_t^{cn} = \gamma_t^{hn}$. Standard errors are in parentheses. Significance level (t-test for testing if each parameter=0): ***1%, ** 5%, *10%.

Figure B.2: Average Partial, Joint and Cross Returns from Sequential Program Participation, Expanding the Measurement System and the Outcome to Include WJIII Applied Problems Test



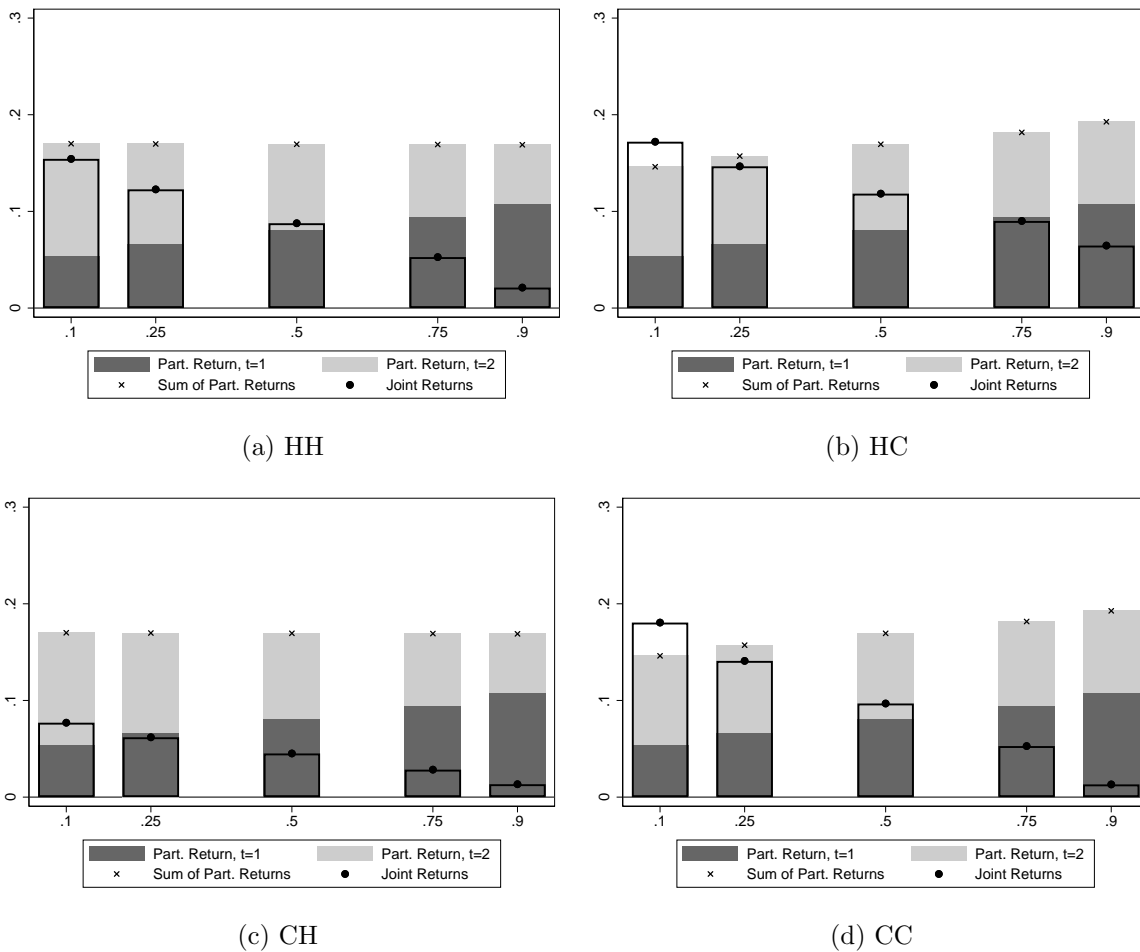
(a) Partial and Joint Returns



(b) Cross>Returns (with 95% C.I.)

Notes : The outcome analyzed is the average of end-of-age-4 scores of PPVT and three WJIII subtests (Letter-Word, Spelling and Applied Problems). We expand the measurement system to include baseline WJIII Applied Problems test score at randomization as an additional measurement of the latent factor. Each bar represents the joint and partial return of different program sequences. In Panel (a) each bar represents the joint and partial return of different program sequences. In Panel (b) each bar represents the cross-returns, i.e. the difference between the joint return and the sum of the partial returns at $t = 1$ and $t = 2$, with the respective 95% confidence intervals. The first bar (HH) shows the return of receiving HS both at $t = 1$ and at $t = 2$. The second bar (HC) shows the return of receiving HS at $t = 1$ and other center care at $t = 2$. The third bar (CH) shows the return of receiving other center care at $t = 1$ and HS at $t = 2$. The fourth bar (CC) shows the return of receiving other center care in both periods. All these returns use two periods in home care as the baseline program sequence for comparison (see Section 5.2 for the exact definition of partial, joint and cross-returns).

