

The Age 14 Follow-up of the Preparing for Life Study*

By ORLA DOYLE

University College Dublin

UCD School of Economics & UCD Geary Institute for Public Policy

Technical Report

May 2025

Summary

This report examines the impact of *Preparing for Life (PFL)*, an Irish prenatally commencing home visiting programme, ten years after the intervention ended. The intervention involved regular visits from trained home visitors from pregnancy until school entry to support parents around child development and parenting. Previous reports of the *PFL* trial found that the programme was effective in boosting children's cognitive skills, with smaller effects on some dimensions of health and socio-emotional skills. The Age 14 Follow-up found that the programme has a sustained and long-term effect on children's cognitive development, with large effect sizes of 0.70 standard deviations (SDs). Significant effects were also found on working memory, attention, and educational expectations, however there were relatively few effects on health or socio-emotional outcomes. There was some evidence that the programme reduced children's waist-to-height ratio and improved the parent-child relationship. All results were estimated using permutation-based hypothesis testing which accounts for attrition using inverse probability weighting and multiple hypothesis testing using the stepdown procedure. While 43% of the original sample participated at the Age 14 Follow-up, the high and low treatment groups were still balanced on all key baseline characteristics. This is one of the few experimental home visiting programmes that tracks participants into adolescence and finds evidence that the intervention continues to have a significant impact on important dimensions of children's skills.

* The evaluation of the *Preparing for Life* programme was funded by the Northside Partnership through the Department of Children, Equality, Disability, Integration and Youth and The Atlantic Philanthropies. Additional support was provided by Tusla and the HSE for the Age 14 Follow-up. We would like to thank all those who supported this research, especially the participating families and community organisations and the *PFL* intervention staff. Thanks to the Early Childhood Research Team at the UCD Geary Institute for Public Policy who significantly contributed to this study. The trial was registered with controlled-trials.com (ISRCTN04631728) and the AEA RCT Registry (AEARCTR-0000066). All study procedures were approved by the UCD Human Research Ethics Committee, the Rotunda Hospital's ethics committee, and the National Maternity Hospital's ethics committee. E-mail: orla.doyle@ucd.ie

1. Introduction

Much of the evidence base that is used to justify investment in the early years is based on a handful of experiments conducted in the U.S. in the 1960/70's and have tracked children throughout their lives. The results of these experiments, such as the Perry Preschool Program and the Abecedarian Program, have found that providing high quality preschool education from ages three to four leads to significantly better life outcomes in adulthood (e.g., Heckman *et al.*, 2010; Heckman, Pinto, and Savelyev, 2013). These findings have been used to justify the roll-out of preschool programmes globally. However, research from developmental neuroscience indicates that brain plasticity and neurogenesis are most pronounced even earlier in the lifecycle, particularly from pregnancy to age three (Thompson and Nelson, 2001; Knudsen *et al.*, 2006), and that parental stimulation and sensitive caregiving is associated with improved child development (Gertler *et al.*, 2014; Miller *et al.*, 2011). Thus, interventions that begin earlier and work with parents may be a more effective strategy. While a number of studies have examined the short-term effects of adopting such an approach, there are relatively few experimental studies that have tested the sustained impact of early life parenting programmes on children's later skills. This report addresses this gap by assessing the impact of *Preparing for Life (PFL)*, a prenatal to age five home visiting programme, ten years after the intervention ended.

Previous reports of the *PFL* programme identified some important effects at earlier ages.¹ Doyle (2020), based on data collected at the end of the trial, found that the programme had a large impact on children's cognitive, social, and behavioural development. The programme raised general conceptual ability, which is a proxy for IQ, by 0.77 of a standard deviation (SD), indicating the malleability of IQ in the early years. Gains were found across all dimensions of cognitive skill including spatial ability, pictorial reasoning, and language ability.

¹ Doyle (2013) describes the design of the *PFL* evaluation. Doyle *et al.* (2014) examine the impact of the *PFL* programme on birth outcomes utilising hospital data and identify a significant treatment effect regarding a reduction in caesarean sections, yet no impact on neonatal outcomes. Doyle *et al.* (2017a) examine the impact of the programme on parent reported cognitive and non-cognitive skills at 6, 12, and 18 months, and find no evidence of effects, yet there are significant improvements in the quality of the home environment. Doyle *et al.* (2015) examine the impact of the programme on child health measured at 6, 12, 18, 24, and 36 months and identify a number of significant treatment effects at 24 months in terms of reducing the incidence of asthma, chest infections, and health problems. O'Sullivan, Fitzpatrick and Doyle (2017) examine the impact of the programme on dietary intake at 12, 18, 24, and 36 months, and its mediating effect on cognitive development at 24 and 36 months; and find evidence of improved nutrition at 24 months in terms of increased protein intake. Doyle *et al.* (2017b) examine the impact of the programme on maternal wellbeing using daily data collected over a 24-hour period using the Day Reconstruction Method, finding little evidence of effects on maternal wellbeing. Cote *et al.* (2018) investigate whether the impact of the programme varied according to children's developmental trajectories and find a positive impact on trajectories of cognitive development and number of health clinic visits for all children, whereas positive impacts on externalizing behaviour problems are restricted to children with the most severe problems.

Although weaker, the programme also impacted several dimensions of children's non-cognitive skills including externalising problems such as aggressive behaviour, and prosocial behaviour such as helping other children. The programme also had an impact on child health by reducing the amount of hospital services the children used, as well as improving how families used these services (Coy and Doyle, 2024).

The first follow-up of the *PFL* cohort at age nine found that the programme continued to have an impact on cognitive scores, with effect sizes of 0.55 SD on general conceptual ability and 0.30 SD and 0.54 SD on achievement tests of reading and math respectively (Doyle, 2024). These results could not be attributed to differences in school quality across the high and low treatment groups as there was no evidence that the programme impacted school choice. The programme, however, had no impact on absenteeism rates or the use of school resources, and the significant treatment effects observed for children's socio-emotional skills at age four were no longer present at age nine. A mediation analysis found that between 35-38% of the treatment effects on age nine test scores was explained by improvements in early parental beliefs, stimulation, and health investments. The size of the treatment effects identified at the end of the trial and at age nine generally exceed findings from meta-analyses of other home visiting programmes (e.g., Sweet and Appelbaum, 2004; Gomby, 2005; Filene *et al.*, 2013).

The sustained impact of the *PFL* programme beyond the lifetime of the intervention may be attributed to both its duration and intensity. By providing a five-year intervention — covering the critical first 2,000 days of a child's life — the programme offered early and continual support to families during a pivotal period of development. This approach aligns with the "technology of skill formation" framework proposed by Cunha and Heckman (2007), which suggests that early skills serve as a foundation for the development of more advanced skills through a process of self-productivity. In turn, this enhances the effectiveness of later investments through dynamic complementarity (Cunha, Heckman, and Schennach, 2010; Heckman and Mosso, 2014). While genetic factors play a role in skill development (Nisbett *et al.*, 2012), research indicates that environmental conditions, such as the quality of the home environment and parenting, can shape and enhance these skills (Weaver *et al.*, 2004).

Empirical research highlights several key aspects of the home environment that predict children's skills, including the quality of the home setting (Todd and Wolpin, 2007), parenting skills (Dooley and Stewart, 2007; Fiorini and Keane, 2014), and parental stimulation (Miller *et al.*, 2014). Yet, socio-economic inequalities in the quality of the home environment exist, with disadvantaged families often facing financial constraints that limit their capacity to invest in their children. Evidence suggests that parents from low socioeconomic backgrounds are more

likely to adopt less effective parenting styles and behaviours (Cunha, Elo, and Culhane, 2013), such as permissive or harsh parenting (Bradley and Corwyn, 2002), and to provide fewer stimulating materials and experiences for their children (Bradley *et al.*, 1989). This may, in part, stem from a knowledge gap regarding optimal parenting practices. Cunha *et al.* (2013) point to differing beliefs about the importance of parenting among parents from low socioeconomic backgrounds. Disadvantaged homes are also less likely to offer pre-academic stimulation, such as reading to children or helping them recognise letters (Miller *et al.*, 2014). To address these issues, the *PFL* programme sought to enhance parenting knowledge and promote developmentally appropriate activities, thereby mitigating the negative effects of socioeconomic disadvantage on children's development. Thus, the programme targeted both behavioural frictions and informational frictions.

This report examines the long-term impact of the *PFL* programme now that the cohort have reached adolescence. It focuses on outcomes measured using directly assessed tests and a self-completion questionnaire. In addition, using pre-existing data from Ireland (Growing Up in Ireland; Health Behaviour of School Age Children) and the UK (Millennium Cohort Study), also collected during adolescence, it situates the *PFL* cohort alongside these representative samples.

The remainder of the report is structured as follows. Section 2 reviews the literature assessing the impact of home visiting programmes into adolescence. Section 3 describes the design of the original study and the Age 14 Follow-up. Section 4 outlines the statistical methods that are used to estimate the results. Section 5 presents the main results. Finally, Section 6 concludes.

2. Literature

The effectiveness of home visiting programmes in the short term is well-documented. A meta-analysis of 60 home visiting programmes by Sweet and Appelbaum (2004) reports an average effect size (ES)² of 0.18 for cognitive skill improvements and 0.10 for non-cognitive skills. Subsequent reviews found similar results, with Miller, Maguire, and Macdonald (2011) reporting an average ES of 0.30 for cognitive skills across seven studies, and Filene *et al.* (2013) reporting an ES of 0.25 based on 51 studies. Collectively, these findings suggest that

² The effect size (ES) represents the magnitude or the size of the difference between the treatment and controls group. While the p-value allows the reader to determine whether or not there is a statistically significant difference between the groups, it does not indicate the strength of the difference. Effect sizes are usually expressed in terms of standard deviation of the outcome variable. Effect sizes are calculated using Cohen's d, where effect sizes of 0.0 to 0.2 are considered small, 0.2 to 0.8 medium, and greater than 0.8 large.

home visiting programmes generally produce small to modest improvements in both cognitive and non-cognitive skills (Gomby, 2005; Peacock *et al.*, 2013; Avellar *et al.*, 2016).

Research regarding their effectiveness beyond the lifetime of the intervention has found mixed results. A study by Bailey *et al.* (2017) examined 67 high-quality early intervention programmes in the U.S., including some home visiting programmes, and identified a general pattern of declining effect sizes. While the average effect size at the end of the intervention was 0.23, this dropped to 0.10 by the end of the first year post-intervention and to 0.05 within one to two years after the programme concluded. Focusing specifically on home visiting programmes, findings on medium- and long-term effects are varied. Bierman *et al.* (2017) found that children who participated in Early Head Start (EHS) demonstrated improved cognitive ability as well as reading and language skills at ages seven to nine. Studies on the Healthy Families America (HFA) programme also revealed significant impacts, with children more likely to be enrolled in gifted programmes, less likely to require special education, and more likely to excel academically at ages six to seven (DuMont *et al.*, 2010; Kirkland and Mitchell-Herzfield, 2012). In terms of socio-emotional outcomes, two EHS studies found evidence of positive effects on children's behaviour, perceived competence, and learning approaches at ages five, seven, and nine (Bierman *et al.*, 2017; Chazan-Cohen, Raikes, and Vogel, 2013). Additionally, the Nurse-Family Partnership (NFP) programme reported a reduction in internalising disorders at age 12 (Kitzman *et al.*, 2010).

Research on the long-term impact of home visiting programmes provides more compelling evidence. The Jamaica Study, which involved weekly home visits for children aged nine to 24 months, found an initial IQ effect size of 0.88 at the end of the intervention. Although this effect diminished by age 7, it re-emerged at ages 11, 17, and 22, with effect sizes ranging from 0.40 to 0.60 (Grantham-McGregor and Smith, 2016). The Abecedarian Programme, which provided centre-based care and home visits from infancy to age five, recorded an initial IQ effect size of 0.74. While this declined to 0.37 on average at ages eight, 12, 15, and 21, the effects were still sustained (Campbell *et al.*, 2001). The NFP programme has also demonstrated significant long-term impacts. Cognitive effect sizes of 0.22 to 0.27 were observed at age six (for both boys and girls) and at age 12 (for boys only) (Heckman *et al.*, 2017). At age 18, children of mothers with low psychological resources showed higher receptive language (ES = 0.24) and math achievement (ES = 0.38) (Kitzman *et al.*, 2019). By age 19, girls in the treatment group had fewer children, were less likely to receive Medicaid, and were less involved in crime (Eckenrode *et al.*, 2010). Even in cases where cognitive effects faded over

time, other long-term benefits, such as reductions in criminal behaviour and receipt of social welfare, were observed (Campbell *et al.*, 2014; Heckman *et al.*, 2017).

2.1 Studies with follow-ups during adolescence

Table 1 summarizes the literature on home visiting programmes with follow-ups conducted during the same period of adolescence as in the current report (between 12-16 years). Overall, there are very few studies that test for the sustained effect of home visiting programmes during this period. Most studies either stop collecting data directly after the intervention ends; or only revisit families during adulthood (often using administrative data); or the interventions are still in infancy thus long-term follow up is not yet possible. The only programme that has conducted assessments during adolescence specifically is the NFP programme. Follow-ups were conducted at ages 12 -16 for the Memphis trial and age 15 for the Elmira trial in 11 separate papers. With few exceptions, there were no impacts on child outcomes.

Of the five papers included in Table 1, three found no effects on child outcomes, and only two found effects on a small number of the outcomes. Kitzman *et al.* (2010) found effects in the Memphis trial at age 12 on the substance use and internalising disorders, but no effects on cognitive scores, achievement tests, other behavioural problems, or educational outcomes. Olds *et al.* (1998) found effects in the Elmira trial at age 15 on convictions and probation violations, but no effects on substance use, risky behaviours, behavioural problems, anti-social behaviour, or school behaviour outcomes. Note that only two of the studies (Kitman *et al.*, 2010; Sidora-Arcoleo *et al.*, 2010) assessed cognitive skills, and no effects were found. The majority of the studies assessed socio-emotional outcomes (focusing on behavioural problems), and significant effects were only identified in one. Thus, the main takeaway from the sparse home visiting literature with assessments during adolescence is that effects on child outcomes are minimal.

Table 1 *Impact of Home Visiting Programmes on Child Outcomes from Ages 12-16*

Author	Sample Size	Programme	Measures	Significant Finding	Effect	Age (years)
Kitzman <i>et al.</i> (2010)	635	Nurse Family Partnership (Memphis)	GPA, Peabody Individual Achievement Tests, Leiter-R sustained attention test, Group achievement test scores, Placement in special education, ever retained in a grade, conduct grades, externalising disorders, internalising disorders, total problems, Days of substance use in the last 30 days	Incidence of substance use, used cigarettes, alcohol or marijuana in the last 30 days, internalising disorders	Favourable	12
Sidora-Arcoleo <i>et al.</i> (2010)	721	Nurse Family Partnership (Memphis)	Peabody Picture Vocabulary Test-Revised, physical aggression (CBCL)	None	None	6-12 years
Enoch <i>et al.</i> (2016)	559	Nurse Family Partnership (Memphis)	Composite externalising disorders continuous total scores (CBCL)	None	None	12
Eckenrode <i>et al.</i> (2001)	228	Nurse Family Partnership (Elmira)	Number of early onset of problem behaviors, Percentage abused or neglected.	None	None	15
Olds <i>et al.</i> (1998)	245	Nurse Family Partnership (Elmira)	Alcohol and drug impairment, Ever pregnant or made someone pregnant, Incidence of sex partners, cigarettes smoked per day, days drank alcohol, days used drugs, times ran away, Number of acting out problems, Number of externalising problems , Number of internalising problems , Number of minor antisocial acts, Ever was person in-need of supervision, Incidence of arrests, Incidence–convictions and probation violations, Incidence–long-term school suspensions, Incidence–sent to youth corrections , Incidence–short-term school suspensions, Number of major delinquent acts	Incidence–convictions and probation violations	Favourable	15

3. *PFL* Study Description

In an effort to break the intergenerational cycle of disadvantage, *PFL* was developed as part of the Irish Government's and The Atlantic Philanthropies' Prevention and Early Intervention Programme. *PFL* was developed by 28 local agencies and community groups who collaborated to design an evidence-based intervention tailored to meet the needs of the local community. The study took place between 2008 and 2015 in a highly disadvantaged Dublin community with the aim of reducing socioeconomic inequalities in children's skills by working directly with parents.

3.1 Initial recruitment and randomisation

Recruitment into the *PFL* programme took place between the 29th of January 2008 and the 4th of August 2010 through two maternity hospitals and/or self-referral using a community-based marketing campaign. The inclusion criteria included all pregnant women residing in the designated *PFL* catchment area during this period, regardless of their social or family circumstances. Based on estimates of a two to five point difference on cognitive development scores (i.e., average standardised effect size of 0.18) from a meta-analysis of home visiting programmes (Sweet and Appelbaum, 2004), a sample size of approximately 117 in each group was required to power the study.

In total, 233 participants were recruited by the *PFL* recruitment officers. This represents a recruitment rate of 52% based on the number of live births during the recruitment period. Of those who joined the programme, an unconditional probability randomisation procedure, with no stratification, assigned 115 to a high treatment group and 118 to a low treatment group. Baseline data from 205 participants (representing 90% of the high treatment group and 86% of the low treatment group) was collected after randomisation yet prior to treatment delivery.³ The baseline variables include 117 measures of socio-demographics, physical and mental health, IQ, parenting attitudes, and self-control, among others. To assess the effectiveness of the randomisation procedure, the baseline characteristics of the high and low treatment groups were compared using permutation tests. At the 10% significance level, the two groups differed on 7.7% (9/117) of measures, which is consistent with pure chance and indicated the success of the randomisation process (see Doyle and *PFL* Evaluation Team, 2010).

