

# Disrupting Education? Experimental Evidence on Technology-Aided Instruction in India

Karthik Muralidharan and Abhijeet Singh and Alejandro J.Ganimian\*

February 22, 2019

## Abstract

We study the impact of a personalized technology-aided after-school instruction program in middle-school grades in urban India using a lottery that provided winners free access to the program. Lottery winners scored  $0.37\sigma$  higher in math and  $0.23\sigma$  higher in Hindi over just a 4.5-month period. IV estimates suggest that attending the program for 90 days would increase math and Hindi test scores by  $0.6\sigma$  and  $0.39\sigma$  respectively. We find similar absolute test score gains for all students, but much greater relative gains for academically-weaker students. Our results suggest that well-designed technology-aided instruction programs can sharply improve productivity in delivering education.

JEL codes: C93, I21, J24, O15

Keywords: computer-aided learning, productivity in education, personalized learning, teaching at the right level, post-primary education, middle school, secondary school

---

\*Muralidharan: Department of Economics, University of California San Diego, 9500 Gilman Drive, La Jolla CA; NBER; J-PAL. E-mail: kamurali@ucsd.edu. Singh: Department of Economics, Stockholm School of Economics, Sveavagen 65, Stockholm, Sweden. E-mail: abhijeet.singh@hhs.se. Ganimian: NYU Steinhardt School of Culture, Education, and Human Development, 246 Greene St, New York, NY. E-mail: alejandro.ganimian@nyu.edu. We thank Esther Duflo (the editor), Abhijit Banerjee, James Berry, Peter Bergman, Prashant Bharadwaj, Gordon Dahl, Roger Gordon, Heather Hill, Priya Mukherjee, Chris Walters and several seminar participants for comments. We thank the staff at Educational Initiatives (EI)—especially, Pranav Kothari, Smita Bardhan, Anurima Chatterjee, and Prasad Sreepakash—for their support of the evaluation. We also thank Maya Escueta, Smit Gade, Riddhima Mishra, and Rama Murthy Sripada for excellent research assistance and field support. Finally, we thank J-PAL’s Post-Primary Education initiative for funding this study. The study was registered with the AEA Trial Registry (RCT ID: AEARCTR-0000980). The operation of Mindspark centers by EI was funded by the Central Square Foundation, Tech Mahindra Foundation and Porticus. All views expressed are those of the authors and not of any of the institutions with which they are affiliated.

## References

- Andrabi, Tahir, Jishnu Das, Asim I. Khwaja, and Tristan Zajonc.** 2011. “Do value-added estimates add value? Accounting for learning dynamics.” *American Economic Journal: Applied Economics*, 3(3): 29–54.
- Angrist, Joshua, and Victor Lavy.** 2002. “New evidence on classroom computers and pupil learning.” *The Economic Journal*, 112(482): 735–765.
- Banerjee, Abhijit V, Shawn Cole, Esther Duflo, and Leigh Linden.** 2007. “Remedying Education: Evidence from Two Randomized Experiments in India.” *The Quarterly Journal of Economics*, 122(3): 1235–1264.
- Barrera-Osorio, Felipe, and Leigh L Linden.** 2009. “The use and misuse of computers in education: evidence from a randomized experiment in Colombia.” (World Bank Policy Research Working Paper No. 4836.) Washington, DC: The World Bank.
- Barrow, Lisa, Lisa Markman, and Cecilia Elena Rouse.** 2009. “Technology’s edge: The educational benefits of computer-aided instruction.” *American Economic Journal: Economic Policy*, 1(1): 52–74.
- Beuermann, Diether W, Julian Cristia, Santiago Cueto, Ofer Malamud, and Yyannu Cruz-Aguayo.** 2015. “One Laptop per Child at home: Short-term impacts from a randomized experiment in Peru.” *American Economic Journal: Applied Economics*, 7(2): 53–80.
- Bhattacharjea, S., W. Wadhwa, and R. Banerji.** 2011. *Inside primary schools: A study of teaching and learning in rural India.* ASER Centre, New Delhi.
- Bold, Tessa, Deon P. Filmer, Gayle Martin, Ezequiel Molina, Christophe Rockmore, Brian William Stacy, Jakob Svensson, and Waly Wane.** 2017. “What do teachers know and do? Does it matter? Evidence from primary schools in Africa.” The World Bank Policy Research Working Paper Series 7956.
- Borman, G. D., J. G. Benson, and L. Overman.** 2009. “A randomized field trial of the Fast ForWord Language computer-based training program.” *Educational Evaluation and Policy Analysis*, 31(1): 82–106.
- Buswell, Guy Thomas, and Charles Hubbard Judd.** 1925. *Summary of educational investigations relating to arithmetic.* University of Chicago.
- Campuzano, L., M. Dynarski, R. Agodini, K. Rall, and A. Pendleton.** 2009. “Effectiveness of reading and mathematics software products: Findings from two student cohorts.” *Unpublished manuscript.* Washington, DC: Mathematica Policy Research.
- Carrillo, Paul E, Mercedes Onofa, and Juan Ponce.** 2010. “Information technology and student achievement: Evidence from a randomized experiment in Ecuador.” (IDB Working Paper No. IDB-WP-223). Washington, DC: Inter-American Development Bank.

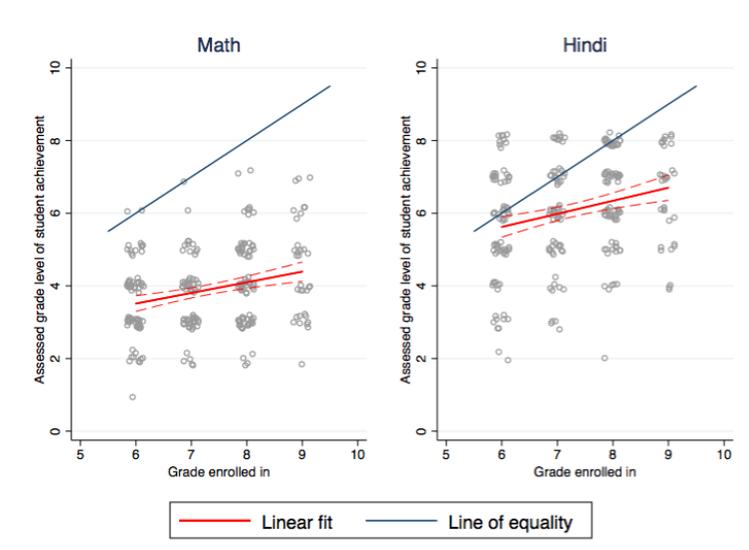
- Cristia, Julian, Pablo Ibararán, Santiago Cueto, Ana Santiago, and Eugenio Severín.** 2012. “Technology and child development: Evidence from the One Laptop per Child program.” (IDB Working Paper No. IDB-WP-304). Washington, DC: Inter-American Development Bank.
- Das, Jishnu, and Tristan Zajonc.** 2010. “India shining and Bharat drowning: Comparing two Indian states to the worldwide distribution in mathematics achievement.” *Journal of Development Economics*, 92(2): 175–187.
- Dewan, Hridaykant, Namrita Batra, and Inder Singh Chabra.** 2012. “Transforming the Elementary Mathematics Curriculum: Issues and Challenges.” In *Mathematics Education in India: Status and Outlook.*, ed. R. Ramanujan and K. Subramaniam. Mumbai, India:Homi Bhabha Centre for Science Education, Tata Institute for Fundamental Research.
- Duflo, E., P. Dupas, and M. Kremer.** 2011. “Peer effects, teacher incentives, and the impact of tracking: Evidence from a randomized evaluation in Kenya.” *American Economic Review*, 101: 1739–1774.
- Dynarski, M., R. Agodini, S. Heaviside, T. Novak, N. Carey, L. Campuzano, B. Means, R. Murphy, W. Penuel, H. Javitz, D. Emery, and W. Sussex.** 2007. “Effectiveness of reading and mathematics software products: Findings from the first student cohort.” *Unpublished manuscript*. Washington, DC: Mathematica Policy Research.
- Fairlie, R. W., and J. Robinson.** 2013. “Experimental Evidence on the Effects of Home Computers on Academic Achievement among Schoolchildren.” *American Economic Journal: Applied Economics*, 5(3): 211–240.
- Goolsbee, Austan, and Jonathan Guryan.** 2006. “The impact of Internet subsidies in public schools.” *The Review of Economics and Statistics*, 88(2): 336–347.
- Kothari, Brij, Avinash Pandey, and Amita R Chudgar.** 2004. “Reading out of the “idiot box”: Same-language subtitling on television in India.” *Information Technologies & International Development*, 2(1): pp–23.
- Kothari, Brij, Joe Takeda, Ashok Joshi, and Avinash Pandey.** 2002. “Same language subtitling: a butterfly for literacy?” *International Journal of Lifelong Education*, 21(1): 55–66.
- Kumar, Ruchi S., Hridaykant Dewan, and K.Subramaniam.** 2012. “The preparation and professional development of mathematics teachers.” In *Mathematics Education in India: Status and Outlook.*, ed. R. Ramanujan and K. Subramaniam. Mumbai, India:Homi Bhabha Centre for Science Education, Tata Institute for Fundamental Research.
- Lai, Fang, Linxiu Zhang, Qinghe Qu, Xiao Hu, Yaojiang Shi, Matthew Boswell, and Scott Rozelle.** 2012. “Does computer-assisted learning improve learning outcomes? Evidence from a randomized experiment in public schools in rural minority areas in Qinghai, China.” (REAP Working Paper No. 237). Rural Education Action Program (REAP). Stanford, CA.

- Lai, Fang, Linxiu Zhang, Xiao Hu, Qinghe Qu, Yaojiang Shi, Yajie Qiao, Matthew Boswell, and Scott Rozelle.** 2013. “Computer assisted learning as extracurricular tutor? Evidence from a randomised experiment in rural boarding schools in Shaanxi.” *Journal of Development Effectiveness*, 52(2): 208–231.
- Lai, Fang, Renfu Luo, Linxiu Zhang, and Scott Huang, Xinzhe Rozelle.** 2015. “Does computer-assisted learning improve learning outcomes? Evidence from a randomized experiment in migrant schools in Beijing.” *Economics of Education*, 47: 34–48.
- Leuven, Edwin, Mikael Lindahl, Hessel Oosterbeek, and Dinand Webbink.** 2007. “The effect of extra funding for disadvantaged pupils on achievement.” *The Review of Economics and Statistics*, 89(4): 721–736.
- Linden, L. L.** 2008. “Complement or substitute? The effect of technology on student achievement in India.” Unpublished manuscript. Abdul Latif Jameel Poverty Action Lab (J-PAL). Cambridge, MA.
- Machin, Stephen, Sandra McNally, and Olmo Silva.** 2007. “New technology in schools: Is there a payoff?” *The Economic Journal*, 117(522): 1145–1167.
- Malamud, Ofer, and C. Pop-Eleches.** 2011. “Home computer use and the development of human capital.” *The Quarterly Journal of Economics*, 126: 987–1027.
- Mo, Di, Johan Swinnen, Linxiu Zhang, Hongmei Yi, Qinghe Qu, Matthew Boswell, and Scott Rozelle.** 2013. “Can one-to-one computing narrow the digital divide and the educational gap in China? The case of Beijing migrant schools.” *World Development*, 46: 14–29.
- Mo, Di, Linxiu Zhang, Renfu Luo, Qinghe Qu, Weiming Huang, Jiafu Wang, Yajie Qiao, Matthew Boswell, and Scott Rozelle.** 2014*a*. “Integrating computer-assisted learning into a regular curriculum: Evidence from a randomised experiment in rural schools in Shaanxi.” *Journal of Development Effectiveness*, 6: 300–323.
- Mo, Di, L. Zhang, J. Wang, W. Huang, Y. Shi, M. Boswell, and S. Rozelle.** 2014*b*. “The persistence of gains in learning from computer assisted learning (CAL): Evidence from a randomized experiment in rural schools in Shaanxi province in China.” *Unpublished manuscript*. Stanford, CA: Rural Education Action Program (REAP).
- Mo, Di, Yu Bai, Matthew Boswell, and Scott Rozelle.** 2016. “Evaluating the effectiveness of computers as tutors in China.”
- Morgan, P., and S. Ritter.** 2002. “An experimental study of the effects of Cognitive Tutor Algebra I on student knowledge and attitude.” Pittsburg, PA: Carnegie Learning.
- Muralidharan, Karthik, and Abhijeet Singh.** 2018. “Improving Public Sector Governance at Scale: Experimental Evidence from a Large-Scale School Governance Improvement Program in India.” University of California San Diego mimeo., San Diego, CA.
- Muralidharan, Karthik, Jishnu Das, Alaka Holla, and Aakash Mohpal.** 2017. “The fiscal cost of weak governance: Evidence from teacher absence in India.” *Journal of Public Economics*, 145: 116–135.

- Murphy, R., W. Penuel, B. Means, C. Korbak, and A. Whaley.** 2001. “E-DESK: A review of recent evidence on the effectiveness of discrete educational software.” *Unpublished manuscript*. Menlo Park, CA: SRI International.
- NCERT.** 2006. *Position Paper of the National Focus Group on Curriculum, Syllabus and Textbooks*. National Council of Educational Research and Training, New Delhi.
- PASEC.** 2015. *Programme d’Analyse des Systèmes éducatifs de la Confédération (PASEC) 2014: Education System Performance in Francophone Africa, Competencies and Learning Factors in Primary Education*. PASEC, Dakar, Senegal.
- Pearson, P.D., R.E. Ferdig, R.L. Blomeyer Jr., and J. Moran.** 2005. “The effects of technology on reading performance in the middle-school grades: A meta-analysis with recommendations for policy.” *Unpublished manuscript*. Naperville, IL: Learning Point Associates.
- Pratham.** 2016. *Annual Status of Education Report 2015*. Pratham, New Delhi.
- Pritchett, Lant, and Amanda Beatty.** 2015. “Slow down, you’re going too fast: Matching curricula to student skill levels.” *International Journal of Educational Development*, 40: 276–288.
- Radatz, Hendrik.** 1979. “Error analysis in mathematics education.” *Journal for Research in mathematics Education*, 163–172.
- Rampal, Anita, and Jayasree Subramaniam.** 2012. “Transforming the Elementary Mathematics Curriculum: Issues and Challenges.” In *Mathematics Education in India: Status and Outlook*, ed. R. Ramanujan and K. Subramaniam. Mumbai, India: Homi Bhabha Centre for Science Education, Tata Institute for Fundamental Research.
- Rockoff, Jonah E.** 2015. “Evaluation report on the School of One i3 expansion.” *Unpublished manuscript*. New York, NY: Columbia University.
- Rouse, Cecilia Elena, and Alan B Krueger.** 2004. “Putting computerized instruction to the test: A randomized evaluation of a “scientifically based” reading program.” *Economics of Education Review*, 23(4): 323–338.
- SACMEQ.** 2007. *Southern and Eastern Africa Consortium for Monitoring Educational Quality (SACMEQ), Various years*. University of Botswana, Gaborone. <http://www.sacmeq.org/>.
- SAFED.** 2017. *Annual Status of Education Report (ASER-Pakistan) 2016*. South Asia Forum for Education Development, Lahore.
- Sankar, Deepa, and Toby Linden.** 2014. “How much and what kind of teaching is there in elementary education in India? Evidence from three states.” (South Asia Human Development Sector Report No. 67). Washington, DC: The World Bank.
- Singh, Abhijeet.** 2015. “Private school effects in urban and rural India: Panel estimates at primary and secondary school ages.” *Journal of Development Economics*, 113: 16–32.

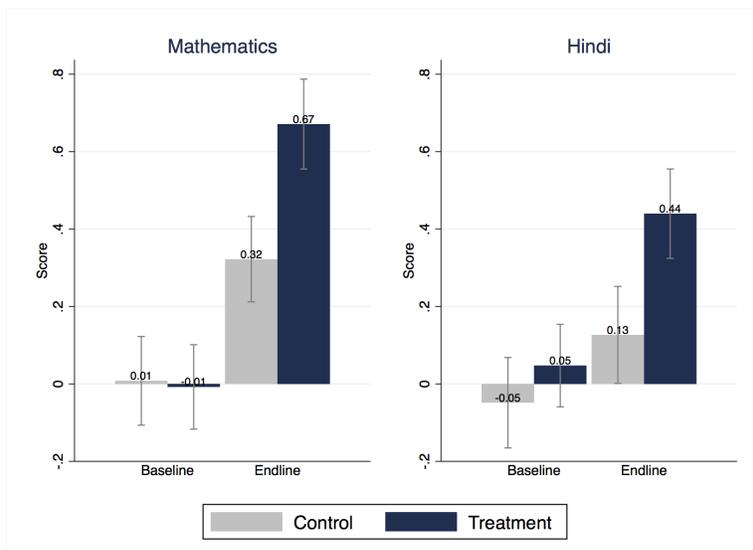
- Sinha, S., R. Banerji, and W. Wadhwa.** 2016. *Teacher performance in Bihar, India: Implications for education*. The World Bank, Washington D.C.
- Uwezo.** 2016. *Are Our Children Learning? Uwezo Uganda 6th Learning Assessment Report*. Twaweza East Africa, Kampala.
- van der Linden, Wim J, and Ronald K Hambleton.** 2013. *Handbook of modern item response theory*. Springer Science & Business Media.
- Waxman, H.C., M.-F. Lin, and G.M. Michko.** 2003. "A meta-analysis of the effectiveness of teaching and learning with technology on student outcomes." *Unpublished manuscript*. CambridgeNaperville, IL: Learning Point Associates.
- Wise, B. W., and R. K. Olson.** 1995. "Computer-based phonological awareness and reading instruction." *Annals of Dyslexia*, 45: 99–122.
- World Bank.** 2016. *What is happening inside classrooms in Indian secondary schools? A time on task study in Madhya Pradesh and Tamil Nadu*. The World Bank, Washington D.C.
- World Bank.** 2018. *World Development Report 2018: Learning to Realize Education's Promise*. World Bank, Washington, DC.