Figure B.3: Partial, Joint and Cross Returns from Sequential Program Participation, Expanding the Measurement System and the Outcome to Include WJIII Applied Problems Test, Different Values of θ



Notes : The outcome analyzed is the average of end-of-age-4 scores of PPVT and three WJIII subtests (Letter-Word, Spelling and Applied Problems). We expand the measurement system to include baseline WJIII Applied Problems test score at randomization as an additional measurement of the latent factor. Within each panel, the different bars represent the joint and partial returns of the same program sequence, but for different levels of baseline ability θ (“Part. Return” in label stands for “Partial Returns”). In particular, we report the returns for the 10th, the 25th, the 50th, the 75th and the 90th percentile of the θ distribution. Panel (a) is the return for two periods in HS (HH), panel (b) for one period in HS followed by other center care (HC), panel (c) for one period in other center care followed by HS (CH) and panel (d) for two periods in other center care (CC).

C Appendix: Further Note on Identification

Measurement equations. Let $M_i^{pp*} \equiv M_i^{pp} - \alpha^{pp} = \theta_i + e_i^{pp}$ and $M_i^{ws*} \equiv \frac{M_i^{ws} - \alpha^{ws}}{\gamma^{ws}} = \theta_i + \frac{e_i^{ws}}{\gamma^{ws}}$. By Kotlarski's lemma, the characteristic function of θ_i is

$$\varphi_\theta(s) = \exp \int_0^s \frac{E[\mathbf{i}M^{ws*} e^{\mathbf{i}\zeta M^{pp*}}]}{E[e^{\mathbf{i}\zeta M^{pp*}}]} d\zeta \quad (\text{C.1})$$

where $\mathbf{i} = \sqrt{-1}$. By inverse Fourier transform, the density of θ_i is $f_\theta(\theta) = (2\pi)^{-1} \int \varphi_\theta(s) e^{-\mathbf{i}s\theta} ds$. Due to independence of θ and e^{pp} , we have $\varphi_{M^{pp*}}(s) = \varphi_\theta(s) \varphi_{e^{pp}}(s)$. We obtain the density of e^{pp} by deconvolution:

$$f_{e^{pp}}(e^{pp}) = (2\pi)^{-1} \int \frac{\varphi_{M^{pp*}}(s)}{\varphi_\theta(s)} e^{-\mathbf{i}s e^{pp}} ds. \quad (\text{C.2})$$

The density of e^{ws} and e^{wl} can be obtained by similar deconvolution procedures.

Choice equations at $\mathbf{t}=2$. The following moment can be expressed as a function of loadings λ_2^h, λ_2^c and average derivatives:

$$E \left[\mathbf{1}\{K_2 = n\} \frac{-\partial \ln f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0)}{\partial M^{pp}} \right] \quad (\text{C.3})$$

$$= E \left[E(\mathbf{1}\{K_2 = n\} | \theta, M^{pp}, K_1 = n, \tilde{\mathbf{X}} = 0) \frac{-\partial \ln f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0)}{\partial M^{pp}} \right] \quad (\text{C.4})$$

$$= - \iint_{-\infty}^{\infty} G_{\nu^h, \nu^c}(-\psi_2^h - \lambda_2^h \theta, -\psi_2^c - \lambda_2^c \theta) \frac{\partial f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0)}{\partial M^{pp}} d\theta dM^{pp} \quad (\text{C.5})$$

$$= \iint_{-\infty}^{\infty} G_{\nu^h, \nu^c}(-\psi_2^h - \lambda_2^h(M^{pp} - e^{pp}), -\psi_2^c - \lambda_2^c(M^{pp} - e^{pp})) \times \frac{\partial f(M^{pp} - e^{pp}, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0)}{\partial M^{pp}} \Big|_{-1} de^{pp} dM^{pp} \quad (\text{C.6})$$

$$= - \iint_{-\infty}^{\infty} \frac{\partial G_{\nu^h, \nu^c}(-\psi_2^h - \lambda_2^h(M^{pp} - e^{pp}), -\psi_2^c - \lambda_2^c(M^{pp} - e^{pp}))}{\partial M^{pp}} \times f(M^{pp} - e^{pp}, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0) de^{pp} dM^{pp} \quad (\text{C.7})$$

$$= \lambda_2^h E \left[\frac{\partial G_{\nu^h, \nu^c}(s_1, s_2)}{\partial s_1} \right] + \lambda_2^c E \left[\frac{\partial G_{\nu^h, \nu^c}(s_1, s_2)}{\partial s_2} \right], \quad (\text{C.8})$$

where $f(\cdot)$ is the joint density of θ and M^{pp} conditional on K_1 and $\tilde{\mathbf{X}}$, the third equality is due to change of variable $\theta = M^{pp} - e^{pp}$, and the fourth equality is due to integration by parts assuming $f(\cdot)$