³ Of the 233 randomly assigned participants, two (high=1; low=1) miscarried, 19 (high=6; low=13) withdrew from the programme before the baseline assessment, and seven (high=4; low=3) did not participate in the baseline but participated in subsequent waves. An analysis of a subset (n = 12) of this group on whom recruitment data but no baseline data are available, implies they do not differ on age, education, employment, and financial status from those who did complete a baseline assessment, however the limited sample size should be noted.

3.2 Treatment

Figure 1 describes the supports provided to the high and low treatment groups. The high treatment consisted of three primary components - a five year home visiting programme, a baby massage course, and the Triple P Positive Parenting Programme. The treatments were built upon the theories of human attachment (Bowlby, 1969), socio-ecological development (Bronfenbrenner, 1979), and social-learning (Bandura, 1977). The home visiting programme aimed to promote children's health and development by building a strong home visitor-parent relationship and focusing on the identification of developmental milestones, appropriate parenting practices, and encouraging enhanced stimulation. The visits started in the prenatal period and continued until school entry. Twice monthly home visits of approximately one hour were prescribed with home visitors from different professional backgrounds including education, social care, and youth studies. The visitors were hired to deliver the programme on a full-time basis and they received extensive training prior to treatment delivery. Supervision took place on a monthly basis to ensure fidelity to the programme model, and families were allocated the same home visitor over the course of the intervention where possible.

Each home visit was structured around *PFL*-developed 'Tip Sheets' which included information on pregnancy, parenting, health, and development. The 210 Tip Sheets were developed by the *PFL* implementation team based on pre-existing and publicly available information. The home visitors could choose when to deliver the Tip Sheets based on the age of the child and the needs of the family, yet the full set of Tip Sheets must have been delivered by the end of the programme. The intervention was delivered using techniques such as role modelling, coaching, discussion, encouragement, and feedback, as well as directly interacting with the *PFL* child. Each home visit began with an update on the family's situation and a discussion of whether the goals agreed at the previous visit were achieved. The home visitor would then guide the parent through the Tip Sheet(s) selected for that visit and following this, new goals would be agreed. The Tip Sheets typically targeted multiple aspects of development. The programme adopted a strengths-based approach to empower parenting in their parenting role and acknowledge that they are experts in their children's lives.

Participants in the high treatment group were also encouraged to take part in a baby massage course in the first year, which consisted of five two-hour individual or group sessions delivered by the home visitors. The purpose of these classes was to equip parents with baby massage skills and to emphasise the importance of early reciprocal interactions and communication between parents and infants.

When the *PFL* children were between two and three years old, the high treatment group was invited to participate in the *Triple P Positive Parenting Programme* (Sanders, Markie-Dadds, and Turner, 2003) which was delivered by the home visitors. The goal of *Triple P* is to encourage positive, effective parenting practices in order to prevent problems in children's development. The programme is based on five principles including providing a safe, engaging environment, the home as a positive place to learn, setting of rules and boundaries, realistic expectations of children, and parental self-care (Sanders, 2012). Meta-analysis of the impact of *Triple P* has identified improved parenting practices and child social, emotional, and behavioural outcomes (Sanders *et al.* 2014). The high treatment participants were encouraged to take part in five two-hour group discussion sessions and three phone calls. The home visitors also used the *Triple P* principles and techniques when delivering the home visits to ensure consistent messaging across the programme components.

In addition to the standard services available to pregnant women and young children, both the high and low treatment groups received a supply of developmental toys annually (to the value of ~€100 per year) including a baby gym, safety items, and developmental toys. They also received four book packs containing six to eight developmentally appropriate books. The groups were also encouraged to attend community-based public health workshops on stress management and healthy eating, as well as social events such as coffee mornings and Christmas parties organized by the *PFL* staff. Programme newsletters and birthday cards were sent annually to each family, in addition to two framed professional photographs of the child. The low treatment group also had access to a *PFL* support worker who could help them avail of community services if needed, and this function was provided by the home visitor for the high treatment group. Note that the low treatment group did not receive the home visiting programme, Tip Sheets, baby massage classes, or the *Triple P* programme.

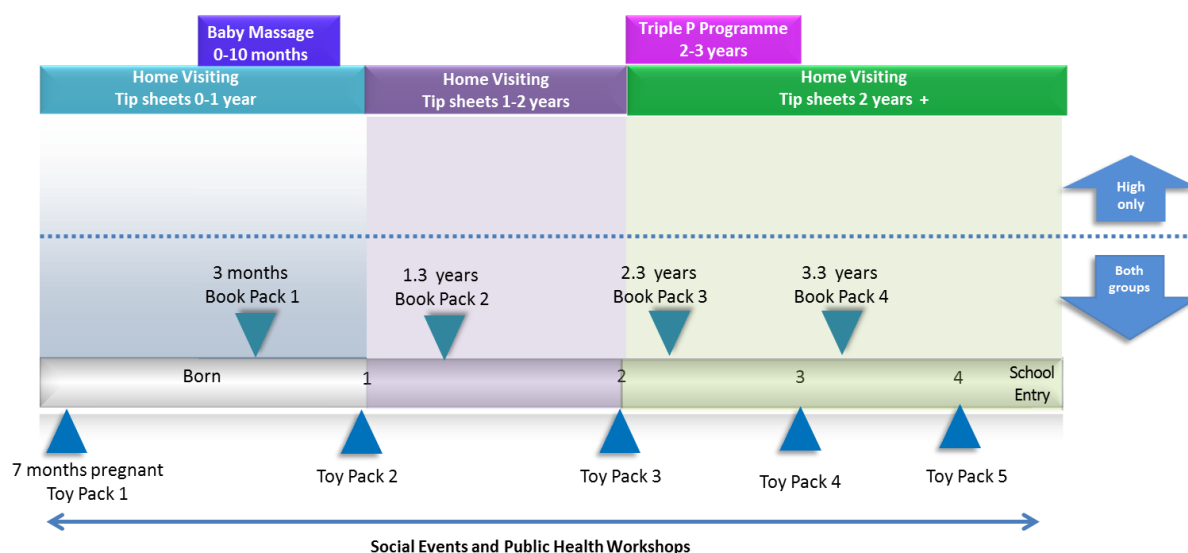


Figure 1 Timing of PFL treatments

3.3 Design of the Age 14 Follow-up⁴

In Summer 2023, the *PFL* evaluation team began designing the Age 14 Follow-up study in partnership with the *PFL* implementation team who were instrumental in helping to contact and re-consent participants into the study.

Recruitment

The Age 14 Follow-up study sought to include as many of the original *PFL* participants as possible. As such, all families that were recruited and randomised into the original *PFL* study between 2008 and 2010 were eligible to take part in the follow-up. However, participants who had officially dropped out or left the study due to death or miscarriage were not contacted. Multiple steps were taken to make contact and reconsent as many families as possible. First, six months prior to formally starting the recruitment process, the *PFL* implementation team attempted to make contact with all families who had not formally withdrawn from the study during previous waves. This informal contact was to inform the families of the upcoming Age 14 Follow-up study and to update their contact details. Second, all families were invited to attend a *PFL* ‘Age 14 Birthday Celebration’ on the 8th October 2023. At this event, which involved talks, games, activities, and food, families were asked if they were interested in hearing about the Age 14 Follow-up study and if they were happy to speak to a researcher about

⁴ Ethical approval to conduct the Age 14 Follow-up study was provided by the UCD Human Research Ethics Committee on 1st Sept 2023.

participating. If a family agreed to meet with the researcher, they were informed about the upcoming study and their consent to participate was sought. A total of 25 participants were recruited at this event. Third, families who did not attend the event, yet agreed to be contacted by the evaluation team at previous assessments, were attempted to be contacted using all available contact data available (e.g., phone, email, house address etc.). Fourth, the *PFL* implementation team also followed up with all other families to inform them of the study and put them in touch with the evaluation team. Families who left the original study area were still invited to participate. 74 participants were recruited using a combination of these methods. In total, 99 participants were recruited between October 2023 and September 2024.

During the recruitment process, eligible participants were provided with information about the follow-up study. Researchers explained the procedures, articulating all relevant study information, consent details, and participant rights. At this time, we sought consent from parents for their child to take part in the study, as well as assent from the child to participate. If a parent provided consent, but the child did not, the assessment did not take place. If a child provided consent, but the parent did not, the assessment did not take place. Parents were provided with an information sheet detailing the study, and children were provided with an age-appropriate information booklet.

Age 14 Sample: Age and Gender

The average age of the *PFL* cohort (high and low treatment groups) at the time of data collection was 14.3 years, with the youngest participant being 12.11 and the oldest 16.2. Importantly, there were no statistically significant differences in the age of the participants in the high and low treatment groups at the time of the interview (high treatment=14.4 years, low treatment=14.3 years; p -value=0.620). There were also no differences regarding participant gender – within the high treatment group, 56% of the sample were girls and 44% were boys, and within the low treatment group, 64% of the sample were girls and 36% were boys; p -value=0.666). Thus, the sample was balanced in terms of age and gender.

Data Collection Procedure

Data for the Age 14 Follow-up were collected between January and September 2024. The majority of the assessments took place in the participant's secondary school (75%), and the remainder took place either in the participant's home or the village centre. To minimize detection bias, all assessments were conducted by trained researchers who were blind to the

treatment condition. There were four assessment components – direct assessment of cognitive skills and executive functioning, a self-completion questionnaire, height and waist measurements, and a saliva sample. All these, bar the saliva sample, are described in more detail below.⁵ When the assessments took place in schools, they were divided into two sessions to reduce fatigue, one before lunch and the other after lunch. Each session lasted about one hour. The participants received a €20 One4All voucher for each session to compensate them for their time and effort.

Session 1

Participants' cognitive skills were assessed using the British Ability Scales III: School Age Battery (BAS III; Elliott, Smith and McCulloch, 2011) (the same assessment used at the Age 9 Follow-up). The BAS III yields an overall score reflecting general cognitive ability (General Conceptual Ability, GCA), as well as three standardised scores for Verbal Ability, Pictorial Reasoning Ability, and Spatial Ability. The participants also conducted several tasks assessing self-regulation/executive functions using the National Institutes of Health Toolbox for Assessment of Neurological and Behavioral Function Cognition Battery (NIH Toolbox; Zelazo and Bauer, 2013) (e.g., the Flanker task to assess inhibitory control, the Dimensional Change Card Sort task to assess attention flexibility, and the List Sorting task to assess working memory).

Session 2

Participants were invited to complete a self-completion questionnaire on an iPad. The questionnaire was programmed in Qualtrics. The participant completed the survey on their own, however a researcher was present in the room, and the participant was informed that they could ask the researcher any questions they had about the survey or if they did not understand a question. To ensure that participants did not lose interest, at two points during the questionnaire, they were prompted on the screen to hand the iPad to the researcher. At the first 'break', the researcher measured the participants' height and waist circumference. At the second 'break', a saliva sample was taken. The questionnaire included a set of standardised instruments and single-item questions capturing measures of socio-emotional development, mental and physical health, health behaviours (diet, substance use), puberty development, self-esteem, attitudes towards school, school absences, educational expectations, life satisfaction,

⁵ A separate report on the results of the biological aging study using the saliva sample will be produced.

relationship with parents, attitudes towards antisocial behaviour, daily activities, and risk and time preferences.

3.4 Age 14 Follow-up Sample and Attrition

Figure 2 depicts the families' participation in the study between programme entry and the Age 14 Follow-up. At the follow-up, data were collected for 99 of the original 233 randomly assigned participants, representing an overall retention rate of 43%.⁶ Despite the substantial time and effort invested in recruitment efforts by both the evaluation and implementation teams, the participation rate was lower than that achieved at Age 9 (50%). At Age 9, a significantly higher proportion of the high treatment group participated compared to the low treatment group (59% vs 42%). However, at Age 14, the level of attrition was relatively similar in both groups (45% vs 40%) and importantly there were no statistically significant differences in the likelihood of participating in the follow-up based on treatment status ($p=0.406$). An analysis of attrition between Age 9 and 14 showed that a higher proportion of the high treatment group dropped out of the study than the low treatment group at this time point.

It is important to note that the composition of the samples was not the same at both waves e.g., some participants at Age 14 did not participate at Age 9, and vice-versa. For example, 117 participants took part in the Age 9 assessment and 99 took part in the Age 14 assessment. While the majority of participants took part in both assessments ($n=82$), 34 participants who took part at Age 9 did not participate at Age 14, and 17 participants who took part at Age 14 did not participate at Age 9. In order to test whether participants with certain characteristics were more likely to drop out, we compared the Age 9 cognitive scores of those who participated at Age 9 and at Age 14 (GCA mean=84.8) to those who did not participate at Age 14 (GCA mean=83.9), and we found no significant difference ($p=0.747$) across the groups. In addition, we compared the Age 14 cognitive scores of those who participated at Age 9 and at Age 14 (GCA mean=83.1) to those who did not participate at Age 9 (GCA mean=78.5) and also found no significant difference ($p=0.190$), however, those who re-joined the study at Age 14 had somewhat lower scores. Overall, this implies that the type of participants who took part in the Age 14 assessment were similar to those who took part at Age 9. Therefore, any

⁶ Note that 99 participants completed the direct assessments and 99 completed the self-completion questionnaire, however the samples are not identical. There was one participant who completed the direct assessment, but did not complete the questionnaire, and another participant who completed the questionnaire but not the direct assessment. Therefore, separate weights are used for both samples in the Inverse Probability Weighting procedure.

differences in the results across waves (if found) may not be attributed to differences in the types of participants who took part.

Figure 2 Participant flow



A re-examination of the comparability of the high and low treatment groups using the Age 14 estimation sample showed that the two groups differed on 11% (13/117) of baseline measures. Using the 10% cut-off level, we would expect 10% of the measures to be statistically significant at random, thus these results are largely consistent with pure chance and indicate that the groups remained largely balanced at the Age14 Follow-up. Table 2 compares the Age 14 participants in the high and low treatment group for a selection of baseline measures. It shows that there are no statistically significant differences across the two groups on all but one of the key socio-demographic and health factors assessed. Specifically, high treatment children who participated in the Age 14 Follow-up are more likely to be part of families where the mother was married during pregnancy.

Table 2 Baseline comparison of the Age 14 high and low treatment groups

	M_{HIGH} (SD)	M_{LOW} (SD)	p^1
Age	26.69 (5.78)	25.81 (6.03)	0.462
Married	0.23 (0.43)	0.10 (0.31)	0.090
No. of children	1.94 (1.19)	1.98 (1.17)	0.891
First time mother	0.48 (0.50)	0.47 (0.50)	0.921
Low education (left ≤ age 16)	0.27 (0.45)	0.23 (0.43)	0.779
Weschler Abbreviated Scale of Intelligence (WASI)	85.25 (11.75)	81.61 (14.49)	0.178
Employed	0.48 (0.50)	0.45 (0.50)	0.727
Resides in social housing	0.54 (0.50)	0.55 (0.50)	0.952
Medical card	0.58 (0.50)	0.64 (0.49)	0.462
Prior physical health condition	0.77 (0.43)	0.70 (0.46)	0.453
Prior mental health condition	0.31 (0.47)	0.28 (0.45)	0.767
Smoking during pregnancy	0.40 (0.50)	0.47 (0.50)	0.484
Drinking alcohol during pregnancy	0.33 (0.47)	0.32 (0.47)	0.922
<i>N</i>	52	47	

Note: All baseline measures were assessed during pregnancy prior to treatment delivery except for WASI which was assessed at 3 months postpartum. ¹ two-tailed p -values calculated from permutation tests with 100,000 replications.

Although the estimation samples were largely balanced in terms of baseline characteristics, it is important to test for differential attrition in the high and low treatment groups. To investigate this, the factors predicting participation in the Age 14 assessment were tested using bivariate tests with 49 baseline measures. Analyses were conducted separately for

the high and low treatment groups to allow for differential attrition processes. In general, there was some evidence of differential attrition, with 11 (22%) baseline measures predicting attrition from the high treatment group, and nine (18%) predicting attrition from the low treatment group (in two-tailed tests, using the 10% significance level).⁷

The factors predicting attrition from both groups differ somewhat, however in both cases we find that, consistent with the home visiting literature (see Roggman *et al.*, 2008), families who did not take part in the Age14 assessment had more risk factors at baseline. Figure 3 shows the baseline measures that were significant predictors of retention in the high and low treatment groups.

Figure 3 Baseline measures significantly predicting retention at Age 14

High Treatment Group (11/49)	Low Treatment Group (9/49)
<ul style="list-style-type: none"> • WASI IQ Score (+) • Married (+) • Use of health services (-) • Employed (+) • Smoke during pregnancy (-) • Age (+) • Vulnerable attachment style questionnaire score (-) • Drink during pregnancy (+) • Birth control (-) • Satisfied with neighbourhood (+) • WHO5 well-being score (+) 	<ul style="list-style-type: none"> • Low education (<16) (-) • TIPI Openness personality (+) • Regular exercise (+) • Non traveller (+) • Impaired by illness (+) • Married (-) • Birth control (+) • Domestic risk factors (+) • Consideration of future consequence scale score (-)

Note: (+) means that the measure is positively associated with participating, while (-) means that the measure is negatively associated with participating.

Table 3 compares a selection of baseline characteristics of those who participated in the Age 14 assessment - ‘stayers’ - to those who did not - ‘non-stayers’. It shows that high treatment participants who completed the Age 14 assessment had parents who were older, and more likely to be married, employed, and drink at baseline, however they were less likely to smoke. Low treatment participants who completed the Age 14 assessment had parents who were less likely to be married and to have low education at baseline. In order to account for differential attrition across the high and low treatment groups, treatment effects were estimated using the Inverse Probability Weighting procedure detailed in the Methods section.

⁷ This analysis is based on the direct assessment sample. For the self-completion questionnaire sample, 12 (24%) baseline measures predicted attrition from the high treatment group, and 9 (18%) predicted attrition from the low treatment group.

