

# Using matched administrative data to design and improve public policy

David Figlio, Northwestern University & NBER

Keynote address, (Ce)<sup>2</sup> workshop 2014

Warsaw, July 10, 2014

# The issue

---

- ▶ **Good information improves policymaking**
  - ▶ Explicit program evaluations
  - ▶ Basic research that informs policymaking
- ▶ **Multiple ways to obtain this information**
  - ▶ Experimental, quasi-experimental, structural, descriptive analyses
  - ▶ Survey data vs. administrative data



# Survey data/new data collections are often great, but they have real limitations...

---

- ▶ Small  $n$  → frequently not the power to detect modest-sized effects
  - ▶ Policymakers should make decisions based on cost-effectiveness, not effect sizes!
  - ▶ Many important questions require looking at small subgroups, either because they are the population of interest or because we expect/suspect heterogeneous effects
- ▶ Purpose-built new data collections will take years to be able to study longer-term effects
- ▶ Self-reports frequently lead to recall bias
- ▶ Surveys/new data collections are very expensive



# The promise of matched administrative data

---

- ▶ Population-level datasets permit analyses heretofore impossible
- ▶ Possibility of retrospective analysis
- ▶ Reduced likelihood of measurement error
- ▶ Orders of magnitude less expensive
  - ▶ For instance, the US Early Childhood Longitudinal Study- Birth Cohort has 13,500 participants followed from birth to first grade, total cost >\$20 million
  - ▶ In contrast, the value of total staff time (and my time) plus equipment, etc. to create a longitudinal file of >1.6 million children from birth through tenth grade and beyond cost <\$500,000
  - ▶ Total operating costs of US Census Data Centers <\$2 million/year
  - ▶ Clearly there are tradeoffs, but the matched data can address the vast majority of questions that can be answered by new data collections, and many more that cannot!



# Examples

---

- ▶ Linking birth records to school records and/or adult labor market records
- ▶ Linking health data with labor market or schooling data
- ▶ Linking program participation data with criminal justice data
- ▶ Linking parents' administrative records with their children's records
- ▶ And so on...



# Challenges

---

- ▶ Politics
- ▶ Privacy
- ▶ Technical challenges
- ▶ ~~Computing capacity~~



# Ways to overcome privacy issues

---

- ▶ Establish data use agreements with qualified users
- ▶ Establish secure data facilities
- ▶ Create de-identified merged files that are not potentially identifiable
- ▶ Suppress or merge small cells in public use data files
- ▶ Add small amounts of noise to public use data files → “analytically valid synthetic data”
- ▶ Create model servers in which users log in to estimate models using actual data but obtain output that is unidentifiable



# A new data resource: Florida “registry” data

---

- ▶ Jeff Roth at the University of Florida and I have built, in conjunction with the Florida Departments of Education and Health, the first, to our knowledge, large-scale dataset that
  - ▶ Links birth records to school records in a highly developed context
  - ▶ Includes annual assessment data, behavioral data, and (as children age) high school completion and postsecondary outcomes – plus early childhood program participation
  - ▶ To date, children born from 1992-2002 matched to school records
  - ▶ >14,000 twin pairs, >1.3 million singletons old enough for test scores
- ▶ Florida is a location with many desirable characteristics for study:
  - ▶ **Large:** Florida’s population of ~17M and ~200K births/year compares to Norway, Denmark, and Sweden combined
  - ▶ **Heterogeneous:** 45% of moms racial/ethnic minorities; 25% of moms foreign born
  - ▶ Politically and socially **representative** of the United States
  - ▶ **Excellent institutional conditions** for matching birth and school data





# Florida matched data

---

- ▶ Only observe school history in Florida if a child
  - ▶ Remains in Florida until school age
  - ▶ Attends a Florida public school
  - ▶ Is successfully matched between birth and school records
- ▶ How good is the match?
  - ▶ Match based on name (with some fuzziness), date of birth, and social security number – checked with other shared variables
  - ▶ American Community Survey: 80.9% of children born in Florida live in Florida at age 5 and attend public school – this is an overstatement
  - ▶ Our match: 80.7% of all births (79.5% for twins)
  - ▶ Therefore, nearly all potentially matchable children are matched
- ▶ In other settings, the potential match can be much higher...



# Attributes of all Florida births and Florida-born kids attending Florida public schools

---

<b>Maternal attribute</b>	<b>Full population of births</b>	<b>Population of kids matched to Florida school records</b>	<b>Population of kids with a third-grade test score</b>
Black	22.6	24.8	25.7
Hispanic	23.0	23.3	23.9
Foreign-born	23.5	22.9	23.2
Married at time of birth	64.8	62.2	60.9
High school dropout	20.9	22.5	23.3
College graduate	20.5	17.5	16.2
Age 21 or below	22.0	23.6	24.2
Age 36 or above	9.8	9.3	9.2



# Example #1: Explicit policy evaluation

---

- ▶ The question: **Does early intervention for autism spectrum disorders (ASD) improve autistic children's life chances?** (with Currie, Goodman, Persico, Roth)
- ▶ Social and language impairments
- ▶ Restricted, stereotyped, repetitive patterns of behavior
- ▶ Symptoms generally apparent by the age of 3 years
- ▶ Affected individuals often require constant care from family members and professionals
- ▶ In school, affected children must frequently be educated in special settings and even mainstreamed children perform very poorly
- ▶ Estimated per-capita lifetime societal cost of \$3.2 million
- ▶ Affects ~1.5 percent of boys and ~0.4 percent of girls → benefit of using large-scale data to study!



# Early intervention for autism

---

- ▶ Autism awareness organizations and many governments argue that early intervention can dramatically improve the lives of individuals with autism spectrum disorders
- ▶ Argument: Early intervention when neuroplasticity is highest can allow for adapted patterns of interaction between child and environment, more typical development of neural circuitry
- ▶ Billions of dollars spent annually in USA alone on early diagnosis and intervention programs



# The available evidence

---

- ▶ There have been many studies of the effects of early intervention A vs. early intervention B
  - ▶ BUT all of these studies are for kids already diagnosed!
  - ▶ And the largest of these studies had only 44 children in the “treatment” group
- ▶ The first order policy question is whether early identification and intervention works per se
  - ▶ Randomized trials are impossible!
  - ▶ Limited by identification strategy
  - ▶ Limited by data – need large number of children with ASD
- ▶ Administrative data offer a potential opportunity to address this question



# New opportunities with administrative data

---

- ▶ Florida has 8,433 children with ASDs in the matched birth-school records – as well as information on early intervention
- ▶ Possible to control for a wide range of initial conditions at birth (e.g., neonatal health, congenital anomalies, complications of labor and delivery)
- ▶ But can Florida provide an opportunity to study this question in a quasi-experimental manner?