vanishes to zero at the boundary (see Stoker (1986)). Similarly:⁶¹

$$E \left[\mathbf{1}\{K_2 = c\} \frac{-\partial \ln f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0)}{\partial M^{pp}} \right] \quad (\text{C.9})$$

$$= - \int \int_{-\infty}^{\infty} G_{\nu^h - \nu^c, -\nu^c}((\psi_2^c - \psi_2^h) + (\lambda_2^c - \lambda_2^h)\theta, \psi_2^c + \lambda_2^c\theta) \frac{\partial f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0)}{\partial M^{pp}} d\theta dM^{pp} \quad (\text{C.10})$$

$$= \lambda_2^h E \left[\frac{\partial G_{\nu^h - \nu^c, -\nu^c}(s_1, s_2)}{\partial s_1} \right] - \lambda_2^c \left[E \left[\frac{\partial G_{\nu^h - \nu^c, -\nu^c}(s_1, s_2)}{\partial s_1} + \frac{\partial G_{\nu^h - \nu^c, -\nu^c}(s_1, s_2)}{\partial s_2} \right] \right]. \quad (\text{C.11})$$

To identify the intercepts ψ^h , ψ^c , two moments are needed to fix the location:

$$E[\mathbf{1}\{K_2 = n\}] = \int \int_{-\infty}^{\infty} G_{\nu^h, \nu^c}(-\psi_2^h - \lambda_2^h\theta, -\psi_2^c - \lambda_2^c\theta) f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0) d\theta dM^{pp} \quad (\text{C.12})$$

$$E[\mathbf{1}\{K_2 = c\}] = \int \int_{-\infty}^{\infty} G_{\nu^h - \nu^c, -\nu^c}((\psi_2^c - \psi_2^h) + (\lambda_2^c - \lambda_2^h)\theta, \psi_2^c + \lambda_2^c\theta) f(\theta, M^{pp} | K_1 = n, \tilde{\mathbf{X}} = 0) d\theta dM^{pp} \quad (\text{C.13})$$

Identification with Only Two Baseline Measurements. Suppose that only two types of baseline test scores are available in the data, including PPVT ($M = pp$) and WJIII Spelling test ($M = ws$). First, consider the (centered) measurement equations:

$$\epsilon_i^{pp} \equiv \theta_i + e_i^{pp} \quad (\text{C.14})$$

$$\epsilon_i^{ws} \equiv \gamma^{ws}\theta_i + e_i^{ws} \quad (\text{C.15})$$

Following Fruehwirth, Navarro, and Takahashi (2016), higher-order cross moments can be used to identify γ^{ws} (see also Bonhomme and Robin (2009)):

- If θ_i is asymmetric: use $E((\epsilon_i^{pp})^2 \epsilon_i^{ws}) = \gamma^{ws} E(\theta_i^3)$ and $E(\epsilon_i^{pp} (\epsilon_i^{ws})^2) = (\gamma^{ws})^2 E(\theta_i^3)$ to identify $\gamma^{ws} = \frac{E(\epsilon_i^{pp} (\epsilon_i^{ws})^2)}{E((\epsilon_i^{pp})^2 \epsilon_i^{ws})}$;
- If θ_i is symmetric and kurtotic (i.e., $E(\theta_i^4) \neq 3[E(\theta_i^2)]^2$): obtain $(\gamma^{ws})^2 = \frac{E(\epsilon_i^{pp} (\epsilon_i^{ws})^3) - 3E(\epsilon_i^{pp} \epsilon_i^{ws}) E((\epsilon_i^{ws})^2)}{E((\epsilon_i^{pp})^3 \epsilon_i^{ws}) - 3E(\epsilon_i^{pp} \epsilon_i^{ws}) E((\epsilon_i^{pp})^2)}$ where the sign of γ^{ws} is the same as that of $E(\epsilon_i^{pp} \epsilon_i^{ws})$ (see Fruehwirth, Navarro, and Takahashi (2016), Appendix B3);

and then the distributions of θ_i , e_i^{pp} and e_i^{ws} are nonparametrically identified via Kotlarski (1967). A special case is when θ_i , e_i^{pp} and e_i^{ws} are all normally distributed. In this case, γ^{ws} is not identified from equations (C.14) and (C.15) alone (see e.g., Reiersøl (1950)).⁶²

⁶¹The characteristic function $\varphi_{(\nu^h - \nu^c, -\nu^c)}(s, w) = E(e^{i[s(\nu^h - \nu^c) + w(-\nu^c)]}) = E(e^{i[s\nu^h - (s-w)\nu^c]}) = \varphi_{(\nu^h, \nu^c)}(s, -(s+w))$.

⁶²If θ_i is normal but e_i^{pp} and e_i^{ws} are nonnormal, γ^{ws} is still identified from equations (C.14) and (C.15) alone (Reiersøl (1950)).