Table 3 Baseline characteristics predicting attrition from the Age 14 sample

	High Treatment Group			Low Treatment Group		
	M_{STAYER}	$M_{\text{NON-STAYER}}$	p^1	M_{STAYER}	$M_{\text{NON-STAYER}}$	p^1
	(SD)	(SD)		(SD)	(SD)	
Age	26.69 (5.78)	24.23 (5.71)	0.031	25.81 (6.03)	24.85 (5.97)	0.425
Married	0.23 (0.43)	0.06 (0.24)	0.011	0.11 (0.31)	0.24 (0.43)	0.067
First time mother	0.48 (0.50)	0.60 (0.50)	0.236	0.47 (0.50)	0.52 (0.50)	0.612
No. of children	1.94 (1.19)	1.94 (1.42)	1.000	1.98 (1.17)	1.85 (1.14)	0.593
Low education (left \leq age 16)	0.27 (0.45)	0.40 (0.50)	0.159	0.23 (0.43)	0.54 (0.50)	0.001
Weschler Abbreviated Scale of Intelligence (WASI)	80.90 (11.32)	77.38 (11.41)	0.116	79.04 (12.78)	77.41 (10.46)	0.476
Employed	0.48 (0.50)	0.25 (0.44)	0.012	0.45 (0.50)	0.35 (0.48)	0.327
Medical card	0.58 (0.50)	0.62 (0.49)	0.693	0.64 (0.49)	0.69 (0.47)	0.617
Prior physical health condition	0.77 (0.43)	0.73 (0.45)	0.656	0.70 (0.46)	0.56 (0.50)	0.137
Prior mental health condition	0.31 (0.47)	0.25 (0.44)	0.514	0.28 (0.45)	0.20 (0.41)	0.387
Smoking during pregnancy	0.40 (0.50)	0.62 (0.49)	0.029	0.47 (0.50)	0.48 (0.50)	0.919
Drinking alcohol during pregnancy	0.33 (0.47)	0.17 (0.38)	0.070	0.30 (0.46)	0.24 (0.43)	0.538
<i>N</i>	52			47		

Note: All baseline measures were assessed during pregnancy prior to treatment delivery except for WASI which was assessed at 3 months postpartum. Baseline data are missing for four participants in the age 9 assessment.

¹ two-tailed p -values calculated from permutation tests with 100,000 replications.

3.5 Power

Attrition has two consequences. First, as discussed above, it can reduce the comparability of the high and low treatment groups which means that the assumption of baseline equivalence no longer holds i.e., the results could be attributed to underlying differences between the two groups rather than differences being ‘caused’ by the programme. However, another consequence of attrition is reduced statistical power. Power is the ability to identify statistically significant differences between the high and low treatment groups if the programme really does have an impact. More formally, it is the probability of not rejecting the null hypothesis even though the null hypothesis is not true. A Type II Error is the probability of (falsely) concluding that there is no treatment effect, even when there is, and power is the probability of avoiding Type II errors. Traditionally, we aim for 80% power in experiments. This means that there is a 20% chance of making a Type II error i.e., 20% of the time we will not be able to reject the null hypothesis of a zero treatment effect despite there being a significant effect.

One of the key factors in determining the power of an experiment is sample size. In general, for a given expected effect size, power is higher when the sample size is larger. Thus, as the sample size falls over time with attrition, the ability to detect statistically significant effects falls. Thus, assuming a power of 80%, with a sample size of 99 (53 = high treatment, 47 = low treatment), a post-hoc power analysis was conducted to determine what size effects we are powered to detect at Age 14 e.g., how large does the difference in outcomes between the high and low treatment groups need to be in order to identify a statistically significant result. This analysis found that we have the power to detect effects of 0.50 standard deviations. For context, the effect size found at Age 9 on children's cognitive skills was 0.55 standard deviations. Thus, if we observe a similar effect size on cognitive skills at Age 14, we will be able to conclude that the programme still has a statistically significant effect on cognitive skills. However, for smaller effect sizes, which may be practically meaningful, we may not be able to detect statistically significant differences. In addition, for effects of 0.20 SD, power is only 26%, which means there is a high likelihood of making a Type II error. Thus, when interpreting the results, we will refer to both statistically significant results (concerning the p value) and clinically significant results (concerning the effect size).

3.6 Statistical methods used to estimate the results

This section describes the statistical methods that were used to estimate the results. Using an intention-to-treat approach, the standard treatment effect framework defines the observed outcome Y_i of participant $i \in I$ by:

$$(1) \quad Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0) \quad i \in I = \{1 \dots N\}$$

where $I = \{1 \dots N\}$ represents the sample space, D_i represents treatment assignment for participant i ($D_i = 1$ for the high treatment group, $D_i = 0$ for the low treatment group) and $(Y_i(0), Y_i(1))$ are the potential outcomes for participant i . The null hypothesis of no treatment effect on outcomes is tested via:

$$(2) \quad Y_i = \beta_0 + \beta_1 D_i + \epsilon_i$$

Given the relatively small sample size, traditional hypothesis testing techniques which are based on large sample assumptions are not appropriate, thus the treatment effects were

estimated using exact permutation-based hypothesis testing (see Good, 2005).⁸ The permutation tests were estimated by calculating the observed t -statistic. The data were then repeatedly shuffled so that the treatment assignment of some participants was switched (100,000 replications were used). The observed t -statistic was then compared to the distribution of t -statistics that resulted from the permutations. The mid- p value was reported and was calculated as follows:

$$(3) \quad MP(t) = P(t^* > t) + 0.5P(t^* = t)$$

where $P(.)$ is the probability distribution, t^* is the randomly permuted t -statistic, and t is the observed t -statistic. Similar to other early childhood intervention studies (e.g., Heckman *et al.* 2010; Campbell *et al.*, 2014; Gertler *et al.*, 2014), one-sided tests with the accepted Type I error rate set at 10% were used given the hypothesis that the high treatment will have a positive effect on children's outcomes.

As there was an imbalance in the proportion of girls and boys in the treatment groups at baseline, and given differential developmental trajectories by gender, all analyses control for gender. As the assumption of exchangeability under the null hypothesis may be violated when controls are included, conditional permutation testing was applied. Using this method, the sample was proportioned into subsets, called orbits, each including participants with common background characteristics, in this case, there was one orbit for boys and one for girls. Under the null of no effect, the outcomes of the high and low treatment groups have the same distributions within an orbit. The exchangeability assumption is thus limited to strata defined by gender.

In order to account for any potential bias due to differential attrition, an Inverse Probability Weighting (IPW) technique (Robins, Rotnitzky, and Zhao, 1994) was applied. First, logistic models were estimated to generate the predicted probability of participation in the Age 14 assessment. The measures which were the significant predictors of attrition in Figure 3 were included in the logistic models modelling the likelihood of participating in the Age 14 assessment. Separate models were estimated for the high and low treatment groups.

⁸ As permutation testing does not depend on the asymptotic behaviour of the test statistic, it is a more appropriate method to use when dealing with non-normal data (Ludbrook and Dudley, 1998). A permutation test is based on the assumption of exchangeability under the null hypothesis, therefore if the null hypothesis is true, taking random permutations of the treatment variable does not change the underlying distribution of outcomes for the high or low treatment groups. Permutation testing has been shown to exhibit power advantages over parametric t tests in simulation studies, particularly when the degree of skewness in the outcome data is correlated with the size of the treatment effect (e.g. Mewhort 2005). Although this method is useful for dealing with non-normal data, it cannot be used to compensate for an under-powered study.

The predicted probabilities from these logistic models were then used as weights in the permutation tests so that a larger weight was given to participants that were underrepresented in the sample due to attrition.

The issue of testing multiple outcomes at multiple time points and thus increasing the likelihood of a Type-I error, was mitigated using the stepdown procedure which controls the Family-Wise Error Rate (Romano and Wolf, 2005). Using this method all outcome measures were placed into a series of stepdown families each representing an underlying construct. The stepdown procedure was conducted by calculating a t -statistic for each null hypothesis in the stepdown family using permutation testing. The results were placed in descending order. The largest t -statistic was then compared with the distribution of maxima permuted t -statistics. If the probability of observing this statistic was $p \geq 0.1$ we failed to reject the joint null hypothesis. If the probability of observing this t -statistic was $p < 0.1$ the joint null hypothesis was rejected, and the most significant outcome was excluded, and the remaining subset of outcomes were tested. This process continued until the resulting subset of hypotheses failed to be rejected or only one outcome remained. By stepping down through the outcomes, the hypothesis that leads to the rejection of the null was found.

The results are discussed using p -values to indicate statistically significant effects, where $p < 0.1$ is considered statistically significant, and Cohen's d effect sizes, where a small effect is 0.2, a medium effect is 0.5, and a large effect is 0.8.

4 Results

4.1 Cognitive outcomes

Cognitive skills were measured using the *British Ability Scales III* (Elliot and Smith, 2011) which consists of six subscales: word definitions, verbal similarities, matrices, quantitative reasoning, recognition of designs, and pattern construction. These sub-scales yield an overall score reflecting general cognitive ability (General Conceptual Ability, GCA), as well as three cluster scores for Verbal Ability, Non-Verbal Ability, and Spatial Ability. The *GCA score* assesses overall cognitive ability such as thinking logically, making decisions, and learning. The *Spatial Ability* score assesses problem solving, spatial visualisation, and short-term visual memory. The *Nonverbal Reasoning* score assesses inductive reasoning. The *Verbal Ability* score assesses verbal reasoning, verbal knowledge, and expressive language. Age-based T scores were calculated for each domain that were standardised to have a mean of 100 and a

standard deviation of 15, as well as cutoff scores indicating whether the participant scores above or below average for the GCA and cluster scores.

Executive function was measured using the *National Institutes of Health Toolbox for Assessment of Neurological and Behavioral Function Cognition Battery* (NIH Toolbox; Zelazo and Bauer, 2013). Executive functions are higher order meta-cognitive processes involved in concentration, reasoning, problem solving, and planning.⁹ The *Flanker Task* was used to assess inhibitory control. Participants were asked to indicate the left-right orientation of a centrally presented stimulus arrow surrounded by congruent or incongruent stimuli arrows. The *Dimensional Change Card Sort* task was used to assess attention flexibility. Participants were asked to match test pictures to a target picture that varied along two dimensions, colour and shape. Finally, the *List Sorting* task was used to assess working memory. Participants were presented with a series of stimuli (either food or animals) on the screen and orally and asked to order the list of items from smallest to largest; and then presented with a series of stimuli and asked to recall the food items in size order followed by the animals in size order from smallest to largest. Age-corrected scores for each of the three NIH toolbox measures was then standardised and summed to create a composite indicator of executive functions.

Table 4 reports the Inverse Probability Weighted (IPW) adjusted means, standard deviations, and *p*-values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size (as measured by the ratio of the treatment effect and the pooled standard deviation), for children's cognitive outcomes.¹⁰ The results indicate that the *PFL* programme had a significant impact on children's skills at Age 14. The treatment increased overall GCA scores by 0.70 SD. The effect size is equivalent to those observed at both the Age 4 and Age 9 assessments. The results also demonstrate that *PFL* had a significant impact on each dimension of cognitive ability including spatial ability (0.51 SD), non-verbal

⁹ Executive functions are comprised of three core abilities. 1. Inhibitory control which involves the ability to override impulse responses and ideally replace them with a more adaptive behaviour. For example, delaying eating a treat to receive a reward. 2. Attention flexibility which involves being able to deliberately focus and maintain attention or to divert attention to a new task if required. For example, blocking out distraction to complete a task or shifting attention from one task to another. And 3, working memory which involves the ability to retain and manipulate information over brief periods of time. Working memory is central to remembering instructions or rules or pieces of information that are necessary to solve a problem.

¹⁰ The non-IPW adjusted results for the BAS outcomes were similar to the results reported here. The number of significant outcomes were the same. In terms of effect sizes, some were slightly larger in the non-IPW results and some were slightly more conservative with somewhat lower effect sizes. For the executive functioning outcomes, the number of significant results were the same in the IPW and non-IPW results, however the effect sizes were smaller and the level of significance higher in the non-IPW results.

reasoning ability (0.53 SD), and verbal ability (0.61 SD). In addition, all four composite scores survive adjustment for multiple hypotheses testing.

Children were classified as scoring above the norm if their score was above 110 points and below the norm if their scores were less than 90 points. Table 4 shows that high treatment children were more likely to score *above* the norm in terms of their overall cognitive ability and their spatial ability; with effect sizes of 0.55 to 0.46 SDs respectively. For example, 17% of children in the high treatment group scored above the norm on their spatial skills compared to 4% in the low treatment group. Non-verbal and verbal ability were not statistically significant, although the effect sizes were of moderate size (~0.40 SD). In addition, the high treatment children were less likely to score *below* the norm across all three cognitive domains, as well as overall ability, results which were robust to multiple hypothesis adjustment. The effect sizes ranged from 0.35 to 0.86 of a standard deviation.

It is important to note that relatively few children, in either the high or low treatment groups, score above the norm, while large proportions of children score below the norm. For example, only 7% of the high treatment group has above the norm GCA scores, while 53% have below the norm scores. The BAS III norms are based on a representative UK sample which includes children across all social groups. The scores identified here thus reflect the disadvantaged nature of the *PFL* cohort, where lower levels of cognitive ability were expected to be observed. Yet the counterfactual (low treatment group) reveals that in the absence of the *PFL* intervention, a significantly greater proportion of the high treatment children would have scored below the norm, thus demonstrating the effectiveness of the programme.

Table 4 Comparison of high and low treatment groups: Cognitive outcomes

	N (HIGH/LOW)	M _{HIGH} (SD)	M _{LOW} (SD)	p ¹	p ²	ES
<i>BAS Composite Scores</i>						
General Conceptual Ability	99 (52/47)	85.41 (14.04)	76.76 (10.74)	0.003	0.010	0.70
Spatial Ability	99 (52/47)	94.09 (17.77)	86.28 (13.02)	0.012	0.012	0.51
Non-Verbal Ability	99 (52/47)	82.94 (12.36)	76.57 (11.62)	0.015	0.023	0.53
Verbal Ability	99 (52/47)	87.39 (12.82)	80.42 (9.89)	0.016	0.023	0.61
<i>BAS Above the Norm %</i>						
General Conceptual Ability	99 (52/47)	0.07 (0.26)	0.00 (0.00)	0.090	0.091	0.55
Spatial Ability	99 (52/47)	0.17 (0.38)	0.04 (0.19)	0.042	0.044	0.46
Non-Verbal Ability	99 (52/47)	0.06 (0.23)	0.00 (0.00)	0.153	0.153	0.48
Verbal Ability	99 (52/47)	0.04 (0.21)	0.00 (0.00)	0.158	0.158	0.42
<i>BAS Below the Norm %</i>						
General Conceptual Ability	99 (52/47)	0.53 (0.50)	0.88 (0.33)	0.000	0.002	0.86
Spatial Ability	99 (52/47)	0.38 (0.49)	0.55 (0.50)	0.063	0.063	0.35
Non-Verbal Ability	99 (52/47)	0.74 (0.44)	0.92 (0.28)	0.045	0.062	0.48
Verbal Ability	99 (52/47)	0.60 (0.49)	0.85 (0.37)	0.020	0.039	0.57
<i>NIH Toolbox Executive Functioning</i>						
Flanker Task - Inhibitory Control	98 (51/47)	95.10 (17.78)	98.27 (17.57)	0.790	0.790	-0.18
Dimensional Change Card Sort Task - Attention Flexibility	99 (52/47)	111.49 (20.52)	113.94 (19.84)	0.624	0.794	-0.12
List Sorting Task - Working Memory	99 (52/47)	102.44 (16.20)	94.68 (13.89)	0.023	0.052	0.52
<i>Other</i>						
Composite Executive Function Score	98 (51/47)	0.07 (0.75)	0.00 (0.78)	0.342	~	0.09

Note: N' indicates the sample size. 'M' indicates the IPW-adjusted mean. 'SD' indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional *p*-value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional *p*-value from IPW-adjusted stepdown permutation test with 100,000 replications. 'Effect size' is the ratio of the treatment effect to the pooled standard deviation.

Figure 3 shows that the distribution of GCA scores for the high treatment group was shifted to the right of the low treatment groups, indicating that the programme impacted children of all ability types – the programme impacted both the average score, and the tails of the distribution. The effects on the BAS scores at Age 14 were similar in magnitude to the BAS results measured at the end of the programme (at approx. 51 months) and at Age 9, demonstrating the sustained impact of the programme almost ten years after the treatment ended.

Figure 3 *Distribution of BAS GCA cognitive scores at age 14*

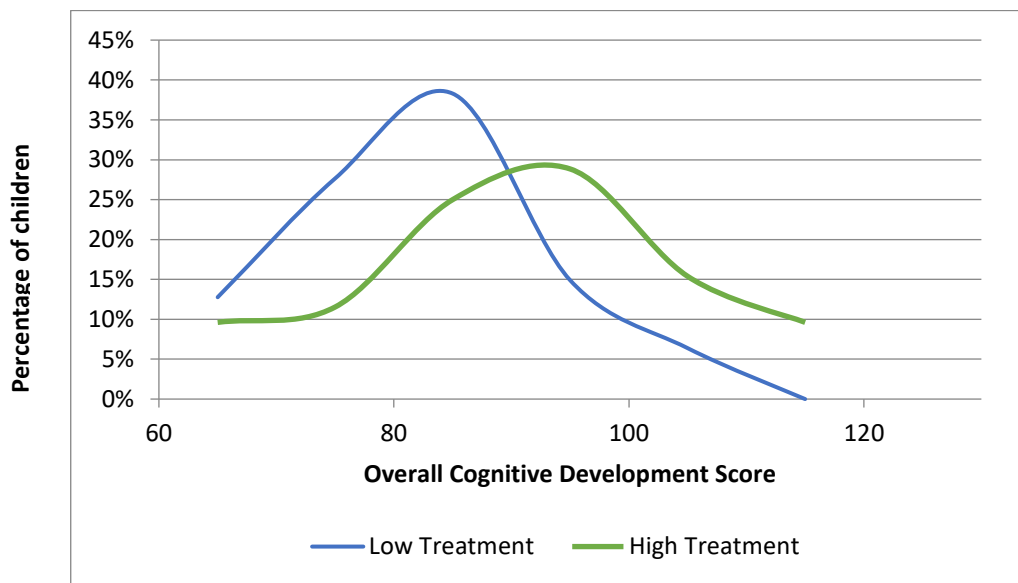


Table 4 also reports the results for executive functioning. At Age 9, the programme impacted all three forms of executive functioning, however at Age 14, only one significant effect on working memory (the ability to retain and manipulate or use information over brief periods of time), was found. Children in the high treatment group have significantly better working memory skills than those in the low treatment group, with a large effect size of 0.52 SD. This result survived multiple hypothesis adjustment. There were no effects on inhibitory control, the ability to override impulse responses, or attention flexibility, the ability to deliberately focus and maintain attention. Indeed, the low treatment group appear slightly better on these domains, however the effects were not significant, and the effect sizes were small.