Figure 1: Assessed levels of student achievement vs. current grade enrolled in school



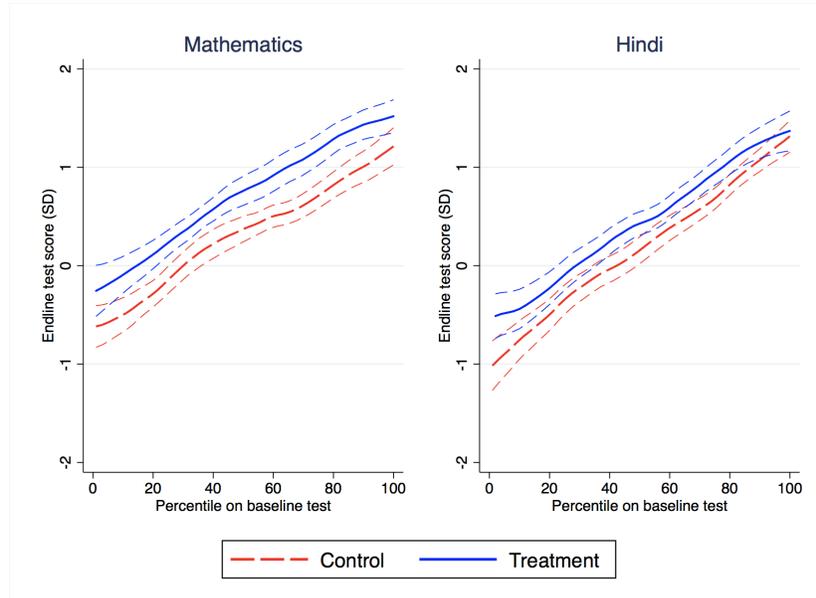
*Note:* This figure shows, for treatment group, the estimated level of student achievement (determined by the Mindspark CAL program) plotted against the grade they are enrolled in. These data are from the *initial* diagnostic test, and do not reflect any instruction provided by Mindspark. In both subjects, we find three main patterns: (a) there is a general deficit between average attainment and grade-expected norms; (b) this deficit is larger in later grades and (c) within each grade, there is a wide dispersion of student achievement.

Figure 2: Mean difference in test scores between lottery winners and losers



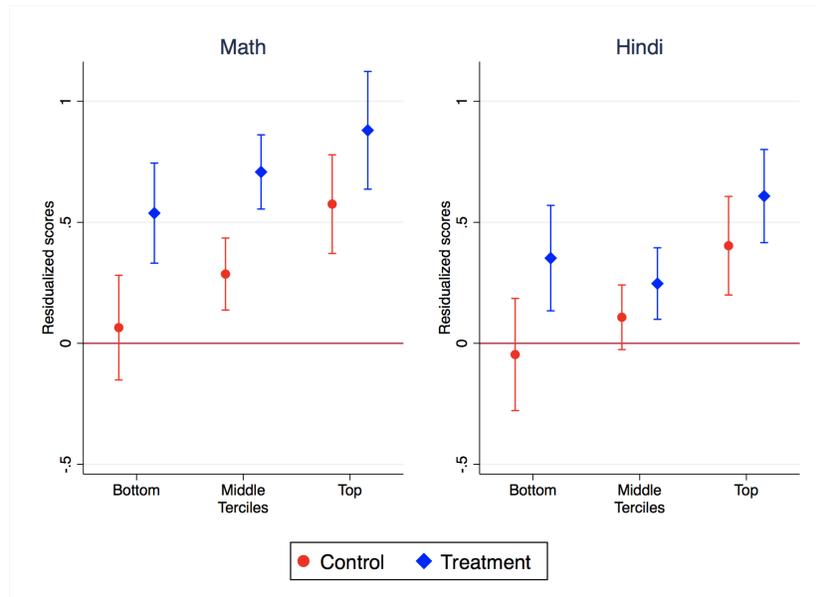
*Note:* This figure shows mean of test scores, normalized with reference to baseline, across treatment and control groups in the two rounds of testing with 95% confidence intervals. Test scores were linked within-subject through IRT models, pooling across grades and across baseline and endline, and are normalized to have a mean of zero and a standard deviation of one in the baseline. Whereas baseline test scores were balanced between lottery-winners and lottery-losers, endline scores are significantly higher for the treatment group.

Figure 3: Non-parametric investigation of treatment effects by baseline percentiles



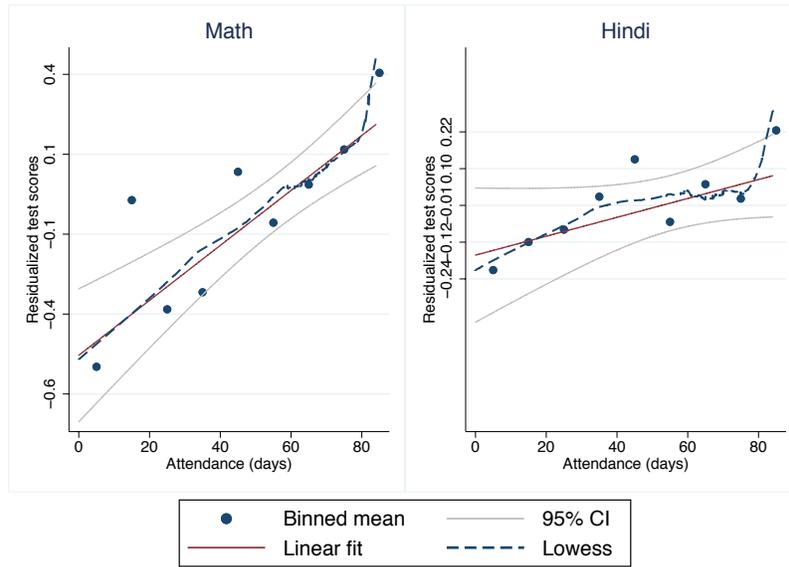
*Note:* The figures present kernel-weighted local mean smoothed plots which relate endline test scores to percentiles in the baseline achievement, separately for the treatment and control groups, alongside 95% confidence intervals. At all percentiles of baseline achievement, treatment group students score higher in the endline test than the control group, with no strong evidence of differential absolute magnitudes of gains across the distribution.

Figure 4: Growth in achievement in treatment and control groups



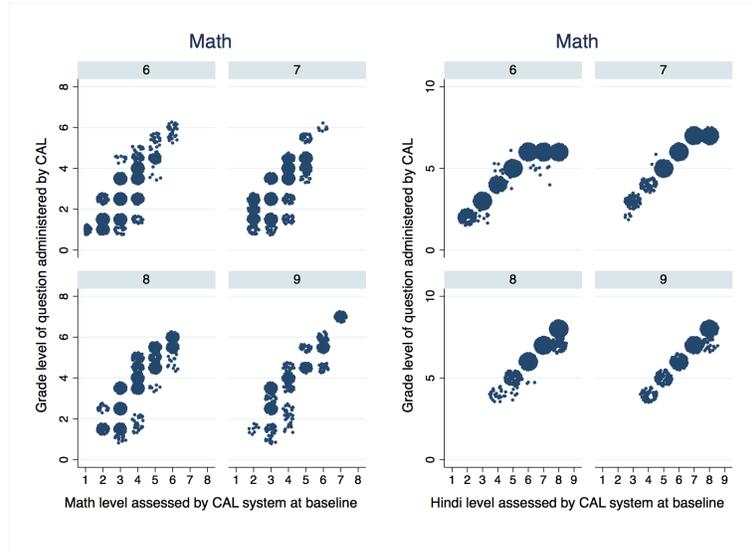
*Note:* This figure shows the growth in student achievement in the treatment and control groups in math and Hindi, as in Table 5. Students in the treatment group see positive value-added in all terciles whereas we cannot reject the null of no academic progress for students in the bottom tercile in the control group.

Figure 5: Dose response relationship



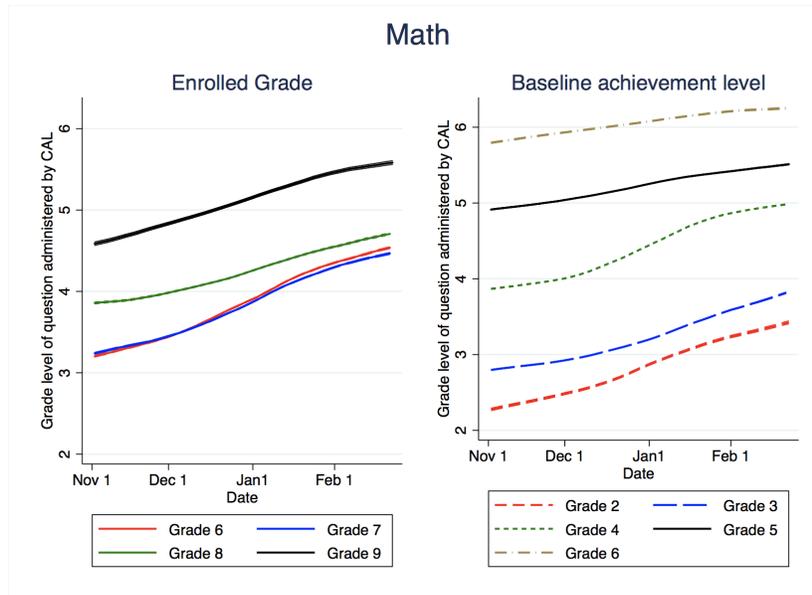
*Note:* This figure explores the relationship between value-added and attendance in the Mindspark program among the lottery-winners. It presents the mean value-added in bins of attendance along with a linear fit and a lowess smoothed non-parametric plot.

Figure 6: Precise customization of instruction by the Mindspark CAL program



*Note:* This figure shows, for treatment group, the grade level of questions administered by the computer adaptive system to students on a single day near the beginning of the intervention. In each grade of enrolment, actual level of student attainment estimated by the CAL software differs widely; this wide range is covered through the customization of instructional content by the CAL software.

Figure 7: Dynamic updating and individualization of content in Mindspark



*Note:* This figure shows kernel-weighted local mean smoothed lines relating the level of difficulty of the math questions administered to students in the treatment group with the date of administration. The left panel presents separate lines by the actual grade of enrolment. The right panel presents separate lines by the level of achievement assessed at baseline by the CAL software. Note that 95% confidence intervals are plotted as well but, given the large data at our disposal, estimates are very precise and the confidence intervals are narrow enough to not be visually discernible.

Table 1: Sample descriptives and balance on observables

	Mean (treatment)	Mean (control)	Difference	SE	N (treatment)	N (control)
<u>Panel A: All students in the baseline sample</u>						
<i>Demographic characteristics</i>						
Female	0.76	0.76	0.004	0.034	314	305
Age (years)	12.67	12.41	0.267	0.143	230	231
SES index	-0.03	0.04	-0.070	0.137	314	305
<i>Grade in school</i>						
Grade 4	0.01	0.01	-0.003	0.007	305	299
Grade 5	0.01	0.02	-0.007	0.010	305	299
Grade 6	0.27	0.30	-0.035	0.037	305	299
Grade 7	0.26	0.26	0.005	0.036	305	299
Grade 8	0.30	0.28	0.017	0.037	305	299
Grade 9	0.15	0.13	0.024	0.028	305	299
<i>Baseline test scores</i>						
Math	-0.01	0.01	-0.016	0.081	313	304
Hindi	0.05	-0.05	0.096	0.080	312	305
Present at endline	0.85	0.90	-0.048	0.027	314	305
<u>Panel B: Only students present in Endline</u>						
<i>Demographic characteristics</i>						
Female	0.77	0.76	0.013	0.036	266	273
Age (years)	12.61	12.37	0.243	0.156	196	203
SES index	-0.17	0.03	-0.193	0.142	266	273
<i>Grade in school</i>						
Grade 4	0.01	0.01	-0.003	0.008	258	269
Grade 5	0.01	0.02	-0.011	0.011	258	269
Grade 6	0.28	0.30	-0.022	0.040	258	269
Grade 7	0.26	0.26	-0.001	0.038	258	269
Grade 8	0.30	0.28	0.020	0.040	258	269
Grade 9	0.14	0.12	0.017	0.029	258	269
<i>Baseline test scores</i>						
Math	-0.03	-0.00	-0.031	0.086	265	272
Hindi	0.05	-0.07	0.124	0.084	266	273

*Note:* Treatment and control here refer to groups who were randomly assigned to receive an offer of Mindspark voucher till March 2016. Variables used in this table are from the baseline data collection in September 2015. The data collection consisted of two parts: (a) a self-administered student survey, from which demographic characteristics are taken and (b) assessment of skills in math and Hindi, administered using pen-and-paper tests. Tests were designed to cover wide ranges of achievement and to be linked across grades, as well as between baseline and endline assessments, using common items. Scores are scaled here using Item Response theory models and standardized to have a mean of zero and standard deviation of one in the baseline. The SES index refers to a wealth index generated using the first factor from a Principal Components Analysis consisting of indicators for ownership of various consumer durables and services in the household.

Table 2: Intent-to-treat (ITT) Effects in a regression framework

	(1)	(2)	(3)	(4)
	Dep var: Standardized IRT scores (endline)			
	Math	Hindi	Math	Hindi
Treatment	0.37 (0.064)	0.23 (0.062)	0.37 (0.064)	0.24 (0.071)
Baseline score	0.58 (0.042)	0.71 (0.040)	0.57 (0.051)	0.68 (0.033)
Constant	0.33 (0.044)	0.17 (0.044)	0.32 (0.031)	0.17 (0.035)
Strata fixed effects	Y	Y	N	N
Observations	535	537	535	537
R-squared	0.403	0.493	0.397	0.473

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher till March 2016. Tests in both math and Hindi were designed to cover wide ranges of achievement and to be linked across grades, as well as between baseline and endline assessments, using common items. Scores are scaled here using Item Response theory models and standardized to have a mean of zero and standard deviation of one in the baseline.

Table 3: Treatment effect by specific competence assessed

(a) Mathematics

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
<i>Dep var:</i> Proportion of questions answered correctly							
	Arithmetic computation	Word problems - computation	Data interpretation	Fractions and decimals	Geometry and Measurement	Numbers	Pattern recognition
Treatment	0.078 (0.016)	0.072 (0.016)	0.042 (0.021)	0.071 (0.020)	0.15 (0.024)	0.15 (0.022)	0.11 (0.028)
Baseline math score	0.13 (0.0080)	0.11 (0.010)	0.082 (0.015)	0.093 (0.012)	0.052 (0.014)	0.068 (0.012)	0.099 (0.016)
Constant	0.66 (0.0079)	0.50 (0.0076)	0.38 (0.010)	0.33 (0.010)	0.39 (0.012)	0.45 (0.011)	0.36 (0.014)
Observations	537	537	537	537	537	537	537
R-squared	0.357	0.229	0.097	0.157	0.097	0.135	0.112

(b) Hindi

	(1)	(2)	(3)	(4)
<i>Dep var:</i> Proportion of questions answered correctly				
	Sentence completion	Retrieve explicitly stated information	Make straightforward inferences	Interpret and integrate ideas and information
Treatment	0.046 (0.022)	0.045 (0.016)	0.065 (0.022)	0.053 (0.015)
Baseline Hindi score	0.13 (0.017)	0.14 (0.0075)	0.15 (0.011)	0.067 (0.013)
Constant	0.72 (0.011)	0.59 (0.0078)	0.51 (0.011)	0.31 (0.0077)
Observations	539	539	539	539
R-squared	0.182	0.380	0.309	0.136

*Note:* Robust standard errors in parentheses. The tables above show the impact of the treatment on specific competences. The dependent variable in each regression is the proportion of questions related to the competence that a student answered correctly. All test questions were multiple choice items with four choices. Baseline scores are IRT scores in the relevant subject from the baseline assessment. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher till March 2016. All regressions include randomization strata fixed effects.