# Florida's Early Steps Program

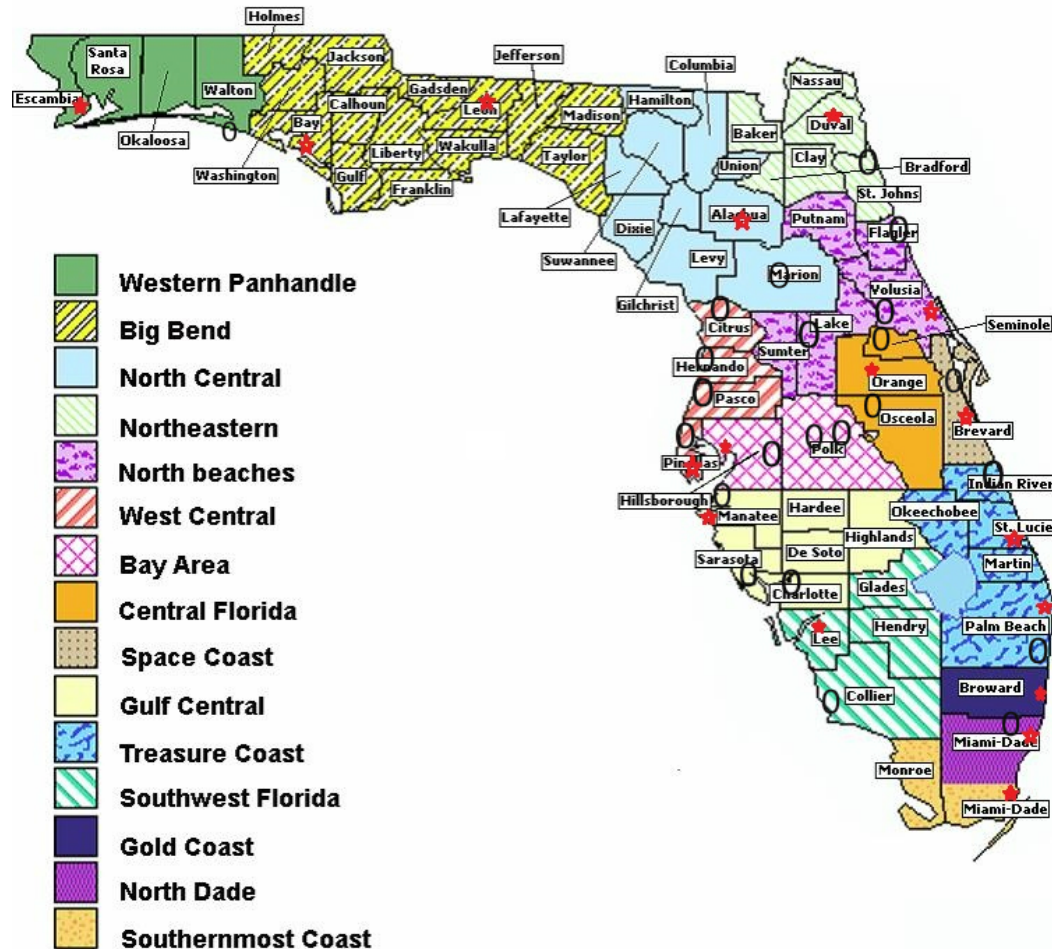
---

- ▶ Florida provides a variety of therapeutic services for children with autism, first through the Early Steps Program (birth to three) and then through the Florida Department of Education
- ▶ Children are screened in one of 16 Early Steps centers located throughout the state
- ▶ Children born as early as 1992-1993 could have been served
- ▶ Once children are determined to be eligible, Early Steps puts together a team of service providers to address individual children's needs; providers are in child's local communities
- ▶ Most services are provided without charge if they are not covered by Medicaid or private insurance