These results are somewhat in contrast to the results for executive functioning reported at Age 9, but they are more in line with the effects reported at the end of the programme where a significant treatment effect was found for children’s ability to control their attention, but not for their ability to delay gratification. As the same tests were administered at Age 9 and 14,

these differences cannot be attributed to the use of different tests. It will be informative to measure the cohort's executive functioning skills later in adolescence to determine whether these Age 14 results are an artefact of the timing of data collection or a true fade out of effects on inhibitory control and attention flexibility.

4.2 Socio-emotional outcomes

Socio-emotional skills were measured using a range of different instruments including the Brief Problems Monitor (Achenbach *et al.*, 2011) and the Strengths and Difficulties Questionnaire (Goodman, 1997), both of which were used at the Age 9 assessment. The *Brief Problems Monitor* (BPM) yields scores across three subscales: *internalising* ($\alpha = 0.83$), *externalising* ($\alpha = 0.77$), and *attention* ($\alpha = 0.81$) problems. The scores for each of the three problems subscales were summed to create a *Total Problems* ($\alpha = 0.87$) score. Scores were then converted to standard scores based on the child's age and gender, and binary indicators of concerning problem behaviour were created. Higher scores are indicative of more behavioural problems.

The *Strengths and Difficulties Questionnaire* (SDQ; Goodman, 1997) is a 25-item questionnaire assessing behaviours, emotions, and relationships. The instrument yields scores across five subdomains: *conduct problems* ($\alpha = 0.71$), *emotional symptoms* ($\alpha = 0.73$), *hyperactivity* ($\alpha = 0.79$), *peer problems* ($\alpha = 0.62$), and *pro-social behaviour* ($\alpha = 0.63$). The five items for each subscale were summed giving a total score of 0 to 10 for each subscale ($\alpha = 0.81$). Cutoff scores were also created to indicate scores that were of clinical concern. In all cases, apart from the prosocial behaviour, higher scores are indicative of more problems.

Three new socio-emotional measures were also included in the Age 14 Follow-up. First, the *Short Mood and Feelings Questionnaire* (Angold and Costello, 1987), which was also used in the *Growing Up in Ireland Age 13 study* and the *Millennium Cohort Study Age 14 study*, measures cognitive, affective, and behavioural-related symptoms of depression during the last two weeks. It includes 13 items which were used to create a summative score ($\alpha = 0.93$), with higher values indicating more negative feelings, as well as a cutoff score (≥ 12) indicating that the participant is at risk of depression. Second, the *Rosenberg Self-Esteem Scale* (Rosenberg, 1965), includes 10 items which were summed to create a continuous score ($\alpha = 0.89$), whereby higher values are indicative of higher levels of self-esteem. This measure was also administered to the *PFL* parents at baseline. Third, a single item life satisfaction question, which asked participants “Here is a picture of a ladder. The top of the ladder “10” is the best possible life

for you and the bottom “0” is the worst possible life for you. In general, where on the ladder do you feel you stand at the moment?”. Higher scores are indicative of greater life satisfaction. The question was also used in the *Health Behaviour of School Age Children* (HBSC) survey.

Table 5 reports the IPW-adjusted means, standard deviations, and *p*-values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size, for socio-emotional outcomes.¹¹ Consistent with Age 9, the results indicate that the *PFL* programme did not have a significant impact on children’s socio-emotional development at Age 14. Only one of the measures was statistically significant and survived adjustment for multiple hypothesis testing. Children in the high treatment group have less attention problems as measured by the *Brief Problems Monitor* scale compared to children in the low treatment group. The results were significant for both the continuous and cutoff scores, with moderate-large effect sizes of 0.61 SD and 0.46 SD respectively. For example, 63% of low treatment children were classified as having significant attention problems compared to 41% of high treatment children. For the other outcomes, the results were mainly in the right direction, e.g., high treatment children have better socio-emotional skills, but the results were not statistically significant, and the effect sizes were very small which is indicative of a true null effect rather than an underpowered study.

¹¹ The non-IPW adjusted results were mostly similar. There were few differences in the number (or level) of statistically significant results between the IPW and non-IPW adjusted models and the effect sizes were largely similar. There were two exceptions, the SDQ hyperactivity score reached statistical significance at the 10% level in the non-IPW results and the BPM attention cutoff score did not reach significance in the non-IPW results.

Table 5 Comparison of high and low treatment groups: Socio-emotional outcomes

	N (HIGH/LOW)	M_{HIGH} (SD)	M_{LOW} (SD)	p^1	p^2	ES
<i>Brief Problem Monitor Scores</i>						
BPM Internalising problems	98 (51/47)	58.80 (7.60)	58.56 (7.45)	0.655	0.655	-0.03
BPM Externalising problems	98 (51/47)	55.14 (6.12)	55.55 (6.56)	0.392	0.625	0.06
BPM Attention problems	98 (51/47)	60.75 (7.67)	65.31 (7.26)	0.007	0.015	0.61
<i>Brief Problem Monitor Cutoff Scores %</i>						
BPM Internalising problems	98 (51/47)	0.29 (0.46)	0.23 (0.43)	0.776	0.776	-0.13
BPM Externalising problems	98 (51/47)	0.09 (0.29)	0.12 (0.33)	0.280	0.528	0.11
BPM Attention problems	98 (51/47)	0.41 (0.50)	0.63 (0.49)	0.033	0.067	0.46
<i>SDQ Scores</i>						
SDQ Conduct Problems	98 (51/47)	2.04 (1.98)	2.56 (2.03)	0.108	0.421	0.26
SDQ Emotional Problems	98 (51/47)	4.16 (2.71)	4.35 (2.73)	0.525	0.702	0.07
SDQ Hyperactivity	98 (51/47)	5.57 (3.09)	6.20 (2.53)	0.196	0.472	0.22
SDQ Peer Problems	98 (51/47)	2.09 (2.03)	2.05 (1.83)	0.548	0.548	-0.02
SDQ Prosocial behaviour (+)	98 (51/47)	8.23 (1.57)	8.18 (1.71)	0.286	0.600	0.03
<i>SDQ Cutoff Scores</i>						
SDQ Conduct Problems %	98 (51/47)	0.10 (0.30)	0.08 (0.27)	0.675	0.675	-0.08
SDQ Emotional Problems %	98 (51/47)	0.31 (0.47)	0.34 (0.48)	0.517	0.828	0.08
SDQ Hyperactivity %	98 (51/47)	0.39 (0.49)	0.49 (0.51)	0.212	0.640	0.21
SDQ Peer Problems %	98 (51/47)	0.21 (0.41)	0.24 (0.43)	0.443	0.739	0.08
SDQ Prosocial behaviour %	98 (51/47)	0.06 (0.23)	0.07 (0.25)	0.343	0.633	0.04
<i>Other Socio-emotional outcomes</i>						
Short Mood and Feelings Questionnaire	98 (51/47)	6.95 (6.65)	7.84 (6.56)	0.413	0.508	0.13
Rosenberg Self-esteem Scale	98 (51/47)	18.76 (5.47)	17.60 (5.86)	0.277	0.455	0.21
Life Satisfaction (1-10)	98 (51/47)	7.48 (2.16)	7.30 (2.06)	0.435	0.435	0.08
<i>Non Stepdown Outcomes</i>						
BPM Total problems standardised score	98 (51/47)	59.62 (7.52)	61.56 (7.60)	0.147	~	0.26
BPM Total problems cutoff %	98 (51/47)	0.37 (0.49)	0.44 (0.50)	0.263	~	0.14
SDQ Total score	98 (51/47)	13.86 (7.04)	15.16 (7.25)	0.279	~	0.18
SDQ Total cutoff %	98 (51/47)	0.36 (0.49)	0.38 (0.49)	0.489	~	0.04
SMFQ Cutoff %	98 (51/47)	0.19 (0.40)	0.27 (0.45)	0.282	~	0.20

Note: N' indicates the sample size. 'M' indicates the IPW-adjusted mean. 'SD' indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional p -value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional p -value from IPW-adjusted stepdown permutation test with 100,000 replications. 'Effect size' is the ratio of the treatment effect to the pooled standard deviation.

In both the high and low treatment groups, over one-third were classified as having high or very high socio-emotional problems using the SDQ instrument. As shown in Table 6, this compares with only 6% in the nationally representative Growing Up in Ireland (GUI) cohort assessed at Age 13. In all cases, the *PFL* cohort report much higher conduct problems, emotional problems, hyperactivity, peer problems, and lower prosocial behaviour than the GUI cohort. A caveat to these results is that the *PFL* measure of SDQ was based on self-reports while the GUI measure was based on parent-reports. However, there was one socio-emotional measure that is common across both studies – the *Short Mood and Feelings Questionnaire* (SMFQ). As shown in Table 6, the *PFL* cohort have considerably higher SMFQ scores than the GUI cohort, indicating poorer mental health. In particular, over 40% of the *PFL* cohort reach the cutoff indicative of depression, compared to only 16% in the GUI cohort. Thus, the *PFL* cohort have significantly poorer socio-emotional skills than the average Irish 13-year-old.

It is also possible to compare the average life satisfaction of the *PFL* cohort to a national representative sample of 15-year-olds who took part in Health Behaviour of School Age Children (HBSC) survey and used the same instrument. The average life satisfaction of girls and boys in the *PFL* cohort was 6.8 and 7.6 respectively, compared to 5.9 and 6.7 among the HBSC sample and 7.8 and 8.3 among the GUI sample. Thus, the *PFL* cohort reported lower life satisfaction compared to the GUI cohort, but higher life satisfaction compared to HBSC cohort. Across all samples, boys reported higher life satisfaction than girls.

Table 6 Comparison *PFL* cohort at age 14 & Growing up in Ireland cohort at age 13

	<i>PFL</i> _{HIGH}	<i>PFL</i> _{LOW}	<i>GUI</i>
<i>SDQ Scores</i>			
SDQ Conduct Problems	2.04 (1.98)	2.56 (2.03)	0.93 (1.30)
SDQ Emotional Problems	4.16 (2.71)	4.35 (2.73)	2.24 (2.29)
SDQ Hyperactivity	5.57 (3.09)	6.20 (2.53)	2.61 (2.42)
SDQ Peer Problems	2.09 (2.03)	2.05 (1.83)	1.25 (1.58)
SDQ Prosocial behaviour (+)	8.23 (1.57)	8.18 (1.71)	8.74 (1.59)
SDQ Total score	13.86 (7.04)	15.16 (7.25)	6.98 (5.42)
SDQ Total cutoff %	36%	38%	6.4%
<i>Other Socio-emotional outcomes</i>			
Short Mood and Feelings Questionnaire	6.95 (6.65)	7.84 (6.56)	3.86
SMFQ Cutoff % > 8	40.38%	44.68%	15.9%
N	51	47	6650

4.3 Health outcomes

A number of different measures were used to assess health and health behaviours at Age 14. Many of the instruments were used in other cohort studies including the GUI and HBSC studies in Ireland and the MCS in the UK.

- **Self-assessed health** – Assessed using a single item “*Would you say your health in general is...Excellent, Very good, Good, Fair, Poor*”. A binary variable was created where 0 = fair/poor and 1 = excellent/very good/good.
- **Diet** – Assessed using 4 items asking “*How often do you eat breakfast/fruit/veg/fast food over a week*”. Four binary variables were created where 0 = sometimes and 1 = never.
- **Puberty** – Assessed using the *Pubertal Development Scale* (Petersen, Crockett, Richards, and Boxer, 1988). The instrument is based on three common items for girls and boys (growth spurts, skin changes, and body hair), and two additional questions for boys (voice changes, facial hair) and girls (breast development, menstruation). Response options on all items (apart from menstruation) are not yet started (1 point); barely started (2 points); definitely started (3 points); seems complete (4 points); I don’t know (missing). For the menstruation item, yes = 4 points; no = 1 point. The point values were averaged to create a Pubertal Development Scale (PDS) score, whereby higher values indicate the participant is more developed.
- **Waist-to-Height (WTH) ratio** – Assessed by measuring the participant’s height and waist circumference (by the fieldworker). The waist-to-height (WTH) ratio was calculated by waist size (cm) divided by height (cm). Higher values are associated with more health problems. A binary risk score was also created where 1 = WTH score > 0.5¹² (moderate/high risk) and 0 if < 0.5 (low risk) of obesity.
- **Substance use** – Assessed using four items asking whether the participant ever smoked cigarettes, vaped, drank alcohol, smoked cannabis, or took illegal drugs [0 = never; 1 = ever].

Table 7 reports the IPW-adjusted means, standard deviations, and *p*-values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size, for health outcomes. The results indicate that the *PFL* programme had few effects

¹² A number of studies state that a WHT ratio cutoff value of 0.5 is a suitable marker for screening of central obesity in children and adolescents of all genders (Browning *et al.*, 2010).

on participants' health at Age 14. Four of the individual measures were statistically significant in the individual tests, however none survived adjustment for multiple hypothesis correction. In addition, in the non-IPW weighted results, only two of the measures were significant. The two measures are related – waist circumference and the waist-to-height ratio (WTH). Overall, the high treatment group had a lower WTH ratio compared to the low treatment group with a moderate effect size of 0.41 SD. In addition, there was almost a six cm difference in the waist circumference of the high and low treatment groups. Although the result was not significant in the more conservative stepdown test, it is suggestive that the programme has had a long-term impact in reducing the waist size of the high treatment group. Note, that this result is consistent with findings from the Age 4 assessment whereby the high treatment group were less likely to be obese/overweight, however at Age 9, no such effect was found.

For the Age 14 assessment waist circumference (to derive the WTH ratio) was used instead of weight (to derive BMI) for two reasons. First, for young adolescents, measuring waist circumference is less sensitive than measuring body weight using a weighing scales. Second, evidence suggests that the WTH ratio is a better measure of central obesity which is a risk factor for cardiometabolic disease in both adults and children (Eslami *et al.*, 2023). Within the sample, 38% of the high treatment group and 44% of the low treatment group were classified as having a high WTH ratio, indicating that a large proportion of the cohort are at risk of obesity. The average WTH ratio was 0.47 and 0.50 for the high and low treatment group respectively. This compares with an average of 0.40 found in a representative sample of adolescences (aged 13-18) from the *Irish National Nutrition Survey* in 2020 (Moore Heslin *et al.*, 2023).

Table 7 Comparison of high and low treatment groups: Health outcomes

	N (HIGH/LOW)	M_{HIGH} (SD)	M_{LOW} (SD)	p^1	p^2	ES
<i>Health outcomes</i>						
Self-rated health %	99 (52/47)	0.79 (0.41)	0.77 (0.42)	0.517	0.764	0.04
Never eats breakfast %	99 (52/47)	0.14 (0.35)	0.25 (0.44)	0.098	0.440	0.29
Never eats fruit %	99 (52/47)	0.19 (0.40)	0.25 (0.44)	0.303	0.740	0.13
Never eats vegetables %	99 (52/47)	0.19 (0.40)	0.23 (0.42)	0.250	0.758	0.09
Eats fast food up to 6 times per week %	99 (52/47)	0.19 (0.39)	0.19 (0.40)	0.486	0.486	0.01
Puberty Development Scale Score	76 (38/38)	3.10 (0.42)	3.07 (0.51)	0.097	0.758	0.06
Waist-to-height ratio	92 (47/45)	0.47 (0.06)	0.50 (0.08)	0.053	0.271	0.41
<i>Substance use %</i>						
Ever drank alcohol	93 (47/46)	0.23 (0.43)	0.25 (0.44)	0.547	0.671	0.04
Ever smoked cigarettes	95 (49/46)	0.05 (0.22)	0.06 (0.25)	0.398	0.713	0.06
Ever vaped	98 (51/47)	0.33 (0.47)	0.23 (0.43)	0.845	0.845	-0.21
Ever tried cannabis	97 (51/46)	0.05 (0.22)	0.10 (0.31)	0.140	0.432	0.22
Ever took illegal drugs	96 (50/46)	0.00 (0.00)	0.02 (0.13)	0.106	0.514	0.25
<i>Non-stepdown measures</i>						
Waist measurement (cms)	92 (47/45)	78.82 (11.32)	84.60 (14.17)	0.024	~	0.45
Height (cms)	97 (51/46)	166.63 (8.04)	167.96 (7.45)	0.938	~	-0.17
High WTH ratio %	92 (47/45)	0.38 (0.49)	0.44 (0.50)	0.405	~	0.12
Age at first drink	25 (13/12)	13.19 (1.03)	13.04 (1.19)	0.460	~	0.14

Note: N' indicates the sample size. 'M' indicates the IPW-adjusted mean. 'SD' indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional p -value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional p -value from IPW-adjusted stepdown permutation test with 100,000 replications. 'Effect size' is the ratio of the treatment effect to the pooled standard deviation.

Regarding the other health outcomes, in most cases, the high treatment group had better health than the low treatment group, but none of these differences were statistically significant in either the individual or stepdown tests. Among the cohort, 79% and 77% of the high and low treatment groups respectively reported their health to be good or very good. This is comparable to the 72% of parents in the GUI sample who reported their Age 13 children to be 'very healthy, with no problems' (GUI, 2023).