Table 4: Heterogeneity in treatment effect by gender, socio-economic status and baseline score

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Dep var: Standardized IRT scores (endline)</i>						
COVARIATES	<u>Female</u>		<u>SES</u>		<u>Baseline score</u>	
	Math	Hindi	Math	Hindi	Math	Hindi
Treatment	0.47 (0.14)	0.27 (0.095)	0.38 (0.065)	0.26 (0.062)	0.37 (0.064)	0.24 (0.070)
Covariate	-0.050 (0.14)	0.21 (0.15)	-0.0028 (0.035)	0.099 (0.021)	0.53 (0.076)	0.70 (0.047)
<b>Interaction</b>	<b>-0.13</b> <b>(0.14)</b>	<b>-0.046</b> <b>(0.12)</b>	<b>0.023</b> <b>(0.050)</b>	<b>-0.0041</b> <b>(0.041)</b>	<b>0.081</b> <b>(0.087)</b>	<b>-0.047</b> <b>(0.071)</b>
Observations	535	537	535	537	535	537
R-squared	0.399	0.474	0.398	0.494	0.399	0.473

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher till March 2016. The SES index and test scores are defined as in Tables 1 and 2 respectively. All regressions include strata fixed effects and control for baseline subject scores.

Table 5: Heterogeneity in treatment effect by within-grade terciles

	(1)	(2)
<i>Dep var: Standardized IRT scores (endline)</i>		
VARIABLES	Math	Hindi
Bottom Tercile	0.13 (0.098)	-0.072 (0.10)
Middle Tercile	0.30 (0.073)	0.14 (0.068)
Top Tercile	0.53 (0.092)	0.46 (0.085)
Treatment	0.33 (0.12)	0.41 (0.12)
Treatment*Middle Tercile	0.083 (0.16)	-0.30 (0.16)
Treatment*Top Tercile	0.068 (0.16)	-0.24 (0.15)
Baseline test score	0.44 (0.066)	0.58 (0.062)
Observations	535	537
R-squared	0.545	0.545

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher till March 2016. Test scores are scaled as in Table 2.

Table 6: Treatment effect on items linked to grade levels

	(1)	(2)	(3)	(4)
VARIABLES	Dep var: Proportion of questions answered correctly			
	Math		Hindi	
	At or above grade level	Below grade level	At or above grade level	Below grade level
Treatment	0.0089 (0.032)	0.081 (0.013)	0.063 (0.027)	0.050 (0.014)
Baseline subject score	0.047 (0.022)	0.099 (0.0069)	0.13 (0.016)	0.13 (0.0068)
Constant	0.31 (0.022)	0.49 (0.0089)	0.45 (0.019)	0.58 (0.0100)
Observations	291	511	292	513
R-squared	0.029	0.346	0.250	0.399

*Note:* Robust standard errors in parentheses. The table shows the impact of the treatment (winning a randomly-assigned voucher) on questions below or at/above grade levels for individual students. The dependent variable is the proportion of questions that a student answered correctly. All test questions were multiple choice items with four choices. Our endline assessments, designed to be informative at students' actual levels of achievement, did not include many items at grade 8 level and above. Therefore students in Grades 8 and 9 are not included in regressions on items at/above grade level. Baseline scores are IRT scores in the relevant subject from the baseline assessment. All regressions include randomization strata fixed effects.

Table 7: Treatment effect on school exams

	(1)	(2)	(3)	(4)	(5)	(6)
VARIABLES	Dep var: Standardized test scores					
	Hindi	Math	Science	Social Sciences	English	Aggregate
Treatment	0.196 (0.088)	0.059 (0.076)	0.077 (0.092)	0.108 (0.110)	0.081 (0.105)	0.100 (0.080)
Baseline Hindi score	0.487 (0.092)		0.292 (0.064)	0.414 (0.096)	0.305 (0.067)	0.336 (0.058)
Baseline math score		0.303 (0.041)	0.097 (0.036)	0.262 (0.058)	0.120 (0.052)	0.167 (0.039)
Constant	1.006 (1.103)	0.142 (0.423)	0.931 (0.347)	1.062 (0.724)	1.487 (0.740)	0.977 (0.600)
Observations	597	596	595	594	597	597
R-squared	0.190	0.073	0.121	0.177	0.144	0.210

*Note:* Robust standard errors in parentheses. This table shows the effect of receiving the Mindspark voucher on the final school exams, held in March 2016 after the completion of the intervention. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher till March 2016. Test scores in the school exams are normalized within school\*grade to have a mean of zero and a standard deviation of one in the control group. All regressions include grade and school fixed effects.

Table 8: Heterogeneous effects on school tests, by terciles of baseline achievement

VARIABLES	(1)	(2)	(3)	(4)	(5)	(6)
	Hindi	Math	Science	Soc. Sc.	English	Aggregate
Treatment	0.058 (0.14)	-0.40 (0.11)	-0.15 (0.16)	-0.17 (0.16)	0.14 (0.11)	-0.052 (0.099)
Treatment*Tercile 2	0.11 (0.23)	0.55 (0.20)	0.31 (0.18)	0.15 (0.24)	-0.30 (0.14)	0.063 (0.16)
Treatment*Tercile 3	0.29 (0.18)	0.82 (0.27)	0.36 (0.19)	0.65 (0.24)	0.14 (0.15)	0.38 (0.13)
Tercile 2	-0.35 (0.27)	-0.27 (0.23)	-0.39 (0.18)	-0.61 (0.29)	0.14 (0.17)	-0.29 (0.19)
Tercile 3	-0.23 (0.31)	-0.48 (0.21)	-0.32 (0.21)	-1.02 (0.38)	0.096 (0.20)	-0.37 (0.21)
Baseline Hindi score	0.53 (0.17)		0.35 (0.083)	0.67 (0.19)	0.25 (0.11)	0.40 (0.10)
Baseline Math score		0.33 (0.072)	0.096 (0.033)	0.27 (0.058)	0.11 (0.051)	0.16 (0.039)
Constant	1.28 (1.09)	0.47 (0.40)	1.27 (0.39)	1.76 (0.76)	1.29 (0.74)	1.24 (0.60)
Observations	597	596	595	594	597	597
R-squared	0.201	0.098	0.132	0.203	0.155	0.226

Treatment Effect by tercile (p-values in brackets)

Tercile 1	0.058 [0.67]	-0.40 [0.002]	-0.15 [0.36]	-0.17 [0.31]	0.14 [0.23]	-0.052 [0.61]
Tercile 2	0.17 [0.27]	0.15 [0.28]	0.16 [0.13]	-0.02 [0.94]	-0.16 [0.25]	0.01 [0.92]
Tercile 3	0.348 [0.04]	0.42 [0.07]	0.21 [0.16]	0.48 [0.04]	0.28 [0.08]	0.33 [0.03]

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher till March 2016. Test scores are scaled as in Table 7.

Table 9: Dose-response of Mindspark attendance

	(1)	(2)	(3)	(4)	(5)	(6)
	<i>Dep var:</i> Standardized IRT scores (endline)					
VARIABLES	IV estimates		OLS VA (full sample)		OLS VA (Treatment group)	
	Math	Hindi	Math	Hindi	Math	Hindi
Attendance (days)	0.0067 (0.0011)	0.0043 (0.0011)	0.0072 (0.00090)	0.0037 (0.00091)	0.0086 (0.0018)	0.0030 (0.0018)
Baseline score	0.56 (0.038)	0.68 (0.036)	0.58 (0.042)	0.71 (0.040)	0.62 (0.061)	0.68 (0.052)
Constant			0.31 (0.041)	0.18 (0.041)	0.22 (0.12)	0.24 (0.11)
Observations	535	537	535	537	264	265
R-squared	0.431	0.479	0.429	0.495	0.446	0.445
Angrist-Pischke F-statistic for weak instrument	1207	1244				
Diff-in-Sargan statistic for exogeneity (p-value)	0.14	0.92				
Extrapolated estimates of 90 days' treatment (SD)	0.603	0.39	0.648	0.333	0.77	0.27

*Note:* Robust standard errors in parentheses. Treatment group students who were randomly-selected for the Mindspark voucher offer but who did not take up the offer have been marked as having 0% attendance, as have all students in the control group. Columns (1) and (2) instrument attendance in Mindspark with the randomized allocation of a scholarship and include randomization strata fixed effects, Columns (3) and (4) present OLS value-added models in the full sample, Columns (5) and (6) present OLS value-added models using only data on the lottery-winners. Scores are scaled here as in Table 2.

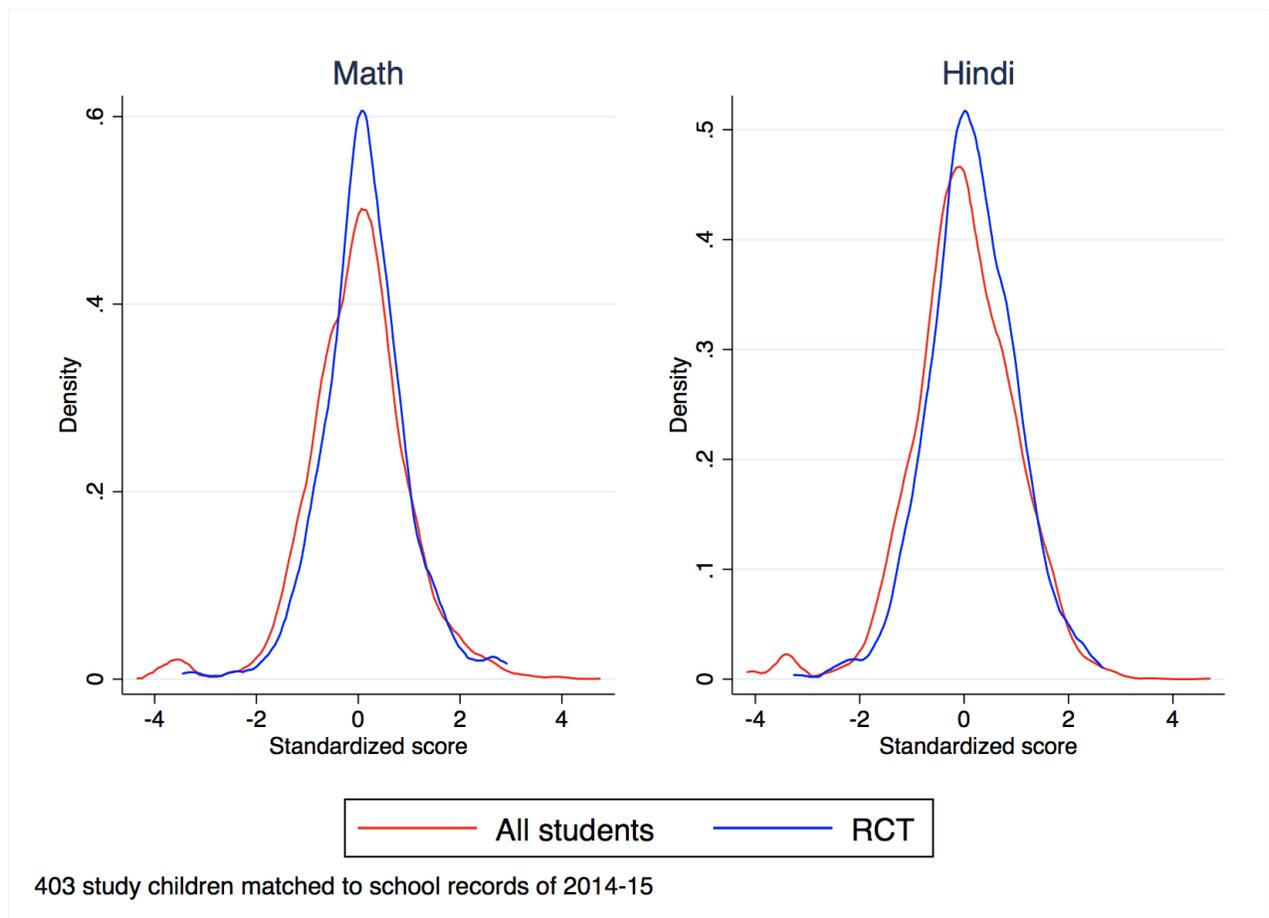
# Online Appendix

## Disrupting Education? Experimental Evidence on Technology-Aided Instruction in India

Karthik Muralidharan and Abhijeet Singh and Alejandro J. Ganimian

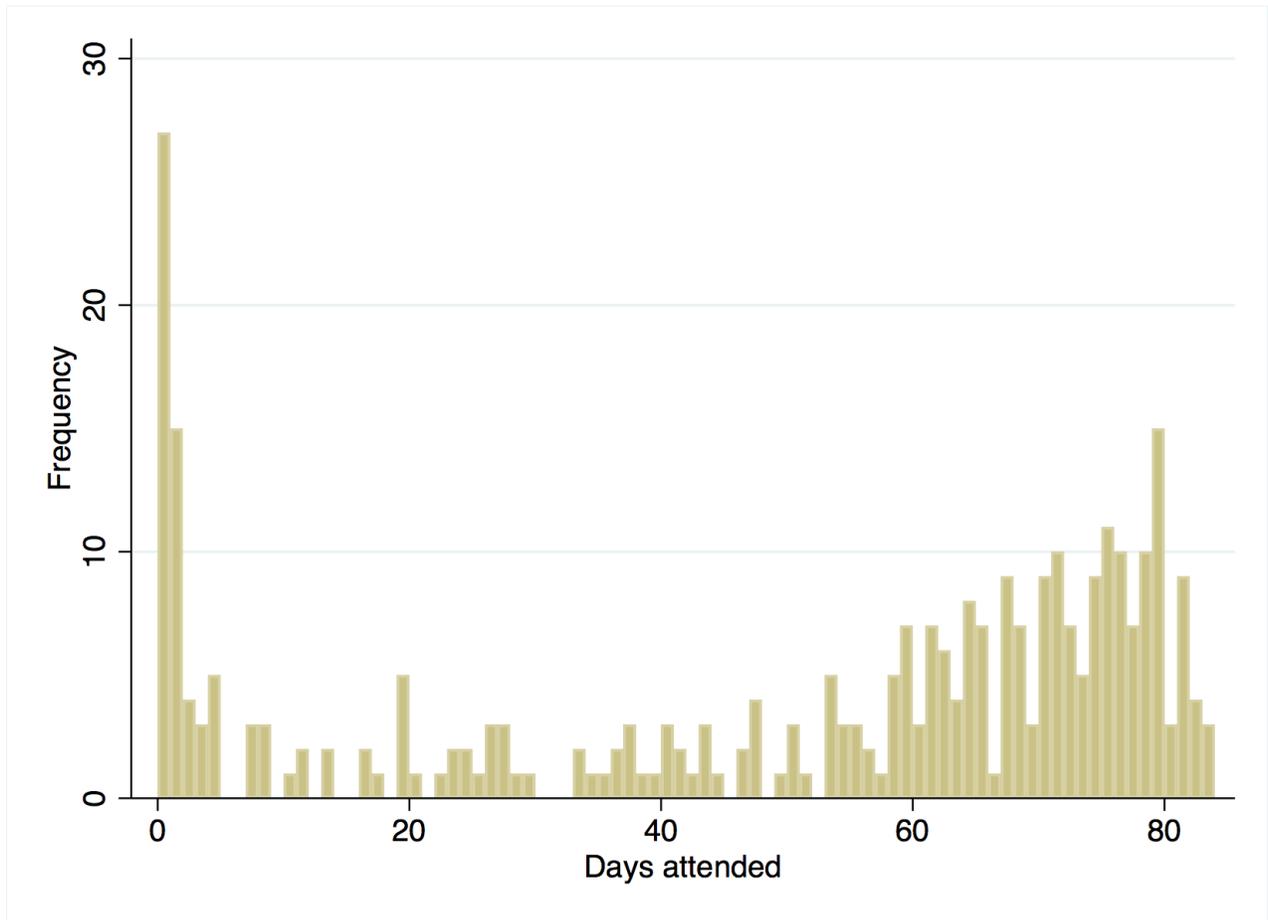
### Appendix A Additional figures and tables

Figure A1: Comparing pre-program achievement of study participants and non-participants