# Locations of Early Steps Centers

## Florida's Early Steps Service Areas

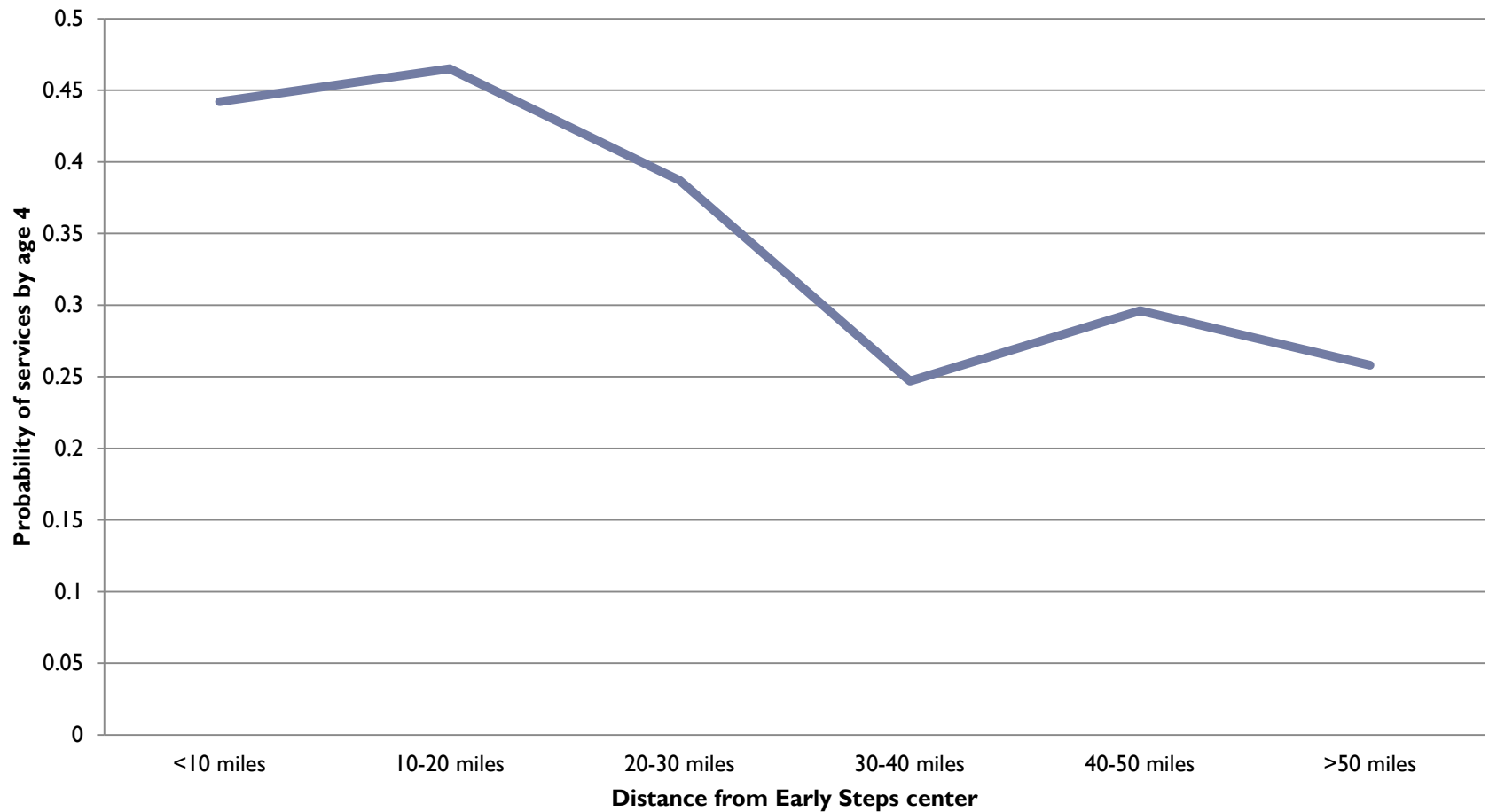




# Children living near Early Steps centers are more likely to receive early intervention

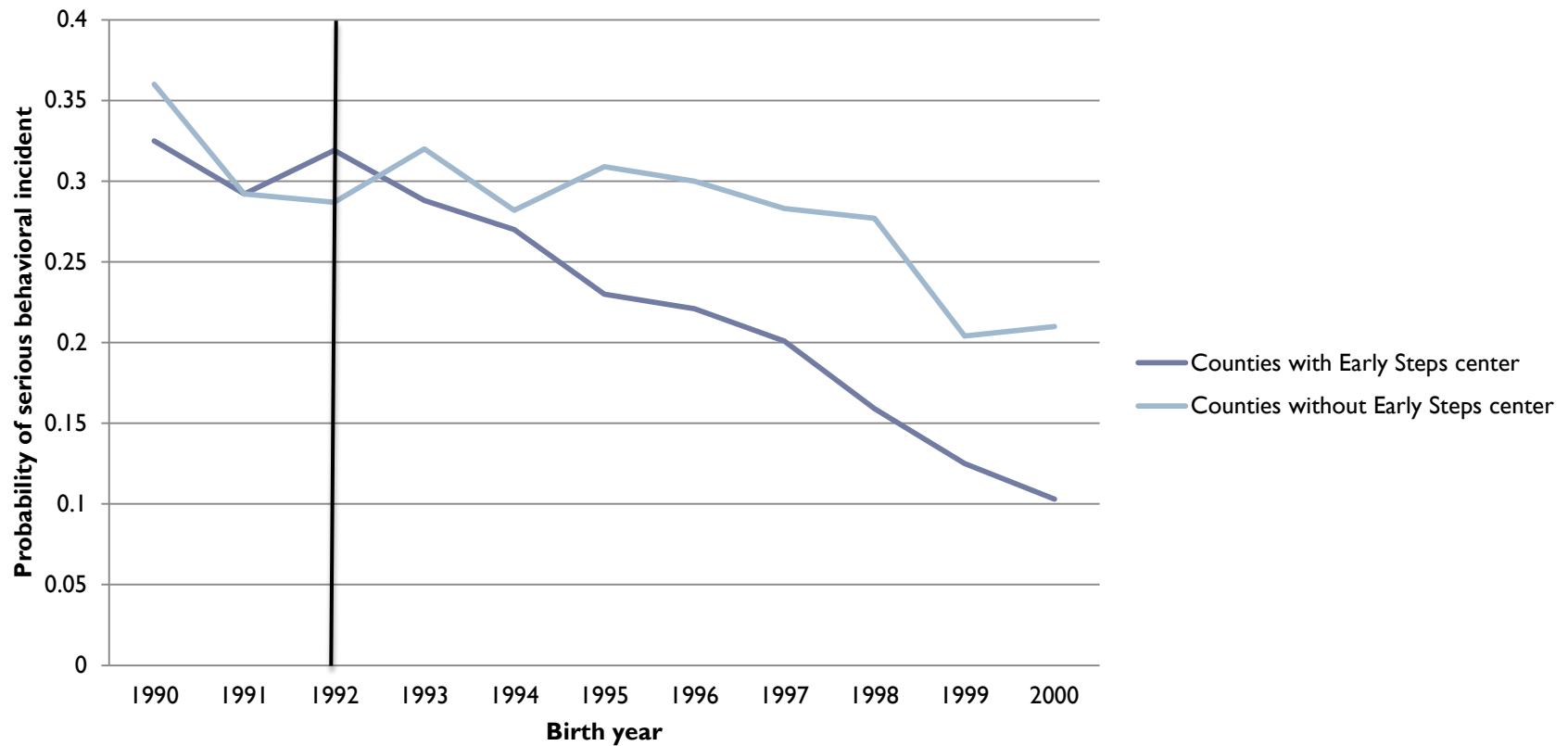
---

**Probability of being served by 4**



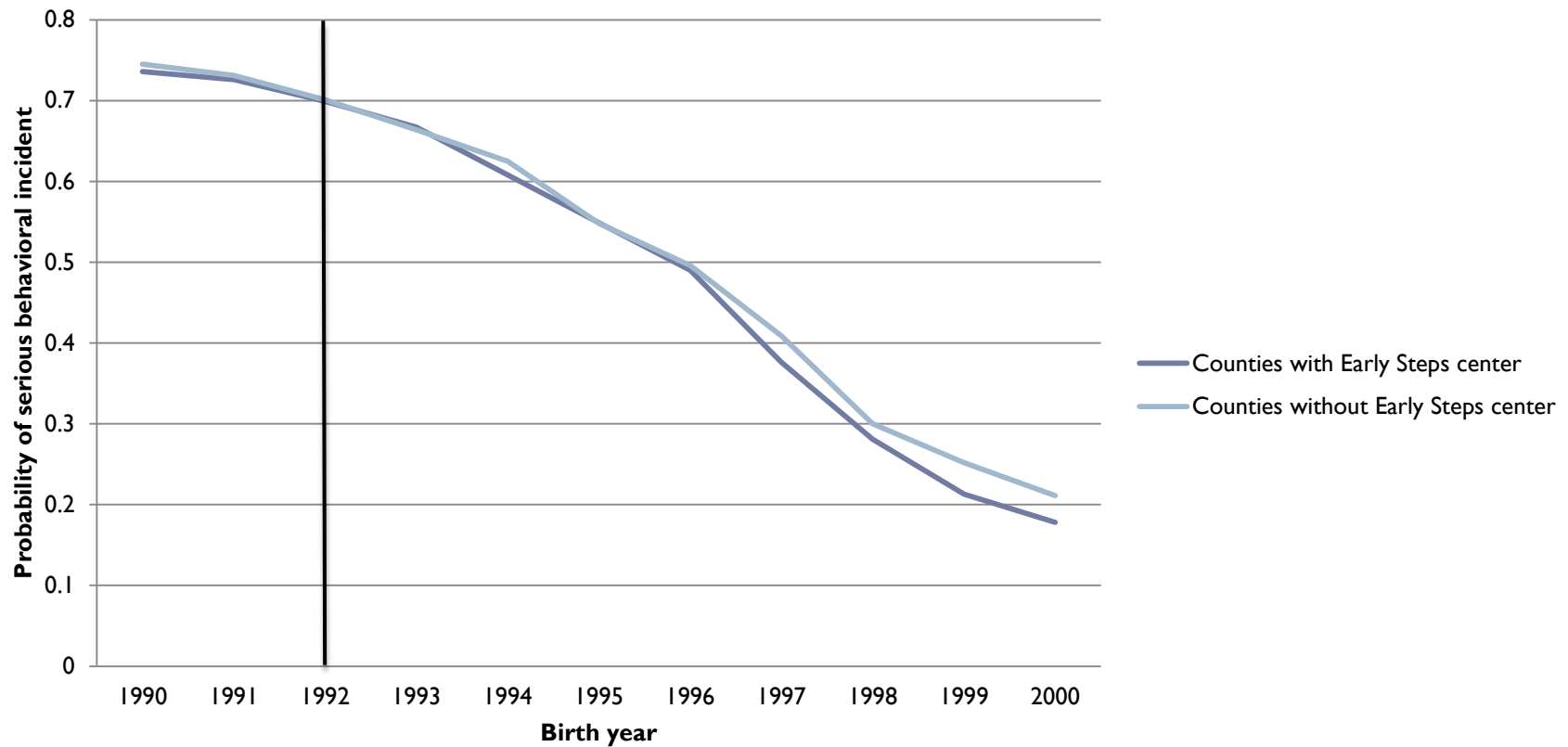
# Cohort-level evidence

**Probability that an autistic child will have a serious behavioral incident, by whether county has an Early Steps center**



# No evidence for other disabilities very rarely served by Early Steps

**Probability that a child with specific learning disabilities will have a serious behavioral incident, by whether county has an Early Steps center**



Note: autistic children are not included in this analysis

## Example #2: Basic research with policy implications

---

- ▶ The question: **What are the effects of neonatal health (specifically, birth weight) on children's cognitive development?** (*American Economic Review*, forthcoming, with Guryan, Karbownik, Roth)
- ▶ Long literature on the effects of neonatal health on many adult outcomes
  - ▶ Wages, disability, adult chronic conditions, human capital accumulation
- ▶ While the existing literature makes clear that there appears to be a permanent effect of poor neonatal health on socio-economic and health outcomes, it is important to know how neonatal health affects child development, whether public policies might be beneficial, and whether parental inputs and neonatal health are complements or substitutes
- ▶ To date, we know little about how neonatal health's effects vary at different stages of development, or whether public policies (e.g., school quality) can help mitigate the relationship