Some of the health measures included in the Age 14 assessment were also present in the Health Behaviour of School-Aged Children (HBSC) Age 15 survey. As shown in Table 8, *PFL* boys and girls are less likely to eat breakfast, fruit, and vegetables every day compared to the HBSC cohort, illustrating that the *PFL* cohort have a poorer diet than the average Irish teenager.

In particular, 31% of the *PFL* cohort eat breakfast every day, compared with 54% in the HBSC cohort and 74% in the GUI cohort. However, they exhibit somewhat better behaviour regarding substance use. The *PFL* cohort were less likely to have smoked cigarettes and drunk alcohol compared to the HBSC cohort, however rates of vaping and cannabis use were similar in the *PFL* and HBSC cohorts. Approximately 20% of boys and 30% of girls have vaped in both cohorts, and 8% of both cohorts have used cannabis. In general, cigarette use was low, and about 40% of the samples have ever drunk alcohol. Interestingly, girls engaged in substance use more frequently than boys in both samples.

Table 8 Comparison of *PFL* cohort at age 14 & HBSC cohort at age 15

	<i>PFL</i> Boys	<i>HBSC</i> Boys	<i>PFL</i> Girls	<i>HBSC</i> Girls
Eats breakfast daily	43%	62%	24%	46%
Eats fruit daily	8%	39%	14%	42%
Eats vegetables daily	18%	42%	15%	44%
Ever smoked cigarettes	5%	12%	9%	15%
Ever vaped	23%	22%	32%	31%
Ever drank alcohol	14%	39%	39%	45%
Ever used cannabis	8%	8%	8%	8%

Note: HBSC data is taken from Health Behaviour in School-aged Children study (2023), Data browser (findings from the 2021/22 international HBSC survey): <https://data-browser.hbsc.org>.

Within the *PFL* cohort, the average puberty development score is 2.67 for boys and 3.37 for girls (with a maximum score of 4). This indicates that boys in the *PFL* cohort were at an earlier stage of pubertal development than girls at Age 14, which is the norm. As the Puberty Development Scale was also included in the Millennium Cohort Study, it is possible to compare the development of the *PFL* cohort and the MCS cohort as the average age in both cohorts at the time of interview was 14.3 years. Table 9 shows the proportion of both groups who started or completed various stages of puberty. In almost all cases, the *PFL* cohort were further along in their pubertal development than the MCS cohort. For example, 86% of the *PFL* cohort had started their growth spurt compared to 65% of the MCS cohort. In addition, the average age to begin menstruation was 11.4 years in the *PFL* cohort compared to 12.1 years in the MCS cohort. There is some evidence that lower levels of socio-economic status is associated with an earlier onset of puberty and earlier age of menarche in developed countries (e.g., Arim *et al.*, 2007; Deardorff *et al.*, 2014; James-Todd *et al.*, 2010; Sun *et al.*, 2017). Thus, the results

reported here align with this literature as the MCS is a nationally representative cohort, while the *PFL* cohort reside in disadvantaged communities. The mechanisms through which low SES may predict earlier pubertal onset are still unknown. However, one possible explanation is that early exposure to stress (initiated by low SES) may impact the epigenome and the regulation of hormones (Manotas *et al.*, 2022).

Table 9 Comparison of *PFL* cohort and *MCS* cohort at Age 14

% started/completed	<i>PFL</i>	<i>MCS</i>
Growth spurt	86%	65%
Body hair	84%	85%
Skin changes	77%	66%
Voice change (males)	71%	63%
Facial hair (males)	29%	37%
Breast growth (females)	92%	79%
Menstruation began (females) (% yes)	96%	91%
Age menstruation began (females)	11.4 yrs	12.1 yrs
N	99	11,000
Average age at time of interview	14.3yrs	14.3yrs

4.4 Educational Engagement & Time-Use outcomes

Several measures were used to assess educational engagement and time use at Age 14. As above, many of these instruments were used in GUI and MSC studies.

- **School liking** – Assessed using a single item “*How do you feel about school in general?*” (on a scale of 1-5). A binary variable was created where 0 = Dislike school and 1= Like school.
- **School engagement** – Assessed using the 4-item *Classroom Climate Measure* (Rowe *et al.*, 2010) which included items such as “*I look forward to going to school*” (on a scale of 1-5). A standardised summative score ($\alpha = 0.83$) of the 4 items was created, whereby higher values are equal to more positive school engagement.
- **School belonging** – Assessed using the 9-item *School Belonging Scale* (used in the German Socio-Economic Panel) which included items such as “*I feel like an outsider in school*” (on a scale of 1-4). A standardised factor score was created whereby higher values are equal to a greater sense of belonging ($\alpha = 0.80$).
- **School absences** – Assessed using a single item “*In the last 12 months, how often did you miss school without your parents’ permission? (most days...never)*”. A binary variable was created where 0 = Ever missed school and 1= Never missed school.
- **Expectations** – Assessed using a single item “*How likely do you think it is that you will go to third level education e.g. university/college?*” Participants were asked to select a

number from 0 to 100 using a slider. Higher scores indicate a higher expectation that the participant will attend university.

- **Time Use** – Assessed using three items asking “*On a normal weekday during term time, how many hours do you spend doing homework / using the internet / on social networking or messaging sites or Apps?*”. There were eight response options ranging from none to more than 7 hours. Continuous measures, whereby higher values mean more time spent on the activity, were used. In addition, binary variables were created whereby 1 = None and 0 = Any for homework, and 1 > 7hrs and 0 = < 7hrs for internet and social media use. Thus, in all cases, higher values correspond to more negative outcomes.

Table 10 reports the IPW-adjusted means, standard deviations, and *p*-values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size, for education engagement and time use outcomes.¹³ It shows that the programme had little impact on educational engagement and no impact on time use. For most outcomes, the high treatment group had better scores than the low treatment group, but only one outcome was statistically significant in the individual test and almost reached significance in the stepdown test (and reached significance in the non-IPW results). In particular, 66% of the high treatment group stated that they intend to go to university in the future compared to 51% in the low treatment group, with an effect size of 0.51 SD.

There were some items which were common or similar in the *PFL* assessment and the GUI and MCS assessments. For example, 66% and 51% of the high and low treatment groups respectively intend to go to university. This compares to 76% in the GUI cohort and 70% in the MCS cohort. If we compare the MCS and *PFL* cohorts, there was no statistically significant difference in the expectation of the high treatment group and the MCS cohort to go to university ($p=0.331$), however the low treatment group was significantly less likely to state that they

¹³ Similar results were found in the non-IPW results.

expect to go to university ($p=0.000$). This suggests that the programme raised the aspirations of the high treatment group to the national average in the UK.

In terms of time use, 55% of the *PFL* cohort reported doing less than 30 minutes of homework on a normal weekday. This compares to 12% in the GUI cohort and 7.9% in the MCS cohort. Also, 16% and 15% of the high and low treatment groups respectively spend more than 7 hours a day on social networking sites, compared to 9.56% in MCS. Finally, 47% of the *PFL* cohort missed school without permission in the last 12 months, compared to only 2.5% in the GUI cohort and 9.26% in the MCS cohort. Thus, the *PFL* cohort have lower levels of engagement with education compared to the national cohorts.

Table 10 Comparison of high and low treatment groups: Educational Engagement Outcomes & Time Use

	N (HIGH/LOW)	M_{HIGH} (SD)	M_{LOW} (SD)	p^1	p^2	ES
<i>Educational outcomes</i>						
Likes school %	99 (52/47)	0.76 (0.43)	0.68 (0.47)	0.301	0.610	0.18
School engagement score (std)	99 (52/47)	0.14 (0.92)	0.13 (1.14)	0.188	0.468	0.27
School belonging score (std)	99 (52/47)	0.14 (0.96)	0.00 (0.93)	0.280	0.280	0.15
School absence (% never)	99 (52/47)	0.57 (0.50)	0.48 (0.50)	0.304	0.474	0.17
Intention of going to university/college %	99 (52/47)	66.34 (28.55)	50.73 (32.58)	0.012	0.100	0.51
<i>Time use continuous measures</i>						
Time spent on homework per day (1-8)	99 (52/47)	2.43 (1.37)	2.17 (1.24)	0.173	0.392	0.20
Time spent on social media per day (1-8)	99 (52/47)	5.15 (2.21)	5.65 (1.72)	0.239	0.393	0.26
Time spent on internet (1-8)	99 (52/47)	6.27 (1.97)	6.55 (1.55)	0.321	0.321	0.16
<i>Time use binary outcomes</i>						
No homework %	99 (52/47)	0.33 (0.47)	0.42 (0.50)	0.192	0.538	0.19
>7hrs social media %	99 (52/47)	0.16 (0.37)	0.15 (0.37)	0.604	0.604	-0.01
>7hrs internet %	99 (52/47)	0.34 (0.48)	0.37 (0.49)	0.487	0.649	0.05

Note: N' indicates the sample size. 'M' indicates the IPW-adjusted mean. 'SD' indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional p -value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional p -value from IPW-adjusted stepdown permutation test with 100,000 replications. 'Effect size' is the ratio of the treatment effect to the pooled standard deviation.

4.5 Antisocial Beliefs and Attitudes outcomes

Antisocial beliefs were measured using the *Antisocial Beliefs and Attitudes Scale* (Butler *et al.*, 2015) which has 28 items measured on a 4-point scale from ‘Strongly disagree’ to ‘Strongly agree’. The instrument contains 3 sub-domains 1) Rule non-compliance ($\alpha = 0.73$) e.g., “*I don’t like having to obey all the rules at home and in school*”, 2) Peer conflict ($\alpha = 0.72$) e.g., “*It’s ok to walk away from a fight*”, and 3) Aggression ($\alpha = 0.72$) e.g., “*It’s ok to hit my mother as long as I don’t hurt her*”, and a total score ($\alpha = 0.82$). Summative scores were created whereby higher scores indicate greater acceptance of antisocial behaviours.

Table 11 reports the IPW-adjusted means, standard deviations, and p -values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size, for the participants’ antisocial beliefs and attitudes.¹⁴ It shows that there were no significant effects on antisocial beliefs and attitudes at Age 14. In some cases, the high treatment group reported higher (i.e., worse) scores, but these were not statistically significant, and the effect sizes were small. Thus, overall, the programme had no impact on beliefs around antisocial behaviour.

Table 11 Comparison of high and low treatment groups: Antisocial Beliefs

<i>Antisocial Beliefs and Attitudes Scale</i>	N (HIGH/LOW)	M_{HIGH} (SD)	M_{LOW} (SD)	p^1	p^2	ES
Rule non-compliance	99 (52/47)	12.25 (3.66)	12.98 (4.31)	0.260	0.483	0.18
Peer conflict	99 (52/47)	9.87 (3.72)	9.21 (4.52)	0.702	0.702	-0.16
Aggression	99 (52/47)	4.70 (3.06)	4.44 (3.53)	0.661	0.838	-0.08
<i>Non stepdown measure</i>						
Total score	99 (52/47)	26.82 (7.62)	26.63 (9.53)	0.549	0.929	-0.02

Notes: N’ indicates the sample size. ‘M’ indicates the IPW-adjusted mean. ‘SD’ indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional p -value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional p -value from IPW-adjusted stepdown permutation test with 100,000 replications. ‘Effect size’ is the ratio of the treatment effect to the pooled standard deviation.

4.6 Parent-Child Relationship outcomes

Two instruments were used to assess the child’s perception of their relationship with their parents. First, the *Inventory of Parents and Peer Attachment-Revised* (IPPA; Gullone and

¹⁴ In the non-IPW results, there was a significant treatment effect on the Rule non-compliance sub-domain, with the high treatment group reporting better compliance than the low treatment group.

Robinson, 2005). This instrument included 50 items measured on a 5-point scale from ‘Always true’ to ‘Never true’ - 25 pertaining to the participant’s relationship with their mother (or the person acting as their mother), and 25 pertaining to their father (or the person acting as their father). The instrument contains three sub-domains 1) Trust ($\alpha_M = 0.89$; $\alpha_F = 0.93$) e.g., “*My mother/father respects my feelings*”, 2) Communication ($\alpha_M = 0.89$; $\alpha_F = 0.91$) e.g., “*My mother/father can tell when I’m upset about something*”, and 3) Alienation ($\alpha_M = 0.81$; $\alpha_F = 0.86$) e.g., “*I don’t get much attention from my mother/father*”, and a Total Score for each parent. For the Trust and Communication sub-domains, higher scores represent more positive outcomes, while higher scores on the Alienation sub-domain represents more negative outcomes. The total score is the sum of the Trust and Communication sub-domains, minus the Alienation sub-domain.

Second, the parent-child relationship was assessed using 10 items (five for mothers, and five for fathers) used in the German Socio-economic Panel (SEOP) survey. Each item asks the participant how often 1) *they turn to their mother/father when worried about something*, 2) *their mother/father encourages or helps them when something is important*, 3) *their mother/father orders them around*, 4) *their mother/father tells them it’s important to do well in school and to study a lot*, and 5) *they argue with their mother and father*. Each item was assessed on a 5-point scale from ‘Very often’ to ‘Never’. Binary variables were created whereby 1 = Often/Very often and 0 = Less often/seldom/never.

Table 12 reports the IPW-adjusted means, standard deviations, and p -values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size, for parent-child relationship outcomes.¹⁵ It shows that the programme had an impact on some dimensions of the parent-child relationship. For the IPPA instrument, the high treatment group reported significantly better communication with their mothers and trust with their fathers, with effect sizes of 0.36 and 0.45 SDs respectively. Neither result, however, survived adjustment for multiple testing. The total IPPA score for mothers was also statistically significant with an effect size of 0.37 SD. For the SOEP items, the high treatment group reported a significantly higher likelihood of turning to their mothers when they are worried, with an effect size of 0.70 SD. They also reported that their mothers and fathers were more likely to encourage or help them with something important, with effect sizes of 0.48, and 0.48

¹⁵ In the non-IPW results, the results were similar to the IPW-adjusted results, however the IPPA Communication with mothers result was not statistically significant in the non-IPW results.

SDs respectively. Importantly, the result for ‘*turning to mother when worried*’ survived multiple hypothesis adjustment. Specifically, 75% of the high treatment group stated that they often or very often turn to their mothers when worried compared to 42% in the low treatment group. Interestingly, there was one treatment effect in a non-hypothesized direction. The low treatment group was more likely to report that their mothers tells them it is important to study (79% in the low treatment group vs 63% in the high treatment group). It is possible that the parents of high treatment children do not need to remind their children to study as, given their higher cognitive skills, they are likely to study without the need for reminders or indeed, they need to study less due to their higher abilities. This is consistent with the finding that they spend more time on homework as shown in Table 9 (although this result was not statistically significant).

Table 12 Comparison of high and low treatment groups: Parenting Relationship outcomes

	N (HIGH/LOW)	M _{HIGH} (SD)	M _{LOW} (SD)	p ¹	p ²	ES
<i>IPPA</i>						
Communication Mother	98 (51/47)	35.93 (7.64)	33.12 (7.91)	0.058	0.244	0.36
Trust Mother	98 (51/47)	37.94 (5.58)	38.40 (7.15)	0.690	0.690	-0.07
Alienation Mother	98 (51/47)	12.59 (5.05)	13.77 (4.86)	0.211	0.410	0.24
Communication Father	98 (51/47)	31.17 (9.11)	28.49 (10.78)	0.142	0.392	0.27
Trust Father	98 (51/47)	37.10 (7.64)	32.85 (11.09)	0.029	0.150	0.45
Alienation Father	98 (51/47)	12.51 (5.53)	13.76 (6.26)	0.313	0.429	0.21
<i>SOEP</i>						
Turn to when worried (Mother) %	98 (51/47)	0.75 (0.44)	0.42 (0.50)	0.001	0.027	0.70
Encourages/helps you with something important (Mother) %	98 (51/47)	0.90 (0.30)	0.72 (0.45)	0.030	0.242	0.48
Tells you it's important to study Mother %	97 (50/47)	0.63 (0.49)	0.79 (0.41)	0.972 [^]	0.972	-0.36
Argues with mother %	96 (50/46)	0.22 (0.42)	0.15 (0.36)	0.199	0.729	0.19
Orders you around (Mother) %	96 (51/45)	0.29 (0.46)	0.37 (0.49)	0.810	0.979	-0.17
Turn to when worried (Father) %	88 (49/39)	0.40 (0.49)	0.39 (0.50)	0.518	0.942	0.00
Encourages/helps you with something important (Father) %	87 (48/39)	0.81 (0.40)	0.60 (0.50)	0.030	0.268	0.48
Tells you it's important to study (Father) %	85 (47/38)	0.56 (0.50)	0.69 (0.47)	0.885	0.984	-0.27
Argues with father %	85 (48/37)	0.25 (0.44)	0.13 (0.34)	0.138	0.583	0.29
Orders you around (Father) %	84 (48/36)	0.17 (0.38)	0.30 (0.47)	0.817	0.979	-0.31
<i>Non stepdown measures</i>						

IPPA Total Score Mothers	98 (51/47)	55.76 (20.06)	47.58 (23.71)	0.077	0.527	0.37
IPPA Total Score Fathers	98 (51/47)	61.28 (16.63)	57.75 (18.21)	0.228	0.894	0.20

Notes: N' indicates the sample size. 'M' indicates the IPW-adjusted mean. 'SD' indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional *p*-value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional *p*-value from IPW-adjusted stepdown permutation test with 100,000 replications. 'Effect size' is the ratio of the treatment effect to the pooled standard deviation. ^ result in non-hypothesised direction.

4.7 Time and Risk Preference outcomes

Time and risk preferences were assessed in two ways, First, using two self-assessed single item questions which were also asked in the MCS study. Second, using task-based elicitation across a series of games, which are the standard means of measuring preferences in an experimental setting.

