

WORKING PAPER # 10
PRINCETON UNIVERSITY
EDUCATION RESEARCH SECTION
DECEMBER 2004
<http://www.ers.princeton.edu>

Does Competition Among Public Schools Benefit Students and Taxpayers? A Comment on Hoxby (2000)

Jesse Rothstein*
Princeton University

December 15, 2004

* Industrial Relations Section, Firestone Library, Princeton, NJ 08544; jrothst@princeton.edu. Phone (609) 258-4045; fax (609) 258-2907. I thank Orley Ashenfelter, David Card, Tom Davidoff, Olivier Deschenes, Hank Farber, Alan Krueger, Steve Levitt, Robert McMillan, Cecilia Rouse, Lori Taylor, Miguel Urquiola, Jacob Vigdor, Diane Whitmore, three anonymous referees, and seminar participants at Berkeley, Princeton, and the Southern Economics Association meetings for helpful discussions. Anna Huang provided excellent research assistance. I am grateful to the National Center for Education Statistics for access to confidential data and to Caroline Hoxby for generously making available her data and program code. I also gratefully acknowledge financial support from the Princeton Industrial Relations Section and the National Science Foundation. Any errors, opinions, findings, conclusions, or recommendations are mine alone and do not necessarily reflect the views of the NSF or NCES.

Does Competition Among Public Schools Benefit Students and Taxpayers? A Comment on Hoxby (2000)

I. Introduction

School choice policies promise to align the incentives of school administrators with the demands of parents, and may therefore lead to more efficient educational production (Friedman, 1962; Brennan and Buchanan, 1980; Chubb and Moe, 1990). Absent a large-scale school voucher program in the United States, however, this prediction has been difficult to test. Several authors (e.g. Borland and Howsen, 1992; Belfield and Levin, 2002) have suggested studying the effects of “Tiebout choice,” the use of the residential location decision to select among local monopoly education providers. The idea here is that fragmented governance induces competition among school districts analogous to that which would occur among schools with non-residential choice.

In an influential paper, Hoxby (2000) points out that current governance structures are potentially endogenous to school productivity, and proposes that variation in topography, which may have influenced optimal jurisdiction size before modern transportation technologies, provides a source of exogenous variation. She estimates instrumental variables regressions of individual test scores and school spending on a metropolitan-level Tiebout choice index, defined as one minus a Herfindahl concentration index with districts’ enrollments as their “market shares,” using as excluded instruments the number of larger and smaller streams in the area. She reports substantial positive effects of district fragmentation on student test scores and negative effects on spending.

This comment presents a reanalysis of Hoxby’s test score results, which form the core of her empirical analysis. These results turn out to be quite sensitive to plausible alterations to Hoxby’s specification. In particular, the large, significant effect of choice on achievement obtains only with Hoxby’s particular streams variables. When I substitute alternative and arguably better constructions of the same variables, I obtain smaller estimates that are never significant. There is also some

evidence of sample selection bias, deriving from Hoxby's decision to exclude private school students from the analysis. I conclude that Hoxby's positive estimated effect of interdistrict competition on student achievement is not robust, and that a fair reading of the evidence does not support claims of a large or significant effect. Similarly, I find little compelling evidence of endogeneity of the choice index to school quality, suggesting that the more precise OLS estimate of zero choice effect on test scores should be preferred to less precise IV estimates. The evidence that competition among schools will improve academic outcomes is thus substantially weaker than it might have appeared.

Section II focuses on replication. Despite several requests, Hoxby has not provided the precise data set from which her published results were derived. She has, however, made available a corrected data set (Hoxby, 2004a). The new data generate results that deviate in important ways from those that were published. In particular, the first stage coefficients, and even basic summary statistics for the streams variables, are substantially different. Moreover, there appear to be errors remaining in Hoxby's data and computer programs, causing some students to be assigned to the wrong metropolitan statistical areas (MSAs) and some others to be randomly assigned to districts and MSAs. When I correct these errors, I obtain somewhat weaker results. In what I consider the best replication sample, Hoxby's specification and instruments indicate an insignificant or marginally significant effect of choice (i.e., district fragmentation) on student achievement.

In Section III, I consider the sensitivity of the results to the particular instrumental variables used. Hoxby's discussion does not make clear precisely how her larger and smaller streams counts are defined. In particular, though Hoxby writes that the source of her smaller streams variable provides "the longitude and latitude of [each stream's] origin and destination" (2000, p. 1222), she actually uses only streams' destinations to assign them to MSAs. A stream that flows through an MSA but ends elsewhere is not included in the MSA's count. I present results using an alternate variable that counts all streams flowing through each MSA, regardless of where they end. I also

demonstrate that Hoxby's larger streams variable is key to the results, and that it plays a substantially different role in the first stage to the individual-level IV model than in the MSA-level model that Hoxby presents as "the implied first-stage regression" (2000, p. 1224-5).¹ The choice coefficient shrinks by 45 to 85% and ceases to be significant when the larger streams variable is excluded. I obtain similarly small and insignificant coefficients when I substitute alternative larger streams counts that, unlike Hoxby's subjectively coded variable, are readily replicable using public-use data.

Finally, Section IV explores the implications of Hoxby's exclusion of private school students from her sample. Hoxby documents a negative relationship between the Tiebout choice index and the metropolitan private enrollment rate. This may produce selection bias in specifications, like Hoxby's, that are estimated only on public sector students (Hsieh and Urquiola, 2003). Estimates from samples that include both public and private school students are free of this potential sample selection bias, and are notably smaller than those from public-sector samples. None are significantly different from zero, even with Hoxby's instruments.

II. Replication

Table 1 presents IV estimates of the district fragmentation effect on each of two test scores, using Hoxby's streams variables as instruments.² The first column reproduces the estimates from Hoxby's Tables 3 and 4. Hoxby's preferred specification is that for 12th grade reading scores in Panel A, although I analyze 8th grade scores as well (in Panel B) because the sample sizes are so much larger.³ Hoxby assumes that the student-level error term is composed of three homoskedastic

¹ The IV model could be estimated at the MSA level as well, as both the endogenous variable (choice) and instruments (streams) vary only across MSAs. Hoxby (2000, p. 1219) claims that her specification "is most efficiently estimated at the individual level." I follow this decision throughout, though I present MSA-level estimates in the appendix.

² The student test score data are drawn from the National Educational Longitudinal Study (NELS). Details of the data set construction, along with summary statistics, control variable coefficients, and alternative specifications, are in an appendix available from the author.

³ I prefer the 8th grade sample, as its design is much more straightforward than in later waves. Students were randomly sampled from within their schools in the 8th grade, then followed across schools in successive waves. As a result, the follow-up samples are not representative of the schools their students attend, nor of their districts or metropolitan areas, though they remain nationally representative. Also, as with any panel data, sample attrition is a potential problem in later survey waves.

components, one common to all students in the same metropolitan area, another common within the district, and the last specific to the student. She computes standard errors using an FGLS estimator, due to Moulton (1986), that accounts for the implied student-level serial correlation. The estimated choice effect is positive and significant in each panel.

An earlier version of this comment discussed several alternative algorithms for assigning students in the NELS data to school districts and metropolitan areas (MSAs), as Hoxby's (2000) discussion did not specify her approach. In response to that draft, Hoxby re-evaluated her assignment algorithm and discovered some errors (Hoxby, 2004c). She has made available, via the National Center for Education Statistics (NCES), a corrected data set that uses a new crosswalk.⁴ Column 2 reports estimates from the Hoxby/NCES data, which provide substantially smaller samples than were used in the published results. Hoxby's computer program, also provided (Hoxby, 2004b), does not compute the "Moulton" standard errors that were used in the published paper, but instead uses Stata's "cluster" option to generate standard errors which are consistent in the presence of arbitrary heteroskedasticity and within-MSA serial correlation. I have implemented the Moulton estimator, and I report both Moulton and clustered standard errors for each specification in Table 1.⁵ Estimates from Hoxby's corrected data (hereafter, the "Hoxby/NCES" data) have somewhat larger standard errors than did those in the published paper, and the 12th grade coefficient ceases to be significant (at the 5% level) when clustered standard errors are used.

In examining the Hoxby/NCES data and code, I have found several remaining glitches. First, some errors remain in the new district-MSA crosswalk: Several Ohio school districts are

⁴ The corrected data set and the programs used to construct it are available from NCES to researchers who are licensed for access to the restricted-use NELS data.

⁵ Hoxby writes that "Robust [clustered] standard errors are larger than standard errors calculated using the Moulton method" (Hoxby 2004b). Both estimators are consistent (with asymptotics in the number of MSAs) under the error components model, and there is no model in which the Moulton estimator is consistent but the cluster estimator is not. A difference between the two estimators may indicate that the error components assumption is incorrect; in that case, cluster is consistent but the Moulton estimator is not. Further discussion of the two estimators, and of my implementation of the Moulton estimator, is in the appendix.

assigned to the Raleigh-Durham MSA; several additional districts have incorrect, invalid or obsolete MSA codes; and over one quarter of metropolitan districts are missing MSA codes. Second, though the clear intent is to use all three waves of the NELS survey to assign students to districts, due to an apparent coding error information about students' second- and third-wave schools is ignored.⁶

Finally, students with missing school IDs from the first wave of the NELS survey—the sample was freshened in later waves—are randomly assigned to schools that entered the survey in later waves. This occurs because Hoxby's program fails to exclude observations with missing IDs when merging the student and school files. Stata's sort algorithm breaks ties randomly when, as here, a unique sort order is not specified. Stata's merge procedure then assigns the first observation with a missing ID from the "master" data set to the first similar observation from the "using" data set, the second to the second, and so on. Because ties among students and schools with missing IDs are broken differently every time the sort command is run, each execution of Hoxby's program produces a different data set, and different estimated choice effects.⁷ To gauge the severity of this unintended stochasticity, I executed Hoxby's data construction program 10,000 times, tabulating the estimated choice effect from each resulting data set. The histogram is available as Appendix Figure A1. The mean choice effect for 12th grade scores is 5.39, quite close to the 5.30 computed from the Hoxby/NCES data. The standard deviation across iterations (0.47) is not particularly large, but the range is quite wide: I obtained estimates as small as 2.17 and as large as 8.15.

After discovering these anomalies, I re-wrote Hoxby's data assembly program, fixing errors in the district-MSA crosswalk and taking care to correctly match students, schools, districts, and

⁶ Hoxby merges the NELS student file to the NELS school file three times in succession, using school ID variables from each of the three survey waves. By the second merge, all variables from the school file exist on the student file. Without specific instruction (which is not provided), the merge command in Stata does not overwrite variables that already exist on the "master" file, so nothing on the student file is altered by the second and third merges.

⁷ Hoxby's program also fails to account for Stata's tie-breaking procedure when creating the MSA-level data set used for her first stage model, and her program thus assigns the Raleigh MSA to the East North Central division (which contains Ohio; see above) 36% of the times it is executed; the Hoxby/NCES data set is one such draw from the distribution.

metropolitan areas. I attempted to follow Hoxby’s algorithm as closely as possible.⁸ I did not at this point attempt to reproduce the “larger streams” variable, but simply relied on the MSA-level count that Hoxby provided and discarded MSAs that were excluded from her tabulation.⁹ Results are presented in column 3 of Table 1. Sample sizes are somewhat larger—correctly assigning districts that were previously classified as non-metropolitan more than offsets the loss of students who are reclassified to an MSA with a missing larger streams value—and approach those seen in Hoxby’s Table 4. Coefficients resemble those found in the Hoxby/NCES data, somewhat smaller for 12th grade scores and somewhat larger for 8th grade scores, with similar patterns of significance.

Column 4 represents a somewhat more expansive interpretation of replication. I retain Hoxby’s specification, but I follow my own judgment in sample and covariate construction rather than directly following her algorithm. Where Hoxby assigns each student to a single district for all three waves even if the student moved between waves, for this sample I use only contemporaneous information to construct distinct assignments for each wave. There are also minor differences in variable definitions.¹⁰ Choice effect estimates are smaller with this sample. For 12th grade scores, the choice effect is insignificant regardless of the standard error computation; for 8th grade scores, it is insignificant with the random effects standard errors but significant when the errors are clustered.

⁸ There were some ambiguities. In particular, each student has nine potential district codes, as each student may have a school code in each of three waves and each school may have different district codes in each wave. Hoxby attempts to assign a single district code for each student, to be used with data from all three waves, but the aforementioned coding errors mean that only the three district codes from the first-wave school are considered. It is not clear how she would resolve discrepancies among the larger set. I assign each student to a separate district for each wave, using only contemporaneous information from the student and school files, then use Hoxby’s majority rule algorithm to select among the three resulting assignments.

⁹ Hoxby uses 1990 MSA definitions. Puzzlingly, she does not provide counts of larger streams for all of the MSAs included in these definitions, but does provide counts for some obsolete MSA codes—from the 1983 or 1981 MSA definitions—that appear in her faulty crosswalk. For example, 19 larger streams are reported for MSA number 3755, which corresponded to the Kansas City, KS PMSA in 1983 but was included in the Kansas City MO-KS MSA (number 3760) in 1990; there is also an entry of 37 larger streams in MSA 3760. It is not clear what algorithm might have produced this redundancy, nor whether the latter count includes the streams attributed to the former.

¹⁰ The largest difference is in what Hoxby calls the “mean of log(income) of metropolitan area” variable. She uses an arithmetic weighted average of the log of each district’s mean income; I use instead the log of the MSA mean income. There are also minor differences in the Gini coefficient and the racial composition variables. Finally, I compute the choice index over 8th grade enrollment, where Hoxby uses total enrollment, reasoning that parents cannot be said to choose between overlapping elementary and secondary districts (Urquiola, 1999). Further details are in the appendix.

Panel A of Table 2 reports mean values of the streams variables. Column 1 is from Hoxby's Table 2, while columns 2 and 3 are computed from the Hoxby/NCES data set and from my replication sample, respectively. There are substantial differences between columns 1 and 2. For some reason, the mean of the larger streams variable is more than five times larger than that reported in the published paper, while the average MSA has only two thirds as many total—larger plus smaller—streams as is indicated by Hoxby's (2000) Table 2.

Both the streams variables and the potentially endogenous choice measure vary only at the MSA level. Though Hoxby's IV estimates are computed at the student level, Hoxby reports only an MSA-level "implied first-stage regression." I reproduce this specification in Panel B, with the published estimates in column 1, those from the Hoxby/NCES data in Column 2, and those from the replication samples in 3 and 4.¹¹ All of the replication estimates are substantially different from those in the published paper. Comparing the Hoxby/NCES estimates to the published results, the larger streams coefficient has fallen by more than 80% and is no longer remotely significant, while the smaller streams coefficient has tripled. Though both of these findings are somewhat attenuated in the replication data sets, they remain worrisome: The logic of the argument for Hoxby's instruments is that streams once represented impediments to travel, and one would expect this to be far more true for larger than for smaller streams, particularly when the threshold for being a "larger" stream is set low enough to include over 40 streams from the average MSA (rather than the 8 indicated in the published paper).

As noted above, the MSA-level estimates are not the actual first stages for the individual-level models in Table 1. The actual first stages are reported in Panel C (for the 12th grade samples) and D (for the 8th grade samples). The streams coefficients are dramatically different: Larger

¹¹ The replication data sample sizes are somewhat smaller, as several invalid MSA codes that were on the Common Core of Data file from which Hoxby took her district-MSA assignments are no longer present and some newly added MSA codes must be excluded for lack of the larger streams variable.

streams are now *negatively* related to choice in five of the six samples, once significantly and once nearly so.¹² Again, this is difficult to reconcile with the story behind the identification strategy.

III. Counting Streams

There are several reasons to worry about the validity of Hoxby's larger streams variable: It derives from Hoxby's subjective count from printed maps—she describes counting streams “of a certain width on the map,” (2000, p. 1222), but does not elaborate; it is missing for several MSAs that were inadvertently excluded from Hoxby's sample;¹³ and, as Hoxby writes, “one has more a priori confidence in the exogeneity of the smaller streams variable because smaller streams are too small to affect modern life,” (2000, p. 1230). Given the evident differences between the larger streams variable described in the published paper and the one included in the Hoxby/NCES data, it is unclear whether the discussion in Hoxby's text even applies to the latter variable.

These concerns cannot be addressed by using the smaller streams variable as the sole instrument, however. Hoxby uses the U.S. Geologic Survey's Geographic Names Information System (GNIS) to count total streams, and defines smaller streams as the number of total streams less the count of larger streams. As a result, any errors in the larger streams variable appear as errors of the opposite sign in the smaller streams count. To avoid reliance on Hoxby's larger streams count, I present estimates that use the total streams count—which can be produced using Hoxby's code from the public-use GNIS data set—as the single instrument.

¹² The divergence between the MSA-level results in Panel B and the individual-level results in Panels C and D appears to derive from differences in the set of MSAs included. Hoxby's first stage estimates and those that I report in Panel B include all MSAs, regardless of whether they contain NELS sample students. When I restrict the sample to those in the NELS data (Appendix Table D5), coefficients are similar to those in Panels C and D. Efficiency can be improved with two-sample IV, using the full sample of MSAs to estimate the first stage. In the Hoxby/NCES data, this yields choice coefficients of 3.68 for 8th grade scores and 2.14 for 12th grade scores, both substantially shrunken from the estimates in Table 1 and neither significant (Appendix Table D6).

¹³ One indication that there may be problems with Hoxby's larger streams count is that when I correct Hoxby's code to correctly assign total streams to MSAs—her incorrect district-MSA crosswalk is used here as well—there are several MSAs with fewer total streams than larger streams. Hoxby writes that the hand counts were “checked against” the GNIS data (2000, p. 1222), but appears not to have caught all discrepancies. Though I argue below that Hoxby systematically undercounts total streams, my correction of this problem reduces but does not eliminate the discrepancies.

I also explore an alternative specification for the “total streams” variable. Despite her reference to GNIS variables describing the longitude and latitude of streams’ origins and destinations, Hoxby’s code uses only a variable indicating the county where a stream’s destination (mouth) is located to assign streams to MSAs. To illustrate the consequences of this, the Mississippi River is attributed only to the non-metropolitan Plaquemines Parish, Louisiana, and not to any of the eight metropolitan areas along its banks.¹⁴ There is little reason to think that a stream’s destination is the key to either its past effects on travel costs or to its current effects on district structure. The USGS distributes an alternative version of the GNIS data that codes each county through which each stream flows, from origin to destination. Using this data file, I construct a “total streams” measure that counts toward an MSA’s total any stream flowing through it.¹⁵

Finally, I explore alternative classifications of streams into “larger” and “smaller” groups. First, following Hoxby (1994b), I compute separate counts of inter-county and intra-county streams and enter them as separate instruments. I also categorize streams based on their lengths, computed as the distance between their sources and mouths, following Hoxby (2000) in requiring a larger stream to exceed 3.5 miles. Each is a crude measure for the variation of interest, but it is difficult to see how either might be endogenous; as a result, either should provide consistent IV estimates of the choice effect.¹⁶ These estimates provide a check on the robustness of the earlier estimates, and have the virtue of being easily replicable using the public-use GNIS data.

¹⁴ This is not documented in the published paper. It does not automatically mean that inland cities lack streams, as a smaller stream’s mouth might be located where it feeds into a larger river. Note that the Mississippi may be included in the *larger* streams counts for the relevant MSA’s, though it is not counted toward the *total* streams. This appears to account for some but not all of the negative smaller streams counts discussed in footnote 13.