D Policy Implications: Additional Details on Optimal Pathways

Suppose that the goal of a social planner is to maximize the population-mean test score at the end of age 4, by assigning program sequence d for each individual at the start of age 3.

We consider the solution to this problem under two cases. In the first instance, we assume that the planner observes θ_i , the latent baseline ability for each individual perfectly. Conditional on the value of θ_i , he then assigns the optimal program sequence \hat{d} in order to maximize the expected potential outcome:

$$\hat{d} = \arg \max \beta_d(\theta) = \arg \max E(Y_i^d | \theta_i = \theta)$$

where the expectation operator integrates out the outcome shocks. $\beta_d(\theta)$ can be evaluated easily given that the outcome shocks are additive and separable in the potential outcome equation. \hat{d} is an implicit function of θ . The resulting average test score in the population is given by

$$\int \beta_{\hat{d}}(\theta) dF(\theta) \tag{D.1}$$

We call \hat{d} the first-best assignment which maximizes the average expected potential outcome in the population at the end of age 4.

In the second instance, suppose that the social planner does not observe θ_i but only observe the baseline test scores. In this case, the optimal program assignment is based on a specific baseline measurement M_i . Define

$$\tilde{d} = \arg \max \beta_d(m) = \arg \max E(Y_i^d | M_i = m)$$

where \tilde{d} is the optimal program sequence that maximizes the expected potential outcome conditional on a given measurement (\tilde{d} is an implicit function of m). Then the resulting mean test score in the population is given by

$$\int \beta_{\tilde{d}}(m) dF(m) \tag{D.2}$$

The difference between (D.2) and (D.1) is the efficiency loss from imperfect information about individual's latent ability. The efficiency loss is decreasing with the signal-to-noise ratio of the specific measurement being used.⁶³

⁶³We evaluate $\beta_d(m)$ by simulating 30 draws of θ and measurement errors per person in the sample. We divide the simulated baseline measurements into 100 equal-sized bins. Within each percentile bin, we compute the mean potential outcomes and derive the optimal program path \tilde{d} .

E Structural Model: Flow Utility

The flow utilities in period $t = 1$ are specified as:⁶⁴

$$\tilde{u}_{i1}^h = \tilde{\psi}_1^h + \tilde{\beta}_1^h tp_{ih} + \tilde{f}(Z_i, q_{ih}, tp_{ih}, \theta_i) + \tilde{\lambda}_1^h \theta_i + \tilde{\nu}_{i1}^h \quad (\text{E.1})$$

$$\tilde{u}_{i1}^c = \tilde{\psi}_1^c + \tilde{\beta}_1^c tp_{ic} + \tilde{\lambda}_1^c \theta_i + \tilde{\nu}_{i1}^c \quad (\text{E.2})$$

$$\tilde{u}_{i1}^n = 0 \quad (\text{E.3})$$

with

$$\tilde{f}(Z_i, q_{ih}, tp_{ih}, \theta_i) = \left(\tilde{\beta}_Z + \tilde{\beta}_{Ztp} tp_{ih} + \tilde{\beta}_{Z\theta} \theta_i \right) Z_i + \tilde{\beta}_q q_{ih} (1 - Z_i). \quad (\text{E.4})$$

Compared to the threshold model (equations (2)-(5)), center quality variables (q_{ih}, q_{ic}) are excluded from the utility flows. An exception is $q_{ih}(1 - Z_i)$, i.e., HS center quality can affect the HS utility flow of the control group. As discussed in Section 4, children in the control group tend to enroll in low-quality HS centers, so we expect $\tilde{\beta}_q < 0$. More broadly, the parameter $\tilde{\beta}_q$ captures implicit barrier costs for the control group to crossover to enroll in HS in period-1; excluding it from the period-1 utility may be a strenuous assumption. By comparison, we assume that HS quality drives HS enrollment for the treatment group via expectations of improved future outcomes. These exclusion restrictions identify the discount factor δ .⁶⁵

The flow utilities in period $t = 2$ are given by:

$$\tilde{u}_{i2}^h = \tilde{\psi}_2^h + \tilde{\psi}_2^{hh} \mathbf{1}(K_{i1} = h) + \tilde{\psi}_2^{ch} \mathbf{1}(K_{i1} = c) + \tilde{\beta}_2^h tp_{ih} + \tilde{\lambda}_2^h \theta_i + \tilde{\nu}_{i2}^h \quad (\text{E.5})$$

$$\tilde{u}_{i2}^c = \tilde{\psi}_2^c + \tilde{\psi}_2^{hc} \mathbf{1}(K_{i1} = h) + \tilde{\psi}_2^{cc} \mathbf{1}(K_{i1} = c) + \tilde{\beta}_2^c tp_{ic} + \tilde{\lambda}_2^c \theta_i + \tilde{\nu}_{i2}^c \quad (\text{E.6})$$

$$\tilde{u}_{i2}^n = 0 \quad (\text{E.7})$$

The exclusion of q_{ih} and q_{ic} from the period-2 flow utility identifies the terminal value scaling factor κ .