For the self-assessed items, time preferences were measured by asking “*On a scale of 0-10, where 0 is never and 10 is always, how patient would you say you are?*”. Higher levels of patience is associated with lower time preferences i.e., the participant places a higher value on the future and a lower value on the present. While low levels of patience is associated with higher time preferences i.e., the participant places a higher value on the present and a lower value on the future. In general, low time preferences are associated with more positive/healthy behaviours and high time preferences are associated with more negative/unhealthy behaviours (e.g., smoking, drinking, not studying etc). Higher scores on the measure presented in Table 13 is indicative of lower time preferences (i.e., more patience).

Risk preferences were measured by asking “*On a scale of 0-10, where 0 is never and 10 is always, how willing to take risks would you say you are?*”. Depending on the domain, a person who is willing to take more risks may have more negative outcomes (e.g., driving fast) or more positive outcomes (e.g., financial investments). Higher scores on the measure presented in Table 13 is indicative of being less risky.

Table 13 report the IPW-adjusted means, standard deviations, and *p*-values that result from weighted individual and stepdown permutation tests, controlling for gender, alongside the effect size, for self-reported time and risk preferences.¹⁶ It shows that the high treatment group have lower time preferences (i.e., they place a higher value on the future than the present) compared to the low treatment group. While the result was not statistically significant, the

¹⁶ The non-IPW results were similar to the IPW-adjusted results.

effect size of 0.31 SD was sizeable, which suggests that the study could be underpowered to detect a significant effect on this item. In terms of risk preferences, the low treatment group report being less risky than the high treatment group, but again, the result was not statistically significant, and the effect size was very small (-0.14 SD), suggesting that the programme had no impact on risk preferences.

It is possible to compare the average time and risk preferences of the *PFL* cohort to the MCS cohort. Within the MCS cohort, the average time preference score was 5.69 (on a 0-10 scale), compared to an average of 5.40 in the high treatment group and 4.66 in the low treatment group. This result was as expected as the MCS is a nationally representative sample of children, thus we would expect them to exhibit higher levels of patience than children living in a disadvantaged community. It is interesting that the high treatment group was closer to the average MCS score than the low treatment group. Indeed, there was no statistically significant difference in scores between the MCS cohort and the high treatment group ($p = 0.374$), however there was a significant difference in the scores of the MCS cohort and the low treatment group ($p = 0.004$). This suggests that the programme played a role bringing the time preferences of the high treatment group into line with the national (UK) average.

A similar analysis was conducted for the risk preference measure. Within the MCS cohort, the average risk preference score was 4.90 (on a 0-10 scale) whereby higher values correspond to being less risky. This compares to an average of 5.14 in the high treatment group and 5.38 in the low treatment group. Thus, in both cases, the *PFL* cohort were less likely to take risks than the MCS cohort, but the differences were not statistically significant (MCS v High treatment: $p = 0.446$; MCS v Low treatment: $p = 0.333$).

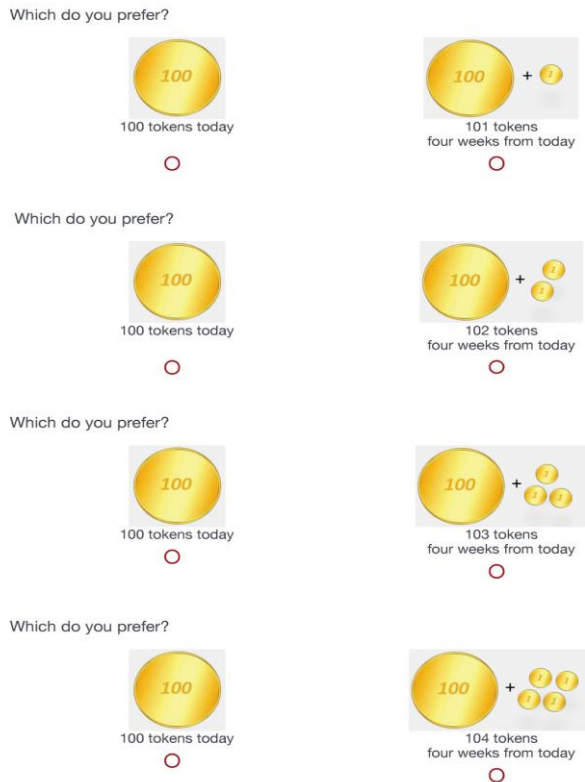
For the task-based elicitation methods, participants were asked to complete three games assessing time preferences (Coller and Williams, 1999) and two games assessing risk preferences (Holt and Laury, 2002) on an iPad.

To measure time preferences, participants were asked to make a choice between receiving 100 tokens today or X tokens four weeks from today (Game 1 – 9 choices), or 100 tokens today or X tokens eight weeks from today (Game 2 – 10 choices), or 100 tokens 4 weeks from today or X tokens eight weeks from today (Game 3 – 9 choices). In each game, the number of tokens increased sequentially from 101 to 125 in Game 1, from 101 to 150 in Game 2, and from 101 to 125 in Game 3. An example is provided below. Each game was analysed

separately. Note, these choices were hypothetical i.e., participants did not receive any monetary payment based on their choice.

Responses to these games were used to create two measures used in the outcome analysis. The first assessed the proportion of times the participant chose the 'later' option over the 'sooner' option e.g., if they chose the 'later' option 3 times and the 'sooner' option 6 times, they received a score of 67% (6/9). If they always chose the 'sooner' option, they received a score of 0%, and if they always chose the 'later' option, they received a score of 100%. Thus, higher values correspond to lower time preferences e.g., the participant is more patient. The second measure assessed the first point at which the participant switched from the 'later' option to the 'sooner' option, e.g., if they chose the 'sooner' option for their first 7 choices, and switched to the 'later' option for their 8th choice, they received a score of 8. If they always chose the 'sooner' option they received a score of 10 and if they always chose the 'later' option, they received a score of 1. If their first choice was the 'later' option, they received a score of 1. For this measure, higher scores correspond to higher time preferences e.g., the participant is less patient.

Table 13 shows that there are no statistically significant differences across the high and low treatment groups on any of the time preference measures assessed. In addition, the effect sizes were very small, suggesting that the programme had no impact on time preferences. For example, the high treatment group chose the 'later' option over the 'sooner' option 31% of the time, compared to 32% of the time for the low treatment group.



Example of task-based elicitation of time preferences

To measure risk preferences, participants were asked to make a choice between two options: receiving X tokens ‘for certain’ or having their reward determined by a coin flip where there was a 50% chance of winning Y tokens and a 50% chance of winning Z tokens. For example, would you prefer option A “*to receive 70 tokens for sure*” or option B “*to have a 50% chance of winning 10 tokens or a 50% chance of winning 100 tokens*”. Participants were informed that each token represents 10 cents, and that they should treat each decision as if it were real money, however, they did not receive a real monetary payout. In the first game, participants made 11 choices whereby the ‘risky’ choice did not change (50% chance of winning 10 tokens and 50% chance of winning 100 tokens), but the value of the ‘certain’ choice descended by 5 tokens (from 70 to 20 tokens) over the course of the 11 questions. For the second game, participants made 11 choices whereby the ‘risky’ choice did not change (50% chance of winning 20 tokens and 50% chance of winning 200 tokens), but the value of the ‘certain’ choice ascended by 5 tokens (from 40 to 140 tokens) over the course of the 11 questions. An example is provided below.

Responses to these games were used to elicit the two measures used in the outcome analysis. The first assessed the proportion of times the participant chose the ‘safe’ option over

the 'risky' option e.g., if they chose the risky option 5 times and the safe option 6 times, their risk preference score was 55% (6/11). If they always chose the safe option, their score was 100%. If they always chose the risky option, their score was 0%. Thus, higher values correspond to less risky choices. The second measure assessed the first point at which the participant switched from choosing the safe option to the risky option e.g., if they chose the safe option for their first 9 choices and then chose the risky option for their 10th choice, they received a score of 10. If they always chose the safe option, they received a score of 12, and if they always chose the risky option, they received a score of 1. If their first choice was the risky option, they received a score of 1. Again, higher values correspond to less risky choices. Each game was analysed separately.

Table 13 shows that the high treatment had consistently lower scores than the low treatment group i.e., they were more risky loving and less risk averse. In these analyses, we test the hypothesis that the high treatment group was less risky than the low treatment group (in a one-tailed test), and we find no evidence that this hypothesis was supported. Indeed, the high *p*-values on Game 1 (for both measures), indicates that the high treatment group was significantly more risky than the low treatment group. For example, the high treatment group chose the safe option over the risky option 35% of the time, while the low treatment group chose the safe option 45% of the time. Thus, the programme may have led to lower risk aversion. As discussed above, the willingness to take risks may have positive or negative outcomes depending on the domain and the level of risk considered. There is some research which shows that individuals with higher levels of cognitive skills are less risk averse than individuals with lower levels of cognitive skills. For example, Andreoni *et al.* (2000) finds that children and adolescents with higher cognitive skills (especially math skills) are more willing to take risks, possibly as it impacts their ability to process information on probability. As the high treatment group have significantly higher cognitive skills than the low treatment group, this may explain this finding on risk preferences. However, an additional analysis showed that there was no correlation between the cognitive scores and risk preferences of the *PFL* cohort (correlation coefficient = 0.04).


Option A




70 tokens for sure

☐

Option B

Heads = 

OR

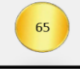
Tails = 

10 tokens with a chance of 50%,
100 tokens with a chance of 50%

☐

Which do you prefer?


Option A




65 tokens for sure

☐

Option B

Heads = 

OR

Tails = 

10 tokens with a chance of 50%,
100 tokens with a chance of 50%

☐

Which do you prefer?

Option A

Option B

Example of task-based elicitation of risk preferences

Table 13 Comparison of high and low treatment groups: Time & Risk Preferences

	N (HIGH/LOW)	M_{HIGH} (SD)	M_{LOW} (SD)	p^1	p^2	ES
<i>Self-assessment</i>						
Time preferences (higher = more patient)	96 (52/44)	5.40 (2.31)	4.66 (2.50)	0.149	0.223	0.31
Risk preferences (higher= less risky)	99 (52/47)	5.14 (1.56)	5.38 (1.87)	0.566	0.566	-0.14
<i>Time preference games: % Later choice (more patient)</i>						
Game 1	98 (51/47)	0.31 (0.32)	0.32 (0.27)	0.582	0.582	-0.01
Game 2	98 (51/47)	0.29 (0.26)	0.27 (0.25)	0.426	0.645	0.09
Game 3	98 (51/47)	0.33 (0.32)	0.31 (0.26)	0.498	0.631	0.07
<i>Time preference games: First switch (less patient)</i>						
Game 1	98 (51/47)	6.89 (3.19)	6.50 (2.98)	0.738	0.738	-0.13
Game 2	98 (51/47)	7.85 (2.86)	7.76 (3.16)	0.579	0.708	-0.03
Game 3	98 (51/47)	6.55 (3.38)	6.72 (2.87)	0.504	0.656	0.06
<i>Risk preference games: % Safe choice (less risky)</i>						
Game 1	98 (51/47)	0.35 (0.25)	0.45 (0.28)	0.932^	0.932^	-0.36
Game 2	98 (51/47)	0.53 (0.26)	0.60 (0.27)	0.794	0.900	-0.25
<i>Risk preference games: First switch (less risky)</i>						
Game1	98 (51/47)	3.48 (2.62)	4.47 (3.34)	0.926^	0.926^	-0.33
Game1	98 (51/47)	1.39 (1.41)	1.55 (1.93)	0.691	0.872	-0.09

Notes: N' indicates the sample size. 'M' indicates the IPW-adjusted mean. 'SD' indicates the IPW-adjusted standard deviation. ¹ one-tailed (right-sided) conditional p -value from individual IPW-adjusted permutation test with 100,000 replications. ² one-tailed (right-sided) conditional p -value from IPW-adjusted stepdown permutation test with 100,000 replications. 'Effect size' is the ratio of the treatment effect to the pooled standard deviation. ^ indicates the result was significant in the non-hypothesized direction.

5 Conclusions

The aim of the Age 14 Follow-up study was to examine whether the large and significant impacts of *PFL* found at the end of the programme and at Age 9 were sustained. Prior evidence on the long-term impact of home visiting programmes into adolescence is inconclusive, as very few studies continue to track children beyond the lifetime of the intervention, and of the few NFP studies, that do, they fail to find effects. In contrast, this study found that *PFL* continues to have a sizeable impact on children's cognitive skills approximately ten years after the participants finished the programme. There was no evidence of cognitive fade-out, with effect

sizes of 0.70 of a standard deviation on overall cognitive ability, and significant effects on some dimensions of executive functioning and health.

Overall, the IQ scores of the *PFL* children were above that of their parents (i.e., the Flynn effect), however the correlation between the high treatment children and their mothers was small and not statistically significant at either age five ($r^{17} = 0.07, p = 0.562$) or age nine ($r = 0.18, p = 0.148$) and significant at age 14 ($r = 0.34, p = 0.015$), compared to the large and significant correlation between the low treatment children and their mothers at age five ($r = 0.31, p = 0.018$), age nine ($r = 0.57, p = 0.001$), and age 14. ($r = 0.54, p = 0.001$).¹⁸ Thus the programme appears to be effective in reducing the intergenerational transmission of IQ scores, however, the correlation between parents and children's IQ is growing over time, even within the high treatment group.

The programme impacted all dimensions of cognitive skill including spatial ability, non-verbal ability, and verbal ability, in addition to reducing the proportion of children scoring below the standardised norm. The magnitude of the cognitive effects at age 14 (0.54 – 70 SD) are similar to those observed at the end of the programme (0.56 - 0.77 SD) and at age 9 (0.39 - 0.76 SD). An additional analysis found that controlling for age five cognitive scores slightly reduces the size of the age 14 treatment effects, however the impact of the programme was still statistically significant.¹⁹ This suggests that *PFL* is continuing to have an impact on children's development beyond the lifetime of the programme. This provides evidence in support of the skill formation model (Cunha and Heckman, 2007) which posits that developing children's skills early in life helps them to develop more advanced skills later in life (a process called self-productivity), and this raises the effectiveness of later investments, such as investments in schooling (a process called dynamic complementarity). If this process continues, and the high treatment group continue to utilise their higher cognitive skills, this is likely to translate into improved outcomes throughout the life cycle.

The size of the cognitive effects were substantially larger than those found in much of the existing literature. For example, the meta-analyses discussed earlier in the report found effect sizes of less than 0.30 for cognitive outcomes (e.g., Layzer *et al.* 2001; Sweet and Appelbaum, 2004; Miller *et al.* 2011; Filene *et al.* 2013; Rayce *et al.* 2017). The effects

¹⁷ r is Pearson's correlation coefficient.

¹⁸ Maternal IQ was measured using the Weschler Abbreviated Scale of Intelligence which assesses cognitive ability across four subscales: vocabulary, similarities of constructs, block design, and matrix reasoning. From this, standardised measures of verbal ability, perceptual reasoning, and a full-scale measure of cognitive functioning, standardised to have a mean of 100 and standard deviation of 15, are generated. The full-scale measure was used in this analysis to correspond with the measure of General Conceptual Ability from the BAS.

¹⁹ Available upon request.

were also larger than the German home visiting programme, *Pro Kind*, which found average effect sizes for cognition of 0.20 - 0.30 SD for girls only at age 2 (Sandner and Jungmann, 2017). The effect sizes were also larger than the NFP Memphis trial which reported effects of 0.13 - 0.27 SD for cognitive skills at age six (Heckman *et al.* 2017). The *PFL* effects were more similar in magnitude to those found in studies of low and middle income countries. The Jamaica study, which was based on weekly home visits for two years starting between nine and 24 months, identified no significant effects on IQ at ages seven to eight, however the cognitive effects re-emerged at the 11, 17, and 22 year follow-ups with effect sizes ranging from 0.40 to 0.60 (Grantham-McGregor and Smith, 2016).

In addition to cognitive skills, another key aspect of children's development is their executive functioning skills, especially self-regulation, as independent of IQ, self-regulation has been shown to predict later academic performance, health, and finances in adulthood (Blair and Raver, 2012; Liew, 2012). However, children from disadvantaged backgrounds typically have poorer self-regulation (Evans and Rosenbaum, 2008). Thus improving children's early skills in these domains could yield cascading benefits into adulthood (Diamond and Lee, 2011). At age nine we found that the programme had a large and substantive impact on all dimensions of the children's executive functioning skills with effect sizes ranging from 0.56 – 0.65 SDs. However, at age 14, we only found effects on working memory (0.52 SD), which means that the high treatment group are better able to retain, manipulate and use information over brief periods of time. There were no effects on their ability to override their automatic impulses (inhibitory control) or maintain and focus their attention (attention flexibility). In addition, the effect sizes on these measures were small and in the non-hypothesized direction, which suggests that the study was not under-powered to detect these effects; rather the effects observed at age nine have dissipated. There is no clear explanation as to why this occurred. It is possible that hormonal fluctuations and an overactive limbic system associated with the teenage period may impede their ability to delay gratification and focus (Arain *et al.* 2013). Somewhat relatedly, the programme had no impact on the participant's attitudes towards anti-social behaviour.

Similarly, the programme had no impact on various dimensions of educational engagement including liking school, a sense of belonging at school, and school absences. Given the finding that the programme raised cognitive scores, it is somewhat surprising that this did not translate into higher levels of school engagement or more positive feelings about school among the high treatment group. However, it is important to note that the majority of the *PFL*

cohort felt positive about school, with 76% and 68% in the high and low treatment groups respectively reporting that they like school, yet over 50% report missing school (without permission) in the last 12 months. Related to this, the results showed that the programme had no impact on changing time use patterns regarding homework or internet/social media use. Over one-third of students reported not doing any homework during the week, which is substantially higher than the 12% reported in the GUI cohort. The cohorts were somewhat more similar regarding social media use, with about 15% of the *PFL* cohort and 10% of the MCS cohort spending more than seven hours a day on social networking sites.