¹⁵ In most of the country, MSAs are composed of whole counties. In New England, however, towns are the basic unit, and some counties are split among several MSAs. Hoxby assigns all of each county’s streams to the MSA containing the plurality of its population. When I reproduce her stream mouths variable, I follow her all-or-nothing rule; my total streams count instead assigns streams fractionally to MSAs in proportion to the MSAs’ shares of the county population.

¹⁶ Measurement error in instruments, so long as it is uncorrelated with the endogenous variable, reduces the precision of IV estimates but does not affect consistency as long as the measures are sufficiently reliable to avoid so-called “weak instruments” problems. As I show below, the first stages are quite strong.

Table 3 presents instrument means (Panel A) and first-stage estimates (Panels B-D, using the close replication sample) for several instrument sets. As before, the first stage is computed at both the MSA and individual levels; corresponding estimates using my alternative sample and covariate definitions are similar and are reported in the appendix. For a benchmark, Column 1 reproduces the estimate from Column 3 of Table 2, using Hoxby's streams variables. Column 2 uses only total streams (by Hoxby's definition, counting only stream mouths), which have positive coefficients at both the MSA and individual levels. Columns 3 and 4 repeat these specifications, using the count of all streams flowing through each MSA in place of the count of stream mouths. This change has little effect on the estimates, with the negative larger streams coefficient still evident in the individual-level model. Columns 5 and 6 use alternative definitions for "larger" streams, first as inter-county streams and second as streams exceeding 3.5 miles in length. Using either definition and in both the MSA and individual samples, the larger streams variable accounts for the full effect of streams on choice, a result that is consistent with the idea that the role of streams derives from their importance as natural barriers to travel.

For each set of instruments, Table 4 reports IV estimates of the choice effect on 12th and 8th grade reading scores, Moulton and clustered standard errors, and p-values for tests of the exogeneity of the choice variable (using the cluster estimator).¹⁷ I also report OLS estimates, each of which indicates a negligible choice effect.

The choice effects are consistently positive and exogeneity of the choice variable is consistently rejected when Hoxby's larger streams count is included as an instrument. Neither of these results holds in any of the specifications that exclude Hoxby's larger streams variable, however. This is partly because the latter estimates are less precise, but this is not the whole story: The

¹⁷ I obtain similar results with Moulton standard errors or when I use the preferred replication sample and covariates.

coefficient estimates are also uniformly smaller, generally less than half as large, when Hoxby's larger streams variable is excluded.

Taking the estimates in Table 4 together, it is clear that Hoxby's conclusions depend critically on her count of larger streams. I attempted my own count for several MSAs that contribute most to the large choice effect estimates, using the same 1/24,000 quadrangle maps that Hoxby reported using. It quickly became apparent that counting streams involves many subjective judgments.¹⁸ Hoxby describes larger streams as those that "were at least 3.5 miles long and of a certain width on the map" (2000, p. 1222), but does not specify what constitutes "a certain width" nor where in a stream's course the width is to be measured. I began with Fort Lauderdale, which may be a particularly difficult case as much of the MSA is swampland and much of the remainder was recovered from swampland by a system of man-made canals. (Even today, airboat trails are more common through much of the MSA than is dry land; it seems unlikely to have been settled by people who viewed water as an obstacle to travel.) I decided not to count canals which ran perfectly straight, generally exactly West to East, but I did count canals which took irregular paths, reasoning that the latter were more likely to correspond to pre-existing rivers. I also counted branches of streams as separate from their parents when they had distinct names (such as the North and South Forks of the Middle River), and counted the intracoastal waterway, which separates the easternmost portion of the Florida coast from the mainland, as a stream for its similar effect on the ease of travel. Where Hoxby reports 5 larger streams in Fort Lauderdale, I counted 12, and a research assistant—working independently—counted 15.

I had a similarly difficult experience with other MSAs, finding that many rivers divide and recombine multiple times, become wider and narrower, and are interrupted by man-made structures

¹⁸ I worked without reference to Hoxby's counts, to prevent being influenced by these. Hoxby's text is confusing about whether linear bodies of water other than streams are included in her count. Her footnote 24 seems to suggest that they are not, but her footnote 16 indicates that she counts "inlets, lakes, ponds, marshes, and swamps" "*if they are roughly curvilinear in form*" (emphasis in original). I followed the latter rule.

throughout their courses. My counts were correlated with Hoxby's, but generally not identical. The exercise makes clear that Hoxby's larger streams variable is subjective and unverifiable without a list of the particular rivers coded as large. In the absence of such a list, which Hoxby has not provided, no two researchers would come up with identical counts. As I have only counted streams for a few MSAs, however, I cannot be certain of the sensitivity of Hoxby's results to the differences that would inevitably arise.

IV. Private Enrollment and Selection Bias

I have concerned myself thus far with replication of Hoxby's primary specification, and with its robustness to plausible alternative decisions about sample and variable construction. In this section, I turn to another issue: Hoxby's specification may not provide consistent estimates of the effect of interest, that of choice on public school productivity, because her sample excludes private school students. In her Table 6, she documents that choice has a significant negative effect on the metropolitan private enrollment share.¹⁹ As a result, Hoxby's specification may be subject to selection bias even with valid instruments (Hsieh and Urquiola, 2003). The reasoning is simple: Suppose that the distribution of student test scores is identical across MSAs when both public and private school students are included, but that MSAs vary in private enrollment patterns. In particular, suppose that some relatively high-scoring students would choose private schools in a low-choice market but would remain in the public sector when Tiebout choice is sufficient to provide public schools with desired characteristics (Rothstein, 2004). Then the average test score among public school students will tend to be higher in high-choice markets purely as a result of differential sample selection.

¹⁹ Using both of her streams instruments in a district-level regression, Hoxby (2000, Table 6) estimates that a one-unit increase in choice leads to a 4.2% (s.e. 1.2%) reduction in private enrollment. Hoxby's SDDDB data set double-counts students in areas served by separate elementary and secondary districts. When I instead estimate the relationship at the MSA level, I estimate a choice effect of -4.8% (s.e. 2.4%), though this result is somewhat sensitive to the sample and covariate construction.

Any resulting bias is present in both OLS and IV estimates, though its sign and magnitude depend on whether the marginal private school student is positively or negatively selected. If the average score is higher among students drawn into the public sector by expansions of choice than among inframarginal public school students, estimates from public school students are (asymptotically) upward-biased; if the average score is lower among marginal students than among the inframarginal, these estimates are downward-biased.²⁰ Hoxby seems to make the former claim when she discusses the consequences of “families with a strong taste for education leav[ing] the public sector by shifting their children into private schools” (2000, p. 1233).

As the NELS survey includes both public and private school students, this potential bias can be easily avoided by simply including both groups in the sample.²¹ The only hurdle is that the CCD cannot be used to assign private schools to school districts and MSAs. As an alternative, I use NELS variables characterizing the demographic composition of the school’s zip code to uniquely assign the vast majority of schools to zip codes, and via these to MSAs.²² As many zip codes span school districts, I cannot use this strategy to assign school districts, and I therefore must exclude district-level covariates from the specification.²³

Panel A of Table 5 reports estimates from public school students who have been matched to MSAs via their schools’ zip codes, using both the “close” and “preferred” covariate definitions.

²⁰ NELS private school students score nearly half a standard deviation higher on the 8th grade reading test than do public school students. This is not particularly informative, however, as the students whose sectoral decision is sensitive to Tiebout choice are likely atypical of the inframarginal private school population.

²¹ Under fairly strong assumptions—including that private schools are not systematically better or worse than public schools; that competition has similar effects on the productivity of public and private schools; and that any peer effects are linear and additive, so that stratification does not have an independent effect on average scores— an unbiased estimate of the choice effect on average school productivity can be obtained by estimating Hoxby’s specification on a pooled sample of public and private school students (Hsieh and Urquiola, 2003). Hoxby (1994a) uses exactly this strategy to test for selection bias from private school enrollment.

²² In the rare cases where a zip code spans multiple MSAs, I assign each student attending school in that zip code to each MSA, with weights proportional to the fraction of the zip code population in each MSA.

²³ Hoxby (2000, Section 7) argues at great length that the inclusion of district-level variables improves the precision but does not affect the coefficients on MSA-level variables as long as MSA-level means are included in the specification. Strictly, this is only true in the limit, as it relies on the assumption that the district-level variables aggregate exactly within the sample to the MSA-level means. In small samples this is not likely to hold, and the choice coefficient is somewhat smaller (more negative) when district-level covariates are excluded from Hoxby’s specification (Appendix Table D3).

Estimates are substantially smaller than those presented earlier, with the divergence due more to the different methods of assigning MSAs than to the exclusion of district-level covariates.²⁴ Panel B adds the private school students to the sample. The choice effect estimates fall notably farther here, and t -statistics are uniformly less than one.

I read the estimates in Table 5 as suggesting, but not conclusively demonstrating, that the students drawn into the public sector by expansions of choice are somewhat positively selected.²⁵ While much of the difference from earlier estimates appears to derive from sensitivity of the results to the exclusion of district-level covariates and to the method by which schools are assigned to MSAs, point estimates do fall even farther when private school students are added to the sample.

V. Discussion

Hoxby's analysis has been very influential, providing what many (e.g. Howell and Peterson, 2002; Maranto, 2001; Bast and Walberg, 2004) have seen as some of the most compelling extant evidence in favor of the proposition that school choice will lead to improvements in the efficiency of educational production. Unfortunately, Hoxby's key results do not seem to be robust to small, reasonable alterations to the sample or to the instrumental variables used. Interested readers are invited to explore alternative specifications beyond those considered here; code to construct both of my replication samples and to perform all analyses is available from my web page, as are all data components that I am at liberty to distribute.

As I document above, there are several problems with the Hoxby/NCES data set. When these are remedied, I estimate somewhat weaker effects of choice on student performance than

²⁴ The declines are largest in the close replication sample, as my zip code matching algorithm, which uses only the contemporaneous school, is more similar to that used in the preferred sample. Students in the close replication sample who were assigned to MSAs based on their 8th or 10th grade school's district code in Panel B are assigned using the 12th grade school's zip code in Panel C.

²⁵ As an alternative test for selection bias, I have estimated a version of Hoxby's specification (using only public school students) that includes a control for an inverse Mill's ratio computed from the MSA private enrollment rate, in the spirit of normal-distribution selection corrections (Gronau, 1974; Heckman, 1979; Card and Payne, 2002). Estimates of the selectivity parameter were extremely imprecisely estimated, and the selection correction had little effect on the estimated choice coefficients.

those that Hoxby reports.²⁶ When I consider slight adjustments to her specification of the streams variables—such as replacing them with plausible, replicable alternative measures—or when I alter the sample to avoid potential selection bias from private enrollment, the significant effect of Tiebout competition on student scores is greatly attenuated and not statistically distinguishable from zero. In my specification including private school students, using my preferred sample, and instrumenting with inter- and intra-county streams (Table 5, Panel B, Column 6), I estimate that a one standard deviation increase in choice raises test scores by just under 0.05 standard deviations, with a standard error somewhat larger than that. This compares unfavorably to, for example, the 0.22 standard deviations that Krueger (1999) estimates as the effect of reducing elementary school class sizes from 22 to 15 students in the Tennessee STAR experiment.

I do not find support, in any of the alternative specifications that I consider, for Hoxby's claim that "naïve estimates (like OLS) that do not account for the endogeneity of school districts are biased toward finding no effects" (2000, p. 1236), nor for her conclusion that "Tiebout choice raises productivity by simultaneously raising achievement and lowering spending" (p. 1236-7). Any relationship between choice and student test scores is too imprecisely estimated to be robustly distinguishable from zero. Hoxby's results for the effect of district fragmentation on school spending, which I examine in the appendix, are only slightly more robust.²⁷

There are only a few hundred metropolitan areas in the United States, and this is evidently too few to precisely estimate any relationship that may exist between jurisdictional fragmentation and either student performance or school spending. One cannot reject large effects of competition, but neither is there strong evidence against a hypothesis of zero effect. It would be premature to

²⁶ The current analysis has not considered Hoxby's analysis of the NLSY, which echoes her NELS analysis in indicating a salutary effect of interdistrict competition on attainment. Hoxby seems to find her NELS estimates the most compelling, however, and focuses her discussion on these.

²⁷ Hoxby (2000, Table 5) reports a choice effect on the log of per pupil spending of -0.076 (Moulton standard error 0.034). The Hoxby/NCES data yield an estimate of -0.074 (0.141); IV estimates in the replication samples similarly fail to reject zero, although OLS estimates are significantly negative.

conclude that schools respond to Tiebout competition by raising productivity, nor to use such a conclusion as justification for policies that expand non-residential forms of school choice.

References

- Bast, Joseph L. and Walberg, Herbert J.** "Can Parents Choose the Best Schools for Their Children?" *Economics of Education Review*, August 2004, 23(4), pp. 431-40.
- Belfield, Clive R. and Levin, Henry M.** "The Effects of Competition between Schools on Educational Outcomes: A Review for the United States." *Review of Educational Research*, Summer 2002, 72(2), pp. 279-341.
- Borland, Melvin V. and Howsen, Roy M.** "Student Academic Achievement and the Degree of Market Concentration in Education." *Economics of Education Review*, March 1992, 11(1), pp. 31-39.
- Brennan, Geoffrey and Buchanan, James.** *The Power to Tax: Analytical Foundation of a Fiscal Constitution*. Cambridge: Cambridge University Press, 1980.
- Card, David and Payne, A. Abigail.** "School Finance Reform, the Distribution of School Spending, and the Distribution of SAT Scores." *Journal of Public Economics*, 2002, 83(1), pp. 49-82.
- Chubb, John and Moe, Terry M.** *Politics, Markets, and America's Schools*. Washington, D.C.: The Brookings Institution, 1990.
- Friedman, Milton.** *Capitalism and Freedom*. Chicago: University of Chicago Press, 1962.
- Gronau, Reuben.** "Wage Comparisons--a Selectivity Bias." *The Journal of Political Economy*, Nov. - Dec. 1974, 82(6), pp. 1119-43.
- Heckman, James J.** "Sample Selection Bias as a Specification Error." *Econometrica*, 1979, 47, pp. 153-61.
- Howell, William G. and Peterson, Peter E.** "Impact of School Voucher Research." *Ps-Political Science & Politics*, December 2002, 35(4), pp. 659-60.
- Hoxby, Caroline M.** "Do Private Schools Provide Competition for Public Schools?" National Bureau of Economic Research Working Paper #4978, December 1994a.
- _____. "Does Competition among Public Schools Benefit Students and Taxpayers?" National Bureau of Economic Research Working Paper 4979, December 1994b.
- _____. "Does Competition among Public Schools Benefit Students and Taxpayers?" *American Economic Review*, December 2000, 90(5), pp. 1209-38.
- _____. "District-Level and Metropolitan-Area Variables Merged with NELS Data." CD, National Center for Education Statistics, September 2, 2004a.
- _____. "Documentation to 'District-Level and Metropolitan-Area Variables Merged with NELS Data': Construct.Do." National Center for Education Statistics, September 2, 2004b.
- _____. "Documentation to 'District-Level and Metropolitan-Area Variables Merged with NELS Data': Runregressions.Do." National Center for Education Statistics, September 2, 2004c.
- Hsieh, Chang-Tai and Urquiola, Miguel.** "When Schools Compete, How Do They Compete? An Assessment of Chile's Nationwide School Voucher Program." Mimeo, Columbia University, August 2003.
- Krueger, Alan B.** "Experimental Estimates of Education Production Functions." *Quarterly Journal of Economics*, May 1999, 114(2), pp. 497-532.
- Maranto, Robert.** "Finishing Touches." *Education Next*, Winter 2001, 2001(4), pp. 20-25.

Moulton, Brent R. "Random Group Effects and the Precision of Regression Estimates." *Journal of Econometrics*, August 1986, 32(3), pp. 385-97.

Rothstein, Jesse M. "Good Principals or Good Peers? Parental Valuation of School Characteristics, Tiebout Equilibrium, and the Incentive Effects of Competition among Jurisdictions." National Bureau of Economic Research Working Paper #10666, August 2004.

Urquiola, Miguel. "Demand Matters: School District Concentration, Composition, and Educational Expenditure." University of California, Berkeley, Center for Labor Economics Working Paper #14, April 1999.

Table 1: IV estimates of choice effect on NELS 8th and 12th grade reading scores in several samples, Hoxby specification

	Published	Hoxby/ NCES data	Close replication sample	Preferred sample and covariates
	(1)	(2)	(3)	(4)
<i>Panel A: 12th grade reading scores</i>				
# of students	6,119	5,475	5,934	6,688
# of MSAs	209	184	194	199
Choice index coefficient	5.77	5.30	4.74	3.29
S.E. (Moulton)	(2.21)	(2.36)	(1.98)	(1.83)
S.E. (Cluster)		(2.94)	(2.42)	(2.56)
P-values, exogeneity test (clustered)		0.02	0.02	0.20
<i>Panel B: 8th grade reading scores</i>				
# of students	10,790	10,175	10,429	11,719
# of MSAs	211	185	186	184
Choice index coefficient	3.82	4.45	5.93	2.93
S.E. (Moulton)	(1.59)	(1.87)	(2.10)	(1.58)
S.E. (Cluster)		(1.99)	(2.32)	(1.40)
P-values, exogeneity test (clustered)		0.00	0.00	0.00

Notes: See Hoxby (2000) and data appendix for description of data, samples, and covariates. Column 1 is from Hoxby (2000), Table 4. Standard error estimators and exogeneity tests are described in the appendix. Following Hoxby, all analyses use NELS sampling weights, adjusted to sum to one within each MSA (though this does not hold exactly in Column 2; see appendix for details). Bold S.E.s indicate that with that S.E., the coefficient is significant at the 5% level.

Table 2: Overview of first stage estimates, different samples
Dependent variable is MSA-level choice index (1- index of concentration across districts)

	Published	Hoxby/ NCES data	Close replication sample	Preferred sample and covariates
	(1)	(2)	(3)	(4)
<i>Panel A: MSA-level sample means</i>				
Larger streams	8	44	45	
Smaller streams	183	84	80	
<i>Panel B: MSA-level first stage estimates</i>				
Larger streams (100s)	0.080 (0.040)	0.012 (0.021)	0.040 (0.021)	0.043 (0.021)
Smaller streams (100s)	0.034 (0.007)	0.096 (0.019)	0.093 (0.018)	0.091 (0.018)
N	316	310	304	304
F statistic (instruments)	24.4	14.8	16.2	16.3
<i>Panel C: Individual-level first stage estimates (12th grade reading sample)</i>				
Larger streams (100s)	nr	-0.043 (0.023)	-0.024 (0.020)	0.015 (0.020)
Smaller streams (100s)	nr	0.133 (0.021)	0.133 (0.017)	0.114 (0.018)
N	nr	5,475	5,934	6,688
F statistic (instruments)	nr	20.5	31.3	28.4
<i>Panel D: Individual-level first stage estimates (8th grade reading sample)</i>				
Larger streams (100s)	nr	-0.045 (0.021)	-0.033 (0.018)	-0.012 (0.018)
Smaller streams (100s)	nr	0.131 (0.022)	0.130 (0.017)	0.132 (0.017)
N	nr	10,175	10,429	11,719
F statistic (instruments)	nr	17.6	30.7	32.1

Notes: "nr"=not reported. Column 1 is from Hoxby (2000), Table 2. Sample sizes in Panels C and D are identical to those in the corresponding columns of Table 1, Panels A and B respectively. Standard errors are clustered in Panels C and D, but are conventionally calculated (under homoskedasticity assumptions) in Panel B.