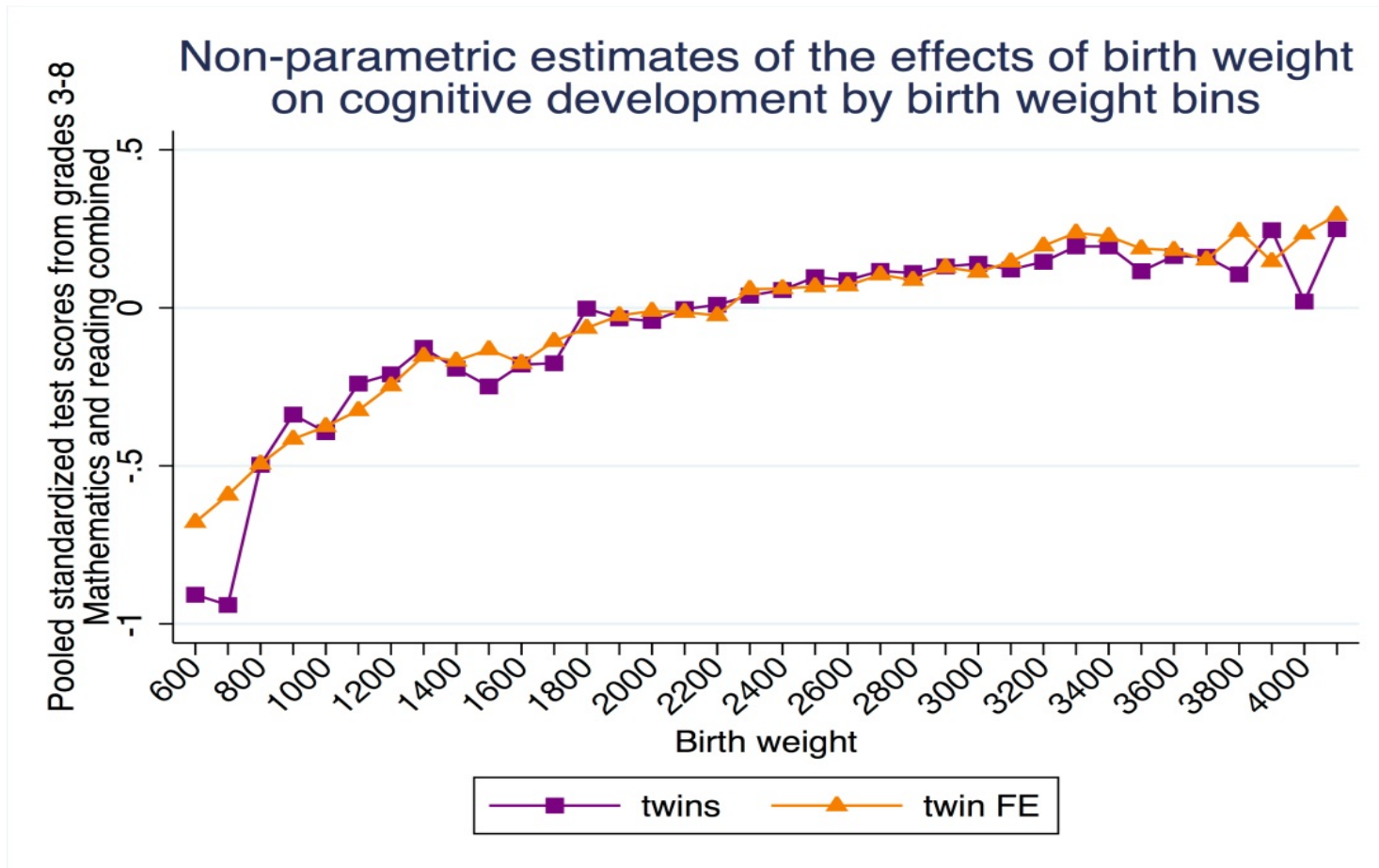
# Opportunities with matched administrative data

---

- ▶ Large scale permits estimating heterogeneous effects, even for twin pairs – allows us to ask whether parental inputs and neonatal health are complements or substitutes
- ▶ Data that include outcomes in childhood allow us to consider more recent births than would otherwise be possible
- ▶ Information about schools attended allow us to observe whether school quality affects the relationship between birth weight and outcomes