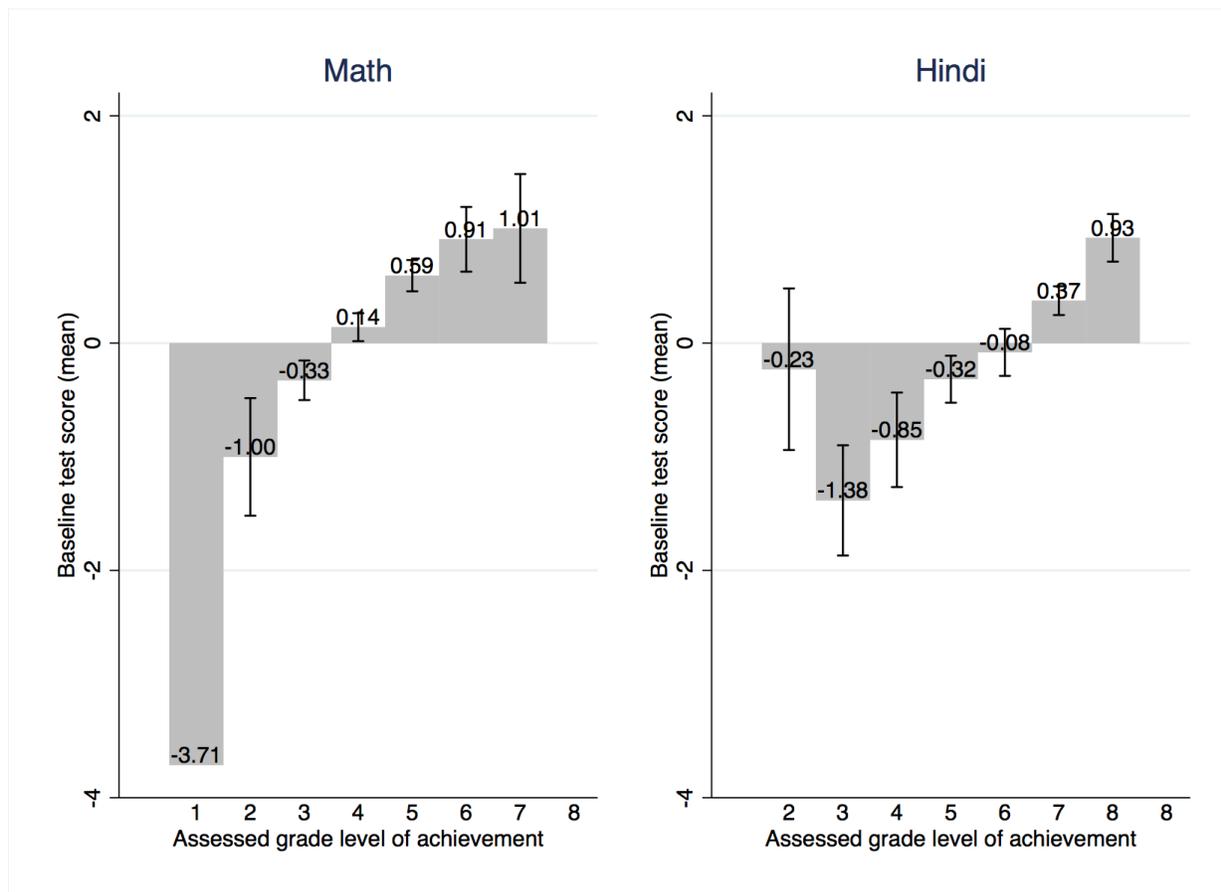
*Note:* The panels compare the final scores for the 2014-15 school year, i.e. the pre-program academic year, for study participants and non-participants. Test scores have been standardized within school\*grade cells. The study participants are positively selected into the RCT in comparison to their peers but the magnitude of selection is modest and there is near-complete common support between the two groups in pre-program academic achievement. See Table A1 for further details.

Figure A2: Distribution of take-up among lottery-winners



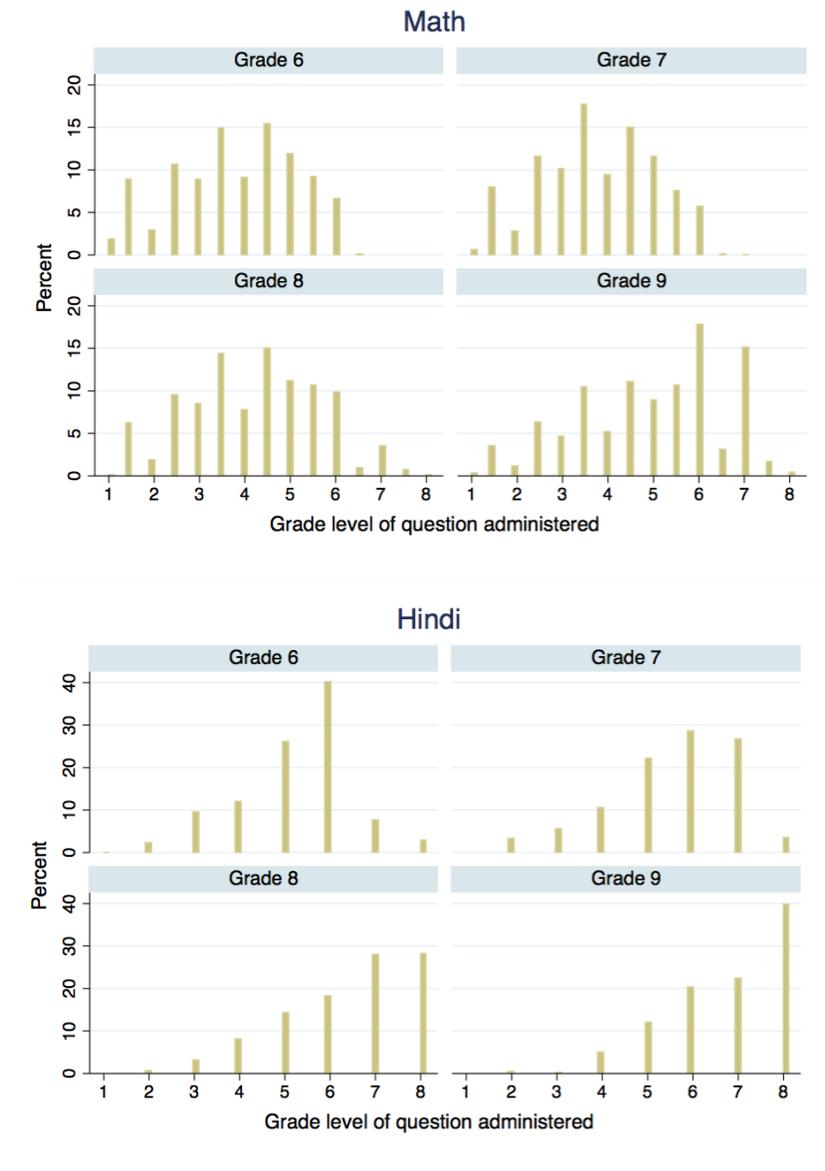
*Note:* This figure shows the distribution of attendance in the Mindspark centers among the lottery-winners. Over the study period, the Mindspark centers were open for 86 working days.

Figure A3: Comparison of Mindspark initial assessment of grade-level of student achievement with (independent) baseline test scores



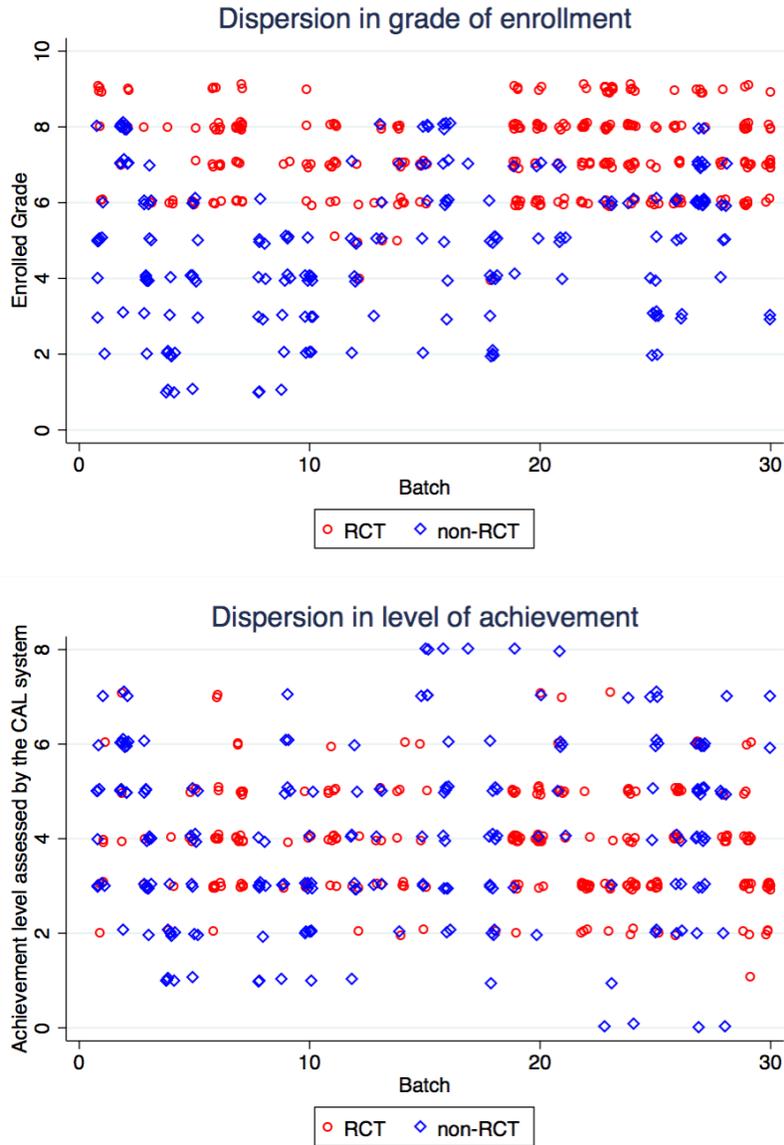
*Note:* The two panels above show mean test scores in Mathematics and Hindi respectively by each level of grade ability as assessed by the Mindspark CAL software at the beginning of the intervention (i.e. soon after the initial baseline) for students in the treatment group. Average test scores on our independently-administered assessments increase with CAL-assessed grade levels of achievement; this serves to validate that the two assessments capture similar variation and that Mindspark assessments of grade ability are meaningful. Only one student was assessed at Grade 1 level in math, and only 10 students at Grade 2 level in Hindi, the lowest categories in our sample in the two subjects. Consequently, scores are very noisy in these categories (and measurement error in the CAL assessments is also likely to be more severe).

Figure A4: Distribution of questions administered by Mindspark CAL system



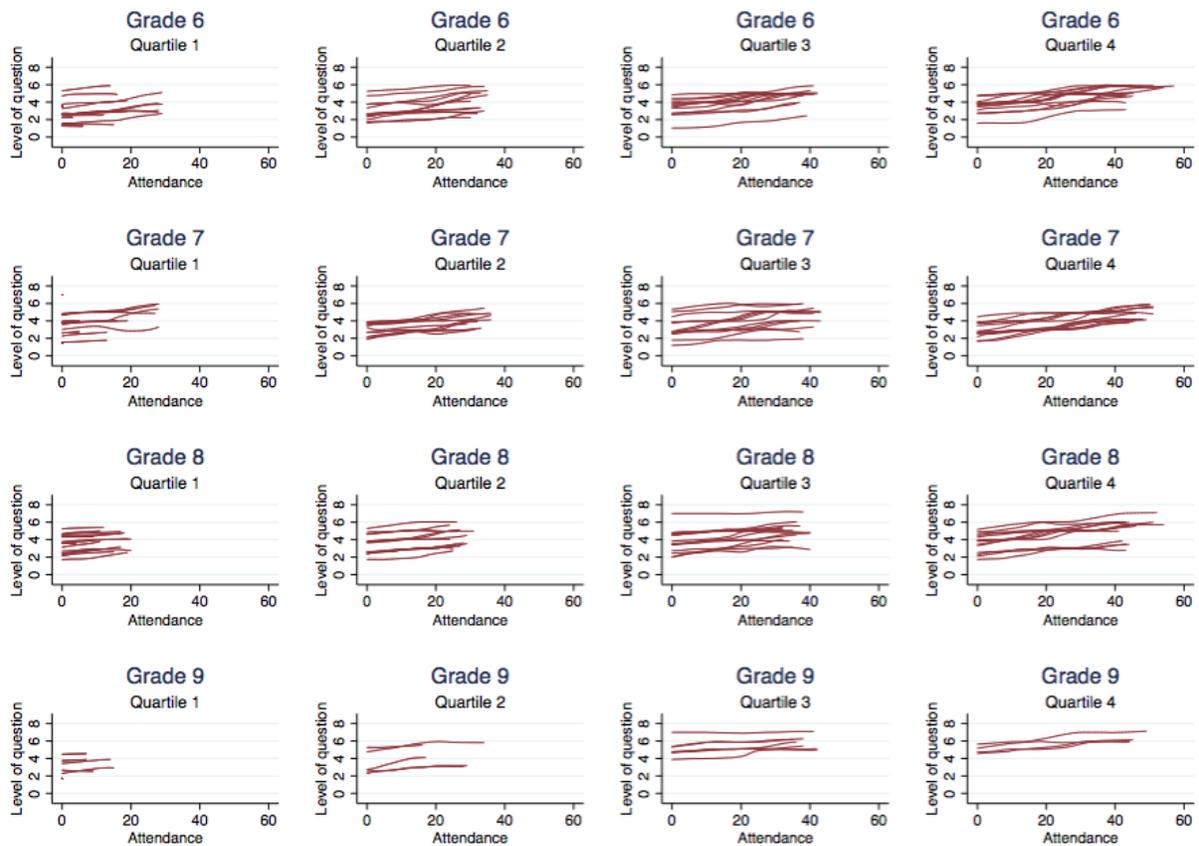
*Note:* The two panels above show the distribution, by grade-level, of the questions that were administered by the Mindspark CAL system over the duration of treatment in both math and Hindi. Note that in math, students received very few questions at the level of the grade they are enrolled in; this reflects the system’s diagnosis of their actual learning levels. In Hindi, by contrast, students received a significant portion of instruction at grade-level competence which is consistent with the initial deficits in achievement in Hindi being substantially smaller than in math (see Figure 1).

Figure A5: Composition of group instruction batches in Mindspark centers



*Note:* The two panels above show the composition of batches in Mindspark centers, by the grade students are enrolled in, and by their level of math achievement, as assessed by the Mindspark CAL system. We separately identify students in the treatment group from fee-paying students who were not part of the study but were part of the small group instruction in each batch. Note that, while our study is focused on students from grades 6-9, the centers cater to students from grades 1-8. Batches are chosen by students based on logistical convenience and hence there is substantial variation in grade levels and student achievement within each batch with little possibility of achievement-based tracking. This confirms that it would not have been possible to customize instruction in the instructor-led small group instruction component of the intervention.

Figure A6: Learning trajectories of individual students in the treatment group



*Note:* Each line in the panels above is a local mean smoothed plot of the grade level of questions administered in Mathematics by the computer adaptive system against the days that the student utilized the Mindspark math software (Attendance). The panels are organized by the grade of enrolment and the within-grade quartile of attendance in Mindspark.

Table A1: Comparing pre-program exam results of study participants and non-participants

	RCT	Non-study	Difference	SE	N(RCT)	N(non-study)
Math	0.13	-0.01	0.14	0.05	409	4067
Hindi	0.16	-0.02	0.17	0.05	409	4067
Science	0.09	-0.01	0.10	0.05	409	4067
Social Science	0.13	-0.01	0.15	0.05	409	4067
English	0.14	-0.01	0.15	0.05	409	4067

*Note:* This table presents the mean scores of study participants and non-participants, standardized within each school\*grade, in the 2014-15 school year. Study participants are, on average, positively selected compared to their peers.

Table A2: Intent-to-treat (ITT) effects with within-grade normalized test scores

	(1)	(2)	(3)	(4)
	Dep var: Standardized IRT scores (endline)			
	Math	Hindi	Math	Hindi
Treatment	0.38 (0.068)	0.23 (0.066)	0.38 (0.069)	0.23 (0.071)
Baseline score	0.59 (0.045)	0.72 (0.039)	0.58 (0.051)	0.70 (0.031)
Constant	0.33 (0.047)	0.20 (0.046)	0.33 (0.034)	0.19 (0.035)
Strata fixed effects	Y	Y	N	N
Observations	523	525	523	525
R-squared	0.384	0.480	0.380	0.470

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher. The SES index refers to a wealth index generated using the first factor from a Principal Components Analysis consisting of indicators for ownership of various consumer durables and services in the household. Tests in both math and Hindi were designed to cover wide ranges of ability and to be linked across grades, as well as between baseline and endline assessments, using common items. Scores are scaled here using Item Response theory models and standardized to have a mean of zero and standard deviation of one in the baseline in each grade.