Table 3: First-stage estimates for alternative instruments, using "close replication" sample and covariates

	(1)	(2)	(3)	(4)	(5)	(6)
Total stream definition:	Stream mouths		All streams			
Larger stream definition:	Hoxby	n/a	Hoxby	n/a	Inter-county >3.5 miles	
<i>Panel A: MSA-level sample means</i>						
Larger streams	45		45		41	70
Smaller streams	80		108		107	75
Total streams		124		148		
<i>Panel B: MSA-level first stage estimates</i>						
Larger streams (100s)	0.040 (0.021)		0.037 (0.021)		0.260 (0.055)	0.177 (0.036)
Smaller streams (100s)	0.093 (0.018)		0.069 (0.013)		0.014 (0.016)	0.013 (0.017)
Total streams (100s)		0.071 (0.013)		0.061 (0.010)		
F statistic (instruments)	16.2	30.9	17.5	36.5	25.8	23.9
<i>Panel C: Individual-level first stage estimates (12th grade reading sample)</i>						
Larger streams (100s)	-0.024 (0.020)		-0.030 (0.019)		0.240 (0.047)	0.190 (0.029)
Smaller streams (100s)	0.133 (0.017)		0.104 (0.013)		0.015 (0.013)	0.001 (0.013)
Total streams (100s)		0.064 (0.011)		0.058 (0.009)		
F statistic (instruments)	31.3	32.0	35.0	37.0	27.5	33.7
<i>Panel D: Individual-level first stage estimates (8th grade reading sample)</i>						
Larger streams (100s)	-0.033 (0.018)		-0.036 (0.017)		0.243 (0.046)	0.177 (0.029)
Smaller streams (100s)	0.130 (0.017)		0.101 (0.012)		0.011 (0.012)	0.001 (0.014)
Total streams (100s)		0.059 (0.011)		0.054 (0.009)		
F statistic (instruments)	30.7	28.9	34.8	34.7	26.5	30.1

Notes: Base samples are those from Column 3 of Tables 1 (individual level) and 2 (Panel B; MSA level), though some observations that were excluded from those samples for missing data on larger streams are included here in Columns 2, 4, 5, and 6. Alternative specifications that use the preferred covariates and sample are in the appendix. In individual-level specifications, standard errors are clustered at the MSA level. "n/a" indicates not applicable: No larger streams instrument is used in this specification, and the only instrument is the "total streams" count.

Table 4: IV estimates of choice effect, using alternative instruments and "close replication" sample

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	OLS	IV					
Total stream definition	n/a	Stream mouths		All streams			
Larger stream definition	n/a	Hoxby	none	Hoxby	none	Inter-county	>3.5 miles
<i>Panel A: 12th grade reading scores</i>							
Choice index coefficient	-0.25	4.74	0.68	4.38	0.87	2.04	1.35
S.E. (Moulton)	(0.79)	(1.98)	(2.79)	(1.98)	(2.59)	(2.36)	(2.30)
S.E. (Cluster)	(0.94)	(2.42)	(3.12)	(2.15)	(2.81)	(2.94)	(2.04)
p-value, exog. test	--	0.02	0.70	0.02	0.66	0.37	0.38
<i>Panel B: 8th grade reading scores</i>							
Choice index coefficient	-0.06	5.93	2.76	5.17	2.78	1.67	0.91
S.E. (Moulton)	(0.70)	(2.10)	(2.54)	(2.01)	(2.33)	(2.09)	(1.93)
S.E. (Cluster)	(0.82)	(2.32)	(3.19)	(2.02)	(2.84)	(1.77)	(1.81)
p-value, exog. test	--	0.00	0.30	0.00	0.24	0.21	0.51

Notes: Base samples are those from Column 3 of Table 1, though some observations that were excluded from that sample for missing data on larger streams are included here in Columns 3 and 5-7. Alternative specifications that use the preferred covariates and sample are in the appendix. Exogeneity tests are based on clustered specification. Bold S.E.s indicate that with that S.E., the coefficient is significant at the 5% level. "n/a" indicates not applicable: No excluded instruments are used in this specification.

Table 5. Exploration of potential bias from exclusion of private school students, 12th grade reading scores

Covariate specification	Close replication			Preferred replication		
	OLS	Hoxby	Inter- and intra-cnty	OLS	Hoxby	Inter- and intra-cnty
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Public school students in zip-code matched sample (no district covariates)</i>						
Choice index coefficient	-0.93	1.40	1.10	-0.76	1.97	2.25
S.E. (Cluster)	(1.05)	(2.44)	(2.66)	(0.97)	(2.20)	(2.30)
N	5,631	5,445	5,631	6,976	6,729	6,976
p-value, exog. test		0.35	0.36		0.22	0.12
<i>Panel B: Public and private school students in zip code-matched sample</i>						
Choice index coefficient	-0.71	0.68	0.84	-0.41	1.35	1.81
S.E. (Cluster)	(0.98)	(2.59)	(2.35)	(0.92)	(2.32)	(2.14)
N	6,900	6,670	6,900	8,553	8,259	8,553
p-value, exog. test		0.63	0.43		0.45	0.22

Notes: Clustered standard errors and test statistics are reported.

Appendices to:

“Does Competition among Public Schools Benefit Students and Taxpayers? A Comment on Hoxby (2000)”

Jesse Rothstein

December 2004

Appendix A: Data

1. The Hoxby/NCES data

Each issue of the *American Economic Review* includes a paragraph describing the “Policy on Data Availability:” “...to publish papers only if the data used in the analysis are clearly and precisely documented and are readily available to any researcher for purposes of replication.” Despite repeated requests over several years, Hoxby has not provided the data used to generate the results in her published paper (Hoxby, 2000). She has, however, provided a corrected data set (Hoxby, 2004a) to the National Center of Educational Statistics, which will distribute it to researchers with licenses for the confidential version of the National Educational Longitudinal Survey. I refer to these data throughout as the “Hoxby/NCES data.” Hoxby (2004c) notes several differences between the data used in her paper and the corrected Hoxby/NCES data. First, where she originally used MSA codes from the School District Data Book (SDDB) to assign school districts to metropolitan areas, she now uses the Common Core of Data (CCD), with a few manual corrections, for this purpose. Second, “[t]he published version of the paper showed a version of the first stage regression that was outdated. The paper went through multiple revisions and Table 2 evidently was not updated so that it corresponded exactly with Tables 3+.” As I document in the main text, the actual first stage regression is substantially different than that indicated in the published paper.

Hoxby constructs her new district-MSA crosswalk from the Common Core of Data agency files, using a field on these files that indicates the MSA in which each district is located. Unfortunately, there are substantial errors in this field. Many of these come from a small number of MSAs whose boundaries have changed over time; the CCD often contains

obsolete MSA codes for districts in these areas.¹ One example is the “Kansas City, MO-KS MSA,” which was coded 3760 in 1990 but had been divided into separate Missouri (code 3760) and Kansas (code 3755) components in 1983. Even the 1991-92 CCD, the most recent used for Hoxby’s crosswalk, reports the obsolete 3755 code for 19 districts. Hoxby corrects six of these in Johnson County, Kansas, but retains the incorrect code for 13 districts in other counties.² Another common occurrence is a metropolitan district for which the CCD is simply missing an MSA code. One such district is the Collier County School District, which serves the entire Naples, Florida MSA (5345); as a result of this miscoding, the Naples MSA does not show up in Hoxby’s analyses. A final sort of error involves MSA codes that are flatly incorrect: The CCD contains MSA codes of 6640, corresponding to the Raleigh-Durham MSA, for several non-metropolitan school districts in Ohio. As a result, four of the 11 districts in the Hoxby/NCES district-level data set that appear to be in Raleigh are actually in Ohio.

In addition to her data, Hoxby provided the Stata program that she wrote to assemble it. This program, entitled “construct.do,” is included on the Hoxby/NCES CD. In examining this program, I discovered several glitches. A few of the most important are enumerated here:

1. Hoxby merges the NELS student-level file to the NELS school-level file, first merging on the base year school ID (*sch_id*), next the first follow-up school ID (*f1sch_id*), and finally the second follow-up ID (*f2sch_id*). The relevant section of her code (lines 25-66³) is reproduced in Table A1. Note that Hoxby does not rename any variables from the school-level data set between merges, nor does she specify the “update” option on her second and third merges (lines 52 and 62). As a result, nothing is changed by these latter two merges: All variables from the “using” data set (the school-level data) already exist on the “master” (student-level) data, and Stata does not alter pre-existing variables on the master data without explicit instruction.⁴

¹ Hoxby writes that “the Common Core files use the 1990 MSA and PMSA codes, even though some files have names that include years that predate 1990. This conclusion is based on inspection of the counties that were not included in the 1980 MSAs and were included in the 1990 MSAs” (Hoxby, 2004b, lines 506-8). In fact, the documentation to the 1987-88 file, from which Hoxby takes many of her MSA codes, includes several pages titled “Alphabetical Listing of Metropolitan Statistical Areas, October 1984...”

(<http://www.nces.ed.gov/ccd/pdf/pau87gen.pdf>), suggesting that the compilers of that file did not intend to be using up-to-date 1990 codes. Inspection of the CCD files indicates that the MSA codes used do not derive from any single set of MSA definitions, but are drawn from several different generations of those definitions.

² Hoxby also reports a “larger streams” count for the 3755 MSA, though the erroneous CCD crosswalk should not have been necessary for the construction of this variable.

³ All line numbers refer to lines in the “construct.do” program distributed with the September 2, 2004 generation of the Hoxby data.

⁴ “[U]pdate varies the action merge takes when an observation is matched. By default, the master dataset is held inviolate—values from the master dataset are retained when variables are found in both datasets. If

2. Another problem appears in the same segment of code. Students and schools that entered the NELS sample in one of the follow-up surveys (11.5% of student observations and 57% of schools) have *sch_id* set to missing. Observations with missing *sch_id* are not excluded from either the master or the using data sets before the first merge is performed on line 42. Stata's merge command does not make an exception for missing values, so student and school observations with missing *sch_id* are assumed to correspond. Because there are many such observations on either side of the merge, Stata handles them in order, assigning the first such student to the first such school, the second to the second, and so on.⁵ Thus, the resulting data set is incorrect—students with missing *sch_id* should not be assigned to any school on this merge—and its precise form depends on the order in which observations with missing *sch_id* appear in the two data sets. As it happens, each data set arrives at this point having been sorted on *sch_id* (lines 29 and 34), without explicit indication of how to sort observations within *sch_id* groups. Stata breaks ties randomly in sorts, so the order of the relevant observations is random, and indeed is different each time the program is run.⁶
3. A third, related glitch only affects the MSA-level data set, which Hoxby uses to estimate the MSA-level version of the first-stage. Hoxby's student-level models include dummy variables for the nine Census divisions, assigned on the basis of the school's location. The MSA-level analogue of this would use division variables that were not dummies, but which represented the fraction of enrollment in the MSA in each division (for those MSAs which span two or more divisions). Hoxby's computer program does not do this. Rather, it assigns a single division to each MSA, hand-coded for MSAs spanning divisions using the division in which the plurality of the MSA's enrollment is located. There are some MSAs, however, which are wholly contained within a single division but which are incorrectly assigned districts in other divisions by the erroneous CCD district-MSA crosswalk. For example, as noted earlier, four Ohio districts are assigned to the Raleigh-Durham metropolitan area. For MSAs that do not span divisions, division codes are assigned based on the first district in each MSA (line 2370), using a district data set sorted by the MSA code (*fipmsa*, line 2369). Again, Stata breaks ties randomly in sorts, so which district happens to come up first within the Raleigh MSA is different each time the program is run, and on 36% (4/11) of iterations Raleigh is assigned to the East North Central division (containing Ohio).⁷

update is specified, however, the values from the using dataset are retained in cases where the master dataset contains missing," (StataCorp, 2003, "Reference G-M," p. 435). Hoxby apparently also intended update-style merges in lines 527-543.

⁵ "When one is performing a match merge, the master and/or using datasets may have multiple observations with the same *varlist* value. These multiple observations are joined sequentially, as in a one-to-one merge. If the datasets have an unequal number of observations with the same *varlist* value, the last such observation in the *shorter* dataset is replicated until the number of observations is equal," (StataCorp, 2003, "Reference G-M," p. 446). Had Hoxby specified the *unique* or *uniquing* options on her line 46, Stata would have refused to perform the match.

⁶ "[S]table" specifies that observations with the same values of the variables in *varlist* are to keep their same relative order in the sorted data as they had previously....Without the *stable* option, the ordering of observations with equal values of *varlist* is randomized" (StataCorp, 2003, "Reference S-Z," p. 88).

⁷ The Hoxby/NCES dataset happens to be one in which Raleigh is assigned to the wrong division.

4. MSA demographic characteristics are computed by summing observations on each district (from the SDDDB) within the MSA. This creates two problems. First, any errors or omissions in the district-MSA crosswalk are reflected in measured MSA characteristics. Second, many people are double-counted: Any person living in an area served by separate elementary and secondary school districts will be counted twice toward the MSA totals.

To gauge the extent of the randomness introduced by errors 2 and 3, I executed Hoxby's program 10,000 times without alteration, tabulating the estimated choice effects (on the 12th grade reading sample) from each iteration. The distribution of estimates is displayed in Figure A1. The distribution is reasonably concentrated around its median (5.41), with a standard deviation of 0.47. However, the range is quite broad: The smallest estimated effect is 2.18, and the largest is 8.15.

2. Replication data sets

Given the above concerns, a complete replication required re-creating the analysis data set. My first step was to create a correct district-MSA crosswalk. Outside of New England, MSAs consist of whole counties. I used the CCD's county codes, which appear to be more reliable than its MSA codes, to assign districts to MSAs. In New England, town boundaries define MSAs and a district's county location is not sufficient to assign it to an MSA. I built the New England portion of the crosswalk by hand, examining each district's name and mailing address and using these to assign each to a town (and therefore an MSA). I have made available the code needed to produce this crosswalk from the public-use CCD data. While I cannot guarantee that the New England portion is completely free of errors and misjudgments, I believe that it is generally accurate.

I created two replication samples, using my repaired district-MSA crosswalk for each. My first replication sample follows Hoxby's algorithm (as expressed in her code) as closely as possible, but for the substitution of the improved crosswalk and the repair of the errors noted above. Repairing #4 required an alternate source of metropolitan demographic characteristics. I used the Summary Tape File 3A from the 1990 census⁸ for this purpose, with one exception that is noted below.

Repairing #1 was also not straightforward. Each school in the NELS school file has up to six codes identifying its district (an "NCES" code and a "QED" code in each of three survey waves). Moreover, each student is assigned to up to three schools, one in each wave.

⁸ This is based on the same raw 1990 long form census data as the SDDDB, but avoids double-counting.

As a result, there are potentially 18 distinct district assignments for each student. By virtue of glitch #1, however, Hoxby observes only the six first-wave codes for each student. I used the following algorithm to assign students to schools:

1. For each school in each wave, use the NCES code in place of the QED code if both are valid district codes; otherwise, use the valid one. (Lines 369-374 of “construct.do” advocate resolving discrepancies between NCES and QED codes in favor of the NCES codes.) This narrows the 18 codes to 9.
2. Match each student to his/her school in each wave, and keep only that wave’s district code from that school. This leaves a maximum of 3 district codes for each student: The BY code from the BY school, the F1 code from the F1 school, and the F2 code from the F2 school.
3. Follow the “majority rule” algorithm implemented in lines 416-491 of “construct.do” to resolve discrepancies among these three district codes, preferring districts appearing in two waves to one and districts appearing earlier to later. The “construct.do” algorithm neglects a few possibilities—e.g. students who attend different districts in two of the waves and are missing a district assignment in the third wave—but I applied the same rule for these.
4. Following lines 432-446 of “construct.do,” assign students who are as yet unassigned to the modal district among successfully-assigned students who attend the same school as them in the BY, F1, or F2 waves. This has the effect of assigning some students who attended private schools in all waves to the public school that some of their private school classmates in, say, the BY wave transferred to in F1. (I omitted private school students from all analyses of this replication sample, however. As a result, this step is not particularly consequential.)

In constructing my second replication, I deviated from Hoxby’s algorithm when what seemed to me more sensible options were available. I refer to this sample as the “preferred replication” sample. The most important difference is in the way that students are assigned to districts. In an effort to prioritize accuracy over maximizing the sample size, I used only contemporaneous information for this purpose, allowing individual respondents’ district assignments to vary between waves. Thus, a student whose F2 school is missing a district code in the F2 wave was omitted from my analysis sample for F2 scores. (Hoxby’s algorithm would use information about the 8th grade school if it were available, and indeed

for all students would use the same district assignment in each wave.) In practice, this meant following steps 1 and 2 of the above algorithm, then constructing distinct student-district matches for each survey wave using only the district code for the school attended in that wave. There are also important differences in the construction of individual variables, as detailed below.

In Section IV, I also extend the preferred replication sample to include private school students. An alternative algorithm is required for these students, as district codes are unavailable for NELS private schools. I make use of information about the demographic characteristics of the zip codes in which schools are located. These variables are drawn from the zip code tabulation of the 1990 Census, so these data (the “STF-3B” file) can be used to assign zip codes to each NELS school. The vast majority of zip codes are contained entirely within a single county; for those that are not, I assign the zip code to the county containing the plurality of its population. The resulting crosswalk from NELS schools to counties is sufficient to assign schools to MSAs outside of New England. In New England, however, counties do not map to unique MSAs. Using mapping software, cartographic boundary files from the 1990 Census, and a third-party vendor’s zip code boundary files (ESRI, 2002), I compute the fraction of each zip code’s land area contained within each MSA, and assign each zip code to the MSA containing the plurality of its area. (Once again, for the vast majority of zip codes there are no ambiguities.) There are a very small number of zip codes containing NELS schools that cannot be assigned in this way, which I code by hand.

Note that this algorithm assigns each NELS school, public and private, to an MSA, but not to a district. For analyses of these data, I exclude the district-level covariates from Hoxby’s specification. For consistency, I use the zip-code based assignments for both public and private schools, even though alternative district-based assignments (which agree in every case where both are defined) are available for the public schools.

All code for the assembly of my replication samples and for estimation of the models presented has been made available in my “makedata_msadist.do” and “makedata_nels.do” program files. Unfortunately, I cannot make the individual-level data available in the same manner, as they derive from a restricted-access version of the NELS data and are available only from NCES and only to licensed researchers. However, researchers with the appropriate licenses should be able to run my computer programs to extract my samples from the restricted-access NELS data. There are a very few lines of code which have been

redacted from my computer programs, to avoid compromising the confidentiality of the NELS data. I have asked NCES to distribute this code—in a file, “confidential.include.zipcode.do,” that is called from the programs that I do distribute—to licensed researchers who request it. In addition, I describe here the covariates included in Hoxby’s specification and how they are constructed in each replication sample.

3. Comparisons of individual covariates

a. Covariate construction, individual level

I begin with the construction of individual-level variables in the NELS data, as these are the most straightforward. For each item, I describe any discrepancies between the Hoxby/NCES construction and my preferred construction. In each case, I followed the Hoxby/NCES algorithm for the “close replication” sample (variable names are prefixed “*ch_*”) and my preferred algorithm for the “preferred replication” sample (prefix “*jr_*”).