While there were no significant differences across the groups about current school engagement, there was a statistically significant difference regarding future educational engagement with a significantly higher proportion of the high treatment group reporting that they intend to attend university compared to the low treatment group (66% v 51%). While lower than the national average found in the GUI cohort (76%), this finding suggests that the programme may have been effective in improving the educational aspirations of the students.

The results also found that the significant effects observed for children's socio-emotional development at age four are no longer present at age 14, which is largely in line with the age nine findings. At earlier time points we found that the programme was effective in reducing the proportion of children within the clinical range of behavioural problems, however, few effects were identified for continuous scores of children's socio-emotional development. These earlier measures were based on parent reports only. At age nine and 14 we assessed socio-emotional skills using child reports. At age 14, only two of the 22 sub-domains considered were statistically significant (and remained significant in the stepdown tests) – the Brief Problem Monitors Attention Problems continuous score and the cutoff score. Specifically, 41% of the high treatment group were classified as having clinically significant attention problems, compared to 63% in the low treatment group. Thus, the programme was effective in reducing the incidence of attention problems, which means the high treatment group were less likely to have problems with concentration or sitting still. However, none of the other measures assessing different dimensions of socio-emotional skills, depression, self-esteem, or life satisfaction were impacted. In addition, a comparison of the *PFL* cohort to other nationally representative samples showed that the *PFL* cohort have poorer socio-emotional skills overall. For example, over 40% of the *PFL* cohort reached the cutoff for depressive symptoms, compared to only 16% in the national sample. This is consistent with findings that

families from low socioeconomic status communities face more mental health challenges than those from other more affluent communities (Kirkbride *et al.*, 2024).

In-line with studies of other home visiting programmes, there was little evidence that the programme continued to have an impact on children's health. While few studies have examined the long-term impact of home visiting programmes on health, those that do, typically find little evidence of effects (Dumont *et al.* 2010; Kitzman *et al.* 2010; Minkowitz *et al.* 2007; Olds *et al.* 2004; Olds *et al.* 2007). At age 14, the *PFL* programme had no impact on children's general health, health behaviours regarding substance use, or dietary intake. The only significant finding with a moderate effect size was on the waist-to-height ratio, which was primarily driven by a reduction in the participant's waist size (a six cm difference in waist circumference between the high and low treatment groups). As discussed earlier, WTH was used instead of BMI as there is evidence that it is a more reliable measure of obesity. That said, at previous waves (age four and nine), we assessed BMI. While no effects were found at age nine, there were significant differences at age four. In particular, a lower proportion of high treatment children (26%) were categorised as overweight or obese compared to the low treatment group (41%). The re-emergence of an effect at age 14 may be a result of using a different instrument or due to natural growth fluctuations during this period. However, as identified at previous waves, a significant proportion of the *PFL* cohort are still at risk of obesity (>40%). This compares with 21% in the GUI cohort.²⁰ This again speaks to evidence that children from lower socioeconomic backgrounds have an increased risk of obesity (Cronin *et al.*, 2022). This is consistent with the finding that the *PFL* cohort experienced an earlier onset of puberty than the MCS cohort (who are of a similar age) in the UK. Earlier onset has been associated with poorer outcomes in adolescence such as more emotional and behavioural problems, as well as in later life such as cardiometabolic diseases (Day *et al.*, 2015; Mensah *et al.*, 2013).

For the first-time, the study measured the quality of the parent-child relationship from the perspective of the child. While only one of the outcomes survived adjustment for multiple hypothesis testing, there was some evidence that the programme improved the child's relationship with both their mothers and fathers (or father figures). The high treatment group reported that they were more likely to turn to their mothers when they were worried and that their parents (mothers and fathers) were more likely to help them with something important. In

²⁰ Note, in GUI at age 13, BMI was measured using self-reported height and weight. In *PFL*, measured waist size was used instead of weight.

addition, they reported better communication with their mothers and more trust with their fathers. The effect sizes of 0.30-0.70 SDs indicate that these are sizeable impacts. These results were in contrast to findings at previous waves where we found no effects on the parent-child relationship when using parent reports (e.g., on the Condon Maternal Attachment Scale, the Maternal Separation Anxiety Scale, or the Parental Acceptance-Rejection Questionnaire). However, we did identify significant effects during the trial on certain dimensions of parenting related to parent behaviour. For example, Doyle *et al.* (2017a) identified significant treatment effects on parenting skills at six and 18 months in terms of improving the quality of the home environment, O'Sullivan *et al.* (2017) found positive treatment effects regarding improved nutrition at 24 months, and Doyle and *PFL* Evaluation Team (2015) reported improved parenting behaviour regarding the use of appropriate disciplinary techniques and increased parental interactions. These practices, interactions, and activities are recognised as key means of stimulating children's development (Farah *et al.* 2008); however, they may have also positively impacted the child's perception of their parents in adolescence.

Also, for the first time, we measured the participant's time and risk preferences. Time preferences were measured by asking participants how patient they are and then conducting a series of games where the participant could choose between a sooner, but smaller payment, or a larger, but later payment. Time preferences are important as previous studies have found that high time preferences (e.g., less patience) is associated with more disciplinary referrals in school, higher school dropout, less saving, and poorer health behaviours in adolescence (Benjamin, Brown, Shapiro, 2013; Castillo *et al.*, 2011; Castillo *et al.*, 2019; Sutter *et al.*, 2013). In general, the literature finds that children from lower socioeconomic families make more impatient choices (Sutter, Zoller, and Glätzle-Rützler, 2019). Indeed, if we compare the time preferences of the *PFL* and MCS cohorts (using the single self-assessed question) we find that the *PFL* cohort as a whole are less patient, however, the high treatment group exhibits higher levels of patience than the low treatment group. While the difference was not statistically significant, the effect size of 0.31 SD suggests that the study may be underpowered to detect the effect. In addition, the difference between the high treatment group and the MCS cohort was not statistically significant which suggests that the programme may have changed the participant's time preferences, indicating the malleability of time preferences to early intervention.

Risk preferences were measured using a single item question asking participants how risky they are, as well as through a series of games that asked participants to decide between a

safe amount of tokens and a lottery that pays either a higher or lower amount of tokens than the safe alternative. In general, higher levels of risk taking are associated with poorer educational outcome (e.g., Castillo *et al.*, 2018). Evidence suggests that children from lower socioeconomic status families are more likely to take risks (Sutter, Zoller, and Glätzle-Rützler, 2019). However, we found that the *PFL* cohort was less risky than the MCS cohort, although the difference was not statistically significant. The findings regarding risk preferences are less straightforward. We hypothesized that the programme would result in higher levels of risk aversion among the high treatment group (e.g., less risky), however we found the opposite. In two of the games, the high treatment group were significantly less likely to take the safe option over the risky option. However, given that the *PFL* cohort has lower risk preferences than a nationally representative sample, the programme's impact on increasing the likelihood of taking a risk, may not necessarily be an issue. While excessive risk taking may have negative consequences, exhibiting moderate level of risk may be an optimal strategy in some cases.

In total, we found that 25 of the 105 outcomes tested (24%) reached conventional levels of statistical significance in the individual tests, and seven of the 20 stepdown tests (35%) were significant. As we used a 10% cutoff level, this indicates that these findings were unlikely to be a result of random Type I errors and that the programme continues to have a significant impact on families. While 43% of the original sample recruited during pregnancy participated at the Age 14 Follow-up, the treatment groups were still balanced on all key baseline characteristics. Thus, the results were unlikely to be subject to bias.

This is one of the few experimental home visiting programmes that has tracked participants into adolescence and found evidence of long-term effects ten years after the families have finished the programme. Although it is difficult to fully compare the results from different home visiting studies due to wide variations in programme goals, target groups, and implementation practices (Gomby *et al.* 1999), the larger effect sizes identified for the *PFL* programme, particularly for the cognitive outcomes, may be attributed to its prenatal start, its longer programme length, its multiple connected treatments, and its inclusive eligibility criteria. In particular, *PFL* both starts earlier and is longer in duration than most other home visiting programmes. Given that *PFL* has had consistent effects on what were traditionally called 'hard' skills (e.g., cognitive skills, working memory, attention), and no lasting effects on 'soft' skills, suggests that the programme may have had a permanent impact on areas of the brain that are particularly malleable to intervention during early childhood (e.g., the frontal lobe) (Knudsen *et al.*, 2006). The *PFL* home visitors worked with participants for a substantial

and critical period of their children's lives; therefore the positive and sizable treatment effects may be a result of the strength and quality of the home visitor-parent relationship which was given an appropriate length of time to build and develop. This is consistent with the home visiting literature which finds that the bond between parents and programme staff is key for understanding programme effects (Wesley, Buysse, and Tyndall 1997).

The larger effects may also be attributed to the extensive and diverse supports offered to the high treatment group. The *PFL* treatment included baby massage classes during infancy and the *Triple P* programme from age two, yet the majority of the other standalone home visiting programmes, such as Nurse Family Partnership and its European equivalents, did not provide such supports. Therefore, a multi-component approach, which offers supports in a variety of formats and settings may help to engage families who favour one form of treatment over another. The larger effect sizes may also be attributed to the nature of the sample. Compared to many other home visiting programmes which include ethnically diverse samples, the *PFL* cohort is relatively homogenous, consisting mainly of ethnically-Irish born participants. This, coupled with the individual-level randomisation in a confined geographical space, reduces variability within the sample, and allows us to uncover treatment effects if indeed they exist.

To conclude, the sizable cognitive advantages generated by the *PFL* programme are likely to have positive impacts on the participant's outcomes throughout life. Thus, it is critical that we continue to track the *PFL* cohort as they progress through secondary school and potentially, into higher level education. One concern as we move forward with the *PFL* evaluation is sample size. While a response rate of 43% was achieved at the Age 14 Follow-up, this figure is likely to reduce further as the participants start attending university or leaving the family home. Therefore, maintaining the *PFL* cohort should be a priority if we are to assess the long-term impact of the programme. This is particularly important given the magnitude of the cognitive effects, especially in comparison to other home visiting programmes, as *PFL* can provide a model for other communities aiming to reduce long-term socioeconomic inequalities.

References

- Achenbach, T.M., McConaughy, S.H., Ivanova, M.Y., & Rescorla, L.A. 2011. *Manual for the ASEBA Brief Problem Monitor (BPM)*. Burlington, VT: ASEBA, 1-33.
- Andreoni, J., Di Girolamo, A., List, J., Mackevicius, C., and Samek, A. 2020. "Risk Preferences of Children and Adolescents in Relation to Gender, Cognitive Skills, Soft Skills, and Executive Functions." *Journal of Economic Behavior and Organization* 179: 729-472.
- Angold, A., Costello, E. J., Messer, S. C., Pickles, A., Winder, F., and Silver, D. 1995. "The development of a short questionnaire for use in epidemiological studies of depression in children and adolescents." *International Journal of Methods in Psychiatric Research* 5, 237 – 249.
- Arain, M., Haque, M., Johal, L., Mathur, P., Nel, W., Rais, A., ...Sharma, S. 2013. "Maturation of the adolescent brain." *Neuropsychiatric Disease and Treatment*, 9: 449–461.
- Arim, R.G., Shapka, J.D., Dahinten, V.S., and Willms, J.D. 2007. "Patterns and correlates of pubertal development in Canadian youth: effects of family context." *Canadian Journal of Public Health* 98(2): 91–6.
- Avellar, S., Paulsell, D., Sama-Miller, E., Del Grosso, P., Akers, L., and Kleinman, R. 2016. *Home Visiting Evidence of Effectiveness Review: Executive Summary*. Office of Planning, Research and Evaluation. Administration for Children and Families, U.S. Department of Health and Human Services. Washington, DC.
- Bailey, D., Duncan, G.J., Odgers, C.L., and Yu, W. 2017. "Persistence and Fadeout in the Impacts of Child and Adolescent Interventions." *Journal of Research on Educational Effectiveness* 10(1): 7–39.
- Bandura, A. 1977. "Self-efficacy: Toward a Unifying Theory of Behavioral Change." *Psychology Review* 84: 191-215.
- Benjamin, D.J., Brown, S.A., and Shapiro, J.M. 2013. "Who is "Behavioral"? Cognitive ability and anomalous preferences." *Journal of the European Economic Association* 11(6): 1231-1255.
- Bierman, K.B., Heinrichs, B.S., Welsh, J.A., Nix, R.L., and Gest, S.G. 2017. "Enriching Preschool Classrooms and Home Visits with Evidence-based Programming: Sustained Benefits for Low-income Children." *Journal of Child Psychology and Psychiatry* 58:129–137.
- Blair, C., and Raver, C.C. 2012. "Child Development in the Context of Adversity: Experiential Canalization of Brain and Behavior." *American Psychologist* 67(4), 309–318.
- Bowlby, J. 1969. *Attachment and Loss, Vol. 1: Attachment*. New York: Basic Books.
- Bradley, R.H., Caldwell, B.M., Rock, S.L., Ramey, C.T., Barnard, K.E., Gray, C. ... Johnson, D.L. 1989. "Home Environment and Cognitive Development in the First 3 Years of Life: A collaborative Study Involving Six Sites and Three Ethnic Groups in North America." *Developmental Psychology* 25: 217–235.
- Bradley, R.H., and Corwyn, R.F. 2002. "Socioeconomic Status and Child Development." *Annual Review of Psychology* 53(1): 371-399.
- Bronfenbrenner, U. 1979. *The Ecology of Human Development: Experiments by Nature and design*. Cambridge, MA: Harvard University Press.
- Browning, L., Hsieh, S., and Ashwell, M. 2010. "A systematic review of waist-to-height ratio as a screening tool for the prediction of cardiovascular disease and diabetes: 0.5 could be a suitable global boundary value." *Nutrition Research Review* 23: 247–69.

- Butler, S. M., Parry, R., and Fearon, R. M. P.** 2015. *Antisocial Beliefs and Attitudes Scale—Revised (ABAS)*. *Psychological Assessment* 27, 291–301.
- Campbell, F., Conti, G., Heckman, J.J., Hyeok Moon, S., Pinto, R., Pungello, E., and Pan, Y.** 2014. “Early Childhood Investments Substantially Boost Adult Health”. *Science* 343, 1478-1485.
- Campbell, F., Pungello, E., Miller-Johnson, S., Burchinal, M., and Ramey, C.T.** 2001. “The development of cognitive and academic abilities: Growth curves from an early childhood educational experiment.” *Developmental Psychology* 37(2):231–242.
- Castillo, J., Jordan, L., and Petrie, R.** 2018. “Children’s rationality, risk attitudes and field behavior”. *European Economic Review* 102: 62-81.
- Castillo, M., Ferraro, P.J., Jordan, J.L., and Petrie, R.** 2011. “The today and tomorrow of kids: time preferences and educational outcomes of children.” *Journal of Public Economics* 95(11): 1377-1385.
- Castillo, M., Jeffrey, J.L., and Petrie, R.** 2019. “Discount Rates of Children and High School Graduation.” *The Economic Journal* 129(619): 1153–1181.
- Chazan-Cohen, R., Raikes, H.H., and Vogel, C.** 2013. “Program Subgroups: Patterns of Impacts for Home-based, Center-based, and mixed-approach programs.” *Monographs of the Society for Research in Child Development* 78(1), 93-109.
- Coller, M., and Williams, M.B.** 1999. “Eliciting individual discount rates.” *Experimental Economics* 2(2), 107-127.
- Côté, S., Orri, M., Tremblay, R.E., and Doyle, O.** 2018. “A Multi-Component Early Intervention Programme and Trajectories of Behavior, Cognition, and Health.” *Pediatrics* 141(5): e20173174.
- Coy, D., and Doyle, O.** 2024. “Should Early Health Investments Work: Evidence from an RCT of a Complex Early Childhood Intervention.” *Journal of Human Capital* 18(1): 74-104.
- Cronin, F.M., Hurley, S.M., Buckley, T. et al.** 2022. “Mediators of socioeconomic differences in overweight and obesity among youth in Ireland and the UK (2011–2021): a systematic review.” *BMC Public Health* 22: 1585.
- Cunha, F., Elo, I., and Culhane, J.** 2013. “Eliciting Maternal Expectations about the Technology of Cognitive Skill Formation.” *NBER Working Paper* No. 19144.
- Cunha, F., and Heckman, J.J.** 2007. “The Technology of Skill Formation.” *American Economic Review* 97 (2): 31-47.
- Cunha, F., Heckman, J.J., and Schennach, S.M.** 2010. “Estimating the Technology of Cognitive and Noncognitive Skill Formation.” *Econometrica* 78: 883–931.
- Day, F.R., Elks, C.E., Murray, A., Ong, K.K., and Perry, J.R.** 2015. “Puberty timing associated with diabetes, cardiovascular disease and also diverse health outcomes in men and women: the UK Biobank study.” *Scientific Reports* 5:11208.
- Deardorff, J., Abrams, B., Ekwaru, J.P., and Rehkopf, D.H.** 2014. “Socioeconomic status and age at menarche: an examination of multiple indicators in an ethnically diverse cohort.” *Annals of Epidemiology* 24(10): 727–33.
- Diamond, A., and Lee, K.** 2011. “Interventions Shown to Aid Executive Function Development in Children 4 to 12 Years Old”. *Science* 333(6045): 959-964.
- Dooley, M., and Stewart, J.** 2007. “Family Income, Parenting Styles and Child Behavioural–Emotional Outcomes.” *Health Economics* 16 (2): 145-162.
- Doyle, O.** 2013. “Breaking the Cycle of Deprivation: An Experimental Evaluation of an Early Childhood Intervention.” *Journal of the Statistical and Social Inquiry Society of Ireland* Vol. XLI, 92-111.
- Doyle, O.** 2020. “The First 2000 Days and Children’s Skills.” *Journal of Political Economy* 128(6) 2067–2122.