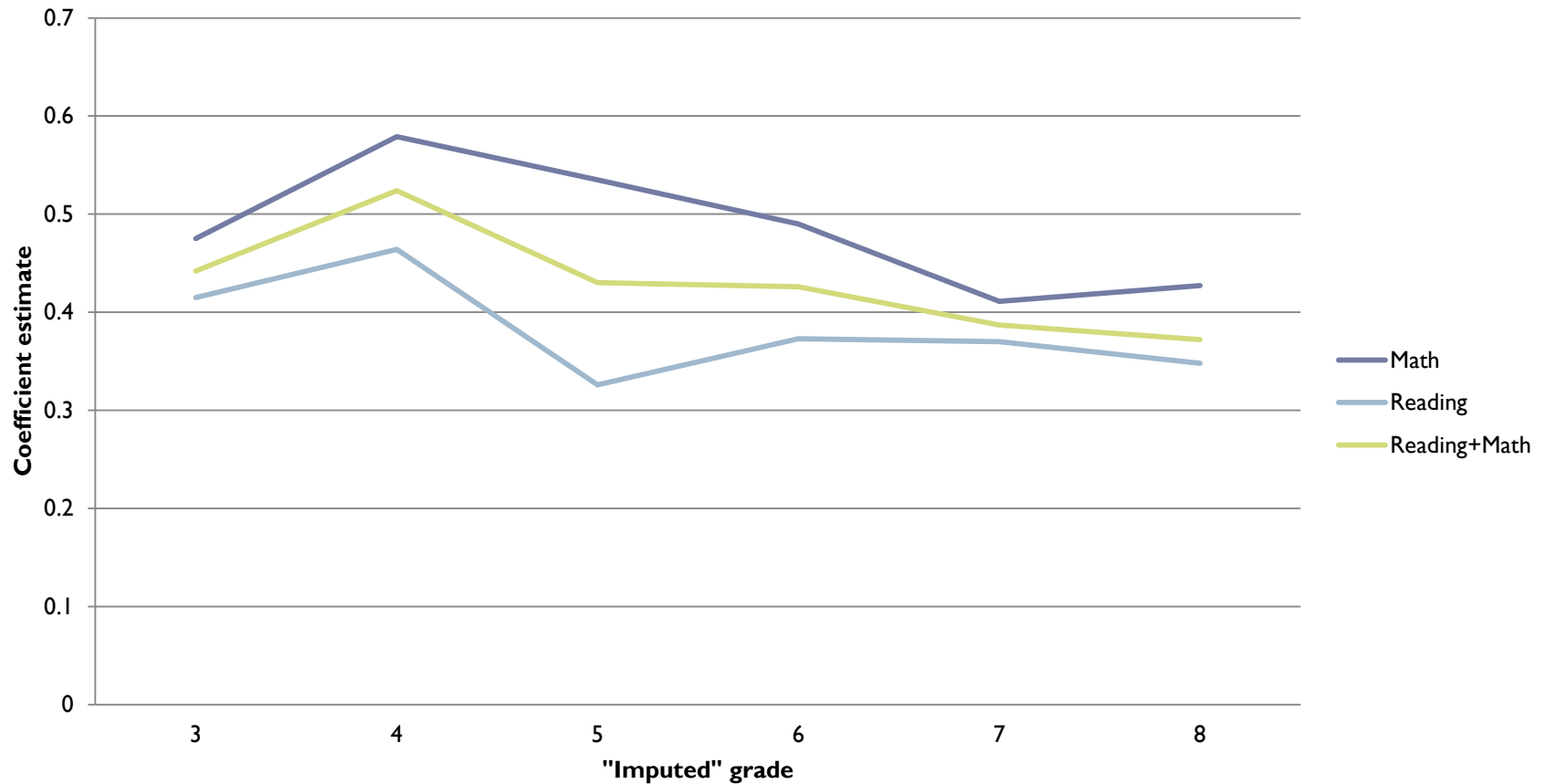


# Non-parametric relationship between birth weight and test scores



# Effects roughly constant over time in twin fixed effect models

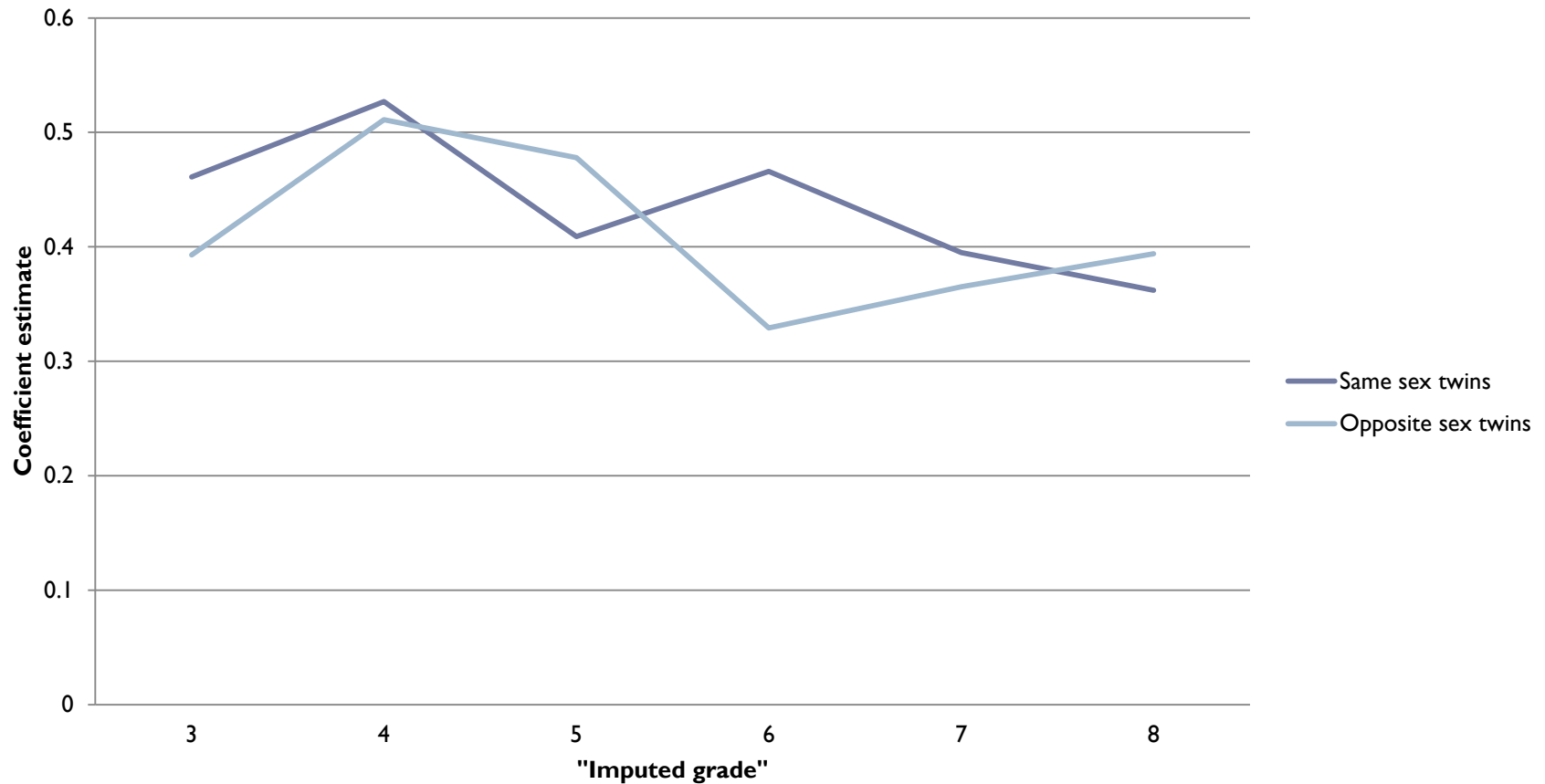
**Estimated effect of log birth weight on test scores, by grade  
(twin FE model)**



# Same-sex versus opposite-sex twin pairs

---

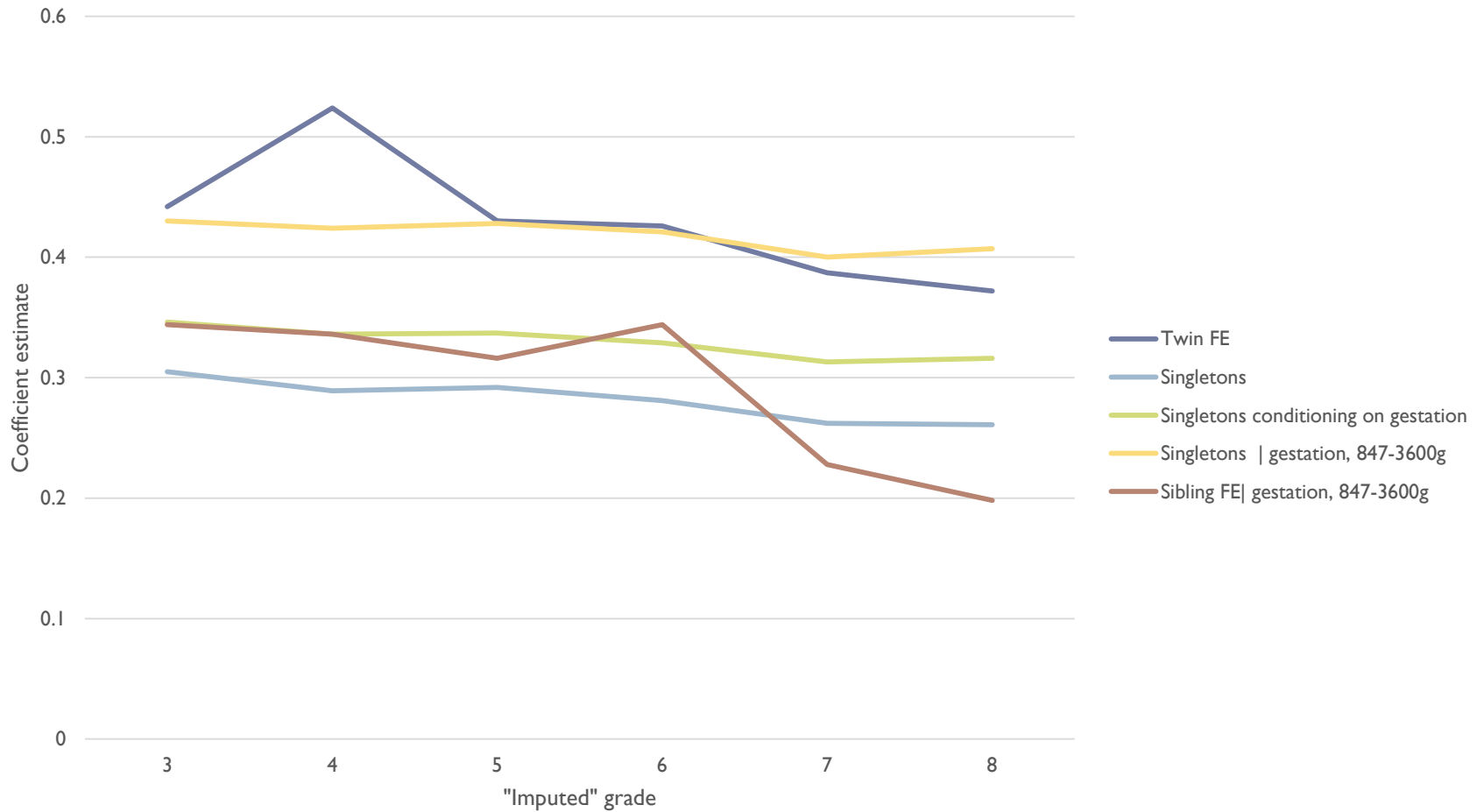
**Estimated effect of log birth weight on test scores (same vs. opposite sex twins (twin FE model))**