- Indicators for attending a public school in each of the three waves (*pubschby*, *pubschf1*, *pubschf2*): Unless otherwise noted, all analyses include only students enrolled in public school in the wave from which the test score data are taken. There are several NELS variables describing the school sector in each wave. In 40 cases, these are not mutually consistent for the base year. Hoxby does not resolve these inconsistencies; I drop these observations from both replication samples for analyses of base year scores.
- The log of family income (*ch_lnfaminc* and *jr_lnfaminc*): The NELS data report family income in bins. Hoxby assigns each family the log of the midpoint of the relevant bin (in thousands). She assigns a “midpoint” of \$800 for the \$1-\$1,000 bin and one of \$220,000 for the \$200,000+ bin. Observations for which the family income is reported as zero are set to missing. *ch_lnfaminc* follows this algorithm. I use a slightly different construction for *jr_lnfaminc*: I assign each family the log of the geometric average of the endpoints of their bin. (That is, a family in the \$1,000-\$3,000 bin is assigned $\ln(2)$ by Hoxby’s algorithm, and $(\ln(1000)+\ln(3000))/2$ in mine.) I assign $\ln(500)$ to families in the \$1-\$1,000 bin, $\ln(250,000)$ for families with incomes above \$200,000, and $\ln(1)$ to families with zero income. Finally, I use the 2nd follow-up survey’s income variable (*f2faminc*) to assign values for students with missing incomes in the base survey variable (*byfaminc*); Hoxby uses only the latter variable.

- Student race (*ch_asian*, *ch_hispanic*, *ch_black*, *jr_asian*, *jr_hispanic*, and *jr_black*): Hoxby uses only the base-year *race* variable to assign these, and sets each to zero for students with *race*=8 (race missing or more than one race reported). Each dummy is set to missing if there is no value for the *race* variable. I supplement this variable with analogues from the follow-up surveys (*f1race*, *f2race1*, *f2rrace1*) when it does not resolve the student's race, and I set the indicators to missing if none of the four survey variables indicate a specific race.
- Student gender (*ch_female*, *jr_female*): Again, Hoxby uses only the base wave *sex* variable, and sets *female* to missing if this variable is missing; I supplement the base wave with analogous variables from the follow-ups (*f1sex*, *f2sex*, *f2rsex*).⁹
- Parental education (*ch_parscol*, *ch_parcolg*, *jr_parscol*, *jr_parcolg*): Hoxby uses the *bys34a* and *bys34b* variables to assign each parent's education, then uses the highest of these to assign her variable. *bys34a* and *bys34b* are student reports of their parents' education in the base year survey. I use instead *bypared*, *f1pared*, and *f2pared*. These use parent reports where they are available, and student reports only when they are not. There are many discrepancies between these variables and the student reports; it seems likely that the parent reports are more accurate.
- Test scores: Hoxby's preferred specification takes the 12th grade reading score, *f22xrstd*, as the dependent variable, though she also reports results for 8th grade reading (*by2xrstd*) and 10th grade mathematics (*f12xmstd*). Each score is normalized to have a mean around 50 and standard deviation around 10. In the main text, I analyze *f22xrstd* and *by2xrstd*, the 12th and 8th grade reading scores. I also present estimates in this appendix for four additional scores: Mathematics in all three waves, and reading in 10th grade (*by2xmstd*, *f12xmstd*, *f22xmstd*, and *f12xrstd*).

Summary statistics for each of these variables are reported in Table A2, for Hoxby's data set and for each of the two replication samples. The rightmost columns of the Table report correlation coefficients between the different samples, computed pairwise over observations that have values in each of two samples. Note that most of these correlations are almost

⁹ One might not expect many missing values for such a basic demographic characteristic. In fact, it is missing for over 10% of observations. These are students who were brought into the NELS sample via "freshenings" in the first and second follow-up surveys. By using only base-year variables, Hoxby excludes freshened observations. Her 12th grade sample is thus not representative of 12th graders in 1992, but only of 1988 8th graders who were in 12th grade in 1992.

exactly one (with the notable exception of the parental education variables). The primary differences between samples are in the number of observations with missing values—note, for example, that the standard deviation of $\ln(\text{family income})$ is higher in the preferred replication sample than in either of the other samples, where it is defined for 3,000 fewer observations.

b. Covariate construction, district level

District covariates are taken from the School District Data Book (SDDB), a tabulation of 1990 Decennial Census data along school district boundaries. There are several extant versions of the SDDB, not all of which are complete.¹⁰ The Hoxby/NCES CD includes a copy of the “Top 100” file (containing most of the variables that are used). There are some necessary variables which are not included in the “Top 100” file. The Hoxby/NCES CD provides extracts of these variables in two separate files. It is not clear what version of the SDDB was used for these, nor does is the code provided that was used to perform the extraction. The supplementary files have somewhat fewer records than does the (presumably complete) “Top 100” file.

In constructing my replication samples, I rely on a version of the SDDB data obtained from the contractor who produced the file for NCES.¹¹ The full data set is too large (1.4 GB zipped) to easily distribute, but I have made available my extraction program and can work with interested readers to help them obtain the data. There do not appear to be major differences between the two versions of the data.¹²

- Racial composition of the district (*ch_d_pop_fas*, *ch_d_pop_fbl*, and *ch_d_pop_fbi*; *jr_d_pop_fas*, *jr_d_pop_fbl*, and *jr_d_pop_fbi*): Hoxby excludes people of “other” race in her computations (so the denominator is the sum of the Hispanic and non-Hispanic black, white, Asian, and American Indian populations). I include the “other race” group in the denominator.

¹⁰ The NBER, for example, has a version acquired from the National Archives at <http://www.nber.org/sddb/>. Documentation on that page reads “Currently we have and have online 156 files, which NARA says is the whole file, or least all they have. They believe there may be records, or parts of records missing from California and Minnesota. We have observed that Minnesota contains 603 undecipherable records.”

¹¹ I am grateful to Cecilia Rouse and Lisa Barrow for providing me with these data.

¹² There are 357 districts that appear in Hoxby’s SDDB extract but not in mine. Hoxby codes only two of these as metropolitan (one apparently incorrectly). Only 11 of the districts in question (and neither of the apparently metropolitan districts) have non-missing enrollment.

- Index of racial homogeneity (*ch_d_berfrace* and *jr_d_berfrace*): This is computed as the sum of each race’s squared population share. As noted above, Hoxby’s population shares exclude “other” race from the denominator, where mine include them. I also include “other” as one of the races over which the sum is taken.
- Educational distribution (*d_ed_scol*, *d_ed_ba*, and *d_berfed*): We use identical algorithms. *d_berfed* is the sum of the squared population shares of the less than high school, high school graduate, some college, and BA+ groups.
- Household income (*d_lnmeanincA*): Hoxby (2000) describes this variable as the “mean of log(household income) in the district.” In fact, she constructs it as the log of the mean household income in the district. I follow this definition.¹³
- Gini coefficient of household income (*ch_d_gini* and *jr_d_gini*): These are constructed from the distribution of district households across 25 income bins. All households in each bin are assumed to have income equal to the midpoint of the bin. Hoxby assigns the bottom bin (\$0-\$5,000) a “midpoint” of \$4,000 and the top bin (\$150,000+) a value of \$175,000. I use \$2500 for the bottom bin. For the top bin, I use a variable describing the total income among families with incomes above \$150,000 (*P81_2*) to construct the actual mean income among families in this bin in the district. Not surprisingly, this has important consequences for the Gini coefficient.
- Index of ethnic homogeneity (*ch_d_berfethn*, *jr_d_berfethn*): This is a modified version of the index of racial homogeneity, mentioned above. Hoxby (2000) does not provide the formula, but cites another paper (Alesina et al., 1999) for its construction. That paper describes homogeneity indices for the Hispanic and white populations, but does not describe how these are to be aggregated into a single ethnic homogeneity index. The formula, as used in the Hoxby/NCES code, is:

¹³ Given the identical construction, one would expect this variable to be perfectly correlated between the Hoxby/NCES data set and the replication data sets (the fact that Hoxby divides average income by \$1,000 before logging notwithstanding). While the correlation is quite high, it is not exactly one. There are two sources of income information on the SDDDB: Variables tabulating the number of households in each of 25 income bins (named *P80*), and variables recording the aggregate income among families with incomes above and below 150,000 (named *P81*). I use the latter variables to construct mean household income, but when I use the former (and Hoxby’s rules, used elsewhere in her code, for assigning households to midpoints of each bin) I replicate her variable to within four decimal places for all but 39 districts. Assuming that the census tabulates each variable correctly, the “aggregate income” approach is almost certainly more accurate than the “bin midpoints” approach.

$$\text{Index} = (\text{FrBl})^2 + (\text{FrAs})^2 + (\text{FrIn})^2 + (\text{FrWh})^{3-\text{WhIndex}} + (\text{FrHi})^{3-\text{HiIndex}} + (\text{FrOt})^2,$$
 where FrBl, FrAs, FrIn, FrWh, FrHi, and FrOt are the population shares black, Asian, American Indian, white, Hispanic, and other race. (As noted above, Hoxby omits the final group). WhIndex and HiIndex are indices of the source-country heterogeneity of the white and Hispanic populations, respectively. The white index is the sum of the squared shares of the white population that are British, Scandinavian, Russian, other Eastern European, Belgian/Dutch, Swiss/Austrian, French, Arab, German, Greek, Hungarian, Irish, Italian, Polish, Portuguese, or Other. Hoxby counts Canadians as British (but French Canadians as French), and includes sub-Saharan African, “USA,” West Indian, and unclassified ancestry in “Other.” I assign both Canadians and French Canadians to the “Other” category, and exclude sub-Saharan Africa, “USA,” West India, and unclassified from the computation on the grounds that these groups are unlikely to be of white race. Hoxby’s Hispanic index is the sum of the squared shares of the Hispanic population with Mexican, Puerto Rican, Cuban, other Central American, or South American ancestry. I add a sixth category, “Other Hispanic;” Hoxby folds this category into the South American group.

Summary statistics for each of these variables are reported in Table A3. All correlate highly across samples, with the only visible deviations appearing in the income variables.

c. Covariate construction, metropolitan level

Hoxby uses the Office of Management and Budget’s June 30, 1990 definitions of Metropolitan Statistical Areas (<http://www.census.gov/population/estimates/metro-city/90mfips.txt>). In larger agglomerations, she uses PMSAs as the metropolitan construct in place of the larger CMSAs.

Hoxby writes, “I derive demographic measures at the metropolitan-area level from the *City and County Data Book*” (2000, p. 1221-2). In fact, her code derives MSA-level demographic characteristics by summing the SDDB observations for all districts attributed to the MSA. As noted previously, this introduces two problems. First, some areas are served by overlapping elementary and secondary school districts; Hoxby’s approach double-counts residents of these areas toward metropolitan averages. Second, there are errors in the crosswalk file that Hoxby uses to assign districts to metropolitan areas. This can lead to

serious misstatement of MSA totals. To take one extreme example, the SDDDB contains records for many subdistricts within the single school district serving the state of Hawaii, while the CCD contains only a single record. As a result, only the administrative headquarters for the Hawaii district is assigned to Honolulu, and the Honolulu MSA is recorded as having only 13,700 residents in 0.77 square miles and. (The true values are 836,000 residents and 600 square miles.)

In place of the aggregated SDDDB data, I use the Summary Tape File 3A of the 1990 Census to construct MSA demographics for both replication samples. I use county-level records for non-New England MSAs and town-level (“county subdivision”) records for MSAs in New England. The source data for the STF-3A file are the same Census long form data as those from which the SDDDB is constructed, but my approach avoids the imperfect match to MSAs and the double-counting problems that arise with the SDDDB data.

- Metropolitan population and land area (*m_pop_n* and *m_arealand*): These are straightforward, but for the errors introduced when they are computed by summing school districts. The Hoxby/NELS data overstate the population by at least 10% in 71 MSAs, largely due to double-counting the populations of overlapping districts.
- Mean income (*ch_m_avglmmeaninc* and *m_lmmeanincA*): Though Hoxby’s paper describes one of her control variables as “Mean of log(income) of metropolitan area,” she does not make clear that this mean is taken over school districts rather than over individuals. That is, she computes the mean income (in levels) for each district, computes the log of this mean, then averages the district log(mean income) across districts in the MSA (weighting by the number of households) to form her MSA variable. This construction cannot be performed using the STF-3A data, and for this one MSA-level variable I follow Hoxby (for the close replication sample) in deriving it from the SDDDB, after correctly assigning districts to MSAs. My preferred replication sample uses the more straightforward log of the MSA mean income.
- Gini coefficient (*ch_m_gini*, *jr_m_gini*): These are constructed identically to the district-level variables, as described above, using the STF-3A data in place of the SDDDB.
- Racial composition of the MSA (*ch_m_pop_fas*, *ch_m_pop_fbl*, and *ch_m_pop_fbi*; *jr_m_pop_fas*, *jr_m_pop_fbl*, and *jr_m_pop_fbi*). Again, the MSA-level construction is the same as the district-level.

- Educational distribution (m_ed_scol , m_ed_ba , and m_herfed): Once more, the same as at the district level.
- Census division (cb_d_div1 - cb_d_div9 , jr_m_divis1 - jr_m_divis9): Hoxby treats division as a district-level characteristic, computing dummy variables based on the division in which the district is located. She uses a somewhat different construction for her MSA-level first stage analysis, hand-coding each MSA that straddles divisions to the division in which the plurality of its population resides. My close replication sample follows Hoxby in each construction. Note, however, that my use of a repaired district-MSA crosswalk avoids the problems (discussed above) that appear in Hoxby’s implementation. For my preferred sample, I treat division as an exclusively metropolitan-level characteristic, and I assign multi-division MSAs fractional values of the division indicators corresponding to the fraction of the MSA population in each division.

Summary statistics for each of these variables are reported in Table A4. They diverge more across samples than did individual or district-level variables, largely because of the differences between what is obtained from aggregated, overlapping SDDDB data and from aggregated Census STF records that do not overlap.

4. Choice index

Hoxby’s choice index is $c_m = 1 - \sum_{jm} (n_{jm} / N_m)^2$, where n_{jm} is the enrollment of district j in MSA m and N_m is the total enrollment in the MSA. Enrollment is drawn from the SDDDB “Top 100” file, which in turn draws the variable from the CCD. For my replication samples, I use the 1989-90 CCD enrollment data directly.

Although I follow Hoxby’s construction exactly for the close replication sample, I make a slight alteration for the preferred replication sample. Where Hoxby constructs her index using enrollment in all grades, I compute it considering only grade-8 enrollment. This makes little difference in MSAs with unified school districts. Where there are separate elementary and secondary districts—which cannot be said to compete against each other—Hoxby’s formula will indicate more competition than actually exists (Urquiola, 1999).

5. Instruments

As instruments for the degree of choice in the MSA, Hoxby uses the number of larger and smaller rivers flowing through the MSA. She describes their construction as follows:

The streams variables are derived from the U.S. Geological Survey's (USGS) 1/24,000 quadrangle maps. It was by using these extremely detailed maps—which allow the viewer to identify even very small streams, buildings, and boundaries—that I initially recognized the relationship between natural barriers and school district boundaries. The measurement of the streams variable was in two stages. Using the physical maps, I first counted all streams that were at least 3.5 miles long and of a certain width on the map. These data were checked against the Geological Survey's *Geographic Names Information System* (GNIS) for accuracy. I derived smaller streams directly from the GNIS. I employ two stream variables: the number of larger streams (measured by hand and often traversing multiple districts, sometimes multiple counties) and the number of smaller streams (from GNIS). (Hoxby, 2000, p. 1222).

Given the vagueness of Hoxby's description of "larger streams" and the subjectivity of the hand measurement, I opted not to try to reproduce Hoxby's counts. This has one unfortunate consequence: Hoxby does not provide counts for all MSAs, and as a result some MSAs must be excluded from analyses that include her larger streams instrument.

Hoxby's description of the genesis of the smaller streams measure is incomplete. She in fact measures *total* streams in the GNIS, and constructs smaller streams as total minus larger streams. As a result, any inaccuracies in the larger streams measure appear (with the opposite sign) in the smaller streams variable. Note, however, that the instrument set can be equivalently formulated as total and larger streams, as these span the same space as do Hoxby's measures.

The Hoxby/NCES program (2004b, lines 2811-2) cites the USGS web site as the source of her GNIS data, with a date of 2004. The GNIS is continually updated as better maps and surveys are completed. It is not clear what might be the impact of changes made between Hoxby's original analysis and her later extraction of the current GNIS data.

In the published paper, Hoxby writes that the GNIS “provides the longitude and latitude of [smaller streams] origin and destination” (2000, note 24, p. 1222). Her code, however, does not make use of these variables, except to construct the streams’ lengths. (She discards streams shorter than one mile.) Another variable, describing the county in which the stream’s “destination” (i.e. mouth) is located, is used instead. The CCD is used to associate counties with MSAs, which means that some areas are mis-assigned.¹⁴ Figure A2 indicates the implications of Hoxby’s assignment algorithm for the Mississippi River. The Mississippi’s mouth is in Plaquemines Parish, Louisiana, which is non-metropolitan. As a result, Hoxby’s algorithm does not include the Mississippi in the total streams count of any of the eight MSAs through which it flows. Of course, the Mississippi is almost certainly included as a “larger” stream, but its exclusion from the total streams count means that smaller streams are necessarily undercounted.

When I reproduce the total streams variable for the replication sample, I use the same GNIS data as Hoxby, continuing to assign streams to counties on the basis of their destinations. I use an accurate county-MSA crosswalk, however, and compute population shares from Census STF data on non-overlapping towns (“county subdivisions”). Following Hoxby, I compute smaller streams by subtracting the number of larger streams that Hoxby counted from the number of total streams that I obtain. This produces a negative number in five MSAs (two in New England), where Hoxby evidently counted more larger streams than the number of total streams (stream mouths) than the GNIS destination variable indicates.¹⁵

The alternative instruments discussed in Section III use a more expansive definition of “total streams,” in which a stream is counted toward any MSA through which it flows, regardless of where it terminates. For this purpose, I use an alternative version of the GNIS database (Geographic Names Office, 1999), which contains a field listing all counties through which each stream flows. (For the Mississippi, this list contains 117 counties, which are shown on Appendix Figure A2.) In New England, where some counties are split into non-metropolitan and metropolitan components, I compute weights for each county

¹⁴ When counties are not wholly contained within a single MSA (in New England, and elsewhere when districts are mis-assigned), all streams in each county are assigned to the MSA in which the plurality of the county’s population resides. The population shares for this computation are taken from the SDDDB, so double-count some areas served by multiple districts.

¹⁵ Note that there are no negative “smaller streams” counts in the Hoxby/NCES data; the problem arises only when I re-create Hoxby’s total streams variable with a corrected county-MSA crosswalk in place of Hoxby’s imperfect district-MSA crosswalk.

corresponding to the share of the population in the county residing in each MSA. Using the county population weights, streams are partially assigned to each MSA with which a county intersects.¹⁶ Outside of New England, I need not approximate MSA boundaries, and each stream gets a weight of one toward each MSA through which it flows. My “total streams” count is the sum of the weights of all streams flowing through each MSA.

I consider three categorizations of the resulting total streams count into larger and smaller streams. First, I use Hoxby’s larger streams variable, defining smaller streams as my total streams count minus Hoxby’s larger streams for those MSAs for which Hoxby’s variable is available. There are six MSAs, five in New England, for which this indicates a negative number of smaller streams.¹⁷ Second, I compute separate counts of inter-county and intra-county streams, where an inter-county stream is one that flows through more than one county. Finally, using the latitude and longitude variables (which exist on both versions of the GNIS data), I compute the distance between the source and mouth of each stream, and compute separate counts of streams that are longer than and shorter than 3.5 miles.¹⁸

Table A5 presents summary statistics for the various streams variables

6. Weights

Hoxby writes that her student-level regressions are “weighted so that each metropolitan area receives equal weight,” (2000, p. 1226). This is operationalized by summing the NELS sample weights for all students in each MSA, then dividing each individual’s weight by the MSA’s total. In the Hoxby/NCES data, however, this normalization is carried out once for each survey wave, without regard to missing values. Many observations that are excluded from the regression analyses because they are missing a test score or one of the covariates are nevertheless included in the weight normalization. As a result, the weights sum to substantially less than one in each MSA when the sum is taken only over observations that appear in the regressions. In the replication samples, I re-

¹⁶ For streams flowing through multiple counties in the same MSA, the maximum population weight is used. (This is analogous to assigning a weight of one to any stream flowing through any metropolitan county in the remainder of the country.)