⁶⁴We use accent $\tilde{\cdot}$ to emphasize that these parameters are different to those in the sequential threshold model.

⁶⁵Because the utility and EV have the same scale in equation (32), δ has a conventional range between 0 and 1. Specifically, the utility scale at $t = 1, 2$, which reflects choices, are normalized with respect to the h alternative ($\text{var}(\tilde{\nu}_{ij}^h) = 1$).

F Computational and estimation procedures

This section provides additional details to the computation of the expected value function in the structural model and estimation procedures used for both the threshold model and the structural model. The discussion is centered around the most general model (see Appendix Section B.1.1 and Figure B.1).

F.1 Computation of the expected value function

For node $j \in \mathcal{J} = \{\mathbf{o}, \mathbf{h}, \mathbf{c}, \mathbf{n}\}$, let \bar{u}_{ij}^k denote the alternative-specific utility flow, *exclusive* of the preference shock, for alternative $k \in \{h, c, n\}$: $\bar{u}_{ij}^h = \tilde{\psi}_j^h + \tilde{f}_j(Z_i, \mathbf{X}_i, \theta_i) + \tilde{\lambda}_j^h \theta_i$; $\bar{u}_{ij}^c = \tilde{\psi}_j^c + \tilde{\lambda}_j^c \theta_i$; $\bar{u}_{ij}^n = 0$. Now consider the decision problem at $t = 2$ for an individual having chosen $K_{i1} = k_1 \in \{h, c, n\}$ at $t = 1$ and is at node $k_1 \in \{\mathbf{h}, \mathbf{c}, \mathbf{n}\}$ (e.g., $k_1 = \mathbf{h}$ if and only if $k_1 = h$). Conditional on the state variables $K_{i1}(= k_1)$, \mathbf{X}_i , θ_i , the probabilities for program choice at $t = 2$, K_{i2} , are (suppress \mathbf{X}_i and θ_i for notational simplicity):

$$\begin{aligned} Pr(K_{i2} = h | K_{i1} = k_1) &= Pr[\tilde{\nu}_{ik_1}^h - \tilde{\nu}_{ik_1}^c > (\bar{u}_{ik_1}^c - \bar{u}_{ik_1}^h) + \kappa(\bar{Y}_{iT}^{k_1, c} - \bar{Y}_{iT}^{k_1, h}); \\ &\quad \tilde{\nu}_{ik_1}^h > -\bar{u}_{ik_1}^h + \kappa(\bar{Y}_{iT}^{k_1, n} - \bar{Y}_{iT}^{k_1, h}) | K_{i1} = k_1] \end{aligned} \quad (\text{F.1})$$

$$\begin{aligned} Pr(K_{i2} = c | K_{i1} = k_1) &= Pr[\tilde{\nu}_{ik_1}^c - \tilde{\nu}_{ik_1}^h > (\bar{u}_{ik_1}^h - \bar{u}_{ik_1}^c) + \kappa(\bar{Y}_{iT}^{k_1, h} - \bar{Y}_{iT}^{k_1, c}); \\ &\quad \tilde{\nu}_{ik_1}^c > -\bar{u}_{ik_1}^c + \kappa(\bar{Y}_{iT}^{k_1, n} - \bar{Y}_{iT}^{k_1, c}) | K_{i1} = k_1] \end{aligned} \quad (\text{F.2})$$

$$Pr(K_{i2} = n | K_{i1} = k_1) = 1 - Pr(K_{i2} = h | K_{i1} = k_1) - Pr(K_{i2} = c | K_{i1} = k_1) \quad (\text{F.3})$$

which are computed by the Geweke-Hajivassiliou-Keane (GHK) simulator (Keane, 1994). For each $k_1 \in \{h, c, n\}$, we then compute the expected value function $E_1 V_{i2}(\cdot)$, which is a probability-weighted average of the value of each future path:

$$\begin{aligned} &E_1 V_{i2}(K_{i1} = k_1; \tilde{\nu}_{ik_1}^h, \tilde{\nu}_{ik_1}^c) \\ &= \sum_{k_2 \in \{h, c\}} Pr(K_{i2} = k_2 | K_{i1} = k_1) E_1 \left(\bar{u}_{ik_1}^{k_2} + \kappa \bar{Y}_{iT}^{k_1, k_2} + \tilde{\nu}_{ik_1}^{k_2} | K_{i2} = k_2, K_{i1} = k_1 \right) + Pr(K_{i2} = n | K_{i1} = k_1) \kappa \bar{Y}_{iT}^{k_1, n} \\ &= \kappa \bar{Y}_{iT}^{k_1, n} + \sum_{k_2 \in \{h, c\}} Pr(K_{i2} = k_2 | K_{i1} = k_1) \left(\bar{u}_{ik_1}^{k_2} + \kappa \left(\bar{Y}_{iT}^{k_1, k_2} - \bar{Y}_{iT}^{k_1, n} \right) \right) \\ &\quad + \sum_{k_2 \in \{h, c\}} Pr(K_{i2} = k_2 | K_{i1} = k_1) E_1(\tilde{\nu}_{ik_1}^{k_2} | K_{i2} = k_2, K_{i1} = k_1) \end{aligned} \quad (\text{F.4})$$