Table A3: Heterogeneous effects on independent tests, by terciles of baseline achievement

VARIABLES	(1)	(2)	(3)	(4)
	Dep var: Proportion correct			
	Math		Hindi	
	At or above grade level	Below grade level	At or above grade level	Below grade level
Treatment	-0.030 (0.054)	0.059 (0.020)	0.095 (0.043)	0.10 (0.026)
Treatment*Tercile 2	0.036 (0.073)	0.056 (0.029)	-0.053 (0.065)	-0.071 (0.037)
Treatment*Tercile 3	0.13 (0.080)	0.023 (0.032)	-0.044 (0.062)	-0.079 (0.033)
Tercile 1	0.24 (0.045)	0.45 (0.017)	0.39 (0.041)	0.49 (0.022)
Tercile 2	0.26 (0.037)	0.46 (0.015)	0.38 (0.030)	0.58 (0.018)
Tercile 3	0.39 (0.042)	0.54 (0.018)	0.55 (0.037)	0.67 (0.019)
Baseline subject score	-0.015 (0.032)	0.069 (0.010)	0.087 (0.023)	0.084 (0.011)
Observations	291	511	292	513
R-squared	0.096	0.371	0.301	0.433
Total Treatment Effect by tercile (p-values in brackets)				
Tercile 1	-0.030 [0.58]	0.059 [0.00]	0.095 [0.03]	0.10 [0.00]
Tercile 2	0.006 [0.91]	0.115 [0.00]	0.042 [0.38]	0.029 [0.24]
Tercile 3	0.10 [0.08]	0.082 [0.00]	0.051 [0.25]	0.021 [0.26]

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher. The total treatment effect by tercile is the sum of the coefficient on treatment and the interaction of the relevant tercile with the treatment. We report, in square brackets below the aggregate treatment effect in each tercile, p-values from an F-test of the hypothesis that this sum of the two coefficients is zero. The dependent variable and baseline scores are scaled as in Table 6

Table A4: Correlates of attendance

VARIABLES	(1)	(2)	(3)	(4)
	Attendance (days)			
Female	3.90 (3.90)	2.65 (3.92)	3.03 (3.88)	4.06 (3.88)
SES index	-3.33 (1.03)	-3.53 (1.05)	-3.47 (1.05)	-3.21 (1.05)
Attends math tuition			-1.83 (4.43)	0.88 (4.55)
Attends Hindi tuition			7.10 (4.40)	5.13 (4.53)
Baseline math score		-0.99 (2.17)	-0.88 (2.24)	-0.81 (2.24)
Baseline Hindi score		3.35 (2.12)	3.83 (2.15)	5.39 (2.14)
Constant	46.6 (3.40)	47.5 (3.42)	45.3 (3.79)	43.7 (3.78)
Grade Fixed Effects	N	N	N	Y
Observations	313	310	310	301
R-squared	0.038	0.046	0.056	0.120

*Note:* Robust standard errors in parentheses. This table shows correlates of days attended in the treatment group i.e. lottery-winners who had been offered a Mindspark voucher. Students from poorer backgrounds, and students with higher baseline achievement in Hindi, appear to have greater attendance but the implied magnitudes of these correlations are small. A standard deviation increase in the SES index is associated with a decline in attendance by about 3 days, and a standard deviation increase in Hindi baseline test scores is associated with an additional 5 days of attendance. We find no evidence of differential attendance by gender or by baseline math score.

Table A5: Quadratic dose-response relationship

	(1)	(2)	(3)	(4)
	Full sample		Treatment group	
	Math	Hindi	Math	Hindi
Attendance (days)	0.0052 (0.0054)	0.0079 (0.0053)	0.0097 (0.0072)	0.0070 (0.0073)
Attendance squared	0.000028 (0.000073)	-0.000058 (0.000072)	-0.000014 (0.000083)	-0.000048 (0.000085)
Baseline subject score	0.58 (0.042)	0.71 (0.040)	0.62 (0.061)	0.68 (0.052)
Constant	0.31 (0.042)	0.18 (0.042)	0.20 (0.14)	0.19 (0.14)
Observations	535	537	264	265
R-squared	0.429	0.496	0.446	0.446

*Note:* Robust standard errors in parentheses. This table models the dose-response relationship between Mindspark attendance and value-added quadratically. Results are estimated using OLS in the full sample and the treatment group only.

Table A6: Dose-response of subject-specific Mindspark attendance

	(1)	(2)	(3)	(4)	(5)	(6)
<i>Dep var: Standardized IRT scores (endline)</i>						
VARIABLES	IV estimates		OLS VA (full sample)		OLS VA (Treatment group)	
	Math	Hindi	Math	Hindi	Math	Hindi
Days of math instruction	0.018 (0.0029)		0.019 (0.0024)		0.022 (0.0047)	
Days of Hindi instruction		0.012 (0.0031)		0.011 (0.0026)		0.0084 (0.0050)
Baseline score	0.56 (0.038)	0.68 (0.036)	0.58 (0.041)	0.71 (0.039)	0.61 (0.060)	0.68 (0.052)
Constant			0.31 (0.041)	0.18 (0.041)	0.22 (0.11)	0.24 (0.11)
Observations	535	537	535	537	264	265
R-squared	0.432	0.478	0.428	0.495	0.445	0.446
	19	19				
Angrist-Pischke F-statistic for weak instrument	1211	1093				
Diff-in-Sargan statistic for exogeneity (p-value)	0.12	0.80				
Extrapolated estimates of 45 days' treatment (SD)	0.81	0.54	0.855	0.495	0.99	0.378

*Note:* Robust standard errors in parentheses. Treatment group students who were randomly-selected for the Mindspark voucher offer but who did not take up the offer have been marked as having 0% attendance, as have all students in the control group. Days attended in Math/Hindi are defined as the number of sessions of either CAL or small group instruction attended in that subject, divided by two. Columns (1) and (2) present IV regressions which instrument attendance with the randomized allocation of a voucher and include fixed effects for randomization strata, Columns (3) and (4) present OLS value-added models for the full sample, and Columns (5) and (6) present OLS value-added models using only data on the lottery-winners. Scores are scaled here as in Table 2.

Table A7: ITT estimates with inverse probability weighting

	(1)	(2)	(3)	(4)
	Dep var: Standardized IRT scores (endline)			
	Math	Hindi	Math	Hindi
Treatment	0.37 (0.063)	0.23 (0.062)	0.38 (0.062)	0.24 (0.061)
Baseline score	0.59 (0.041)	0.71 (0.040)	0.57 (0.038)	0.68 (0.037)
Constant	0.32 (0.044)	0.18 (0.044)	0.32 (0.043)	0.17 (0.042)
Strata fixed effects	N	N	Y	Y
Observations	535	535	535	535
R-squared	0.405	0.487	0.454	0.535

*Note:* Robust standard errors in parentheses. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$  Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher. Results in this table are weighted by the inverse of the predicted probability of having scores in both math and Hindi in the endline; the probability is predicted using a probit model with baseline subject scores, sex of the child, SES index and dummies for individual Mindspark centers as predictors. Tests in both math and Hindi were designed to cover wide ranges of ability and to be linked across grades, as well as between baseline and endline assessments, using common items. Scores are scaled here as in Table 2.

Table A8: Lee bounds estimates of ITT effects

	(1)	(2)
	Math	Hindi
Lower	0.309 (0.092)	0.183 (0.102)
Upper	0.447 (0.085)	0.294 (0.082)
Lower 95% CI	0.157	0.012
Upper 95% CI	0.587	0.43

*Note:* Analytic standard errors in parentheses. This table presents Lee(2009) bounds on the ITT effects of winning a voucher in both math and Hindi. We use residuals from a regression of endline test scores on baseline test scores (value-added) as the dependent variable, and scale scores as in Table 2, to keep our analysis of bounds analogous to the main ITT effects. The bounds are tightened using dummy variables for the Mindspark centres.

Table A9: ITT estimates, by source of test item

VARIABLES	(1)	(2)	(3)	(4)
	Math	Math	Hindi	Hindi
	EI items	non-EI items	EI items	non-EI items
Treatment	0.11 (0.013)	0.075 (0.011)	0.055 (0.017)	0.044 (0.011)
Baseline score	0.092 (0.011)	0.096 (0.0084)	0.14 (0.0093)	0.12 (0.0052)
Constant	0.46 (0.0064)	0.47 (0.0055)	0.61 (0.0082)	0.48 (0.0056)
Observations	537	537	539	539
R-squared	0.226	0.358	0.308	0.416

*Note:* Robust standard errors in parentheses. Treatment is a dummy variable indicating a randomly-assigned offer of a Mindspark voucher. Tests in both math and Hindi were assembled using items from different international and Indian assessments, some of which were developed by EI. EI developed assessments include the Student Learning Survey, the Quality Education Study and the Andhra Pradesh Randomized Studies in Education. The dependent variables are defined as the proportion correct on items taken from assessments developed by EI and on other non-EI items. All test questions were multiple choice items with four choices. Baseline scores are IRT scores normalized to have a mean of zero and a standard deviation of one.

Table A10: Treatment effect on take-up of other private tutoring

VARIABLES	(1) Math	(2) Hindi	(3) English	(4) Science	(5) Social Science
Post Sept-2015	0.019 (0.011)	0.018 (0.0096)	0.026 (0.0098)	0.018 (0.0080)	0.014 (0.0071)
Post * Treatment	0.013 (0.016)	-0.010 (0.012)	-0.0039 (0.013)	0.0017 (0.012)	-0.0056 (0.0086)
Constant	0.21 (0.0053)	0.13 (0.0040)	0.18 (0.0044)	0.14 (0.0041)	0.098 (0.0029)
Observations	3,735	3,735	3,735	3,735	3,735
R-squared	0.009	0.004	0.010	0.007	0.005
Number of students	415	415	415	415	415

*Note:* Robust standard errors in parentheses. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . This table shows individual fixed-effects estimates of receiving the Mindspark voucher on the take-up in other private tutoring in various subjects. The dependent variable is whether a child was attending extra tutoring in a given month between July 2015 and March 2016 in the particular subject. This was collected using telephonic interviews with the parents of study students. Observations are at the month\*child level. Treatment is a dummy variable indicating a randomly-assigned offer of Mindspark voucher.

## Appendix B Classroom Heterogeneity and Curricular Mismatch

As discussed in Sections ?? and ??, we conjecture that an important reason for the large effects we find is that the CAL software was able to accommodate the large heterogeneity in student learning levels within the same grade by personalizing instruction and teaching “at the right level” for all students. In this Appendix, we (a) provide evidence that the patterns in Figure 1 (a large fraction of students being behind grade-level standards and wide variation in academic preparation of students enrolled in the same grade) are present in other developing country settings as well, and (b) discuss qualitative evidence on pedagogical practice to show that the default instructional practice in these settings is to teach to the curriculum and textbook, which is likely to be above the learning levels of most students.

### B.1 Comparing the distribution of achievement in our study sample with other samples

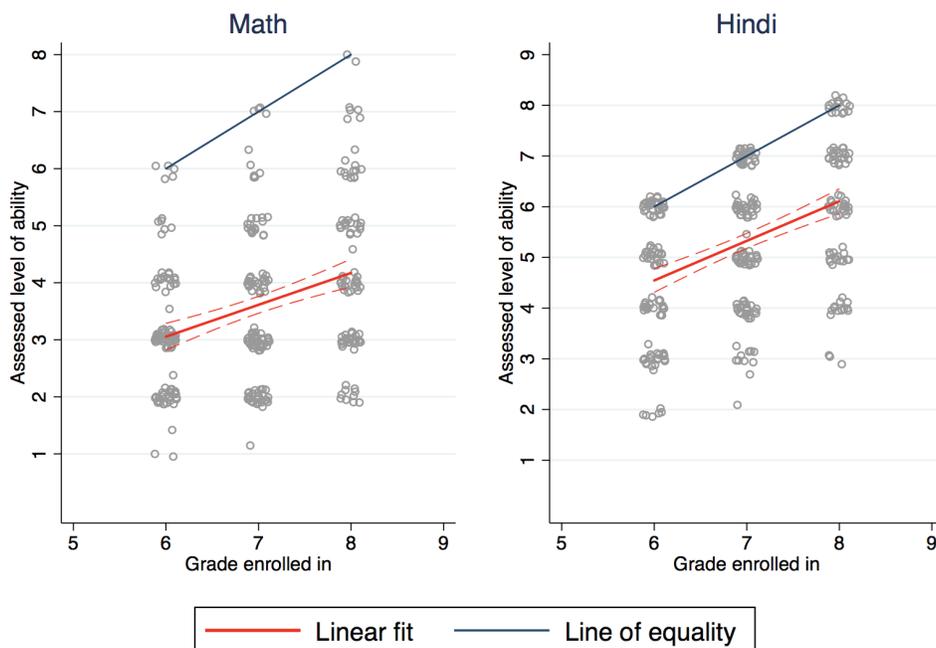
As mentioned in Section ??, an important advantage of the CAL data is the ability to characterize the mean and variance in grade-level preparation of students to produce the description shown in Figure 1. Yet, a limitation of the data in Figure 1 is that it comes from a self-selected sample of around 300 students in Delhi (though these students are quite similar to the other students in their school as seen in Figure A1). We show now that these patterns are replicated in much larger and representative data sets of student learning in India.

#### B.1.1 Rajasthan

In September 2017, subsequent to our study, Educational Initiatives signed an agreement with the Government of the Indian state of Rajasthan to introduce the Mindspark software in 40 government schools in the state. This deployment was spread across urban and rural areas in 4 districts (Churu, Jhunjunun, Udaipur, and Dungarpur) spanning the northern and southern ends of the state of Rajasthan, and covered 3276 students across grades 6-8. A similar diagnostic exercise that informed Figure 1 was conducted for all these students and the data is presented in Figure B1.

The patterns observed in Figure 1 are completely replicated in this larger and more representative (there was no student self-selection here) sample from a different state. Similar to the Delhi RCT sample, we see large absolute deficits against curricular standards (that grow in higher grades) and widespread dispersion within a grade. In math, the average Grade 6 student is 2.9 grade levels below curricular standards (compared to 2.5 grade levels below in Delhi), which rises to nearly 4 grade levels below by Grade 8 (similar to the sample in Delhi). In Hindi, the mean deficit in achievement compared to curricular standards is 1.5 grade levels

Figure B1: Assessed achievement level vs. enrolled grade in 40 public schools in Rajasthan



Note: Each dot represents 10 students

in Grade 6, rising to 2 grade levels in Grade 8.<sup>1</sup> Thus, the patterns in the Rajasthan data are nearly identical to those in Delhi.

Since the Rajasthan data covers all students in the enrolled classes, we can also directly examine the within-classroom heterogeneity in learning levels (which we cannot see in Delhi because the sample there only includes students who signed up for the after-school Mindspark program). Using data from 116 unique middle-school classrooms across 40 schools, we see that the median classroom in these schools has a range of about 4 grade levels of achievement in both math and language. Consistent with the Delhi data, the dispersion is greater in higher grades and, at a maximum, we see a spread of up to 6 grade levels in achievement (Table B1).

The Rajasthan data also allows us to decompose the within-grade variation in Figure B1 into between and within classroom variation. Specifically, we find that classroom fixed effects account for 31% (19%) of the variation in grade-6 scores in math (Hindi), 24% (15%) of the variation in grade-7 scores in math (Hindi), and 19% (7%) of the variation in grade-8 scores in math (Hindi). Thus, the vast majority of the dispersion in learning levels in the same

<sup>1</sup>In 2017, Educational Initiatives modified the diagnostic test such that the maximum grade that a student would be assigned is the grade they are enrolled in. Thus, while students could advance to levels beyond curricular standards dynamically through the system, they could not start above grade level. This would understate the spread of achievement in the Rajasthan sample relative to the Delhi sample in Hindi (this is not an issue for math since almost no students are above grade level in math in Delhi).

Table B1: Classroom-level heterogeneity in 40 schools in Rajasthan

Grade		Mathematics		Hindi	
		Range	p90 - p10	Range	p90 - p10
6	Mean	3.2	2.2	3.5	2.8
	Median	3	2	4	3
	Maximum	5	4	5	4
	N	40	40	40	40
7	Mean	4.1	3	3.9	3
	Median	4	3	4	3
	Maximum	6	5	5	4
	N	40	40	40	40
8	Mean	4.2	3	4.2	3.3
	Median	5	3	4.5	3.5
	Maximum	6	5	6	5
	N	36	36	36	36
Total	Mean	3.8	2.7	3.8	3
	Median	4	3	4	3
	Maximum	6	5	6	5
	N	116	116	116	116

grade seen in Figure Table B1 is *within* classrooms and not between them, underscoring the challenge faced by teachers in effectively catering to such variation.