¹⁷ The non-New England MSA is Topeka, Kansas. Hoxby reports 82 larger streams, but I count only 36 stream mouths and 41 total streams. Hoxby’s count of 109 total streams appears to derive from the CCD’s assignment of Jefferson County and Osage County school districts to the Topeka MSA. These counties were part of the MSA in 1981, but have been excluded since the release of the 1983 definitions.

¹⁸ There are a small number of streams missing latitude and longitude information, for which length cannot be computed. These are included in earlier counts, but are excluded from the categorization by length.

normalize for each specification, including only observations used in that specification. Note (Table A2) that the average weight is substantially smaller in the Hoxby/NCES sample than in the replication samples, and that the product of the average weight and the number of observations in the analysis samples equals the number of MSAs for the replication samples but not in the Hoxby/NCES sample.

Appendix B: Econometric Specification and Standard Error Computation

Simplifying notation slightly, Hoxby specifies her model for student achievement as

$$(B1) \quad A_{ikm} = C_m \beta_1 + X_{ikm} \beta_2 + \bar{X}_{km} \beta_3 + \bar{X}_m \beta_4 + e_{ikm},$$

where A_{ikm} is the test score of student i in district k in MSA m ; C_m is the choice index for MSA m ; and X_{ikm} , \bar{X}_{km} , and \bar{X}_m are individual, district mean, and MSA mean covariates.¹⁹

As there are likely omitted variables at all three levels, it is unreasonable to assume that e_{ikm} is independent across individuals within the same district or MSA. Hoxby specifies an error components model,

$$(B2) \quad e_{ikm} = \varepsilon_m + \varepsilon_{km} + \varepsilon_{ikm},$$

assuming that each component is independent and identically distributed across observations at that level. That is,

$$(B3) \quad \text{cov}(e_{ikm}, e_{jdn}) = \begin{cases} 0 & \text{if } m \neq n \\ \sigma_m^2 & \text{if } m = n, k \neq d \\ \sigma_m^2 + \sigma_k^2 & \text{if } m = n, k = d, i \neq j \\ \sigma_m^2 + \sigma_k^2 + \sigma_i^2 & \text{if } m = n, k = d, i = j \end{cases}$$

where $\sigma_m^2 = \text{var}(\varepsilon_m)$, $\sigma_k^2 = \text{var}(\varepsilon_{km})$, and $\sigma_i^2 = \text{var}(\varepsilon_{ikm})$.

¹⁹ Hoxby (2000) devotes considerable discussion to the complication that the district heterogeneity variables do not average to their MSA-level analogues. I neglect this complication; one may simply think that \bar{X}_{km} and \bar{X}_m are arbitrary vectors of district- and MSA-level covariates. One related point is worth mentioning, however. In her footnote 21, Hoxby claims that when the district heterogeneity variables are excluded, the individual-level first stage model (analogous to the outcome equation indicated above) is identical to an MSA-level analogue that excludes X_{ikm} and \bar{X}_{km} , as there can be no correlation between these two vectors and the MSA-level choice variable conditional on \bar{X}_m . This is correct only so long as the sample average of the individual and district-level variables is identical to the MSA means. When, as in this case, the individual variables are observed only for a sample, this will not be true, and there is no guarantee that coefficients on MSA-level covariates are invariant to the level at which the model is estimated.

It is useful to work with matrix notation. Let $W = [C_m \ X_{ikm} \ \bar{X}_{km} \ \bar{X}_{ikm}]$ be the matrix of right-hand-side variables from (B1), and $\beta = [\beta_1' \ \beta_2' \ \beta_3' \ \beta_4']$. (B1) then becomes $A = W\beta + e$. If C_m is endogenous, $E[W'e] \neq 0$. Let R_m be the streams instrumental variables, and let $Z = [R_m \ X_{ikm} \ \bar{X}_{km} \ \bar{X}_{ikm}]$. If the instruments are valid, $E[Z'e] = 0$, and

$$(B4) \quad \hat{\beta}_{IV} = (W'P_ZW)^{-1}W'P_ZA = \beta + (W'P_ZW)^{-1}W'P_Ze,$$

where $P_Z = Z(Z'Z)^{-1}Z'$.

Inference proceeds by noting that

$$(B5) \quad \text{var}(\hat{\beta}_{IV}) = (W'P_ZW)^{-1}W'P_Z'\Gamma P_ZW(W'P_ZW)^{-1},$$

where $\Gamma = E[ee']$ is the variance-covariance matrix with elements specified in (B3). Under conventional assumptions of i.i.d. observations (i.e. with $\sigma_m^2 = \sigma_k^2 = 0$), $\Gamma = \sigma_i^2 I$, and we obtain the conventional variance expression, $\text{var}(\hat{\beta}_{IV}) = \sigma_i^2 (W'P_ZW)^{-1}$. This expression does not apply, however, with non-zero district- and MSA-level error components.

Hoxby's so-called "Moulton" formula for $\text{var}(\hat{\beta}_{IV})$ (developed by Moulton, 1986, for the OLS case) simply forms a consistent estimate of Γ from estimates of the error component variances, σ_m^2 , σ_k^2 , and σ_i^2 . The resulting $\hat{\Gamma}$ is then plugged in for Γ in (B5) to estimate $\text{var}(\hat{\beta}_{IV})$. This approach requires only that the three error component variances be estimated consistently. Several consistent estimators are available, and it is not clear that any one should be preferred to another. The options multiply in the "unbalanced panel" case that is relevant here, where districts and MSAs contain unequal numbers of observations. Hoxby has not provided code for her implementation, and does not specify how she estimates the variance parameters.

My implementation of the Moulton approach estimates the error component variances from contrasts between individual, district-mean, and MSA-mean residual variances, extending Greene's (2000, p. 570-2) discussion of random effects in unbalanced panels to the three-component hierarchy considered here. Let $\bar{e}_{km} = \frac{1}{N_{km}} e_{ikm}$ and

$\bar{e}_m = \frac{1}{N_m} e_{ikm}$ be the district- and MSA-level means of e_{ikm} , where N_{km} is the number of

individual observations at district k in MSA m and $N_m = \sum_k N_{km}$ is the number of individuals in the MSA. Ignoring degrees of freedom corrections, which only complicate the notation,

$$(B6a) \quad \text{var}(e_{ikm} - \bar{e}_{km}) = \text{var}(\varepsilon_{ikm}) = \sigma_i^2;$$

$$(B6b) \quad \text{var}(\bar{e}_{km} - \bar{e}_m) = \text{var}\left(\varepsilon_{km} + \frac{1}{N_{km}}\varepsilon_{ikm}\right) = \sigma_k^2 + \frac{1}{N_{km}}\sigma_i^2; \text{ and}$$

$$(B6c) \quad \text{var}(\bar{e}_m) = \text{var}\left(\varepsilon_m + \frac{1}{P_m}\varepsilon_{km} + \frac{1}{N_m}\varepsilon_{ikm}\right) = \sigma_m^2 + \frac{1}{P_m}\sigma_k^2 + \frac{1}{N_m}\sigma_i^2,$$

where P_m is the number of districts in MSA m . Note that with unbalanced data, N_{km} can vary across districts, as can P_m and N_m across MSAs. Solving the sample analogues of (B6a), (B6b), and (B6c) for the variance components thus involves the sample averages of N_{km}^{-1} , N_m^{-1} , and P_m^{-1} , but these are readily estimated from the data. Code for my implementation may be downloaded from my web site.

I also present so-called ‘‘clustered’’ standard error estimates. These proceed from the observation that we do not require a consistent estimate of Γ , but only of $Z'\Gamma Z$. This can be consistently estimated under substantially weaker assumptions than (B3). In particular, we can allow for arbitrary heteroskedasticity and within-MSA correlation patterns. Let Z^m be the sub-matrix of Z corresponding to MSA m , and let e^m be the analogous subvector of e . So long as $\text{cov}(e_{ikm}, e_{jdn}) = 0$ whenever $m \neq n$, $G = \sum_m Z^m e^m e^{m'} Z^m$ is consistent for $Z'\Gamma Z$. The cluster variance estimator, then, is

$$(B7) \quad \text{var}(\hat{\beta}_{IV}) = (W' P_Z W)^{-1} W' Z (Z' Z)^{-1} G (Z' Z)^{-1} Z' W (W' P_Z W)^{-1}.$$

As with the Moulton estimator, asymptotic consistency is achieved as the number of MSAs goes to infinity. An important reason to prefer the ‘‘cluster’’ estimator is that it is fully automated, where the Moulton estimator requires the researcher to choose among several estimators of the error component variances and is therefore more difficult to replicate.

1. Exogeneity and overidentification tests

It is useful to test for the endogeneity of the choice index, C_m . With i.i.d. errors, one might use a Hausman test: Under the null hypothesis that C_m is exogenous, the OLS estimator of β is efficient but the IV estimator is consistent; under the alternative of

endogeneity, OLS is inconsistent but IV remains consistent. With the serial correlation implied by the error components model, however, OLS is no longer efficient even under the null hypothesis. Thus, the exogeneity tests reported in the text use an alternative, “artificial regression” test that can be made consistent to clustering (Davidson and MacKinnon, 1993, p. 237-42). I form \hat{C}_m , the fitted values of C_m from the first stage regression—using the same sample used for IV—and estimate the regression A_{ikm} on C_m , \hat{C}_m , X_{ikm} , \bar{X}_{km} , and \bar{X}_m . Under the null hypothesis of exogeneity, the coefficient on \hat{C}_m should be zero; under the alternative, it should be non-zero. I report tests of the significance of this coefficient that use clustered standard errors, although one might equally well use the Moulton standard errors for this purpose (and indeed they give similar results).

In specifications involving multiple instruments, one might also like to report overidentification tests of the mutual consistency of the different instruments. Again, conventional Hausman test-based approaches do not work with serially correlated errors, as the regular two-stage least squares estimator is not efficient. Hoxby and Paserman (1998) propose a Hausman test for the error components model that is based on the GMM estimate of the overidentified specification. Following Hoxby (2000), I present only 2SLS coefficient estimates and not the more efficient GMM estimates. Accordingly, I choose not to present GMM-based overidentification tests, but instead simply report specifications that omit the suspect instruments.²⁰

Appendix C: Full coefficient vectors

Tables C1 through C5 present the full coefficient vectors from the models in Tables 1 through 5 of the main text.

Appendix D: Additional Specifications

The remaining appendix tables present extensions and alternative specifications. Table D1 presents estimates of the first-stage regression (analogous to those in Table 3) using the “preferred” replication sample. As before, Hoxby’s larger streams variable plays a substantially different role in the MSA-level “implied first stage” than in the actual first stages to the individual-level models. Table D2 presents the corresponding IV estimates

²⁰ A second reason to avoid overidentification tests in this context is that they are likely to have little power when, as here, only one of the instruments plays an important role in the first stage.

from the preferred replication sample. As in the “close” sample (Table 4), only specifications that include Hoxby’s larger streams variable as an instrument show any indication that choice is endogenous or that it has a significant positive effect on outcomes. OLS estimates are negative (though far from significant) for both 8th and 12th grade reading scores.

Table D3 extends the sample selection analysis from Table 5 to 8th grade reading scores. Choice coefficients fall somewhat when district-level covariates are dropped from Hoxby’s basic specification, and fall substantially farther when public schools are assigned to MSAs on the basis of their zip codes. There is an additional downward effect of adding private schools to the sample, and *t*-statistics are uniformly below one in this specification.

Table D4 extends the analysis to the four NELS test scores that have not yet been considered. There is no indication that choice has a different effect on mathematics scores or on 10th grade reading scores than it does on the 8th and 12th grade reading scores considered thus far.

Table D5 presents a version of the MSA-level first stage estimated only on MSAs in the NELS sample. Comparing this to Table 3, it is clear that the divergence between Hoxby’s MSA-level first stage (including all MSAs) and the individual-level MSAs derives from the sample of MSAs included rather than from the level at which the model is estimated.

Table D6 presents “two sample” IV estimates of the basic model, in which the first stage is estimated at the MSA level, predicted values are formed, and these predicted values are used in the individual-level second stage regression. In each panel, the first row presents estimates that use only NELS MSAs for the first stage, while the second row presents estimates in which the full set of MSAs are used to estimate the first stage model. This approach introduces two complications. First, the individual- and district-level covariates must be excluded from the MSA-level first stages. This does not affect the IV estimates of the choice coefficient, so long as all first-stage controls are included in the second stage as well, but changes the interpretation of the second-stage coefficients on the control variables (which are not shown in any case). Second, special formulas are needed to compute correct standard errors. I do not compute these, but merely present clustered S.E.s from the second stage regression; these will tend to understate the true variability of the indicated coefficients.

As with conventional IV, the two-sample estimates indicate significant choice effects when Hoxby’s instruments are used and when only the NELS MSAs are included in the first stage. When the first stage sample is broadened to include all MSAs, the choice coefficient falls in the specifications using Hoxby’s instruments but not in those using inter- and intra-county streams as instruments. Even with the incorrect standard errors, none of the specifications indicates a significant choice effect when all MSAs are used to compute the first stage.

Table D7 presents estimates of the reduced form relationship between streams and student test scores. These are quite noisy, and the streams coefficients are never significantly different from zero in models that exclude Hoxby’s larger streams variable. In models including this variable, it always has a negative coefficient, while smaller streams have a significantly positive coefficient. When these variables are added together, however, the “total streams” coefficient is statistically and substantively indistinguishable from zero.

Finally, Table D8 presents IV estimates for MSA average test scores, computed entirely at the MSA level. Very few of these are significant, even when averages of individual- and district-level covariates are included in the specification (in which case the specification is essentially identical to Hoxby’s, but is more conservative about the standard error computation). The two models that are significant both use Hoxby’s streams variables.

Appendix E: Estimates of the choice effect on spending

Much of Hoxby’s discussion, and all of my analysis thus far, focuses on the relationship between district fragmentation and student achievement. However, her conclusions relate to the choice effect on school productivity, and she presents analyses of the relationship between choice and school spending in her Tables 5 and 6. The relevant coefficients from her preferred specification are reproduced in Panel A of Table E1. Panel B reproduces this specification in the Hoxby/NCES data. Standard errors are much larger here, and indeed the published standard error from the IV specification much more closely resembles the Hoxby/NCES classical standard error—which assumes that observations are i.i.d.—than it does either the Moulton or clustered errors from that sample. There is no indication in the Hoxby/NCES sample of a significant IV effect of choice on per-pupil spending, though the OLS estimate is weakly significant. Panels C and D repeat the same specification in the replication samples, and also present estimates that use inter- and intra-

county streams as instruments. Again, there is no indication of a significant effect in IV, but OLS estimates are significantly negative.

Hoxby's model, however, may be misspecified. She estimates it at the school district level, and takes as her dependent variable the log of per pupil spending in the district. The convexity of the log transformation is the source of the problem. Consider an MSA with two equally-sized schools, A and B, with different levels of per-pupil spending, y_A and y_B . If these schools are in the same district, Hoxby's dependent variable will be $\ln(0.5*(y_A + y_B))$. If the schools are divided into two districts, however, the average of Hoxby's dependent variable in the MSA will be $0.5*(\ln(y_A) + \ln(y_B))$, which is smaller. Thus, without any behavioral effect of choice at all, one would estimate a negative coefficient in Hoxby's regression.

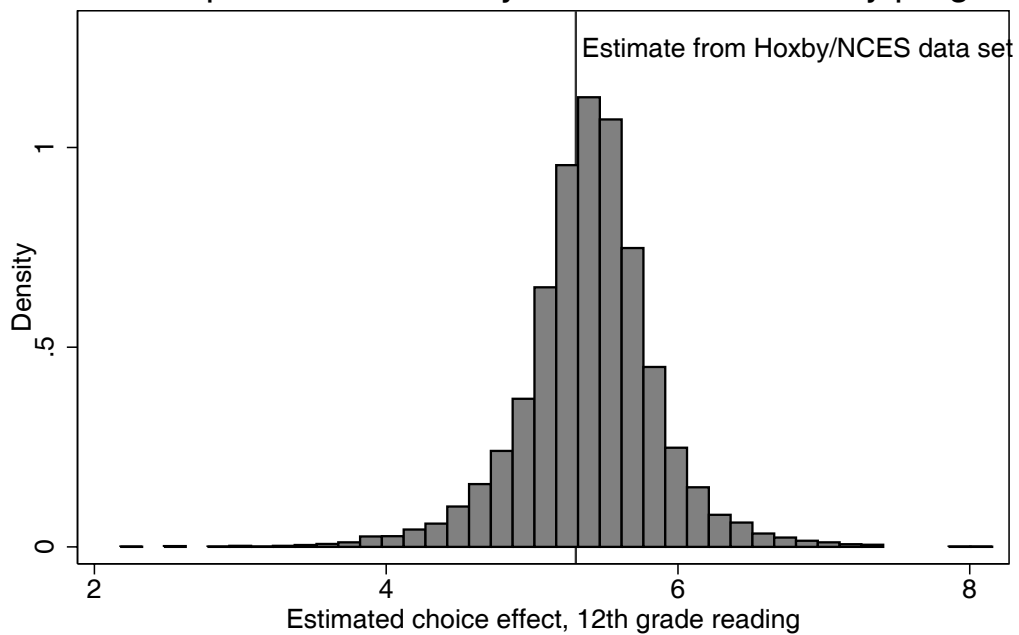
One way to avoid any bias that this mechanical relationship might introduce is to estimate the specification at the MSA level, taking as the dependent variable the log of average per pupil spending in the MSA. This is done in columns 4-6 of Table E1. Evidently, the mechanical bias is not an issue—estimates are nearly identical (with very similar standard errors) to those obtained at the district level).