The first two terms reflect the weighted average of the *deterministic* value of each path. The last term contains expectations of truncated distributions of preference shocks, i.e., selection into alternatives that have larger realizations of shocks. To compute these expectations we invoke an analytical result in Rosenbaum (1961) and Tallis (1961):

$$\begin{aligned} Pr(X > a, Y > b) E(Y | X > a, Y > b) &= \left[r \phi(a) \Phi \left(-\frac{b - ra}{\sqrt{1 - r^2}} \right) + \phi(b) \Phi \left(-\frac{a - rb}{\sqrt{1 - r^2}} \right) \right] \\ &\equiv \bar{\Lambda}(a, b, r) \end{aligned} \quad (\text{F.5})$$

where X and Y follow a standard bivariate normal distribution with zero means, unit variances and correlation coefficient r ; $\phi(\cdot)$ and $\Phi(\cdot)$ are the PDF and CDF of the standard normal distribution, respectively. We have:

$$\begin{aligned}
& Pr(K_{i2} = h | K_{i1} = k_1) E_1(\tilde{\nu}_{ik_1}^h | K_{i2} = h, K_{i1} = k_1) \\
&= Pr(K_{i2} = h | K_{i1} = k_1) E_1(\tilde{\nu}_{ik_1}^h | \frac{\tilde{\nu}_{ik_1}^h - \tilde{\nu}_{ik_1}^c}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}} > a_1; \tilde{\nu}_{ik_1}^h > b_1; K_{i1} = k_1) \\
&= \bar{\Lambda}(a_1, b_1, r_1)
\end{aligned} \tag{F.6}$$

where $a_1 = \frac{(\bar{u}_{ik_1}^c - \bar{u}_{ik_1}^h) + \kappa(\bar{Y}_{iT}^{k_1, c} - \bar{Y}_{iT}^{k_1, h})}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}}$, $b_1 = -\bar{u}_{ik_1}^h + \kappa(\bar{Y}_{iT}^{k_1, n} - \bar{Y}_{iT}^{k_1, h})$, and $r_1 = \frac{1 - \tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c}}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}}$ is the correlation coefficient between $\frac{\tilde{\nu}_{ik_1}^h - \tilde{\nu}_{ik_1}^c}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}}$ and $\tilde{\nu}_{ik_1}^h$. Similarly,

$$\begin{aligned}
& Pr(K_{i2} = c | K_{i1} = k_1) E_1(\tilde{\nu}_{ik_1}^c | K_{i2} = c, K_{i1} = k_1) \\
&= \tilde{\sigma}_{\nu^c} Pr(K_{i2} = c | K_{i1} = k_1) E_1(\frac{\tilde{\nu}_{ik_1}^c}{\tilde{\sigma}_{\nu^c}} | \frac{\tilde{\nu}_{ik_1}^c - \tilde{\nu}_{ik_1}^h}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}} > a_2; \frac{\tilde{\nu}_{ik_1}^c}{\tilde{\sigma}_{\nu^c}} > b_2; K_{i1} = k_1) \\
&= \tilde{\sigma}_{\nu^c} \bar{\Lambda}(a_2, b_2, r_2)
\end{aligned} \tag{F.7}$$

where $a_2 = \frac{(\bar{u}_{ik_1}^h - \bar{u}_{ik_1}^c) + \kappa(\bar{Y}_{iT}^{k_1, h} - \bar{Y}_{iT}^{k_1, c})}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}}$, $b_2 = \frac{-\bar{u}_{ik_1}^c + \kappa(\bar{Y}_{iT}^{k_1, n} - \bar{Y}_{iT}^{k_1, c})}{\tilde{\sigma}_{\nu^c}}$, and $r_2 = \frac{\tilde{\sigma}_{\nu^c} - \tilde{\rho}_{hc}}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}}$ is the correlation coefficient between $\frac{\tilde{\nu}_{ik_1}^c - \tilde{\nu}_{ik_1}^h}{\sqrt{\tilde{\sigma}_{\nu^c}^2 - 2\tilde{\rho}_{hc}\tilde{\sigma}_{\nu^c} + 1}}$ and $\frac{\tilde{\nu}_{ik_1}^c}{\tilde{\sigma}_{\nu^c}}$.

This analytical result reduces the computational burden of the expected value function, which now largely involves the computation of choice probabilities in equations (F.1) and (F.2) for each $k_1 \in \{h, c, n\}$. Because Z_i and \mathbf{X}_i differ for each individual, we solve the dynamic programming problem for each individual in the sample separately.