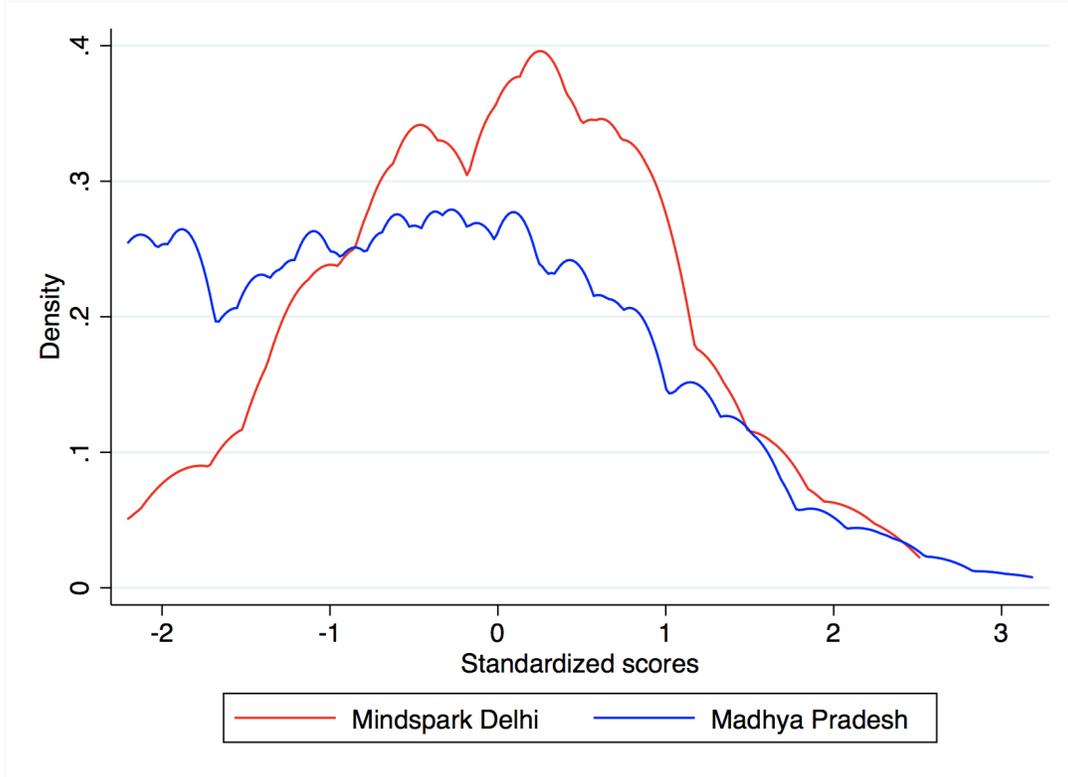
### B.1.2 Madhya Pradesh

While data from the Mindspark CAL system from Rajasthan government schools provides the most direct comparison with the Delhi sample, an alternative comparison is possible using our independent assessments. In a separate contemporaneous study in the Indian state of Madhya Pradesh (MP) on the impact of a school-governance reform (Muralidharan and Singh, 2018), we administered a common subset of items from the Mindspark endline assessments. The MP sample consists of 2760 students in grades 6-8 (who were present on the day of the assessment) in a representative sample of government middle schools in 5 districts of Madhya Pradesh (MP).<sup>2</sup> Both the Delhi and the MP assessments were administered in February 2016. In Figure B2 we present the distribution of achievement in the Madhya Pradesh sample with the control group in the Delhi Mindspark RCT on only the common items across the two studies; scores have been normalized to have a mean of zero and SD of one in the control group in the Delhi Mindspark RCT. The main results are that (a) the mean learning levels in MP are about 0.45 standard deviations below that in the Delhi sample, and (b) the standard deviation of the levels of student learning are about 25% greater than in the Delhi sample. Thus, both the key facts in Figure 1 (from the Delhi) sample of (a) low levels of learning, and

<sup>2</sup>Madhya Pradesh is the fifth-largest state in India by population with over 75 million inhabitants according to the 2011 Census. The state education system consists of over 112,000 schools.

(b) high variation within a grade are replicated in the MP data and appear to be even more pronounced.

Figure B2: Distribution of achievement across the control group in the Mindspark RCT vs a representative sample of schools in Madhya Pradesh



### B.1.3 Other countries and Indian states

There are two challenges in replicating the patterns of Figure 1 in other settings. First, most high-quality datasets on education in developing countries are from primary schools, whereas our focus in this paper is on post-primary grades. Second, while other datasets may allow the fact of variance in learning levels to be documented, the measures of learning are typically not linked to grade-level standards making it difficult to quantify the grade-level equivalent of learning gaps and variation (as we do above). We therefore focus on highlighting one key statistic on learning in developing countries, which is the fraction of students at the end of primary school (fifth or sixth grade) who are not able to read at the second grade level. The main advantage of this statistic is that is available in representative samples in several settings, and is also a meaningful measure of the phenomenon we are interested in – learning gaps (indicating a minimum of a three-year gap) and variation (since these students will be at least three years behind classmates who are at curricular standards). Table B2 presents this number for several Indian states and other countries.

Table B2: Learning standards in Grade 5 in India and selected other countries

State/Country	% Children who cannot read grade 2 level text	% Children who cannot do a division problem with a single-digit divisor	State/Country	% Children who cannot read grade 2 level text	% Children who cannot do a division problem with a single-digit divisor
Andhra Pradesh	55.1	37.2	Odisha	51.6	26.6
Arunachal Pradesh	25.5	19.0	Punjab	69.2	47.9
Assam	38.0	13.6	Rajasthan	54.2	28.2
Bihar	42.0	32.6	Tamil Nadu	45.2	21.4
Chattisgarh	55.9	23.0	Telangana	47.1	30.4
Gujarat	53.0	16.1	Tripura	51.0	19.9
Haryana	68.3	48.9	Uttar Pradesh	43.2	22.6
Himachal Pradesh	70.5	53.7	Uttarakhand	63.7	37.0
Jharkhand	36.4	23.5	West Bengal	50.2	29.0
Karnataka	42.1	19.7	<b>All India (rural)</b>	<b>47.8</b>	<b>25.9</b>
Kerala	69.2	38.6	<b>Pakistan (rural)</b>	<b>52.1</b>	<b>48.4</b>
Madhya Pradesh	38.7	19.4	Balochistan	41.7	39.9
Maharashtra	62.5	20.3	Punjab	65.0	59.6
Manipur	70.7	52.5	Sindh	36.6	24.3
Meghalaya	47.9	10.7	<b>Uganda</b>	<b>40.1</b>	<b>60.8</b>
Mizoram	46.0	27.7			
Nagaland	50.1	21.2			

Sources: Data for Indian states is taken from Pratham (2016), for Pakistan from SAFED (2017) and for Uganda from Uwezo (2016).

Note that students in Rajasthan perform slightly better than the national average for rural India, with several large states (such as Bihar, Madhya Pradesh and Uttar Pradesh) scoring substantially lower indicating that the challenges illustrated in Figure B1 are likely to be even more severe in these settings. Similar patterns are also shown for two other countries (Pakistan, with major states shown separately, and Uganda) in which the grade of testing, the task tested, and the form of reporting is comparable with the ASER tests in India.

The pattern of large learning deficits, with significant heterogeneity within the same grade, is much more general. Table B3 presents data from the World Development Report 2018 (World Bank, 2018) which consolidates data from 24 sub-Saharan countries, across three different assessments, to classify Grade 6 students by levels of competence in Reading and Mathematics. In most countries, a substantial proportion of students are classified as being “not competent” in mathematics.<sup>3</sup> However, there is substantial heterogeneity within the same grade in a country. In Kenya, for instance, about 30-40% of the sample is classified in

<sup>3</sup>For a concrete sense of what “not competent” means, in the PASEC assessment, this implies the inability to perform any but the most basic arithmetic operations with whole numbers (i.e. without demonstrating any knowledge of decimals or fractions or the ability to answer questions involving units of time, length or basic questions in geometry). In reading, it implies the inability to combine two pieces of explicit information in a text to draw simple inferences. In the SACMEQ assessments, “not competent” in reading implies the inability to link and interpret information located in various parts of the text; in math, it implies the inability to translate verbal or graphic information into simple word problems.

Table B3: Heterogeneity in achievement of Grade 6 students in 24 African countries

Country	Mathematics			Reading		
	Not competent	Low competence	High competence	Not competent	Low competence	High competence
All PASEC countries	57.6	24.7	17.7	61.6	25.1	13.3
All SACMEQ countries	36.8	18.4	44.8	63	20.2	16.8
Benin	48.3	29	22.7	60.2	29	10.8
Botswana	24.2	19.2	56.6	56.5	27.2	16.4
Burkina Faso	43.1	35.5	21.4	41.1	36.9	21.9
Burundi	43.5	49.1	7.4	13.2	46.8	39.9
Cameroon	51.2	24.7	24.1	64.6	23.7	11.8
Chad	84.3	12.8	3	80.9	16.1	3
Congo Rep.	59.3	23.5	17.1	71	23.1	5.9
Cote d'Ivoire	52	25.6	22.4	73.1	23.7	3.1
Kenya	19.8	19.6	60.6	38.3	32.1	29.6
Lesotho	52.5	25.5	22	81.1	13.6	5.3
Malawi	73.3	19.9	6.9	91.6	6.6	1.8
Mauritius	21.1	12.1	66.8	26.7	17.9	55.3
Mozambique	43.5	25	31.5	74.1	20.9	5
Namibia	38.7	25.5	35.8	81.7	12.2	6.1
Niger	91.5	6.4	2.1	92.4	6.3	1.4
Senegal	38.8	26.3	34.8	41.2	29.7	29.1
Seychelles	21.9	10.3	67.8	42.3	26.4	31.3
South Africa	48.3	14.7	37	69.2	15.4	15.5
Swaziland	7	20.7	72.2	44.3	37	18.7
Tanzania	10.1	12	77.9	43	25.5	31.5
Togo	61.6	22.6	15.8	52.5	27.9	19.7
Uganda	45.9	23.7	30.5	74.9	18	7.1
Zambia	72.6	14.9	12.4	91.8	6.5	1.7
Zimbabwe	37.2	20.7	42.1	57.2	22.6	20.2

Sources: This table draws upon figures presented in World Bank (2018), based on original data from SACMEQ (2007) , PASEC (2015) and the World Development Indicators.

each of the three bins of competence in mathematics (not competent, low competence, and high competence), highlighting the challenges of delivering a single program of instruction to all students in a classroom.

Taken together, the data presented in this section highlight that the two key patterns we highlight in Figure 1 of (a) large learning deficits relative to curricular standards, and (b) large heterogeneity in learning levels within the same grade, are typical of many developing country education systems.

## **B.2 Teaching to the curriculum**

Inadequate and widely-dispersed academic preparation within a classroom would be a challenge for instruction in any setting. But it is made more severe if curricula and pedagogy are not responsive to this dispersion. Combined with the low general levels of achievement, this leads to substantial mismatch between the instruction delivered in the classroom and students' ability to engage with it. We see strong indirect evidence of this from our data.

First, we see that students scoring in the lowest-tercile of the within-grade baseline achievement distribution (who are at least a few grade levels behind the level of the curriculum) make no progress in absolute learning levels despite being enrolled in school – suggesting that the level of instruction within the classroom was too far ahead (and likely to have been at the level of the curriculum). Second, even though we see no program impact on average for treated students on the grade-level school tests, we see significant positive effects on these tests for students scoring in the top-tercile of the within-grade baseline achievement distribution. Since these students were exposed to Mindspark content that was closer to their grade level, it suggests that the school exams (and instruction in the school) are likely to have adhered to grade-level curricular standards.

In this section, we present additional qualitative evidence to show that classroom instruction in Indian schools closely tracks the textbook and curriculum, regardless of how far behind those standards most students may be. Two main sets of factors contribute to this.

### **B.2.1 Curriculum and syllabi**

The first set relates to the prescribed curricula, syllabi and assessment. The way curricular standards are set and then transmitted in classroom teaching is largely determined by the (high-stakes) examination system, which serves later as a screening mechanism for future educational prospects and, eventually, white-collar jobs. In particular, it is not responsive to contextual factors about students' actual achievement or needs.<sup>4</sup> Although the National

---

<sup>4</sup>The National Focus Group on Curriculum, Syllabus and Textbooks, which underpinned the revised National Curriculum Framework in 2005, summarized the Indian education system as “largely a monolithic system perpetuating a kind of education which has resulted in a set of practices adopted for development

Curriculum Framework in 2005 did recommend unburdening the curriculum and making it more relevant, this has been hard to achieve in practice.<sup>5</sup> This focus on exam-oriented learning is particularly severe in middle and high schools, given major exam-based transition points after Grades 8 and 10. Given that post-primary education relies a great deal on foundational skills having been mastered, this focus means that a significant proportion of students are unable to engage with classroom instruction in a meaningful sense.<sup>6</sup>

### **B.2.2 The lack of differentiated instruction**

The second set of issues relate to the ability and desire of teachers to address low and dispersed achievement in their classrooms of their own accord. While, in theory, it is possible for teachers to provide differentiated instruction to cater to widespread heterogeneity, there is no evidence that they do so. Sinha, Banerji and Wadhwa (2016) report, for instance, that 88% of primary and upper primary school teachers in Bihar believed that their main objective was to “complete the syllabus”, even if nearly half of them agreed with or did not dispute the statement that “the textbooks are too difficult for children” (p. 24). Classroom observations at both primary and post-primary levels find consistently little evidence of differentiated or small-group instruction, with an overwhelming reliance on blackboard teaching and lecturing (Bhattacharjea, Wadhwa and Banerji, 2011; Sankar and Linden, 2014; World Bank, 2016; Sinha, Banerji and Wadhwa, 2016). Remedial instruction is also uncommon, and tracking of students into ability-based sections within school is made impractical in most public school settings in India because schools are small and rarely have more than one section per grade.<sup>7</sup>

In addition to reflecting the overall syllabus-determined orientation of the education system, the lack of remedial or differentiated instruction probably also reflects beliefs among some teachers about students’ ability to learn. As Kumar, Dewan and K.Subramaniam (2012)

---

of curriculum, syllabus and textbooks that is guided by the patterns and requirements of the examination system, rather than by the needs determined by a mix of criteria based on the child’s learning requirement, aims of education and the socio-economic and cultural contexts of learners.” (NCERT, 2006)

<sup>5</sup>See e.g. Dewan and Chabra (2012) on the opposition to revising math curricula: “Even though the NCF is very clear on this issue, state functionaries continue to feel that reducing topics leads to loss of mathematical knowledge and children of their state are being deprived in this process. They also feel that such reductions will make their children unfit for various competitive examinations that they will take at the end of schooling.”

<sup>6</sup>See e.g. Rampal and Subramaniam (2012) for a concrete example: “Mathematics at the upper primary level is premised on the ability to read and write numbers, and make sense of arithmetical expressions, as a starting point towards algebra. As children are not equipped to cope with this, classroom transaction gets reduced to children copying meaningless symbols from the blackboard, or from commercially available guidebooks in which the problems are worked out. Such classrooms where students cannot make sense of arithmetic expressions are not singular but fairly typical of classrooms catering to students from socioeconomically marginalised sections, or from rural backgrounds. They constitute a significant part of the student population.”

<sup>7</sup>If anything, the opposite situation with the same teacher simultaneously teaching multiple grades is more typical. This is because the Indian government has prioritized universal access to school, resulting in several very small schools across rural India. The average enrollment in public schools in rural India is under 100 students across five primary grades, and the majority feature multi-grade teaching (Muralidharan et al., 2017).

discuss: “It is quite common for educators and administrators to believe that children from disadvantaged socio-economic backgrounds are incapable of learning mathematics, either because of an inherent lack of ability or because they do not have the cultural preparation and attitude to learning.” Finally, it is not clear that, even had they wished to, teachers can effectively diagnose student errors and provide appropriate support. In a study of 150 secondary schools in two states (Madhya Pradesh and Tamil Nadu) in the 2014-2015 school year, it was found that language teachers were only able to identify student errors 50% of the time and math teachers were only able to do so 40% of the time (World Bank, 2016, p. 47). These challenges are not unique to India and similar findings of low teacher human capital and ability to support weaker students is also documented elsewhere; Bold et al. (2017), for example, use primary data from seven sub-Saharan African countries to document that “general pedagogical knowledge and the ability to assess students’ learning and respond to that assessment is poor across the seven countries, with roughly only 1 in 10 teachers being classified as having minimum knowledge in general pedagogy and none having minimum knowledge in student assessment.”

In sum, the core challenge of curriculum mismatch is general across the Indian education system. While direct evidence is scarce for other settings, it is likely that this challenge also generalizes to other developing country settings which are beset with low achievement and potentially over-ambitious curricula (see Pritchett and Beatty (2015)). Personalized instruction may also have significant potential for improving learning outcomes in these settings.<sup>8</sup>

---

<sup>8</sup>For experimental evidence, see Duflo, Dupas and Kremer (2011) which finds positive effects of tracking across the initial skill distribution and attributes it to the ability to customize instruction closer to skill levels of students within a classroom.

## Appendix C Prior research on hardware and software

Tables C1 and C2 offer an overview of experimental and quasi-experimental impact evaluations of interventions providing hardware and software to improve children’s learning. The tables only include studies focusing on students in primary and secondary school (not pre-school or higher education) and only report effects in math and language (not on other outcomes assessed in these studies, e.g., familiarity with computers or socio-emotional skills).