References

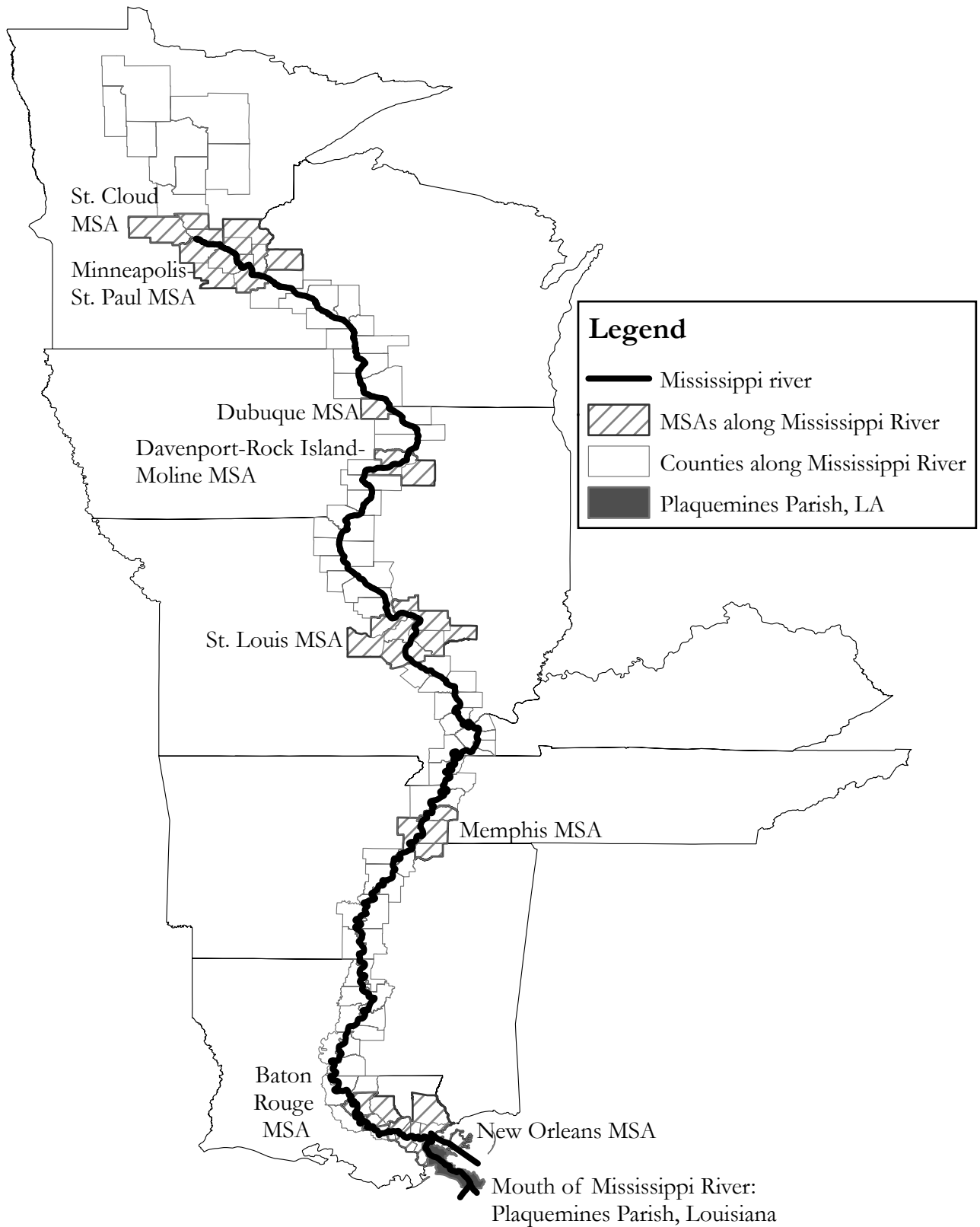
- Alesina, Alberto F.; Baqir, Reza and Hoxby, Caroline M.** "Ethnic Diversity and the Number of Jurisdictions." Working paper, Harvard University, August 1999.
- Davidson, Russell and MacKinnon, James G.** *Estimation and Inference in Econometrics*. New York: Oxford University Press, 1993.
- ESRI.** "Data & Maps 2002." CD, ESRI, 2002.
- Geographic Names Office.** "National Geographic Names Database." GNIS Digital Gazetteer CD, U.S. Geological Survey, 1999.
- Greene, William H.** *Econometric Analysis*. New Jersey: Prentice Hall, 2000.
- Hoxby, Caroline M.** "Does Competition among Public Schools Benefit Students and Taxpayers?" *American Economic Review*, December 2000, 90(5), pp. 1209-38.
- _____. "District-Level and Metropolitan-Area Variables Merged with NELS Data." CD, National Center for Education Statistics, September 2, 2004a.
- _____. "Documentation to 'District-Level and Metropolitan-Area Variables Merged with NELS Data': Construct.Do." National Center for Education Statistics, September 2, 2004b.
- _____. "Documentation to 'District-Level and Metropolitan-Area Variables Merged with NELS Data': Runregressions.Do." National Center for Education Statistics, September 2, 2004c.

- Hoxby, Caroline M. and Paserman, M. Daniele.** "Overidentification Tests with Grouped Data." National Bureau of Economic Research Working paper #T0223, February 1998.
- Moulton, Brent R.** "Random Group Effects and the Precision of Regression Estimates." *Journal of Econometrics*, August 1986, 32(3), pp. 385-97.
- StataCorp.** *Stata Base Reference Manual*. College Station, Texas: Stata Press, 2003.
- Urquiola, Miguel.** "Demand Matters: School District Concentration, Composition, and Educational Expenditure." University of California, Berkeley, Center for Labor Economics Working Paper #14, April 1999.

Appendix Figure A1. Distribution of estimated choice effect across replications of Hoxby/NCES data assembly program



Appendix Figure A2. Alternative algorithms for assigning the Mississippi River to MSAs



Appendix Table A1. Segment of Hoxby's construct.do program (Hoxby, 2004b) that merges NELS student file to NELS school file.

```
25 *Infix data from stmeg.pri and scmeg.pri.;
26
27 cd nels;
28 infix using nels.fromstmeg.infix.dct, using(d:\nels92\stmeg.pri);
29 sort sch_id;
30 compress;
31 save nels.fromstmeg.dta, replace;
32 clear;
33 infix using nels.fromscmeg.infix.dct, using(d:\nels92\scmeg.pri);
34 sort sch_id;
35 compress;
36 save nels.fromscmeg.dta, replace;
37 clear;
38
39 *Merge the variables from the two datasets above.;
40
41 use nels.fromstmeg.dta;
42 merge sch_id using nels.fromscmeg.dta;
43 tab _merge;
44 drop if _merge==2;
45 drop _merge;
46 save nels.dta, replace;
47 use nels.fromscmeg.dta;
48 sort f1sch_id;
49 save nels.fromscmeg.dta, replace;
50 use nels.dta;
51 sort f1sch_id;
52 merge f1sch_id using nels.fromscmeg.dta;
53 tab _merge;
54 drop if _merge==2;
55 drop _merge;
56 save nels.dta, replace;
57 use nels.fromscmeg.dta;
58 sort f2sch_id;
59 save nels.fromscmeg.dta, replace;
60 use nels.dta;
61 sort f2sch_id;
62 merge f2sch_id using nels.fromscmeg.dta;
63 tab _merge;
64 drop if _merge==2;
65 drop _merge;
66 save nels.dta, replace;
```

Appendix Table A2: Summary statistics for individual-level covariates, NELS

	Hoxby/NELS data		Close replication		Preferred replication		Correlations		
	(1)		(2)		(3)		(1) & (2)	(1) & (3)	(2) & (3)
	Mean	S.D.	Mean	S.D.	Mean	S.D.			
Public school (BY)	0.80	(0.40)	0.80	(0.40)	same as close		1.000		
Public school (F2)	0.87	(0.33)	0.87	(0.33)	same as close		1.000		
ln(fam. Income)	3.36	(0.94)	3.36	(0.94)	10.19	(1.21)	1.000	0.999	0.999
Asian	0.06	(0.24)	0.06	(0.24)	0.07	(0.25)	1.000	0.996	0.996
Hispanic	0.13	(0.34)	0.13	(0.34)	0.14	(0.35)	1.000	0.999	0.999
Black	0.12	(0.33)	0.12	(0.33)	0.13	(0.33)	1.000	0.996	0.996
Female	0.50	(0.50)	0.50	(0.50)	0.49	(0.50)	1.000	1.000	1.000
Parents some college	0.23	(0.42)	0.23	(0.42)	0.38	(0.49)	1.000	0.453	0.453
Parents BA+	0.37	(0.48)	0.37	(0.48)	0.30	(0.46)	1.000	0.759	0.759
8th grade reading score	50.5	(10.2)	same as Hoxby/NELS						
12th grade reading score	50.7	(10.0)	same as Hoxby/NELS						
8th gr. weight (*100)	0.74	(1.26)	1.78	(2.04)	1.57	(1.48)	0.743	0.963	0.794
12th gr. weight (*100)	0.79	(3.01)	3.27	(4.97)	2.98	(5.31)	0.640	0.413	0.825
Fr. in MSA									
8th grade	0.56		0.70		0.55				
12th grade	0.53		0.66		0.52				
Analysis sample size									
8th grade reading	10,175		10,429		11,719				
12th grade reading	5,475		5,934		6,688				
Fraction of district assignments that agree									
8th grade							0.65	0.73	0.75
12th grade							0.65	0.41	0.59
Fraction of MSA assignments that agree									
8th grade							0.77	0.94	0.81
12th grade							0.77	0.68	0.66

Notes: All statistics are unweighted. See text for details. N=27,805, though many variables are missing for many observations.

Appendix Table A3: Summary statistics for district-level covariates, SDDB

	Hoxby/NELS data		Close replication		Preferred replication		Correlations		
	(1)		(2)		(3)		(1) & (2)	(1) & (3)	(2) & (3)
	Mean	S.D.	Mean	S.D.	Mean	S.D.			
Fr. Asian	0.02	(0.08)	0.01	(0.02)	0.01	(0.02)	1.000	1.000	1.000
Fr. Hispanic	0.05	(0.12)	0.05	(0.12)	0.05	(0.12)	1.000	1.000	1.000
Fr. Black	0.04	(0.11)	0.04	(0.11)	0.04	(0.11)	1.000	1.000	1.000
Racial homog. index	0.84	(0.17)	0.85	(0.17)	0.85	(0.17)	1.000	1.000	1.000
Fr. some college	0.25	(0.08)	same as Hoxby/NELS						
Fr. BA+	0.15	(0.10)	same as Hoxby/NELS						
Educ. homog. index	0.31	(0.05)	same as Hoxby/NELS						
log(mean HH inc.)	10.37	(0.32)	10.37	(0.34)	same as close		0.986		
Gini coeff.	0.38	(0.05)	0.38	(0.05)	0.40	(0.06)	0.993	0.928	0.934
Ethnic homog. index	0.78	(0.21)	0.79	(0.21)	0.81	(0.21)	1.000	1.000	1.000
Fr. in MSAs	0.39		0.40		same as close				
Fraction of MSA assignments that agree							0.97	0.97	1.00

Notes: All statistics are unweighted. See text for details. N=14,947 in replication samples, 15,304 in Hoxby data.

Appendix Table A4: Summary statistics for MSA-level covariates, SDDDB

	Hoxby/NELS		Close		Preferred		Correlations		
	(1)		(2)		(3)		(1) & (2)	(1) & (3)	(2) & (3)
	Mean	S.D.	Mean	S.D.	Mean	S.D.			
Choice	0.67	(0.27)	0.66	(0.27)	0.66	(0.27)	0.971	0.966	0.995
Population (10 millions)	0.063	(0.112)	0.058	(0.098)	same as close		0.985		
ln(Population)	12.66	(1.12)	12.63	(1.00)	same as close		0.960		
Land Area (100,000s of sq. miles)	0.017	(0.019)	0.017	(0.020)	same as close		0.958		
ln(Land Area)	7.06	(1.01)	7.07	(0.89)	same as close		0.871		
Fr. Asian	0.02	(0.05)	0.02	(0.04)	0.02	(0.04)	0.992	0.992	1.000
Fr. Hispanic	0.07	(0.13)	0.07	(0.13)	0.07	(0.13)	0.999	0.999	1.000
Fr. Black	0.09	(0.09)	0.10	(0.09)	0.10	(0.09)	0.988	0.988	1.000
Racial homog. index	0.72	(0.16)	0.72	(0.16)	0.71	(0.16)	0.992	0.991	1.000
Fr. some college	0.28	(0.05)	0.26	(0.05)	same as close		0.893		
Fr. BA+	0.19	(0.06)	0.20	(0.07)	same as close		0.980		
Educ. homog. index	0.27	(0.02)	0.27	(0.02)	same as close		0.954		
log(mean HH inc.)	3.56	(0.17)	10.48	(0.18)	10.50	(0.19)	0.987	0.976	0.992
Gini coeff.	0.40	(0.02)	0.40	(0.02)	0.43	(0.03)	0.963	0.906	0.937
Ethnic homog. index	0.64	(0.19)	0.63	(0.19)	0.63	(0.19)	0.991	0.991	1.000
# of invalid MSA codes	5								

Notes: N=327 MSAs in Hoxby data, 335 in replication samples. All analyses are unweighted

Table A5: Summary statistics and correlations for streams variables

	Mean	S.D.	Correlation matrix				
			(1)	(2)	(3)	(4)	(5)
<i>Larger streams</i>							
<i>Published</i>	7.9	(14.8)					
(1) Hoxby/NCES data (N=314)	44.3	(64.1)	1.00				
(2) Inter-county streams	41.2	(33.7)	0.50	1.00			
(3) Long streams (>3.5 miles)	69.9	(57.8)	0.55	0.84	1.00		
<i>Total streams</i>							
(1) Hoxby/NCES data (N=314)	128.3	(119.8)	1.00				
(2) Replication stream mouths (N=319)	124.4	(119.2)	0.97	1.00			
(3) Replication total streams	147.8	(149.2)	0.89	0.92	1.00		
<i>Smaller streams</i>							
<i>Published</i>	182.7	(208.8)					
(1) Hoxby/NCES data (N=314)	84.0	(78.0)	1.00				
(2) Replic. total mouths - Hoxby larger (N=304)	80.3	(82.6)	0.94	1.00			
(3) Replic. total streams - Hoxby larger (N=311)	108.1	(118.1)	0.82	0.86	1.00		
(4) Intra-county streams	106.6	(122.4)	0.78	0.78	0.90	1.00	
(5) Short streams (<3.5 miles)	75.3	(102.6)	0.73	0.73	0.86	0.96	1.00

Notes: N=335 except where otherwise indicated. All statistics are unweighted.

Appendix Table C1: Full coefficient vectors for models in Table 1

	12th grade reading score								8th grade reading score							
	Published		Hoxby/ NCES		Close replication		Preferred replication		Hoxby/ NCES		Close replication		Preferred replication			
	(1)	(2)	(3)	(4)	(2)	(3)	(4)	(2)	(3)	(4)						
Choice index	5.77 (2.21)	5.30 (2.36)	4.74 (1.98)	3.29 (1.83)	4.45 (1.87)	5.93 (2.10)	2.93 (1.58)									
<i>Individual covariates</i>																
ln(Family income)	1.54 (0.16)	1.62 (0.18)	1.75 (0.17)	0.85 (0.12)	1.60 (0.12)	1.67 (0.12)	1.03 (0.08)									
Asian	0.28 (0.59)	0.95 (0.51)	0.35 (0.49)	-1.47 (0.43)	0.25 (0.38)	0.17 (0.38)	-0.67 (0.34)									
Hispanic	-2.87 (0.52)	-1.71 (0.47)	-2.44 (0.45)	-3.49 (0.41)	-2.23 (0.32)	-2.58 (0.32)	-2.41 (0.29)									
Black	-5.49 (0.50)	-4.07 (0.52)	-4.58 (0.50)	-5.82 (0.45)	-4.14 (0.34)	-4.46 (0.34)	-4.03 (0.30)									
Female	1.96 (0.23)	2.31 (0.25)	2.18 (0.24)	2.00 (0.22)	2.25 (0.18)	2.27 (0.18)	2.25 (0.17)									
Parents' highest ed is BA+	5.45 (0.30)	5.07 (0.33)	5.14 (0.31)	6.51 (0.34)	5.42 (0.24)	5.13 (0.24)	6.98 (0.26)									
Parents' highest ed is some college	2.31 (0.30)	2.98 (0.33)	2.90 (0.32)	2.52 (0.28)	3.09 (0.23)	3.25 (0.23)	2.49 (0.20)									
<i>District/MSA covariates</i>																
Population (tens of millions)	nr		-1.59 (1.09)	-2.85 (1.38)	-0.10 (1.29)			-0.11 (0.84)		-1.08 (1.35)			0.00 (1.03)			
Land area (hundreds of thousands of sq. mi.)	nr		-13.13 (7.66)	-8.40 (5.85)	2.15 (6.25)			-14.68 (6.12)		-16.55 (6.02)			-6.35 (4.63)			
ln(mean HH income)	nr	-5.42 (5.53)	3.45 (2.86)	-3.66 (3.95)	1.83 (1.99)	-0.52 (3.15)	2.19 (1.56)	-1.84 (2.61)	5.08 (2.06)	-7.41 (3.02)	4.11 (1.68)	-4.73 (3.00)	3.87 (1.28)	-3.10 (2.23)		
Gini coefficient, HH income	nr	-12.77 (12.02)	8.55 (10.67)	-19.12 (21.10)	5.38 (8.60)	8.27 (19.34)	-6.32 (5.51)	5.79 (13.44)	6.85 (7.54)	-15.70 (17.17)	3.41 (7.25)	-4.05 (19.81)	0.63 (4.29)	-3.21 (11.85)		
Pct. Asian	nr	-5.62 (13.07)	6.42 (6.73)	-3.94 (8.18)	30.26 (6.77)	-39.41 (10.72)	8.52 (5.47)	-7.49 (8.96)	10.72 (4.73)	-16.60 (6.01)	16.07 (4.89)	-25.76 (8.67)	5.31 (4.34)	-12.85 (7.52)		
Pct. Black	nr	-0.73 (6.06)	14.38 (3.39)	-9.84 (7.25)	16.08 (3.40)	-14.61 (7.23)	6.46 (2.93)	1.98 (6.31)	-0.59 (2.13)	-3.39 (5.28)	4.24 (2.36)	-6.35 (6.57)	0.17 (2.00)	-3.65 (4.99)		
Pct. Hispanic	nr	0.25 (3.52)	10.84 (3.75)	-7.93 (4.87)	11.62 (3.56)	-8.88 (5.14)	7.18 (3.02)	1.09 (4.60)	-2.21 (2.68)	2.52 (3.67)	1.05 (2.83)	0.70 (4.70)	0.20 (2.46)	0.52 (3.85)		
Racial homogeneity index (0=homog., 1=heterog.)	nr	-9.60 (7.84)	87.42 (17.46)	-42.61 (20.51)	90.55 (17.99)	-51.73 (28.31)	52.72 (15.99)	-0.12 (26.77)	8.13 (11.72)	-8.13 (15.71)	30.07 (13.67)	-35.03 (28.20)	6.27 (11.87)	-13.43 (22.17)		
Ethnic homogeneity index	nr	16.31 (10.70)	-63.75 (13.36)	30.27 (15.88)	-66.22 (13.64)	35.65 (21.53)	-43.63 (12.17)	6.12 (20.34)	-6.87 (9.10)	6.81 (12.22)	-21.48 (10.59)	27.03 (21.54)	-3.65 (9.17)	11.38 (16.95)		
Pct. Some college	nr	5.27 (7.14)	1.42 (7.34)	-2.47 (10.34)	-4.04 (7.01)	16.72 (11.12)	4.19 (6.24)	2.39 (10.11)	-3.57 (5.27)	2.75 (8.16)	-9.25 (5.81)	18.75 (10.73)	-4.37 (4.91)	3.01 (8.57)		
Pct. College graduates	nr	3.16 (5.93)	1.39 (5.68)	2.53 (8.50)	1.11 (4.83)	3.58 (6.44)	1.62 (4.06)	4.59 (5.48)	-1.96 (4.17)	7.96 (6.47)	-0.48 (4.17)	6.59 (6.40)	-0.96 (3.21)	3.83 (4.78)		
Educational homogeneity index	nr	-5.45 (12.95)	-11.50 (8.98)	-8.59 (21.75)	-15.34 (8.25)	12.53 (20.62)	-15.67 (7.18)	-9.45 (19.04)	-9.18 (6.55)	-1.65 (17.53)	-17.01 (7.03)	11.38 (20.85)	-14.41 (5.70)	2.88 (16.36)		
N	6,119		5,475		5,934		6,688		10,175		10,429		11,719			
R ²	nr		0.161		0.173		0.200		0.181		0.180		0.197			

Notes: All models include census division fixed effects. Moulton S.E.s are reported. "nr"=not reported.