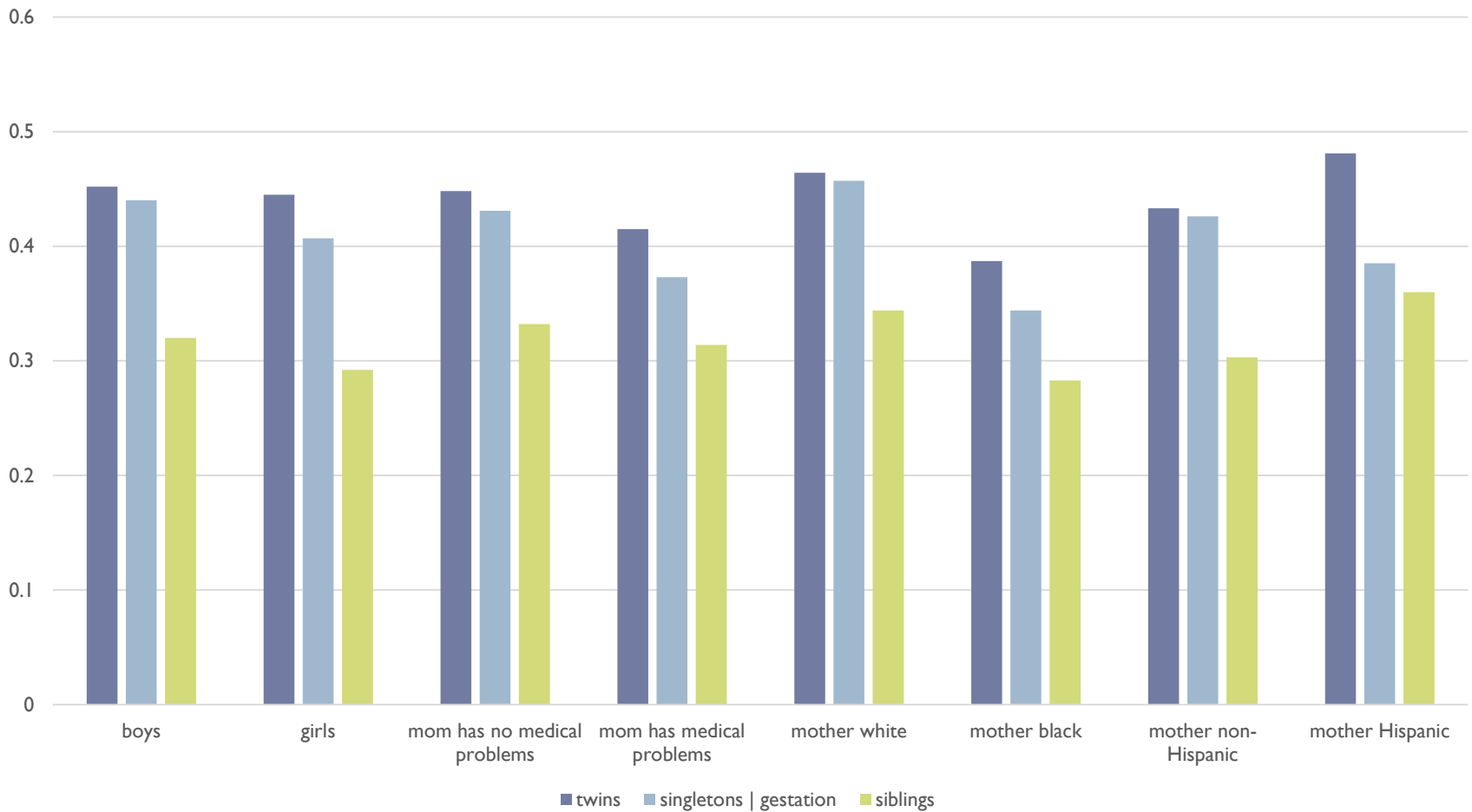
# Twin fixed effects versus singletons

Estimated effect of log birth weight: twin FE, singletons, sibling FEs



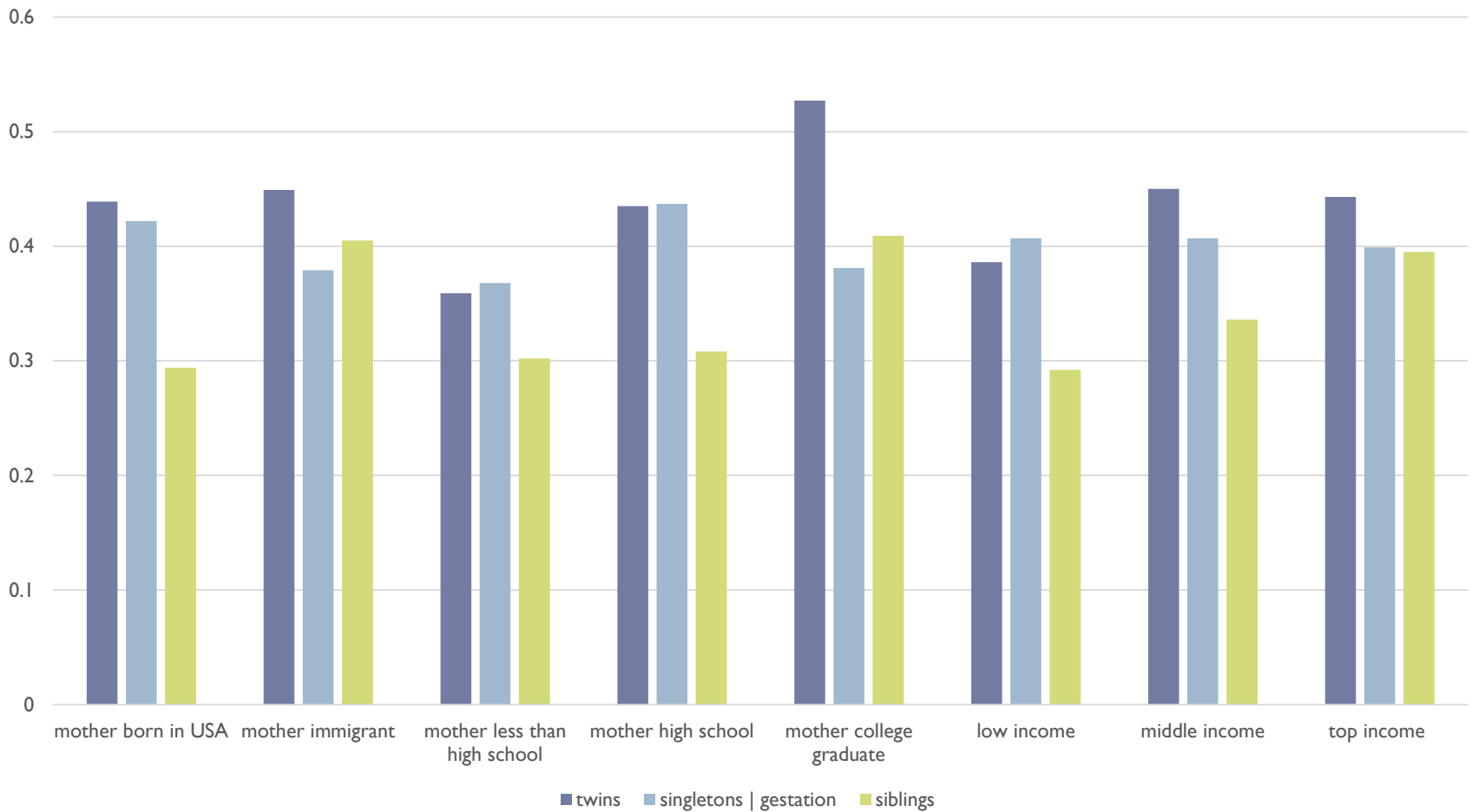
# Differences by groups, part 1

Coefficient on log birth weight, by subgroup, part I



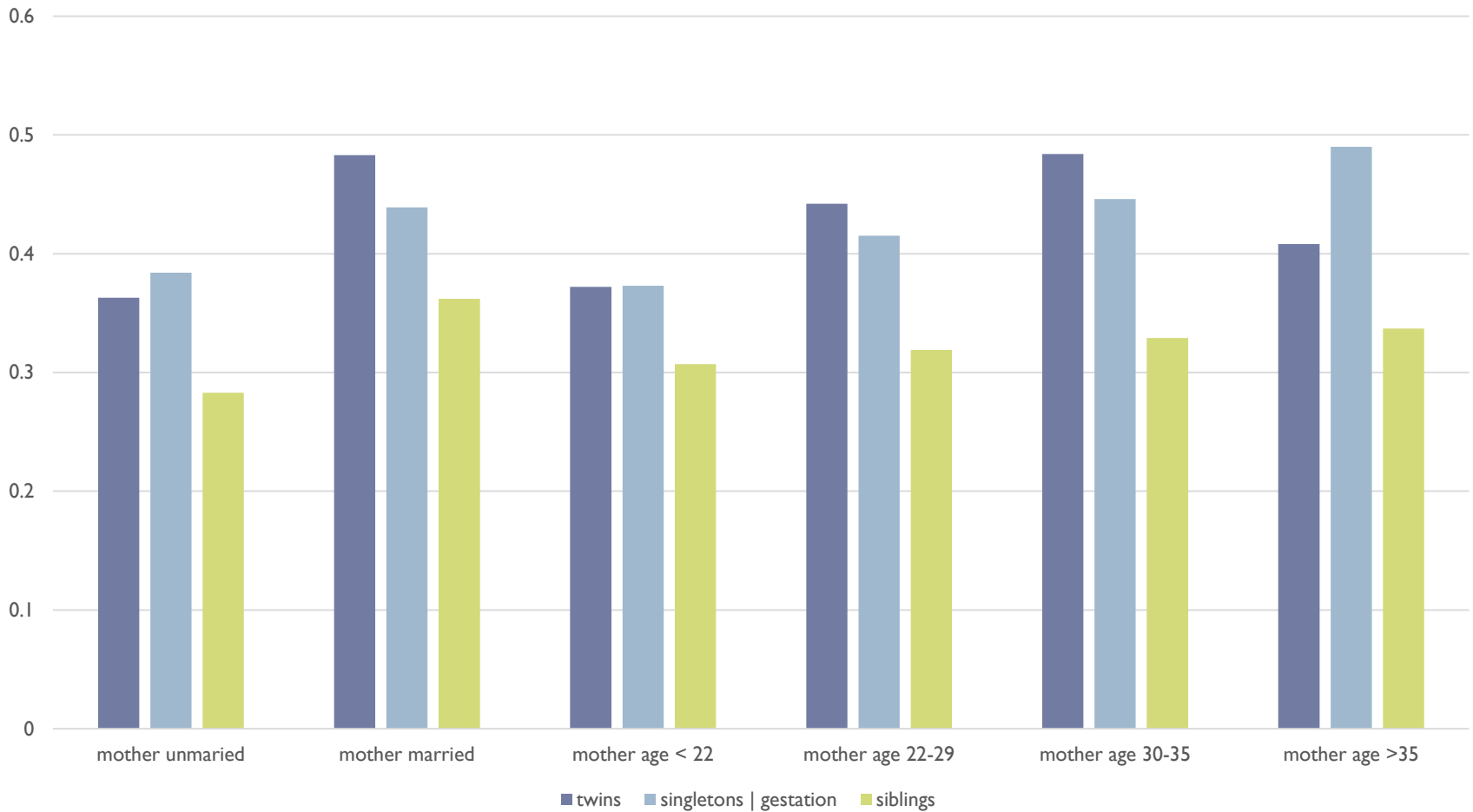
# Differences by groups, part 2

Coefficient on log birth weight, by subgroup, part II

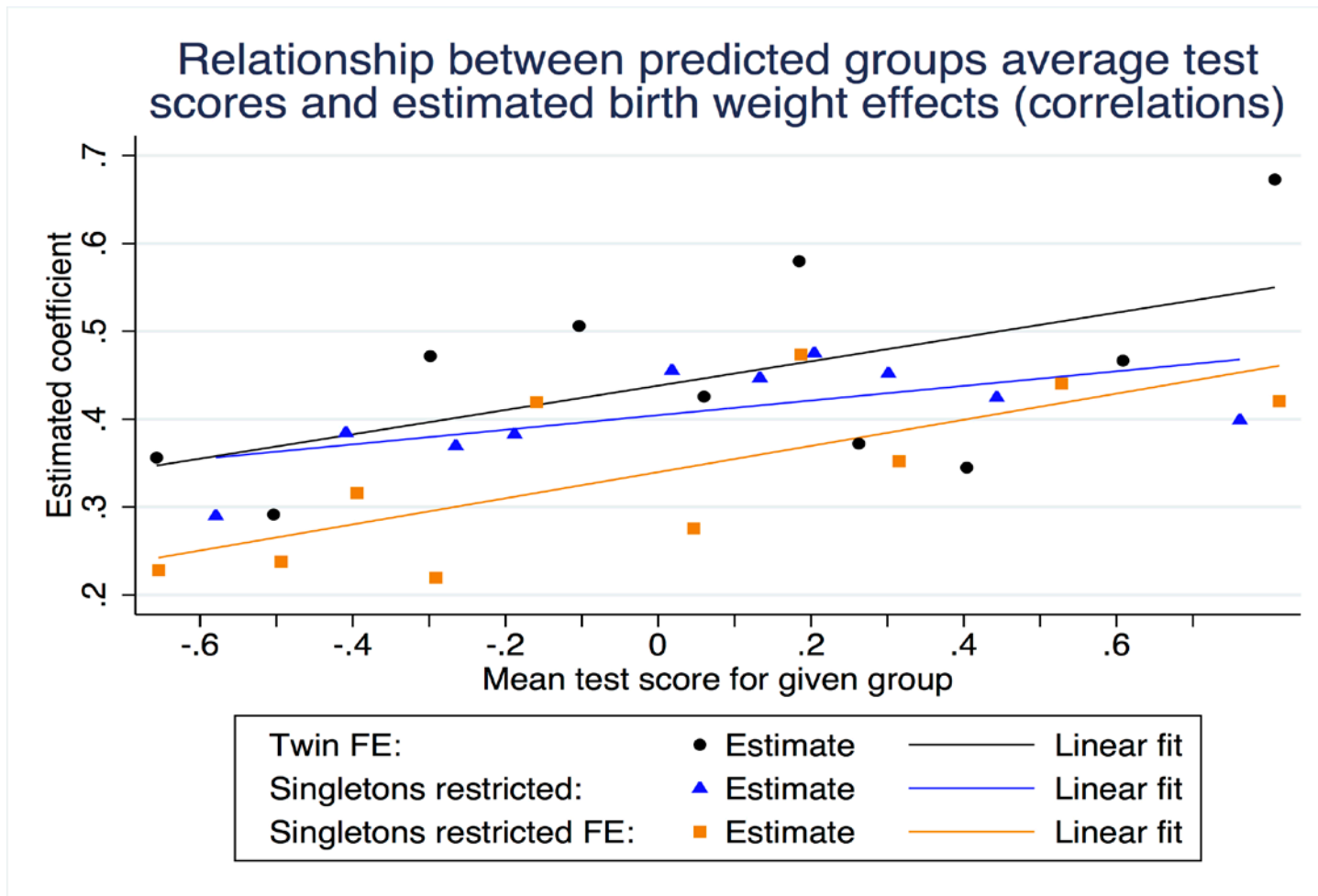


# Differences by groups, part 3

Coefficient on log birth weight, by subgroup, part III



# Test performance and estimated birth weight effects across groups



# Does school quality affect the birth weight gap?

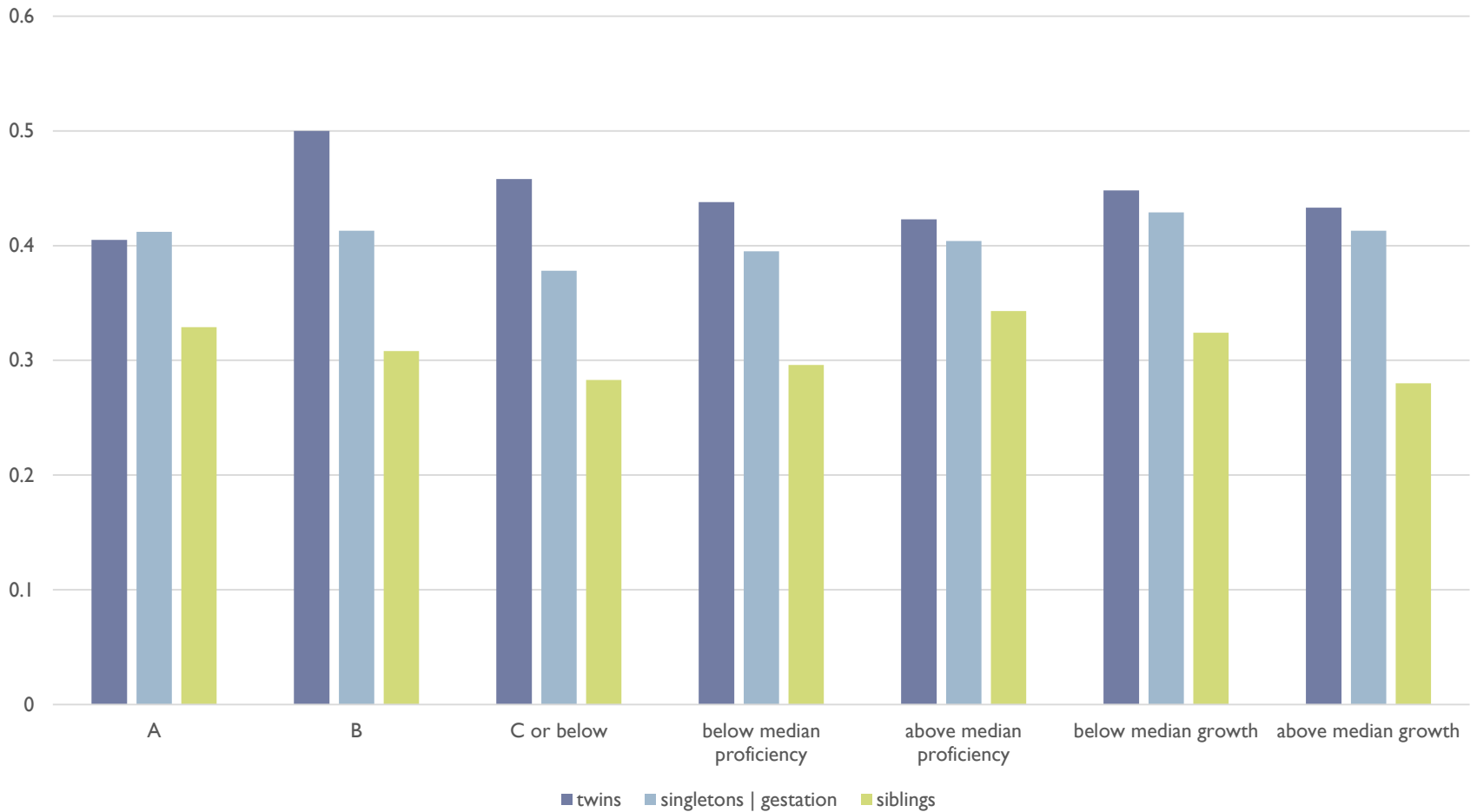
---

- ▶ Since 1999, Florida has graded schools on an A (best) to F (worst) basis