F.2 Model estimation procedure

The model is estimated by the method of maximum likelihood. We first describe the estimation of the structural model, and then describe the estimation of the sequential threshold model as a special case. In both models, the researcher observes the individual's experimental treatment group status Z_i , covariates \mathbf{X}_i , baseline measurements M_i^M ($M \in \{pp, wj\}$), program choice sequence (K_{i1}, K_{i2}) , and outcome Y_{iT} .

The parameters of the structural model are listed as follows: preferences in utility ($\tilde{\psi}_j^h, \tilde{\psi}_j^c, \tilde{\lambda}_j^h, \tilde{\lambda}_j^c, \tilde{\beta}_j^{(\cdot)}$ for $j \in \{o, h, c, n\}$; $\tilde{\sigma}_{\nu^c}, \tilde{\rho}_{hc}$), terminal value scaling factor and discount factor (κ, δ), baseline measurement equations (α^M, σ_e^M for $M \in \{pp, wj\}$; γ^{wj}), outcome equations ($\alpha_T^d, \beta_T^d, \gamma_T^d$ for $d \in \{(h, h), (h, c), (h, n), (c, h), (c, c), (c, n), (n, h), (n, c), (n, n)\}$; σ_ε), and the factor distribution (σ_θ).

For each iteration in the parameter space, computation of the likelihood for individual i consists of three nested loops. The inner loop computes the likelihood for baseline measurements and for choices and outcomes at $t = 1, 2, T$ given the individual's state variables, factor θ_i , and the expected value function obtained from the backward recursion procedure. The middle loop carries out the backward

recursion procedure of the dynamic programming problem given the individual's state variables and factor θ_i . The outer loop integrates out the likelihood with respect to the factor distribution.

For the inner loop, the choice probability at $t = 1$ (node $j = o$) is:

$$Pr(K_{i1} = k_1 | Z_i, \mathbf{X}_i, \theta_i) \tag{F.8}$$

$$= \begin{cases} Pr[\tilde{v}_{io}^h - \tilde{v}_{io}^c > (\bar{u}_{io}^c - \bar{u}_{io}^h) + \delta(E_1 V_{i2}(c, \mathbf{X}_i, \theta_i; \tilde{v}_{ic}^h, \tilde{v}_{ic}^c) - E_1 V_{i2}(h, \mathbf{X}_i, \theta_i; \tilde{v}_{ih}^h, \tilde{v}_{ih}^c)); \\ \quad \tilde{v}_{io}^h > -\bar{u}_{io}^h + \delta(E_1 V_{i2}(n, \mathbf{X}_i, \theta_i; \tilde{v}_{in}^h, \tilde{v}_{in}^c) - E_1 V_{i2}(h, \mathbf{X}_i, \theta_i; \tilde{v}_{ih}^h, \tilde{v}_{ih}^c)) | Z_i, \mathbf{X}_i, \theta_i] & \text{if } k_1 = h \\ Pr[\tilde{v}_{io}^c - \tilde{v}_{io}^h > (\bar{u}_{io}^h - \bar{u}_{io}^c) + \delta(E_1 V_{i2}(h, \mathbf{X}_i, \theta_i; \tilde{v}_{ih}^h, \tilde{v}_{ih}^c) - E_1 V_{i2}(c, \mathbf{X}_i, \theta_i; \tilde{v}_{ic}^h, \tilde{v}_{ic}^c)); \\ \quad \tilde{v}_{io}^c > -\bar{u}_{io}^c + \delta(E_1 V_{i2}(n, \mathbf{X}_i, \theta_i; \tilde{v}_{in}^h, \tilde{v}_{in}^c) - E_1 V_{i2}(c, \mathbf{X}_i, \theta_i; \tilde{v}_{ic}^h, \tilde{v}_{ic}^c)) | Z_i, \mathbf{X}_i, \theta_i] & \text{if } k_1 = c \\ 1 - Pr(K_{i1} = h | Z_i, \mathbf{X}_i, \theta_i) - Pr(K_{i1} = c | Z_i, \mathbf{X}_i, \theta_i) & \text{if } k_1 = n \end{cases}$$

and the choice probability at $t = 2$ (node $j = k_1 \in \{h, c, n\}$) is:

$$Pr(K_{i2} = k_2 | K_{i1} = k_1, Z_i, \mathbf{X}_i, \theta_i) \tag{F.9}$$

$$= Pr(K_{i2} = k_2 | K_{i1} = k_1, \mathbf{X}_i, \theta_i)$$

$$= \begin{cases} Pr[\tilde{v}_{ik_1}^h - \tilde{v}_{ik_1}^c > (\bar{u}_{ik_1}^c - \bar{u}_{ik_1}^h) + \kappa(\bar{Y}_{iT}^{k_1, c}(\mathbf{X}_i, \theta_i) - \bar{Y}_{iT}^{k_1, h}(\mathbf{X}_i, \theta_i)); \\ \quad \tilde{v}_{ik_1}^h > -\bar{u}_{ik_1}^h + \kappa(\bar{Y}_{iT}^{k_1, n}(\mathbf{X}_i, \theta_i) - \bar{Y}_{iT}^{k_1, h}(\mathbf{X}_i, \theta_i)) | K_{i1} = k_1, \mathbf{X}_i, \theta_i] & \text{if } k_2 = h \\ Pr[\tilde{v}_{ik_1}^c - \tilde{v}_{ik_1}^h > (\bar{u}_{ik_1}^h - \bar{u}_{ik_1}^c) + \kappa(\bar{Y}_{iT}^{k_1, h}(\mathbf{X}_i, \theta_i) - \bar{Y}_{iT}^{k_1, c}(\mathbf{X}_i, \theta_i)); \\ \quad \tilde{v}_{ik_1}^c > -\bar{u}_{ik_1}^c + \kappa(\bar{Y}_{iT}^{k_1, n}(\mathbf{X}_i, \theta_i) - \bar{Y}_{iT}^{k_1, c}(\mathbf{X}_i, \theta_i)) | K_{i1} = k_1, \mathbf{X}_i, \theta_i] & \text{if } k_2 = c \\ 1 - Pr(K_{i2} = h | K_{i1} = k_1, \mathbf{X}_i, \theta_i) - Pr(K_{i2} = c | K_{i1} = k_1, \mathbf{X}_i, \theta_i) & \text{if } k_2 = n \end{cases}$$

where we suppress the arguments of $\bar{u}_{ij}^{(\cdot)}$ for notational simplicity and the choice probability at $t = 2$ does not depend on Z_i when conditioned on $K_{i1} = k_1$ (see Section 6 for specifications). We compute each choice probability using the GHK simulator. Define the likelihood contribution of individual i given the observables and unobserved factor θ_i as

$$L_i(\theta) \equiv Pr(K_{i1} = k_1 | Z_i = z, \mathbf{X}_i = x, \theta_i = \theta) Pr(K_{i2} = k_2 | K_{i1} = k_1, \mathbf{X}_i = x, \theta_i = \theta) \times \\ f_{Y_{iT} | K_{i1}, K_{i2}, hsq_i, \theta_i}(y | k_1, k_2, hsq, \theta) f_{M_i^{pp} | \theta_i}(m^{pp} | \theta) f_{M_i^{wj} | \theta_i}(m^{wj} | \theta) \tag{F.10}$$

where $f_{Y_{iT} | K_{i1}, K_{i2}, hsq_i, \theta_i}(\cdot)$, $f_{M_i^{pp} | \theta_i}(\cdot)$, and $f_{M_i^{wj} | \theta_i}(\cdot)$ are the conditional density functions of the outcome and baseline measurement equations, respectively. We then compute $\ell_i \equiv \ln \int_{-\infty}^{\infty} L_i(\theta) f_{\theta_i}(\theta) d\theta$ where $f_{\theta_i}(\cdot)$ is the density function of the unobserved factor and the integration is carried out using the Gauss-Hermite quadrature. The log-likelihood function is $\ell \equiv \sum_{i=1}^N \ell_i$. Denote the parameter vector by Θ and the maximum likelihood estimates by $\hat{\Theta}$. We compute the standard errors using the outer-product-of-gradient (OPG) estimator (sometimes called the BHHH estimator (Berndt, Hall, Hall, and Hausman (1974))): $\widehat{var}(\hat{\Theta}) = [\frac{1}{N} \sum_{i=1}^N (\frac{\partial}{\partial \Theta} \ell_i(\hat{\Theta})) (\frac{\partial}{\partial \Theta} \ell_i(\hat{\Theta}))']^{-1}$, which is consistent for the inverse Fisher information.

The sequential threshold model is estimated separately, but its estimation procedure is the same as in the structural model except that there is no dynamic programming involved (effectively, both κ and δ are zero). The sequential threshold model contains the following parameters in the perceived value

equations (in lieu of flow utility parameters, terminal value scaling factor and discount factor): $\psi_j^h, \psi_j^c, \lambda_j^h, \lambda_j^c, \beta_j^{(\cdot)}$ for $j \in \{\mathbf{o}, \mathbf{h}, \mathbf{c}, \mathbf{n}\}$; $\sigma_{\nu^c}, \rho_{hc}$. The parameters in the baseline measurement equations, outcome equations, and the factor distribution are defined in the same manner as in the structural model. For each iteration in the parameter space, computation of the likelihood for individual i consists of two nested loops only. The inner loop computes the likelihood for baseline measurements and for choices and outcomes at $t = 1, 2, T$ given the individual's state variables and factor θ_i (equations (F.8), (F.9) and (F.10)). The outer loop integrates out the likelihood with respect to the factor distribution.