- Doyle, O.** 2024. “Can Early Intervention have a Sustained Effect on Human Capital?” *Journal of Human Resources* 59(5): 1599-1636.
- Doyle, O., Delaney, L., O’Farrelly, C., Fitzpatrick, N., and Daly, M.** 2017b. “Can Early Intervention Policies Improve Well-being? Evidence from a randomized controlled trial.” *PLoS ONE* 12 (1): e0169829.
- Doyle, O., Fitzpatrick, N., Rawdon, C., and Lovett, J.** 2015. “Early Intervention and Child Health: Evidence from a Dublin-based Trial.” *Economics and Human Biology* 19: 224-245.
- Doyle, O., Harmon, C., Heckman, J.J., and Tremblay, R.** 2009. “Investing in Early Human Development: Timing and Economic Efficiency.” *Economics and Human Biology* 7 (1): 1-6.
- Doyle, O., Harmon, C., Heckman, J.J., Logue, C., and Hyeok Moon, S.** 2017a “Measuring Investment in Human Capital Formation: An Experimental Analysis of Early Life Outcomes.” *Labour Economics* 45: 40-58.
- Doyle, O., McGlanaghy, E., Palamaro Munsell, E., and McAuliffe, F.M.** 2014. “Home Based Educational Intervention to Improve Perinatal Outcomes for a Disadvantaged Community: A Randomised Control Trial.” *European Journal of Obstetrics and Gynaecology* 180: 162-167.
- Doyle, O., and PFL Evaluation Team.** 2010. *Assessing the Impact of Preparing for Life: Baseline Report*. Report to Preparing for Life Programme. Atlantic Philanthropies & Department of Children and Youth Affairs.
- Doyle, O., and PFL Evaluation Team.** 2015. *Assessing the Impact of Preparing for Life at 48 Months*. Report to Preparing for Life Programme. Atlantic Philanthropies & Department of Children and Youth Affairs.
- Doyle, O., and PFL Evaluation Team.** 2016. *Final Report: Did Preparing for Life Improve Children’s School Readiness*. Report to Preparing for Life Programme. Atlantic Philanthropies & Department of Children and Youth Affairs.
- DuMont, K., Kirkland, K., Mitchell-Herzfeld, S., Ehrhard-Dietzel, S., Rodriguez, M. L., Lee, E., ... and Greene, R.** 2010. “A Randomized Trial of Healthy Families New York (HFNY): Does Home Visiting Prevent Child Maltreatment”. *Rensselaer, NY: New York State Office of Children & Family Services and Albany, NY: The University of Albany, State University of New York*.
- Eckenrode, J., Campa, M., Luckey, W.L., et al.** 2010. “Long Term Effects of Prenatal and Infancy Nurse Home Visitation on the Life Course of Youths: 19-Year Follow-Up of a Randomized-Controlled Trial.” *JAMA Pediatrics* 164:9–15.
- Elliott, C.D., Smith, P., and McCulloch, K.** 2011. *British Ability Scales III*. Windsor, UK: NFER-Nelson.
- Enoch, M.A., Kitzman, H., Smith, J.A., Anson, E., Hodgkinson, C.A., Goldman, D., and Olds, D. L.** 2016. “A prospective cohort study of influences on externalizing behaviors across childhood: Results from a nurse home visiting randomized controlled trial.” *Journal of the American Academy of Child and Adolescent Psychiatry* 55(5), 376–382.
- Eslami, M., Pourghazi, F., Khazdouz, M., Tian, J., Pourrostami, K., Esmaeili-Abdar, Z., Ejtahed, H.S., and Qorbani, M.** 2023. “Optimal cut-off value of waist circumference-to-height ratio to predict central obesity in children and adolescents: A systematic review and meta-analysis of diagnostic studies.” *Frontiers in Nutrition* 4(9): 985319.
- Evans, G.W., and Rosenbaum, J.** 2008. “Self-regulation and the Income-Achievement Gap.” *Early Childhood Research Quarterly* 23(4), 504-514.

- Farah, M.J., Betancourt, L., Shera, D.M., Savage, J.H., Giannetta, J.M., Brodsky, N.L., ... and Hurt, H.** 2008. "Environmental Stimulation, Parental Nurturance and Cognitive Development in Humans." *Developmental Science* 11 (5): 793-801.
- Filene, J.H., Kaminski, J.W., Valle, L.A., and Cachat, P.** 2013. "Components Associated with Home Visiting Program Outcomes: A Meta-analysis." *Pediatrics* 132(Supplement): S100–S109.
- Fiorini, M., and Keane, M.** 2014. "How the Allocation of Children's Time Affects Cognitive and Non-cognitive Development." *Journal of Labor Economics* 32 (4): 787-836.
- Gertler, P., Heckman, J.J., Pinto, R., Zanolini, A., Vermeersch, C., Walker, S., Chang, S.M., and Grantham-McGregor, S.** 2014. "Labor Market Returns to an Early Childhood Stimulation Intervention in Jamaica." *Science* 344 (6187): 998-1001.
- Gomby, D.S.** 2005. *Home Visitation in 2005: Outcomes for Children and Parents* (Vol. 7). Invest in Kids Working Paper No. 7. Committee for Economic Development: Invest in Kids Working Group.
- Gomby, D.S., Culross, P.L., and. Behrman, R.E.** 1999. "Home Visiting: Recent Program Evaluations: Analysis and Recommendations." *The Future of Children* 9 (1): 4.
- Good, P.** 2005. *Permutation, Parametric and Bootstrap Tests of Hypotheses* (3rd ed.), New York: Springer.
- Goodman, R.** 1997. "The Strengths and Difficulties Questionnaire: A Research Note." *Journal of Child Psychology and Psychiatry* 38 (5): 581-586.
- Grantham-McGregor, S. and Smith, J.A.** 2016. "Extending The Jamaican Early Childhood Development Intervention." *Journal of Applied Research on Children: Informing Policy for Children at Risk* 7 (2), Article 4.
- Gullone, E., and Robinson, K.** 2005. "The Inventory of Parent and Peer Attachment--Revised (IPPA-R) for Children: A Psychometric Investigation." *Clinical Psychology & Psychotherapy* 12(1), 67–79.
- Heckman, J.J., Holland, M.L., Makino, K.K., Pinto, R., and Rosales-Rueda, M.** 2017. "An Analysis of the Memphis Nurse-Family Partnership Program." NBER Working Paper No. 23610.
- Heckman, J.J., Moon, S., Pinto, R., Savelyev, P.A., and Yavitz, A.** 2010. "Analyzing Social Experiments as Implemented: A Reexamination of the Evidence from the HighScope Perry Preschool Program." *Quantitative Economics* 1 (2): 1-46.
- Heckman, J.J., and Mosso, S.** 2014. "The Economics of Human Development and Social Mobility." *Annual Review of Economics* 6 (1): 689–733.
- Heckman J.J., Pinto, R., and Savelyev, P.A.** 2013. "Understanding the Mechanisms Through Which an Influential Early Childhood Program Boosted Adult Outcomes." *American Economic Review* 103: 2052–86.
- Holt, C.A., and Laury, S.K.** 2002. "Risk Aversion and Incentive Effects." *The American Economic Review* 92(5): 1644–55.
- James-Todd, T., Tehranifar, P., Rich-Edwards, J., Titievsky, L., and Terry, M.B.** 2010. "The impact of socioeconomic status across early life on age at menarche among a racially diverse population of girls." *Annals of Epidemiology* 20(11): 836–42.
- Kirkbride, J.B., Anglin, D.M., Colman, I., Dykxhoorn, J., Jones, P.B., Patalay, P., Pitman, A., Soneson, E., Steare, T., Wright, T., and Griffiths, S.L.** 2024. "The social determinants of mental health and disorder: evidence, prevention and recommendations." *World Psychiatry* 23(1): 58-90.
- Kirkland, K., and Mitchell-Herzfeld, S.** 2012. *Evaluating the effectiveness of home visiting services in promoting children's adjustment in school: Final report to the Pew Center on the States*. Rensselaer, NY: New York State Office of Children and Family Services, Bureau of Evaluation and Research.

- Kitzman, H. J., Olds, D.L., Cole, R.E., Hanks, C.A., Anson, E.A., Arcoleo, K.J., ... and Holmberg, J.R.** 2010. "Enduring Effects of Prenatal and Infancy Home Visiting by Nurses on Children: Follow-up of a Randomized Trial among Children at Age 12 years." *Archives of Pediatrics & Adolescent Medicine* 164(5): 412-418.
- Knudsen, E.I., Heckman, J.J., Cameron, J.L., and Shonkoff, J.P.** 2006. "Economic, Neurobiological, and Behavioral Perspectives on Building America's Future Workforce." *Proceedings of the National Academy of Science USA* 103 (27): 10155–10162.
- Layzer, J.I., Goodson, B.D., Bernstein L, and Price, C.** 2001. *National Evaluation of Family Support Programs. Volume A: The Meta-analysis. Final Report.* Cambridge, Mass: Abt Associates Inc.
- Liew, J.** 2012. "Effortful Control, Executive Functions, and Education: Bringing Self-Regulatory and Social-Emotional Competencies to the Table." *Child Development Perspectives* 6(2): 105-111.
- Ludbrook, J., and Dudley, H.** 1998. "Why Permutation Tests are Superior to t and F Tests in Biomedical Research." *American Statistician* 52 (2): 127-132.
- Manotas, M.C., Mauricio González, D., Céspedes, C., Forero, C., Rojas Moreno, A.P.** 2022. "Genetic and Epigenetic Control of Puberty." *Sexual Development* 16 (1): 1–10.
- Mensah, F.K., Bayer, J.K., Wake, M., Carlin, J.B., Allen, N.B., and Patton, G.C.** 2013. "Early puberty and childhood social and behavioral adjustment." *Journal of Adolescence Health* 53(1): 118–124.
- Mewhort, D.J.K.** 2005. "A Comparison of the Randomization Test with the F test when Error is Skewed." *Behavior Research Methods* 37: 426–435.
- Miller, E.B., Farkas, G., Vandell, D.L., and Duncan, G.** 2014. "Do the Effects of Head Start Vary by Parental Preacademic Stimulation?" *Child Development* 85: 1385–1400.
- Miller, S., Maguire, L.K., and Macdonald, G.** 2011. "Home-based Child Development Interventions for Preschool Children from Socially Disadvantaged Families." *Cochrane Database of Systematic Reviews* 12, CD008131.
- Minkovitz, C.S., Strobino, D., Mistry, K.B., Scharfstein, D.O., Grason, H., Hou, W., Ialongo, N., and Guyer, B.** 2007. "Healthy Steps for Young Children: Sustained Results at 5.5 Years." *Pediatrics* 120 (3): e658-668.
- Moore Heslin, A., O'Donnell, A., Kehoe, L., Walton, J., Flynn, A., Kearney, J., and McNulty, B.** 2020. "Adolescent overweight and obesity in Ireland-Trends and sociodemographic associations between 1990 and 2020." *Pediatrics Obesity* 18(2) :e12988.
- Nisbett, R.E., Aronson, J., Blair, C., Dickens, W., Flynn, J., Halpern, D.F., and Turkheimer, E.** 2012. "Intelligence: New Findings and Theoretical Developments." *American Psychologist* 67(2): 130–159.
- Nixon, E.** 2023. *Growing-up in Ireland: Socio-emotional and behavioural outcomes in early adolescence.* Department of Children, Equality, Disability, Integration and Youth.
- O'Sullivan, A., Fitzpatrick, N., and Doyle, O.** 2017. "Effects of Dietary Recommendations During Early Childhood on Cognitive Functioning: A Randomized Controlled Trial." *Public Health Nutrition* 20 (1): 154-164.
- OECD.** 2016. *Enhancing Child Well-being to Promote Inclusive Growth.* OECD Publishing, Paris.
- Olds, D.L., Kitman, H., Cole, R., Robinson, J., Sidora, K., Luckey, D., Henderson Jr, C.R., et al.** 2004. "Effects of Nurse Home-visiting on Maternal Life Course and Child Development: Age 6 Follow-up Results of a Randomized Trial." *Pediatrics* 114 (6): 1550-1559.

- Olds, D.L., Henderson, C.R., Cole, R., Eckenrode, J., Kitzman, H., Luckey, D., et al.** 1998. "Long-term effects of nurse home visitation on children's criminal and antisocial behavior: 15-year follow-up of a randomized controlled trial." *JAMA: The Journal of the American Medical Association* 280(14), 1238–1244.
- Olds, D.L., Kitzman, H., Hanks, C., Cole, R., Anson, E., Sidora-Arcoleo, K., ... and Stevenson, A. J.** 2007. "Effects of nurse home visiting on maternal and child functioning: age-9 follow-up of a randomized trial." *Pediatrics* 120(4), e832-e845.
- Olds, D.L., Henderson, C.R., Kitzman, H.J., Eckenrode, J.J., Cole, R.E., and Tatelbaum, R.C.** 1999. "Prenatal and Infancy Home Visitation by Nurses: Recent Findings." *The Future of Our Children* 9 (1): 44–65.
- Peacock, S., Konrad, S., Watson, E., Nickel, D., and Muhajarine, N.** 2013. "Effectiveness of Home Visiting Programs on Child Outcomes: A Systematic Review." *BMC Public Health* 13 (1): 17.
- Petersen, A.C., Crockett, L., Richards, M., and Boxer, A.** 1998. "A self-report measure of pubertal status: Reliability, validity, and initial norms." *Journal of Youth Adolescence* 7(2):117-33.
- Rayce, S.B., Rasmussen, I.S., Klest, S.K., Patras, J., and Pontoppidan, M.** 2017. "Effects of Parenting Interventions for At-risk Parents with Infants: A Systematic Review and Meta-analyses." *BMJ Open* e015707.
- Robins, J.M., Rotnitzky, A., and Zhao, L.P.** 1994. "Estimation of Regression Coefficients when Some Regressors are Not Always Observed." *Journal of the American Statistical Association* 89 (427): 846-866.
- Roggman L.A., Cook, G.A., Peterson, C.A., and Raikes, H.H.** 2008. "Who Drops Out of Early Head Start Home Visiting Programs?" *Early Education And Development* 19 (4).
- Romano, J.P., and Wolf, M.** 2005. "Exact and Approximate Stepdown Methods for Multiple Hypothesis Testing." *Journal of the American Statistical Association* 100 (469): 94-108.
- Rosenberg, M.** 1965. *Society and the adolescent self-image*. Princeton, NJ: Princeton University Press.
- Rowe, E.W., Kim, S., Baker, J.A., Kamphaus, R.W., and Horne, A.M.** 2010. "Student Personal Perception of Classroom Climate: Exploratory and Confirmatory Factor Analyses." *Educational and Psychological Measurement* 70(5), 858-879.
- Sanders, M.R.** 2012. "Development, Evaluation, and Multinational Dissemination of the Triple P-Positive Parenting Program." *Annual Review of Clinical Psychology* 8(1): 345–379.
- Sanders, M.R., Markie-Dadds, C., and Turner, K.** 2003. "Theoretical, Scientific and Clinical Foundations of the Triple P-Positive Parenting Program: A Population Approach to the Promotion of Parenting Competence." *Parenting Research and Practice Monograph* 1: 1-21.
- Sanders, M.R., Kirby, J.N., Tellegen, C.L., and Day, J.J.** 2014. "The Triple P-Positive Parenting Program: A Systematic Review and Meta-analysis of a Multi-level System of Parenting Support." *Clinical Psychology Review* 34:337–57.
- Sandner, M., and Jungmann, T.** 2017. "Gender-specific Effects of Early Childhood Intervention: Evidence from a Randomized Controlled Trial." *Labour Economics* 45: 59-78.
- Sidora-Arcoleo, K., Anson, E., Lorber, M., Cole, R., Olds, D., and Kitzman, H.** 2010. "Differential Effects of a Nurse Home-visiting Intervention on Physically Aggressive Behavior in Children," *Journal of Pediatric Nursing* 25(1): 35-45.
- Sun, Y., Mensah, F.K., Azzopardi, P., Patton, G.C., and Wake, M.** 2017. "Childhood social disadvantage and pubertal timing: a national birth cohort from Australia." *Pediatrics* 139(6): e20164099.

- Sutter, M., Zoller, C., and Glätzle-Rützler, D.** 2019. "Economic behavior of children and adolescents – A first survey of experimental economics results." *European Economic Review* 111: 98-121.
- Sutter, M., Kocher, M.G., Glätzle-Rützler, D., and Trautmann, S.T.** 2013. "Impatience and uncertainty: experimental decisions predict adolescents' field behavior." *American Economic Review* 103(1): 510-531.
- Sweet, M.A., and Appelbaum, M.I.** 2004. "Is Home Visiting an Effective Strategy? A Meta-Analytic Review of Home Visiting Programs for Families with Young Children." *Child Development* 75: 1435-1456.
- Thompson, R.A., and Nelson, C.A.** 2001. "Developmental Science and The Media: Early Brain Development." *American Psychologist* 56 (1): 5–15.
- Todd, P.E., and Wolpin, K.I.** 2007. "The Production of Cognitive Achievement in Children: Home, School, and Racial Test Score Gaps." *Journal of Human Capital* 1 (1): 91-136.
- Weaver, I.C.G., Cervoni, N., Champagne, F.A., D'Alessio, A.C., Sharma, S., Seckl, J.R., Dymov, S., Szyf, M., and Meaney, M.J.** 2004. "Epigenetic Programming by Maternal Behavior." *Nature Neuroscience* 7: 847–854.
- Wesley, P.W., Buysse, V., and Tyndall, S.** 1997. "Family and Professional Perspectives on Early Intervention: An Exploration Using Focus Groups." *Topics in Early Childhood Special Education* 17 (4): 435-456.
- Zelazo, P.D., and Bauer, P.J.** 2013. "National Institutes of Health Toolbox—Cognitive Function Battery (NIH Toolbox CFB): Validation for Children Between 3 and 15 Years". *Monograph of the Society for Research in Child Development* 78(4).