### C.1 Selecting studies

This does not intend to be a comprehensive review of the literature. Specifically, we have excluded several impact evaluations of programs (mostly, within education) due to major design flaws (e.g., extremely small sample sizes, having no control group, or dropping attriters from the analysis). These flaws are widely documented in meta-analyses of this literature (see, for example, Murphy et al., 2001; Pearson et al., 2005; Waxman, Lin and Michko, 2003).

We implemented additional exclusions for each table. In Table C1, we excluded DID designs in which identification is questionable and studies evaluating the impact of subsidies for Internet (for example, Goolsbee and Guryan, 2006). In Table C2, we excluded impact evaluations of software products for subjects other than math and language or designed to address specific learning disabilities (e.g., dyslexia, speech impairment).

### C.2 Reporting effects

To report effect sizes, we followed the following procedure: (a) we reported the difference between treatment and control groups adjusted for baseline performance whenever this was available; (b) if this difference was not available, we reported the simple difference between treatment and control groups (without any covariates other than randomization blocks if applicable); and (c) if neither difference was available, we reported the difference between treatment and control groups adjusted for baseline performance and/or any other covariates that the authors included.

In all RCTs, we reported the intent-to-treat (ITT) effect; in all RDDs and IVs, we reported the local average treatment effect (LATE). In all cases, we only reported the magnitude of effect sizes that were statistically significant at the 5% level. These decisions are non-trivial, as the specifications preferred by the authors of some studies (and reported in the abstracts) are only significant at the 10% level or only become significant at the 5% level after the inclusion of multiple covariates. Otherwise, we mentioned that a program had “no effect” on

the respective subject. Again, this decision is non-trivial because some of these studies were under-powered to detect small to moderate effects.

### C.3 Categories in each table

In both tables, we documented the study, the impact evaluation method employed by the authors, the sample, the program, the subject for which the software/hardware was designed to target, and its intensity. Additionally, in Table C1, we documented: (a) whether the hardware provided included pre-installed software; (b) whether the hardware required any participation from the instructor; and (c) whether the hardware was accompanied by training for teachers. In Table C2, we documented: (a) whether the software was linked to an official curriculum (and if so, how); (b) whether the software was adaptive (i.e., whether it could *dynamically* adjust the difficulty of questions and/or activities based on students' performance); and (c) whether the software provided *differentiated* feedback (i.e., whether students saw different messages depending on the incorrect answer that they selected).

Table C1: Impact evaluations of hardware

Study	Method	Sample	Program	Subject	Intensity	Software included?	Instructor's role?	Teacher training?	Effect	Cost
Angrist and Lavy (2002)	IV	Grades 4 and 8, 122 Jewish schools in Israel	Tomorrow-98	Math and language (Hebrew)	Target student-computer ratio of 1:10 in each school	Yes, included educational software from a private company	Not specified	Yes, training for teachers to integrate computers into teaching	Grade 4: $-0.4$ to $-0.3\sigma$ in math and no effect in language	USD 3,000 per machine, including hardware, software, and setup; at 40 computers per school, USD 120,000 per school
Barrera-Osorio and Linden (2009)	RCT	Grades 3-9, 97 public schools in six school districts, Colombia	Computers for Education	Math and language (Spanish)	15 computers per school	Not specified	Use the computers to support children on basic skills (esp. Spanish)	Yes, 20-month training for teachers, provided by a local university	No effect in language or math	Not specified
Malamud and Pop-Eleches (2011)	RDD	Grades 1-12, in six regions, Romania	Euro 200 Program	Math and language (English and Romanian)	One voucher (worth USD 300) towards the purchase of a computer for use at home	Pre-installed software, but educational software provided separately and not always installed	Not specified	Yes, 530 multimedia lessons on the use of computers for educational purposes for students	$-0.44\sigma$ in math GPA, $-0.56\sigma$ in Romanian GPA, and $-0.63\sigma$ in English	Cost of the voucher plus management costs not specified

Cristia et al. (2012)	RCT	319 schools in eight rural areas, Peru	One Laptop per Child	Math and language (Spanish)	One laptop per student and teacher for use at school and home	Yes, 39 applications including: standard applications, educational games, music editing, programming environments, sound and video recording, encyclopedia; also 200 age-appropriate e-books	Not specified	Yes, 40-hour training aimed at facilitating the use of laptops for pedagogical purposes	No effect in math or language	USD 200 per laptop
Mo et al. (2013)	RCT	Grade 3, 13 migrant schools in Beijing, China	One Laptop per Child	Math and language (Chinese)	One laptop per student for use at home	Yes, three sets of software: a commercial, game-based math learning program; a similar program for Chinese; a third program developed by the research team	Not specified	No, but one training session with children and their parents	No effect in math or language	Not specified
Beuermann et al. (2015)	RCT	Grade 2, 28 public schools in Lima, Peru	One Laptop per Child	Math and language (Spanish)	Four laptops (one per student) in each class/section for use at school	Yes, 32 applications including: standard applications, educational games, music editing, programming environments, sound and video recording, encyclopedia	Not specified	No, but weekly training sessions during seven weeks for students	No effect in math or language	USD 188 per laptop

Leuven et al. (2007)	RDD	Grade 8, 150 schools in the Netherlands	Not specified	Math and language (Dutch)	Not specified	Not specified	Not specified	Not specified	-0.08 SDs in language and no effect in math	This study estimates the effect of USD 90 per pupil for hardware and software
Machin, McNally and Silva (2007)	IV	Grade 6, 627 (1999-2001) and 810 (2001-2002) primary and 616 (1999-2000) and 714 (2001-2002) secondary schools in England	Not specified	Math and language (English)	Target student-computer ratio of 1:8 in each primary school and 1:5 in each secondary school	Some schools spent funds for ICT for software	Not specified	Yes, in-service training for teachers and school librarians	2.2 pp. increase in the percentage of children reaching minimally acceptable standards in end-of-year exams	This study estimates the effect of doubling funding for ICT (hardware and software) for a Local Education Authority
Fairlie and Robinson (2013)	RCT	Grades 6-10, 15 middle and high public schools in five school districts in California, United States	Not specified	Math and language (English)	One computer per child for use at home	Yes, Microsoft Windows and Office	No	No	No effect in language or math	Not specified

Table C2: Impact evaluations of software

Study	Method	Sample	Program	Subject	Intensity	Linked to curriculum?	Dynamically adaptive?	Differentiated feedback?	Effect	Cost
Banerjee et al. (2007)	RCT	Grade 4, 100 municipal schools in Gujarat, India	Year 1: off-the-shelf program developed by Pratham; Year 2: program developed by Media-Pro	Math	120 min./week during or before/after school; 2 children per computer	Gujarati curriculum, focus on basic skills	Yes, question difficulty responds to ability	Not specified	Year 1: $0.35\sigma$ on math and no effect in language; Year 2: $0.48\sigma$ on math and no effect in language	INR 722 (USD 15.18) per student per year
Linden (2008)	RCT	Grades 2-3, 60 Gyan Shala schools in Gujarat, India	Gyan Shala Computer Assisted Learning (CAL) program	Math	Version 1: 60 min./day during school; Version 2: 60 min./day after school; Both: 2 children per computer (split screen)	Gujarati curriculum, reinforces material taught that day	Not specified	Not specified	Version 1: no effect in math or language; Version 2: no effect in math or language	USD 5 per student per year
Carrillo, Onofa and Ponce (2010)	RCT	Grades 3-5, 16 public schools in Guayaquil, Ecuador	Personalized Complementary and Interconnected Learning (APCI) program	Math and language (Spanish)	180 min./week during school	Personalized curriculum based on screening test	No, but questions depend on screening test	Not specified	No effect in math or language	Not specified
Lai et al. (2012)	RCT	Grade 3, 57 public rural schools, Qinghai, China	Not specified	Language (Mandarin)	Two 40-min. mandatory sessions/week during lunch breaks or after school; teams of 2 children	National curriculum, reinforces material taught that week	No, same questions for all students	No, if students had a question, they could discuss it with their teammate, but not the teacher	No effect in language and $0.23\sigma$ in math	Not specified
Lai et al. (2013)	RCT	Grades 3 and 5, 72 rural boarding schools, Shaanxi, China	Not specified	Math	Two 40-min. mandatory sessions/week after school; teams of 2 children	National curriculum, reinforces material taught that week	No, same questions for all students	No, if students had a question, they could discuss it with their teammate, but not the teacher	$0.12\sigma$ in language, across both grades	Not specified

Mo et al. (2014a)	RCT	Grades 3 and 5, 72 rural schools, Shaanxi, China	Not specified	Math	Two 40-min. mandatory sessions/week during computer lessons; teams of 2 children	National curriculum, reinforces material taught that week	No, same questions for all students	No, if students had a question, they could discuss it with their teammate, but not the teacher	0.18 $\sigma$ in math	USD 9439 in total for 1 year
Mo et al. (2014b)	RCT	Grades 3 and 5, 72 rural schools, Shaanxi, China	Not specified	Math	Two 40-min. mandatory sessions/week during computer lessons; teams of 2 children	National curriculum, reinforces material taught that week	No, same questions for all students	No, if students had a question, they could discuss it with their teammate, but not the teacher	Phase 1: no effect in math; Phase 2: 0.3 $\sigma$ in math	USD 9439 in total for 1 year
Lai et al. (2015)	RCT	Grade 3, 43 migrant schools, Beijing, China	Not specified	Math	Two 40-min. mandatory sessions/week during lunch breaks or after school	National curriculum, reinforces material taught that week	No, same questions for all students	No, if students had a question, they could discuss it with their teammate, but not the teacher	0.15 $\sigma$ in math and no effect in language	USD 7.9-8.8 per child for 6 months
Mo et al. (2016)	RCT	Grade 5, 120 schools, Qinghai, China	Not specified	Language (English)	Version 1: Two 40-min. mandatory sessions/week during regular computer lessons; Version 2: English lessons (also optional during lunch or other breaks); Both: teams of 2 children	National curriculum, reinforces material taught that week	Version 1: No feedback during regular computer lessons; Version 2: feedback from teachers during English lessons	Version 1: if students had a question, they could discuss it with their teammate, but not the teacher; Version 2: feedback from English teacher	Version 1: 0.16 $\sigma$ in language; Version 2: no effect in language	Version 1: RMB 32.09 (USD 5.09) per year; Version 2: RMB 24.42 (USD 3.87) per year

Wise and Olson (1995)	RCT	Grades 2-5, 4 public schools in Boulder, Colorado, United States	Reading with Orthographic and Segmented Speech (ROSS) programs	Language and reading (English)	Both versions: 420 total min., in 30- and 15-min. sessions; teams of 3 children	Not specified	No, but harder problems introduced only once easier problems solved correctly; also in Version 2, teachers explained questions answered incorrectly	No, but students can request help when they do not understand a word	Positive effect on the Lindamond Test of Auditory Con-ceptualization (LAC), Phoneme Deletion test and Nonword Reading (ESs not reported); no effect on other language and reading domains	Not specified
Morgan and Ritter (2002)	RCT	Grade 9, 4 public schools in Moore Independent School District, Oklahoma, United States	Cognitive Tutor - Algebra I	Math	Not specified	Not specified	Not specified	Not specified	Positive effect (ES not reported) in math	Not specified
Rouse and Krueger (2004)	RCT	Grades 4-6, 4 public schools in urban district in northeast United States	Fast For Word (FFW) programs	Language and reading (English)	90-100 min./day during lessons ("pull-out") or before/after school, 5 days a week, for 6-8 weeks	Not specified	No, but harder problems introduced only once easier problems solved correctly	Not specified	No effect on Reading Edge test, Clinical Evaluation of Language Fundamentals 3rd Edition (CELF-3-RP), Success For All (SFA) test, or State Reading Test	USD 30,000 for a 1-year license for 30 computers, plus USD 100 per site for professional training

Dynanski et al. (2007)	RCT	Grades 4-6, 4 public schools in urban district in northeast United States	Fast For Word (FFW) programs	Language and reading (English)	90-100 min./day during lessons ("pull-out") or before/after school, 5 days a week, for 6-8 weeks	Not specified	No, but harder problems introduced only once easier problems solved correctly	Not specified	No effect on Reading Edge test, Clinical Evaluation of Language Fundamentals 3rd Edition (CELF-3-RP), Success For All (SFA) test, or State Reading Test	USD 30,000 for a 1-year license for 30 computers, plus USD 100 per site for professional training
		Grade 4, 43 public schools in 11 school districts, United States	Leapfrog, Read 180, Academy of Reading, Knowledgebox	Reading (English)	Varies by product, but 70% used them during class time; 25% used them before school, during lunch breaks, or time allotted to other subjects; and 6% of teachers used them during both	Not specified	Not specified, but all four products automatically created individual "learning paths" for each student	Not specified, but all four products provided immediate feedback to students; one provided feedback of mastery; two provided feedback on diagnostics	No effect in reading	USD 18 to USD 184 per student year (depending on the product)
		Grade 6, 28 public schools in 10 school districts, United States	Larson Pre-Algebra, Achieve Now, iLearn Math	Math	Varies by product, but 76% used them during class time; 11% used them before school, during lunch breaks, or time allotted to other subjects; and 13% of teachers used them during both	Not specified	Not specified, but all three products automatically created individual "learning paths" for each student	Not specified, but all three products provided immediate feedback to students; one provided feedback of mastery; two provided feedback on diagnostics	No effect in math	USD 9 to USD 30 per student year (depending on the product)

Algebra I, 23 public schools in 10 school districts, United States	Cognitive Tutor - Algebra I, PLATO Algebra, Larson Algebra	Math	Varies by product, but 94% used them during class time; and 6% of teachers used them during both	Not specified	Not specified, but two products automatically created individual "learning paths" for each student	Not specified, but all three products provided immediate feedback to students; two provided feedback of mastery; two provided feedback on diagnostics	No effect in math	USD 7 to USD 30 per student year year (depending on the product)	
Barrow, Markman and Rouse (2009)	RCT	Grades 8, 10	I Can Learn	Math	Not specified	National Council of Teachers of Mathematics (NCTM) standards and district course objectives	No, but students who do not pass comprehensive tests repeat lessons until they pass them	0.17 $\sigma$ in math	30-seat lab costs USD 100,000, with an additional USD 150,000 for pre-algebra, algebra, and classroom management software
Borman, Benson and Overman (2009)	RCT	Grades 2 and 7, 8 public schools in Baltimore, Maryland, United States	Fast For Word (FFW) Language	Language and reading (English)	100 min./day, five days a week, for four to eight weeks, during lessons ("pull-out")	Not specified	No, all children start at the same basic level and advance only after attaining a pre-determined level of proficiency	Grade 2: no effect in language or reading; Grade 7: no effect in language or reading	Not specified
Cam-puzano et al. (2009)	RCT	Grade 1, 12 public schools in 2 school districts, United States	Destination Reading - Course 1	Reading (English)	20 min./day, twice a week, during school	Not specified	Not specified	No effect in reading	USD 78 per student per year
		Grade 1, 12 public schools in 3 school districts, United States	Headsprout	Reading (English)	30 min./day, three times a week, during school	Not specified	Not specified	0.01 SDs in reading ( $p < 0.05$ )	USD 146 per student per year

Grade 1, 8 public schools in 3 school districts, United States	PLATO Focus	Reading (English)	15-30 min./day (frequency per week not specified)	Not specified	No, but teachers can choose the order and difficulty level for activities	Not specified	No effect in reading	USD 351 per student per year
Grade 1, 13 public schools in 3 school districts, United States	Waterford Early Reading Program - Levels 1-3	Reading (English)	17-30 min./day, three times a week, during school	Not specified	Not specified	Not specified	No effect in reading	USD 223 per student per year
Grade 4, 15 public schools in 4 school districts, United States	Academy of Reading	Reading (English)	25 min./day, three or more days a week, during school	Not specified	Not specified	Not specified	No effect in reading	USD 217 per student per year
Grade 4, 19 public schools in 4 school districts, United States	LeapTrack	Reading (English)	15 min./day, three to five days a week, during school	Not specified	No, but diagnostic assessments determine "learning path" for each student	Not specified	0.09 $\sigma$ in reading	USD 154 per student per year
Grade 6, 13 public schools in 3 school districts, United States	PLATO Achieve Now - Mathematics Series 3	Math	30 min./day, four days a week, for at least 10 weeks, during school	Not specified	No, but diagnostic assessment determines which activities students should attempt	Not specified	No effect in math	USD 36 per student per year
Grade 6, 13 public schools in 5 school districts, United States	Larson Pre-Algebra	Math	Varies according to the number of topics/weeks in the course, but recommended at least one a week	Not specified	Not specified	Not specified	No effect in math	USD 15 per student per year

Algebra I, 11 public schools in 4 school districts, United States	Cognitive Tutor - Algebra I	Math	Two days a week (plus textbook three days a week)	Not specified	Not specified	Not specified	No effect in math	USD 69 per student per year		
Algebra I, 12 public schools in 5 school districts, United States	Larson Algebra I	Math	Varies according to the number of topics/weeks in the course, but recommended at least one a week	Not specified	Not specified	Not specified	No effect in math	USD 13 per student per year		
Rockoff (2015)	RCT	Grades 6-8, 8 public middle schools in New York, NY, United States	School of One (So1)	Math	Not specified	No, activities sourced from publishers, software providers, and other educational groups	Yes, "learning algorithm" draws on students' performance on each lesson and recommends a "playlist" for each student; at the end of the day, students take a "playlist update"	No, but possibility to get feedback from live reinforcement of prior lessons, live tutoring, small group collaboration, virtual live instruction, and virtual live tutoring	No effect on New York State Math Test or Northwest Evaluation Association (NWEA) test	Not specified

## Appendix D Mindspark software

This appendix provides a more detailed description of the working of the Mindspark computer-assisted learning (CAL) software, and specifics of how it was implemented in the after-school Mindspark centers evaluated in our study.