  - ▶ Initially based mainly on average proficiency rates on the criterion-referenced Florida Comprehensive Assessment Test
  - ▶ From 2002 based on a combination of average proficiency rates and average student-level test score gains from year to year
- ▶ We measure “school quality” in three ways:
  - ▶ State-awarded school grade
  - ▶ Average FCAT performance level
  - ▶ Average FCAT gain score



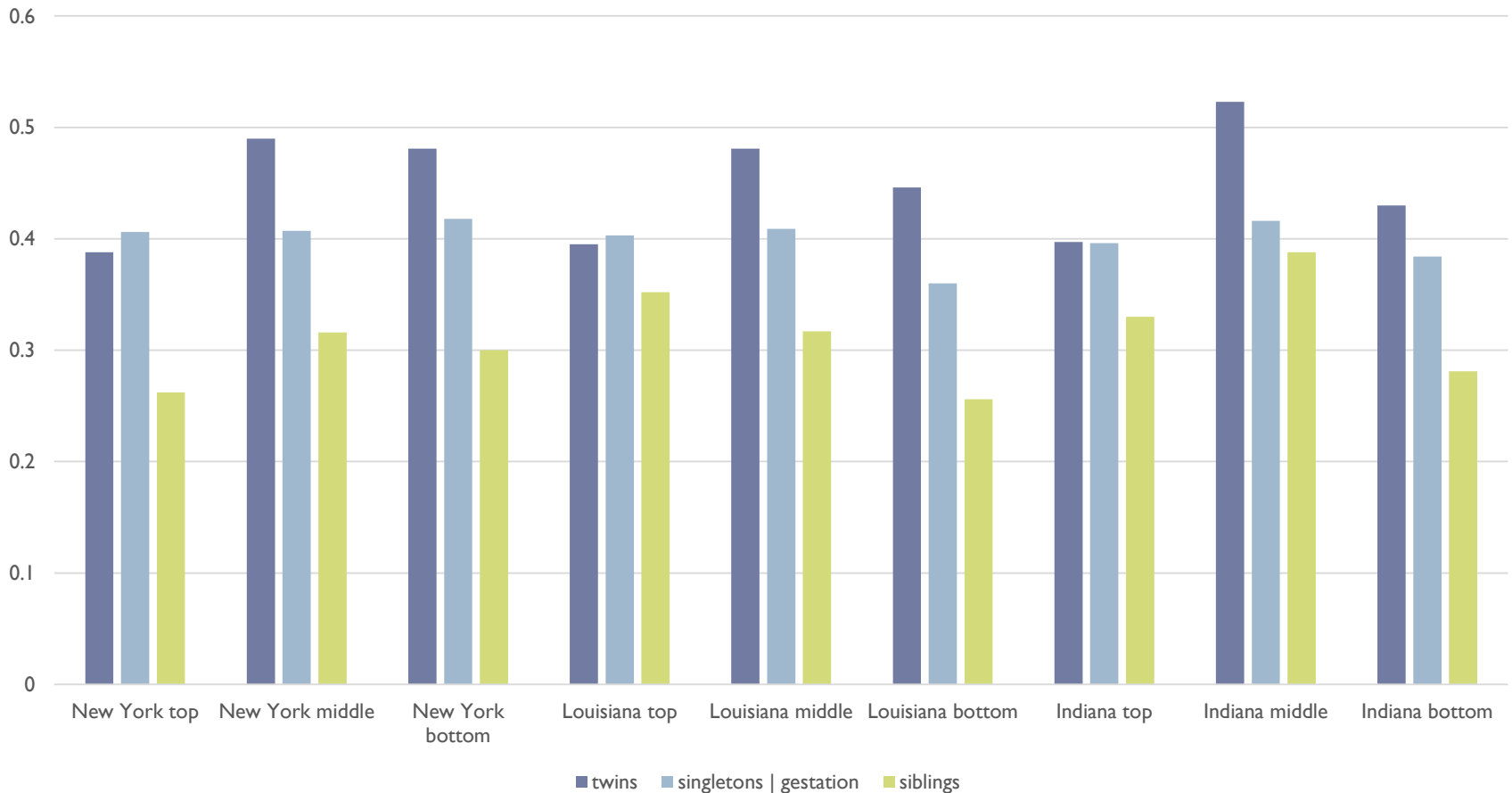
# Does school quality affect the birth weight gap?

Coefficient on log birth weight, by measure of school quality



# Results appear invariant to method of grading schools

Coefficient on log birth weight, running Florida data through other state school grading systems





# Potential implications for policy and practice

---

- ▶ For health policy and practice: The costs of early induction of birth may be greater than previously thought
- ▶ For education policy and practice: There may be educational benefits to greater communication between schools and health care providers – as doing so may help schools target their resources more efficiently



# The bottom line

---

- ▶ These are just two examples of the papers I've written with these matched data (eight projects so far!) – and there will be many more once these data are “democratized”, as is my objective
- ▶ Matched administrative data allow us to, at low cost, study a wide range of research questions – both basic research and direct policy evaluation – that are of immediate concern to policymakers
- ▶ Investing in matched administrative data and providing opportunities to the research community to conduct research and/or collaborate will pay large dividends for science and, one hopes, for policy – at a fraction of the price