Appendix Table C2: Full coefficient vectors for first stage models in Table 2

	MSA level				Indiv. level (12th gr. reading samp.)					
	Published	Hoxby/ NCES	Close replication	Preferred replication	Hoxby/ NCES		Close replication		Preferred replication	
	(1)	(2)	(3)	(4)	(2)	(3)	(3)	(4)	(4)	(4)
Larger streams (100s)	0.080 (0.040)	0.012 (0.021)	0.040 (0.021)	0.043 (0.021)	-0.043 (0.023)	-0.024 (0.020)			0.015 (0.020)	
Smaller streams (100s)	0.034 (0.007)	0.096 (0.019)	0.093 (0.018)	0.091 (0.018)	0.133 (0.021)	0.133 (0.017)			0.114 (0.018)	
<i>District/MSA covariates</i>	<i>MSA</i>	<i>MSA</i>	<i>MSA</i>	<i>MSA</i>	<i>Dist.</i>	<i>MSA</i>	<i>Dist.</i>	<i>MSA</i>	<i>Dist.</i>	<i>MSA</i>
Population (tens of millions)	150.00 (130.00)	0.09 (0.14)	0.07 (0.16)	0.06 (0.16)		0.11 (0.14)		0.03 (0.16)		-0.04 (0.16)
Land area (hundreds of thousands of sq. mi.)	0.50 (0.50)	1.43 (0.79)	0.90 (0.73)	1.00 (0.72)		0.53 (0.55)		0.65 (0.54)		0.60 (0.58)
ln(mean HH income)	-0.25 (0.16)	-0.13 (0.17)	-0.95 (0.44)	-0.91 (0.44)	-0.37 (0.14)	0.13 (0.27)	-0.30 (0.12)	0.21 (0.24)	-0.25 (0.10)	0.36 (0.21)
Gini coefficient, HH income	-3.58 (0.81)	-3.11 (0.90)	0.35 (0.28)	0.22 (0.27)	-1.43 (0.46)	-2.20 (1.30)	-1.18 (0.46)	-2.58 (1.23)	-0.70 (0.36)	-2.23 (0.86)
Pct. Asian	1.88 (1.00)	-0.84 (0.34)	-0.66 (1.08)	-0.65 (1.07)	-0.63 (0.48)	-0.51 (0.54)	-1.59 (0.64)	0.77 (1.07)	-1.55 (0.59)	0.30 (1.07)
Pct. Black	0.83 (0.41)	0.63 (0.31)	-0.01 (0.15)	0.16 (0.13)	-0.23 (0.19)	0.46 (0.42)	-0.52 (0.22)	1.00 (0.52)	-0.54 (0.19)	0.77 (0.47)
Pct. Hispanic	0.09 (0.18)	0.29 (0.18)	-3.32 (0.95)	-2.62 (0.72)	-0.27 (0.26)	0.60 (0.39)	-0.31 (0.28)	0.94 (0.46)	-0.42 (0.28)	0.86 (0.47)
Racial homogeneity index (0=homog., 1=heterog.)	-0.18 (0.49)	0.61 (1.47)	-1.06 (0.50)	-1.09 (0.50)	-3.52 (1.18)	2.83 (1.97)	-5.49 (1.24)	5.33 (2.52)	-5.84 (1.28)	3.36 (2.73)
Ethnic homogeneity index	0.70 (0.70)	-0.27 (1.12)	0.64 (0.35)	0.60 (0.33)	2.84 (0.90)	-2.08 (1.49)	4.25 (0.93)	-3.82 (1.86)	4.57 (0.97)	-2.43 (2.00)
Pct. Some college	-1.00 (0.41)	-0.92 (0.38)	0.34 (0.21)	0.35 (0.21)	0.45 (0.40)	-0.94 (0.60)	0.64 (0.34)	-1.37 (0.63)	0.67 (0.33)	-1.03 (0.64)
Pct. College graduates	0.95 (0.42)	0.67 (0.38)	0.89 (1.80)	0.64 (1.85)	0.68 (0.28)	0.03 (0.47)	0.69 (0.30)	-0.29 (0.39)	0.51 (0.25)	-0.34 (0.33)
Educational homogeneity index	-2.95 (1.03)	-0.82 (0.98)	-0.54 (1.36)	-0.33 (1.39)	0.67 (0.46)	-1.36 (1.36)	0.90 (0.53)	-2.33 (1.74)	1.26 (0.47)	-2.24 (1.77)
<i>Individual covariates</i>										
ln(Family income)					0.003 (0.008)		0.000 (0.009)		0.004 (0.006)	
Asian					0.003 (0.020)		0.015 (0.019)		0.008 (0.023)	
Hispanic					-0.012 (0.020)		-0.019 (0.022)		0.010 (0.022)	
Black					-0.012 (0.019)		-0.011 (0.021)		-0.003 (0.018)	
Female					-0.011 (0.008)		-0.006 (0.010)		0.007 (0.009)	
Parents' highest ed is BA+					-0.015 (0.009)		-0.007 (0.011)		-0.009 (0.013)	
Parents' highest ed is some college					-0.012 (0.009)		-0.011 (0.015)		-0.016 (0.011)	
N	316	310	304	304	5,475		5,934		6,688	
R2	nr	0.516	0.506	0.517	0.575		0.568		0.583	

Notes: All models include census division fixed effects. Cluster S.E.s are reported for indiv.-level models. "nr"=not reported. Published specification (col. 1) also includes controls for the shares of the MSA population that are 0-19 and 65+ years old.

Appendix Table C3: Full coefficient vectors for alternative first stage models in Table 3

Total stream defn. Larger stream defn. Level	Stream mouths						All streams								
	n/a			Hoxby			n/a			Inter-county			>3.5 miles		
	MSA		Indiv	MSA		Indiv	MSA		Indiv	MSA		Indiv	MSA		Indiv
	(2A)	(2B)	(3A)	(3B)	(4A)	(4B)	(5A)	(5B)	(6A)	(6B)					
Larger streams (100s)			0.037	-0.030					0.260	0.240		0.177	0.190		
			(0.021)	(0.019)					(0.055)	(0.047)		(0.036)	(0.029)		
Smaller streams (100s)			0.069	0.104					0.014	0.015		0.013	0.001		
			(0.013)	(0.013)					(0.016)	(0.013)		(0.017)	(0.013)		
Total streams (100s)	0.071	0.064					0.061	0.058							
	(0.013)	(0.011)					(0.010)	(0.009)							
<i>District/MSA covariates</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>
Population (tens of millions)	0.07	0.00	0.07	0.01	0.09	0.02	0.01	-0.04	0.05	-0.01	0.05	-0.01	0.05	-0.01	
	(0.16)	(0.17)	(0.16)	(0.16)	(0.15)	(0.16)	(0.15)	(0.18)	(0.15)	(0.16)	(0.15)	(0.16)	(0.15)	(0.16)	
Land area (100,000s of sq. mi.)	0.88	0.86	0.82	0.64	0.72	0.71	0.70	0.71	-0.53	-0.59					
	(0.72)	(0.72)	(0.73)	(0.54)	(0.71)	(0.67)	(0.69)	(0.60)	(0.79)	(0.52)					
ln(mean HH income)	-0.03	-0.27	0.26	0.01	-0.29	0.22	-0.07	-0.30	0.22	-0.03	-0.32	0.28	-0.02	-0.33	0.33
	(0.14)	(0.13)	(0.24)	(0.15)	(0.12)	(0.24)	(0.13)	(0.12)	(0.22)	(0.12)	(0.12)	(0.22)	(0.13)	(0.12)	(0.22)
Gini coefficient, HH income	-3.46	-0.90	-2.80	-3.30	-1.05	-2.54	-3.43	-0.92	-2.89	-3.03	-0.83	-2.53	-2.95	-0.85	-2.22
	(0.89)	(0.48)	(1.32)	(0.93)	(0.47)	(1.20)	(0.84)	(0.47)	(1.19)	(0.83)	(0.48)	(1.18)	(0.84)	(0.46)	(1.17)
Pct. Asian	-1.14	-1.51	0.51	-0.92	-1.54	0.79	-1.01	-1.51	0.62	-0.96	-1.47	0.71	-0.94	-1.55	0.87
	(0.48)	(0.62)	(1.04)	(0.49)	(0.64)	(1.07)	(0.47)	(0.60)	(1.01)	(0.46)	(0.60)	(1.00)	(0.46)	(0.61)	(0.99)
Pct. Black	0.58	-0.49	0.89	0.70	-0.51	1.04	0.61	-0.48	0.93	0.56	-0.45	0.91	0.49	-0.53	0.93
	(0.33)	(0.22)	(0.52)	(0.34)	(0.21)	(0.50)	(0.32)	(0.21)	(0.51)	(0.32)	(0.21)	(0.49)	(0.32)	(0.21)	(0.48)
Pct. Hispanic	0.30	-0.04	0.49	0.37	-0.30	0.92	0.33	-0.07	0.55	0.34	-0.02	0.53	0.34	-0.01	0.55
	(0.20)	(0.30)	(0.45)	(0.20)	(0.27)	(0.44)	(0.19)	(0.28)	(0.41)	(0.19)	(0.27)	(0.39)	(0.19)	(0.25)	(0.38)
Racial homogeneity index	0.27	-4.89	3.43	1.06	-5.29	5.26	0.42	-4.84	3.84	0.24	-4.53	3.72	0.42	-4.87	4.30
	(1.71)	(1.19)	(2.47)	(1.74)	(1.21)	(2.45)	(1.63)	(1.13)	(2.35)	(1.60)	(1.13)	(2.30)	(1.61)	(1.13)	(2.16)
Ethnic homogeneity index	-0.10	3.86	-2.44	-0.69	4.10	-3.77	-0.24	3.82	-2.76	-0.10	3.60	-2.67	-0.23	3.83	-3.04
	(1.30)	(0.89)	(1.82)	(1.31)	(0.91)	(1.80)	(1.23)	(0.84)	(1.72)	(1.21)	(0.84)	(1.70)	(1.21)	(0.84)	(1.60)
Pct. Some college	-1.04	0.68	-1.43	-0.90	0.63	-1.33	-0.93	0.66	-1.40	-0.77	0.66	-1.16	-0.94	0.75	-1.40
	(0.43)	(0.33)	(0.67)	(0.44)	(0.33)	(0.64)	(0.42)	(0.32)	(0.67)	(0.41)	(0.33)	(0.67)	(0.41)	(0.32)	(0.64)
Pct. College graduates	0.34	0.59	-0.22	0.28	0.68	-0.30	0.30	0.60	-0.25	0.23	0.63	-0.35	0.18	0.61	-0.43
	(0.27)	(0.30)	(0.38)	(0.28)	(0.30)	(0.38)	(0.26)	(0.29)	(0.37)	(0.26)	(0.29)	(0.35)	(0.26)	(0.29)	(0.36)
Educational homogeneity index	-0.64	0.81	-2.12	-0.37	0.85	-2.13	-0.62	0.48	-2.13	-0.27	0.48	-1.53	-0.49	0.56	-1.81
	(1.04)	(0.51)	(1.69)	(1.06)	(0.46)	(1.63)	(0.96)	(0.44)	(1.46)	(0.95)	(0.45)	(1.47)	(0.95)	(0.44)	(1.40)
<i>Individual covariates</i>															
ln(Family income)		0.00		0.00		0.00		0.00		0.00		0.00		0.00	
		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)	
Asian		0.01		0.02		0.02		0.02		0.02		0.02		0.01	
		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)	
Hispanic		-0.01		-0.02		-0.02		-0.02		-0.02		-0.01		-0.01	
		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)	
Black		0.00		-0.01		0.00		0.00		0.00		0.00		-0.01	
		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)		(0.02)	
Female		0.00		-0.01		0.00		0.00		0.00		0.00		0.00	
		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)	
Parents' highest ed is BA+		-0.01		-0.01		-0.01		-0.01		-0.01		-0.01		-0.01	
		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)	
Parents' highest ed is some college		-0.01		-0.01		-0.01		-0.01		-0.01		-0.01		-0.01	
		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)		(0.01)	
N	319	5,987	311	6,014	335	6,139	335	6,139	335	6,139	335	6,139	335	6,139	
R2	0.506	0.531	0.508	0.569	0.505	0.535	0.525	0.557	0.520	0.560					

Notes: All models include census division fixed effects. Cluster S.E.s are reported for indiv.-level models, which use the 12th grade reading score samples.

Appendix Table C4: Full coefficient vectors for alternative-instruments IV models in Table 4, 12th grade reading scores

Model	OLS		IV									
	n/a		Stream mouths		All streams							
Total stream defn.	n/a		n/a		Hoxby		n/a		Inter-county		>3.5 miles	
Larger stream defn.	n/a		n/a		Hoxby		n/a		Inter-county		>3.5 miles	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
Choice index	-0.25		0.68		4.38		0.87		2.04		1.35	
	(0.79)		(2.79)		(1.98)		(2.59)		(2.36)		(2.30)	
<i>District/MSA covariates</i>	<i>Dist</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>	<i>Dist</i>	<i>MSA</i>
Population (tens of millions)	-0.83		-1.98		-2.65		-1.15		-1.48		-1.28	
	(2.01)		(1.42)		(1.36)		(1.40)		(1.37)		(1.38)	
Land area (100,000s of sq. mi.)	1.38		0.29		-7.72		-1.18		-3.84		-2.25	
	(7.90)		(6.44)		(5.85)		(6.25)		(6.08)		(5.99)	
ln(mean HH income)	-0.76	-1.02	0.18	-0.17	1.62	-0.51	-0.41	-1.12	-0.04	-1.21	-0.26	-1.15
	(2.03)	(2.51)	(2.00)	(3.10)	(1.98)	(3.13)	(1.96)	(3.03)	(1.95)	(3.04)	(1.95)	(3.03)
Gini coefficient, HH income	-6.99	-5.20	-2.72	-1.11	3.64	9.07	-5.92	-1.62	-4.82	2.10	-5.48	-0.12
	(8.57)	(15.23)	(8.53)	(21.58)	(8.56)	(19.19)	(8.38)	(20.54)	(8.33)	(19.79)	(8.29)	(19.72)
Pct. Asian	20.31	-32.23	22.77	-34.96	29.64	-38.98	22.23	-33.34	24.22	-34.50	23.03	-33.81
	(11.22)	(15.98)	(6.90)	(10.60)	(6.72)	(10.69)	(6.75)	(10.48)	(6.70)	(10.50)	(6.69)	(10.49)
Pct. Black	13.29	-8.03	13.34	-9.23	16.10	-14.40	13.82	-9.38	14.37	-10.78	14.04	-9.95
	(3.62)	(6.66)	(3.41)	(7.72)	(3.37)	(7.10)	(3.35)	(7.38)	(3.34)	(7.23)	(3.34)	(7.23)
Pct. Hispanic	7.65	-2.31	7.84	-3.72	11.59	-8.70	7.85	-3.05	8.07	-3.82	7.94	-3.36
	(4.17)	(5.42)	(3.47)	(5.20)	(3.53)	(5.01)	(3.41)	(4.98)	(3.41)	(4.92)	(3.40)	(4.91)
Racial homogeneity index	61.54	-16.69	65.47	-25.04	88.45	-49.35	67.11	-22.81	72.90	-29.16	69.45	-25.38
	(19.22)	(28.84)	(18.38)	(29.72)	(17.79)	(28.17)	(17.82)	(29.08)	(17.60)	(28.42)	(17.58)	(28.51)
Ethnic homogeneity index	-44.64	11.16	-47.62	17.02	-64.51	33.82	-49.03	15.42	-53.58	19.85	-50.87	17.21
	(14.35)	(21.30)	(13.99)	(22.26)	(13.49)	(21.35)	(13.53)	(21.75)	(13.35)	(21.33)	(13.33)	(21.39)
Pct. Some college	-4.61	7.60	-4.38	9.35	-4.28	16.36	-5.35	9.91	-6.12	12.29	-5.66	10.87
	(6.92)	(8.96)	(6.94)	(11.27)	(7.00)	(11.12)	(6.85)	(11.09)	(6.86)	(11.03)	(6.85)	(11.00)
Pct. College graduates	7.50	4.42	6.30	4.30	1.59	3.25	6.80	4.26	6.08	4.09	6.51	4.19
	(4.86)	(5.46)	(4.80)	(6.37)	(4.82)	(6.42)	(4.72)	(6.25)	(4.71)	(6.24)	(4.70)	(6.23)
Educational homogeneity index	-13.89	0.07	-7.76	4.83	-17.81	14.05	-14.58	3.66	-15.30	7.38	-14.87	5.16
	(8.93)	(16.05)	(8.32)	(20.97)	(8.18)	(19.70)	(8.09)	(19.65)	(8.07)	(19.57)	(8.05)	(19.58)
<i>Individual covariates</i>												
ln(Family income)	1.75		1.77		1.77		1.75		1.75		1.75	
	(0.17)		(0.17)		(0.17)		(0.17)		(0.17)		(0.17)	
Asian	0.44		0.44		0.34		0.43		0.42		0.43	
	(0.69)		(0.49)		(0.49)		(0.49)		(0.49)		(0.49)	
Hispanic	-2.80		-2.61		-2.51		-2.77		-2.75		-2.76	
	(0.53)		(0.45)		(0.44)		(0.44)		(0.44)		(0.44)	
Black	-4.62		-4.48		-4.60		-4.62		-4.61		-4.62	
	(0.46)		(0.50)		(0.50)		(0.50)		(0.50)		(0.50)	
Female	2.16		2.15		2.18		2.16		2.17		2.17	
	(0.23)		(0.24)		(0.24)		(0.23)		(0.23)		(0.23)	
Parents' highest ed is BA+	5.02		5.05		5.10		5.04		5.06		5.05	
	(0.31)		(0.31)		(0.31)		(0.31)		(0.31)		(0.31)	
Parents' highest ed is some college	2.81		2.82		2.88		2.82		2.84		2.83	
	(0.31)		(0.32)		(0.32)		(0.31)		(0.31)		(0.31)	
N	6,139		5,987		6,014		6,139		6,139		6,139	
R2	0.186		0.186		0.175		0.185		0.184		0.185	

Notes: All models include census division fixed effects. Moulton S.E.s are reported.

Appendix Table C5: Full coefficient vectors for selected models from Table 5

Sample/covariates	Close replication					
	Hoxby sm. & lg. streams			Inter- & intra-cnty. streams		
Instruments	Base samp., no	Zip-code		Base samp., no	Zip-code	
Specification	district-lvl	matched pub.	Pub. & pvt.	district-lvl	matched pub.	Pub. & pvt.
	covariates	schls	schools	covariates	schls	schools
	(1)	(2)	(3)	(4)	(5)	(6)
Choice index	4.61 (2.49)	1.40 (2.44)	0.68 (2.59)	1.76 (2.86)	1.10 (2.66)	0.84 (2.35)
<i>MSA-level covariates</i>						
Population (tens of millions)	-3.55 (2.46)	-0.97 (1.78)	-1.22 (1.68)	-2.02 (2.09)	-1.02 (1.72)	-1.82 (1.61)
Land area (100,000s of sq. mi.)	-9.44 (9.74)	1.20 (7.15)	0.80 (7.28)	-4.11 (8.89)	0.28 (7.09)	-0.32 (6.47)
ln(mean HH income)	0.99 (2.69)	-1.50 (2.37)	-0.88 (2.20)	-0.97 (2.27)	-2.22 (2.34)	-0.57 (2.17)
Gini coefficient, HH income	7.31 (18.32)	-4.28 (17.35)	-0.87 (17.75)	-6.86 (18.15)	-7.00 (18.57)	2.17 (16.99)
Pct. Asian	-5.35 (7.60)	-5.99 (7.10)	-2.48 (6.65)	-8.13 (6.56)	-7.62 (6.93)	-4.53 (6.14)
Pct. Black	-1.62 (7.14)	6.20 (7.47)	7.26 (6.95)	0.54 (7.68)	5.31 (7.89)	6.18 (7.14)
Pct. Hispanic	3.05 (3.67)	7.42 (4.09)	6.84 (3.81)	3.98 (3.71)	7.10 (4.20)	6.35 (3.87)
Racial homogeneity index	6.46 (33.64)	19.17 (29.51)	14.58 (28.27)	16.62 (29.60)	18.57 (28.49)	4.35 (27.39)
Ethnic homogeneity index	-5.13 (25.47)	-11.20 (22.65)	-7.08 (21.87)	-12.58 (21.83)	-11.88 (21.59)	0.60 (20.74)
Pct. Some college	11.91 (10.20)	11.48 (9.37)	7.07 (9.00)	4.87 (9.24)	8.26 (9.09)	6.31 (8.82)
Pct. College graduates	1.63 (6.66)	9.78 (6.12)	8.39 (5.90)	6.48 (6.13)	12.00 (5.83)	9.80 (5.46)
Educational homogeneity index	-8.72 (21.92)	-19.84 (19.44)	-17.92 (18.08)	-10.07 (20.26)	-12.45 (18.51)	-14.18 (17.01)
<i>Individual covariates</i>						
ln(Family income)	1.76 (0.24)	1.07 (0.28)	1.43 (0.27)	1.75 (0.23)	1.12 (0.27)	1.42 (0.25)
Asian	0.67 (0.68)	0.41 (0.79)	0.42 (0.74)	0.75 (0.65)	0.45 (0.77)	0.34 (0.71)
Hispanic	-2.37 (0.84)	-2.93 (0.91)	-2.69 (0.83)	-2.68 (0.82)	-3.01 (0.87)	-2.67 (0.79)
Black	-4.15 (0.87)	-4.85 (0.93)	-5.07 (0.82)	-4.17 (0.83)	-4.88 (0.91)	-5.15 (0.80)
Female	2.18 (0.39)	2.43 (0.43)	2.63 (0.40)	2.15 (0.37)	2.47 (0.41)	2.59 (0.38)
Parents' highest ed is BA+	5.48 (0.48)	5.87 (0.54)	6.14 (0.50)	5.39 (0.47)	5.81 (0.53)	6.15 (0.48)
Parents' highest ed is some college	3.10 (0.52)	3.05 (0.57)	3.08 (0.53)	3.01 (0.49)	3.16 (0.55)	3.21 (0.51)
N	5,939	5,445	6,670	6,144	5,631	6,900
R2	0.164	0.170	0.191	0.178	0.174	0.199

Notes: All models include census division fixed effects. Clustered S.E.s are reported.

Appendix Table D1: First-stage estimates for alternative instruments, using "preferred" replication sample and covariates

	(1)	(2)	(3)	(4)	(5)	(6)
Total stream definition	Stream mouths		All streams			
Larger stream definition	Hoxby	n/a	Hoxby	n/a	Inter-county >3.5 miles	
<i>MSA level</i>						
Larger streams (100s)	0.043 (0.021)		0.040 (0.021)		0.250 (0.054)	0.178 (0.035)
Smaller streams (100s)	0.091 (0.018)		0.068 (0.012)		0.017 (0.016)	0.012 (0.017)
Total streams (100s)		0.071 (0.013)		0.061 (0.010)		
F statistic, instruments	16.3	31.7	17.8	37.9	26.1	25.0
<i>Individual level (12th grade reading sample)</i>						
Larger streams (100s)	0.015 (0.020)		0.001 (0.018)		0.236 (0.044)	0.171 (0.029)
Smaller streams (100s)	0.114 (0.018)		0.099 (0.012)		0.025 (0.014)	0.018 (0.017)
Total streams (100s)		0.072 (0.011)		0.066 (0.009)		
F statistic, instruments	28.4	44.2	38.8	53.3	34.3	36.7
<i>Individual level (8th grade reading sample)</i>						
Larger streams (100s)	-0.012 (0.018)		-0.017 (0.017)		0.227 (0.044)	0.151 (0.029)
Smaller streams (100s)	0.132 (0.017)		0.102 (0.012)		0.021 (0.014)	0.018 (0.017)
Total streams (100s)		0.067 (0.012)		0.060 (0.009)		
F statistic, instruments	32.1	33.2	34.7	40.7	28.2	28.1

Notes: Base samples are those from Column 4 of Tables 1 (individual level) and 2 (Panel B; MSA level), though some observations that were excluded from those samples for missing data on larger streams are included here in Columns 2, 4, 5, and 6. In individual-level specifications, standard errors are clustered at the MSA level.

Appendix Table D2: IV estimates of choice effect, using alternative instruments and "preferred" replication sample

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	OLS	IV					
Total stream definition	n/a	Stream mouths		All streams			
Larger stream definition	n/a	Hoxby	n/a	Hoxby	n/a	Inter-county	>3.5 miles
<i>Panel : 12th grade reading scores</i>							
Choice effect	-0.17	3.29	2.79	4.05	3.41	3.36	2.64
S.E. (Moulton)	(0.62)	(1.83)	(2.44)	(1.79)	(2.13)	(1.96)	(1.92)
S.E. (Cluster)	(0.97)	(2.56)	(2.69)	(2.06)	(2.45)	(2.60)	(1.98)
p-value, exog. test	--	0.20	0.25	0.04	0.10	0.12	0.10
<i>Panel A: 8th grade reading scores</i>							
Choice index	-0.55	2.93	0.58	2.79	1.03	0.95	0.51
S.E. (Moulton)	(0.62)	(1.58)	(1.95)	(1.55)	(1.83)	(1.72)	(1.65)
S.E. (Cluster)	(0.61)	(1.40)	(1.99)	(1.29)	(1.82)	(1.39)	(1.53)
p-value, exog. test	--	0.00	0.55	0.00	0.34	0.20	0.44

Notes: Base samples are those from Column 4 of Table 1, though some observations that were excluded from that sample for missing data on larger streams are included here in Columns 3 and 5-7. Exogeneity tests are based on clustered specification. Bold S.E.s indicate that with that S.E., the coefficient is significant at the 5% level.

Appendix Table D3. Exploration of potential bias from exclusion of private school students, 8th grade reading scores

Covariate specification	Close replication			Preferred replication		
	OLS	Hoxby	Inter- and intra-cnty	OLS	Hoxby	Inter- and intra-cnty
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Base specification, replication sample of public school students</i>						
Choice index coefficient	-0.06 (0.82)	5.93 (2.32)	1.67 (1.77)	-0.55 (0.61)	2.93 (1.40)	0.95 (1.39)
N	10,709	10,429	10,709	12,049	11,719	12,049
p-value, exog. test		0.00	0.21		0.00	0.20
<i>Panel B: Without district-level covariates</i>						
Choice index coefficient	-0.13 (0.80)	5.03 (2.33)	1.98 (1.82)	-0.60 (0.60)	2.06 (1.28)	0.97 (1.40)
N	10,729	10,449	10,729	12,049	11,719	12,049
p-value, exog. test		0.01	0.15		0.03	0.19
<i>Panel C: Public school students in zip-code matched sample (no district covariates)</i>						
Choice index coefficient	-0.87 (0.70)	1.93 (1.55)	0.19 (1.50)	-0.59 (0.60)	1.57 (1.30)	0.20 (1.35)
N	10,394	10,117	10,394	11,992	11,662	11,992
p-value, exog. test		0.04	0.39		0.08	0.50
<i>Panel D: Public and private school students in zip code-matched sample</i>						
Choice index coefficient	-0.37 (0.65)	1.07 (1.56)	0.65 (1.57)	-0.13 (0.59)	0.78 (1.40)	0.45 (1.48)
N	13,879	13,482	13,879	16,026	15,558	16,026
p-value, exog. test		0.30	0.44		0.48	0.66

Notes: Clustered standard errors and test statistics are reported. Bold coefficients are significant at the 5% level.

Appendix Table D4: Choice effect estimates for all six NELS test scores

Sample Model Instruments	Hoxby/NELS		Close replication			Preferred replication		
	OLS	IV	OLS	IV		OLS	IV	
		Hoxby streams		Hoxby streams	Inter- and intra-county streams		Hoxby streams	Inter- and intra-county streams
		(2)		(4)	(5)		(7)	(8)
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	
<i>Panel A: 8th grade</i>								
Reading	-0.65 (0.71)	4.45 (1.99)	-0.06 (0.82)	5.93 (2.32)	1.67 (1.77)	-0.55 (0.61)	2.93 (1.40)	0.95 (1.39)
Mathematics	-0.61 (0.73)	4.23 (1.79)	-0.55 (0.71)	4.52 (1.77)	-0.64 (1.57)	-0.80 (0.63)	3.22 (1.59)	-0.60 (1.41)
<i>Panel B: 10th grade</i>								
Reading	-0.85 (1.01)	6.73 (2.51)	-0.88 (1.06)	6.19 (2.64)	5.34 (2.90)	0.19 (0.95)	3.94 (2.43)	4.39 (2.43)
Mathematics	-0.89 (0.93)	7.91 (2.21)	-1.10 (0.86)	4.95 (2.14)	1.32 (2.32)	0.10 (0.75)	2.67 (1.99)	0.00 (2.08)
<i>Panel C: 12th grade</i>								
Reading	-1.24 (1.15)	5.30 (2.94)	-0.25 (0.94)	4.74 (2.42)	2.04 (2.94)	-0.17 (0.97)	3.29 (2.56)	3.36 (2.60)
Mathematics	-1.21 (0.93)	3.49 (2.29)	-0.34 (0.80)	2.48 (2.08)	1.56 (2.29)	0.60 (1.01)	2.15 (2.06)	2.69 (2.30)

Notes: Cluster S.E.s in parentheses. Bold coefficients are significant at the 5% level.

Appendix Table D5: MSA-level first stage estimates, using only MSAs in the NELS 12th grade sample

	(1)	(2)	(3)	(4)	(5)	(6)
Total stream definition	Stream mouths		All streams			
Larger stream definition	Hoxby	n/a	Hoxby	n/a	Inter-county	>3.5 miles
<i>Panel A: Hoxby/NELS sample</i>						
Larger streams (100s)	-0.044					
	(0.028)					
Smaller streams (100s)	0.143					
	(0.025)					
Total streams (100s)		0.057				
		(0.016)				
F statistic (instruments)	16.0	13.0				
<i>Panel B: Close replication sample</i>						
Larger streams (100s)	-0.013		-0.018		0.242	0.166
	(0.028)		(0.028)		(0.064)	(0.043)
Smaller streams (100s)	0.143		0.113		0.024	0.021
	(0.024)		(0.018)		(0.020)	(0.023)
Total streams (100s)		0.073		0.066		
		(0.016)		(0.013)		
F statistic (instruments)	18.3	21.7	20.2	26.8	17.8	16.4
<i>Panel C: Preferred replication sample</i>						
Larger streams (100s)	0.019		0.005		0.245	0.164
	(0.026)		(0.026)		(0.062)	(0.042)
Smaller streams (100s)	0.121		0.106		0.029	0.029
	(0.022)		(0.017)		(0.019)	(0.022)
Total streams (100s)		0.078		0.071		
		(0.015)		(0.012)		
F statistic (instruments)	16.9	26.6	20.8	33.4	21.4	19.3

Appendix Table D6: Two-sample IV estimates of choice effects on 12th grade reading scores

Sample Streams instruments	Hoxby/NELS	Close replication		Preferred replication	
	Hoxby	Hoxby	Inter- and intra-cnty	Hoxby	Inter- and intra-cnty
	(1)	(2)	(3)	(4)	(5)
<i>Panel A: 8th grade reading scores</i>					
<i>First stage uses only NELS MSAs</i>					
Choice index coefficient	3.99 (1.60)	4.22 (1.75)	0.76 (1.51)	2.70 (1.23)	0.90 (1.30)
<i>First stage uses all MSAs</i>					
Choice index coefficient	3.68 (2.27)	3.46 (2.15)	0.80 (1.51)	2.15 (1.68)	0.90 (1.28)
<i>Panel B: 12th grade reading scores</i>					
<i>First stage uses only NELS MSAs</i>					
Choice index coefficient	4.58 (2.61)	4.64 (2.20)	2.31 (2.24)	3.09 (2.51)	3.15 (2.36)
<i>First stage uses all MSAs</i>					
Choice index coefficient	2.14 (2.98)	4.15 (2.57)	2.71 (2.27)	3.69 (2.75)	3.41 (2.64)

Notes: First stage is estimated at MSA level, on the indicated sample of MSAs. Reported coefficients are those on the first stage fitted value from the second stage regression. Reported standard errors are clustered S.E.s from the second stage regression. These are unadjusted for the presence of a generated regressor, so are not correct, and are likely downward-biased estimates of the true standard errors. Bold coefficients are significant at the 5% level with these S.E.s.

Appendix Table D7: Reduced-form estimates, effect of streams on test scores

Total stream definition Larger stream definition	8th grade reading scores						12th grade reading scores					
	Stream mouths		All streams				Stream mouths		All streams			
	Hoxby	n/a	Hoxby	n/a	Inter- county	>3.5 miles	Hoxby	n/a	Hoxby	n/a	Inter- county	>3.5 miles
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
<i>Panel A: Hoxby/NELS sample</i>												
Larger streams (100s)	-0.50 (0.26)						-1.39 (0.55)					
Smaller streams (100s)	0.60 (0.23)						0.77 (0.38)					
Total streams (100s)		0.09 (0.16)							-0.18 (0.23)			
<i>Panel B: Close replication sample</i>												
Larger streams (100s)	-0.36 (0.20)		-0.37 (0.20)		0.50 (0.41)	-0.21 (0.28)	-1.00 (0.41)		-1.02 (0.40)		0.57 (1.18)	0.52 (0.49)
Smaller streams (100s)	0.36 (0.16)		0.26 (0.12)		-0.10 (0.15)	0.13 (0.18)	0.58 (0.28)		0.44 (0.20)		-0.25 (0.28)	-0.42 (0.32)
Total streams (100s)		0.00 (0.11)		0.01 (0.09)				-0.17 (0.16)		-0.09 (0.14)		
<i>Panel C: Preferred replication sample</i>												
Larger streams (100s)	-0.30 (0.19)		-0.30 (0.18)		0.53 (0.38)	-0.15 (0.27)	-0.54 (0.34)		-0.55 (0.33)		0.97 (1.10)	0.47 (0.53)
Smaller streams (100s)	0.34 (0.14)		0.25 (0.10)		-0.08 (0.13)	0.14 (0.16)	0.60 (0.27)		0.52 (0.19)		-0.08 (0.24)	-0.07 (0.32)
Total streams (100s)		0.02 (0.10)		0.03 (0.08)				0.06 (0.15)		0.12 (0.14)		

Note: Regressions are estimated by OLS, at individual level, with clustered standard errors. Bold coefficients are significantly different from zero at the 5% level.

Appendix Table D8: MSA-level estimates of the choice effect

	12th grade reading scores							8th grade reading scores						
	OLS		IV					OLS		IV				
	Total stream	n/a	Mouths		All streams			Total stream	n/a	Mouths		All streams		
Larger stream	n/a	Hoxby	n/a	Hoxby	n/a	Inter-county	>3.5 miles	n/a	Hoxby	n/a	Hoxby	n/a	Inter-county	>3.5 miles
<i>CD Sample</i>														
All controls	-0.48	4.30	-3.77	4.26	-2.60	-0.92	-0.74	-0.79	3.57	0.42	3.42	0.87	-0.41	-2.08
	(1.20)	(3.15)	(5.13)	(3.39)	(4.82)	(3.71)	(3.88)	(0.75)	(2.18)	(3.54)	(2.31)	(3.28)	(2.28)	(2.52)
MSA controls	-1.33	3.89	-2.56	3.83	-1.84	0.37	0.09	-0.56	3.93	1.43	4.29	2.12	3.37	2.95
	(1.29)	(3.25)	(4.62)	(3.46)	(4.45)	(3.60)	(3.82)	(0.99)	(2.62)	(3.74)	(2.82)	(3.60)	(2.94)	(3.09)
<i>Close replication sample</i>														
All controls	0.45	5.08	0.57	4.72	0.94	2.21	1.81	-0.17	5.50	1.98	4.98	2.53	1.10	-0.43
	(1.14)	(3.03)	(3.79)	(2.87)	(3.45)	(2.96)	(2.86)	(0.86)	(2.60)	(3.48)	(2.43)	(3.14)	(2.39)	(2.36)
MSA controls	-1.26	3.74	1.50	2.58	0.88	2.39	1.64	0.04	5.23	3.62	3.69	2.90	3.22	3.20
	(1.26)	(3.30)	(4.28)	(3.11)	(3.86)	(3.32)	(3.29)	(1.07)	(2.75)	(3.68)	(2.58)	(3.29)	(2.77)	(2.80)
<i>Preferred replication sample</i>														
All controls	0.17	3.38	3.19	4.39	3.89	3.24	3.22	-0.38	2.68	0.46	3.01	1.21	1.04	-0.04
	(1.19)	(3.38)	(3.82)	(3.08)	(3.42)	(2.94)	(2.92)	(0.74)	(2.06)	(2.71)	(1.98)	(2.47)	(2.01)	(2.10)
MSA controls	-0.81	2.73	3.53	2.49	3.20	3.24	2.66	-0.57	2.75	2.05	1.82	1.74	2.54	2.47
	(1.32)	(3.42)	(3.89)	(3.13)	(3.52)	(3.17)	(3.28)	(1.00)	(2.42)	(3.05)	(2.33)	(2.81)	(2.53)	(2.63)

Notes: "All controls" specifications include MSA-level averages of individual- and district-level covariates, computed within the NELS sample. "MSA controls" specifications exclude the individual- and district-level covariates. Bold coefficients are significant at the 5% level.

Appendix Table E1: Estimates of choice effect on ln(per pupil spending)

	District-level analysis			MSA-level analysis		
	OLS	IV, Hoxby streams	IV, Inter- and intra-county streams	OLS	IV, Hoxby streams	IV, Inter- and intra-cnty streams
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A: Published estimates</i>						
Choice index coefficient	-0.072	-0.076				
S.E. (Moulton)	(0.022)	(0.034)				
<i>Panel B: Hoxby/NCES data</i>						
Choice index coefficient	-0.070	-0.074				
S.E. (Classical)	(0.012)	(0.037)				
S.E. (Moulton)	(0.024)	(0.141)				
S.E. (Cluster)	(0.037)	(0.111)				
<i>Panel C: Replication, "close" covariates</i>						
Choice index coefficient	-0.121	-0.201	-0.139	-0.124	-0.187	-0.135
S.E. (Classical)				(0.035)	(0.111)	(0.092)
S.E. (Moulton)	(0.023)	(0.154)	(0.084)			
S.E. (Cluster)	(0.034)	(0.111)	(0.137)			
<i>Panel D: Replication, "preferred" covariates</i>						
Choice index coefficient	-0.120	-0.211	-0.142	-0.119	-0.194	-0.120
S.E. (Classical)				(0.035)	(0.111)	(0.092)
S.E. (Moulton)	(0.023)	(0.153)	(0.084)			
S.E. (Cluster)	(0.034)	(0.112)	(0.131)			

Notes: Dependent variable is log of average per pupil spending in the district (columns 1-3, N=5,336-5,804) or MSA (columns 4-6, N=302-333). "Hoxby streams" instruments are Hoxby's larger and smaller streams variables. All specifications include usual list of MSA-level covariates and division fixed effects; those in columns 1-3 also include usual district-level covariates.