### D.1 Computer training

The first time that students log into the Mindspark software, they are presented with an optional routine (taking 10-15 minutes) designed to familiarize them with the user interface and exercises on math or language.

### D.2 Diagnostic test

After the familiarization routine, students are presented with diagnostic tests in math and Hindi which are used by the Mindspark platform to algorithmically determine their initial achievement level (at which instruction will be targeted). Tests contain four to five questions per grade level in each subject. All students are shown questions from grade 1 up to their grade level. However, if students answer at least 75% of the questions for their corresponding grade level correctly, they can be shown questions up to two grade levels above their own.<sup>9</sup> If they answer 25% or less of the questions for one grade level above their actual grade, the diagnostic test shows no more questions. Initial achievement levels determined by the Mindspark system on the basis of these tests are only used to customize the first set of content that students are provided. Further customization is based on student performance on these content modules and does not depend on their performance on the initial diagnostic test (which is only used for initial calibration of each student’s learning level).

### D.3 Math and Hindi content

Mindspark contains a number of activities that are assigned to specific grade levels, based on analyses of state-level curricula. All of the items are developed by EI’s education specialists. The Mindspark centers focus on a specific subject per day: there are two days assigned to math, two days assigned to Hindi, one day assigned to English, and a “free” day, in which students can choose a subject.

Math and Hindi items are organized differently. In math, “topics” (e.g., whole number operations) are divided into “teacher topics” (e.g., addition), which are divided into “clusters” (e.g., addition in a number line), which are divided into “student difficulty levels” (SDLs) (e.g., moving from one place to another on the number line), which are in turn divided into questions (e.g., the same exercise with slightly different numbers). The Mindspark software

---

<sup>9</sup>For example, a grade 4 student will always see questions from grade 1 up to grade 4. However, if he/she answers over 75% of grade 4 questions correctly, he/she will be shown grade 5 questions; and if he/she answers over 75% of grade 5 questions correctly, he/she will be shown grade 6 questions.

currently has 21 topics, 105 teacher topics and 550 clusters. The organization of math content reflects the mostly linear nature of math learning (e.g., you cannot learn multiplication without understanding addition). This is also why students must pass an SDL to move on to the next one, and SDLs always increase in difficulty.

In Hindi, there are two types of questions: “passages” (i.e., reading comprehension questions) and “non-passages” (i.e., questions not linked to any reading). Passage questions are grouped by grades (1 through 8), which are in turn divided into levels (low, medium, or high). Non-passage questions are grouped into “skills” (e.g., grammar), which are divided into “sub-skills” (e.g., nouns), which are in turn divided into questions (e.g., the same exercise with slightly different words). The Mindspark software currently has around 330 passages (i.e., 20 to 50 per grade) linked to nearly 6,000 questions, and for non-passage questions, 13 skills and 50 sub-skills, linked to roughly 8,200 questions. The Hindi content is organized in this way because language learning is not as linear as math (e.g., a student may still read and comprehend part of a text even if he/she does not understand grammar or all the vocabulary words in it). As a result there are no SDLs in Hindi, and content is not necessarily as linear or clearly mapped into grade-level difficulty as in math.

The pedagogical effectiveness of the language-learning content is increased by using videos with same-language subtitling (SLS). The SLS approach relies on a “karaoke” style and promotes language learning by having text on the screen accompany an audio with on-screen highlighting of the syllable on the screen at the same time that it is heard, and has been shown to be highly effective at promoting adult literacy in India (Kothari et al., 2002; Kothari, Pandey and Chudgar, 2004). In Mindspark, the SLS approach is implemented by showing students animated stories with Hindi audio alongside subtitling in Hindi to help the student read along and improve phonetic recognition, as well as pronunciation.

## **D.4 Personalization**

### **D.4.1 Dynamic adaptation to levels of student achievement**

In math, the questions within a teacher topic progressively increase in difficulty, based on EI’s data analytics and classification by their education specialists. When a child does not pass a learning unit, the learning gap is identified and appropriate remedial action is taken. It could be leading the child through a step-by-step explanation of a concept, a review of the fundamentals of that concept, or simply more questions about the concept.

Figure D1 provides an illustration of how adaptability works. For example, a child could be assigned to the “decimal comparison test”, an exercise in which he/she needs to compare two decimal numbers and indicate which one is greater. If he/she gets most questions in that test correctly, he/she is assigned to the “hidden numbers game”, a slightly harder exercise in which he/she also needs to compare two decimal numbers, but needs to do so with as

little information as possible (i.e., so that children understand that the digit to the left of the decimal is the most important and those to the right of the decimal are in decreasing order of importance). However, if he/she gets most of the questions in the decimal comparison test incorrectly, he/she is assigned to a number of remedial activities seeking to reinforce fundamental concepts about decimals.

In Hindi, in the first part, students start with passages of low difficulty and progress towards higher-difficulty passages. If a child performs poorly on a passage, he/she is assigned to a lower-difficulty passage. In the second part, students start with questions of low difficulty in each skill and progress towards higher-difficulty questions. Thus, a student might be seeing low-difficulty questions on a given skill and medium-difficulty questions on another.

#### **D.4.2 Error analysis**

Beyond adapting the level of difficulty of the content to that of the student, Mindspark also aims to identify specific sources of conceptual misunderstanding for students who may otherwise be at a similar overall level of learning. Thus, while two students may have the same score on a certain topic (say scoring 60% on fractions), the reasons for their missing the remaining questions may be very different, and this may not be easy for a teacher to identify. A distinctive feature of the Mindspark system is the use of detailed data on student responses to each question to analyze and identify *patterns* of errors in student responses to allow for identifying the precise misunderstanding/misconception that a student may have on a given topic, and to target further content accordingly.

The idea that educators can learn as much (or perhaps more) from analyzing patterns of student errors than from their correct answers has a long tradition in education research (for instance, see Buswell and Judd (1925) and Radatz (1979) for discussions of the use of “error analysis” in mathematics education). Yet, implementing this idea in practice is highly non-trivial in a typical classroom setting for individual teachers. The power of ‘big data’ in improving the design and delivery of educational content is especially promising in the area of error analysis, as seen in the example below.

Figure D2 shows three examples of student errors in questions on “decimal comparison”. These patterns of errors were identified by the Mindspark software, and subsequently EI staff interviewed a sample of students who made these errors to understand their underlying misconceptions. In the first example, students get the comparison wrong because they exhibited what EI classifies as “whole number thinking”. Specifically, students believed 3.27 was greater than 3.3 because, given that the integer in both cases was the same (i.e., 3), they compared the numbers to the left of the decimal point (i.e., 27 and 3) and concluded (incorrectly) that since 27 is greater than 3, 3.27 was greater than 3.3.

In the second example, the error cannot be because of the reason above (since 27 is greater than 18). In this case, EI diagnosed the nature of the misconception as “reverse order thinking”. In this case, students know that the ‘hundred’ place value is greater than the ‘ten’ place value, but also believe as a result that the ‘hundred $th$ ’ place value is greater than the ‘tent $h$ ’ place value. Therefore, they compared 81 to 27 and concluded (incorrectly) that 3.18 was greater than 3.27.

Finally, the error in the last example cannot be because of either of the two patterns above (since 27 is less than 39, and 7 is less than 9). In this case, EI diagnosed the nature of the misconception as “reciprocal thinking”. Specifically, students in this case understood that the component of the number to the right of the decimal is a fraction, but they then proceeded to take the reciprocal of the number to the right of the decimal, the way standard fractions are written. Thus, they were comparing  $\frac{1}{27}$  to  $\frac{1}{39}$  as opposed to 0.27 to 0.39 and as a result (incorrectly) classified the former as greater.

It is important to note that the fraction of students making each type of error is quite small (5%, 4%, and 3% respectively), which would make it much more difficult for a teacher to detect these patterns in a typical classroom (since the sample of students in a classroom would be small). The comparative advantage of the computer-based system is clearly apparent in a case like this, since it is able to analyze patterns from thousands of students, with each student attempting a large set of such comparisons. This enables both pattern recognition at the aggregate level and diagnosis at the individual student-level as to whether a given student is exhibiting that pattern. Consistent with this approach, Mindspark then targets follow-up content based on the system’s classification of the patterns of student errors as seen in Figure D1 (which also shows how each student would do 30 comparisons in the initial set of exercises to enable a precise diagnosis of misconceptions).

## D.5 Feedback

The pedagogical approach favoured within the Mindspark system prioritizes active student engagement at all times. Learning is meant to build upon feedback to students on incorrect questions. Also, most questions are preceded by an example and interactive content that provide step-by-step instructions on how students should approach solving the question.

In math, feedback consists of feedback to wrong answers, through animations or text with voice-over. In Hindi, students receive explanations of difficult words and are shown how to use them in a sentence. The degree of personalization of feedback differs by question: (a) in some questions, there is no feedback to incorrect answers; (b) in others, all students get the same feedback to an incorrect answer; and (c) yet in others, students get different types of feedback depending on the wrong answer they selected.

Algorithms for the appropriate feedback and further instruction that follow a particular pattern of errors are informed by data analyses of student errors, student interviews conducted by EI's education specialists to understand misconceptions, and published research on pedagogy. All decisions of the software in terms of what content to provide after classification of errors are 'hard coded' at this point. Mindspark does not currently employ any machine-learning algorithms (although the database offers significant potential for the development of such tools).

In addition to its adaptive nature, the Mindspark software allows the center staff to provide students with an 'injection' of items on a given topic if they believe a student needs to review that topic. However, once the student completes this injection, the software reverts to the item being completed when the injection was given and relies on its adaptive nature.

Figure D1: Mindspark adaptability in math

## Example of Technology Enabling Personalized Learning to Learn Decimals

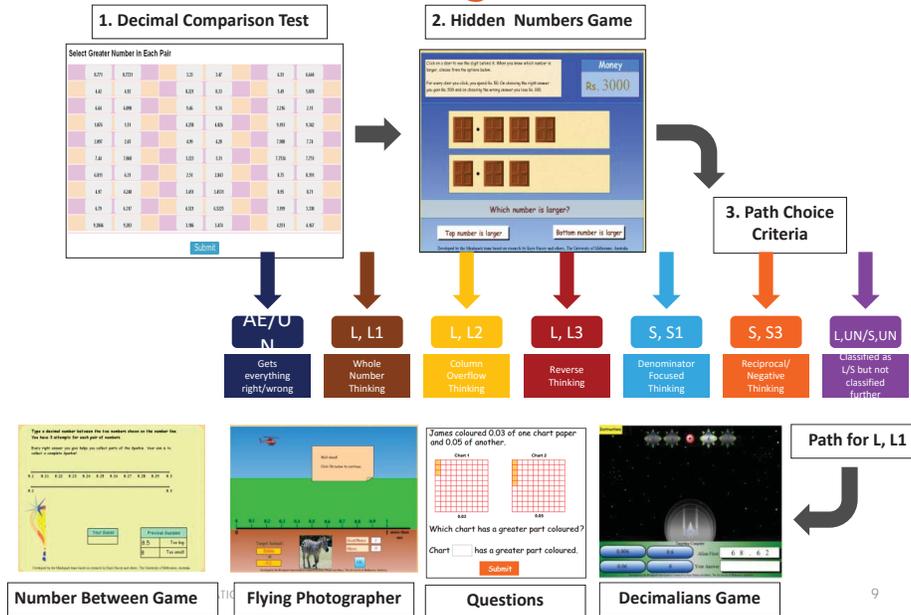


Figure D2: Student errors in math

## Why Would Some Students Think....



## Appendix E Test design

### E.1 Overview

We measured student achievement, which is the main outcome for our evaluation, using independent assessments in math and Hindi. These tests were administered under the supervision of the research team at both baseline and endline. Here we present details about the test content and development, administration, and scoring.

### E.2 Objectives of test design

Our test design was informed by three main objectives. First, was to develop a test which would be informative over a wide range of achievement. Recognizing that students may be much below grade-appropriate levels of achievement, test booklets included items ranging from very basic primary school appropriate competences to harder items which are closer to grade-appropriate standards.

Our secondary objective was to ensure that we measured a broad construct of achievement which included both curricular skills and the ability to apply them in simple problems.

Our third, and related, objective was to ensure that the test would be a fair benchmark to judge the actual skill acquisition of students. Reflecting this need, tests were administered using pen-and-paper rather than on computers so that they do not conflate increments in actual achievement with greater familiarity with computers in the treatment group. Further, the items were taken from a wide range of independent assessments detailed below, and selected by the research team without consultation with Education Initiatives, to ensure that the selection of items was not prone to “teaching to the test” in the intervention.

### E.3 Test content

We aimed to test a wide range of abilities. The math tests range from simple arithmetic computation to more complex interpretation of data from charts and framed examples as in the PISA assessments. The Hindi assessments included some “easy” items such as matching pictures to words or Cloze items requiring students to complete a sentence by supplying the missing word. Most of the focus of the assessment was on reading comprehension, which was assessed by reading passages of varying difficulty and answering questions that may ask students to either retrieve explicitly stated information or to draw more complex inferences based on what they had read. In keeping with our focus on measuring functional abilities, many of the passages were framed as real-life tasks (e.g. a newspaper article, a health immunization poster, or a school notice) to measure the ability of students to complete standard tasks.

In both subjects, we assembled the tests using publicly available items from a wide range of research assessments. In math, the tests drew upon items from the Trends in Mathematics and

Science Study (TIMSS) 4th and 8th grade assessments, OECD’s Programme for International Student Assessment (PISA), the Young Lives student assessments administered in four countries including India, the Andhra Pradesh Randomized Studies in Education (APRESt), the India-based Student Learning Survey (SLS) and Quality Education Study (QES); these are collectively some of the most validated tests internationally and in the Indian context.

In Hindi, the tests used items administered by Progress in International Reading Literacy Study (PIRLS) and from Young Lives, SLS and PISA. These items, available in the public domain only in English, were translated and adapted into Hindi.

#### **E.4 Test booklets**

We developed multiple booklets in both baseline and endline for both subjects. In the baseline assessment, separate booklets were developed for students in grades 4-5, grades 6-7 and grades 8-9. In the endline assessment, given the very low number of grades 4-5 students in our study sample, a single booklet was administered to students in grades 4-7 and a separate booklet for students in grades 8-9. Importantly, there was substantial overlap that was maintained between the booklets for different grades and between the baseline and endline assessments. This overlap was maintained across items of all difficulty levels to allow for robust linking using Item Response Theory (IRT). Table E1 presents a break-up of questions by grade level of difficulty in each of the booklets at baseline and endline.

Test booklets were piloted prior to baseline and items were selected based on their ability to discriminate achievement among students in this context. Further, a detailed Item analysis of all items administered in the baseline was carried out prior to the finalization of the endline test to ensure that the subset of items selected for repetition in the endline performed well in terms of discrimination and were distributed across the ability range in our sample. Table E2 presents the number of common items which were retained across test booklets administered.

#### **E.5 Test scoring**

All items administered were multiple-choice questions, responses to which were marked as correct or incorrect dichotomously. The tests were scored using Item Response Theory (IRT) models.

IRT models specify a relationship between a single underlying latent achievement variable (“ability”) and the probability of answering a particular test question (“item”) correctly. While standard in the international assessments literature for generating comparative test scores, the use of IRT models is much less prevalent in the economics of education literature in developing countries (for notable exceptions, see Das and Zajonc (2010), Andrabi et al. (2011), Singh (2015)). For a detailed introduction to IRT models, please see van der Linden and Hambleton (2013) and Das and Zajonc (2010).

The use of IRT models offers important advantages in an application such as ours, especially in comparison to the usual practice of presenting percentage correct scores or normalized raw scores. First, it allows for items to contribute differentially to the underlying ability measure; this is particularly important in tests such as ours where the hardest items are significantly more complex than the easiest items on the test.

Second, it allows us to robustly link all test scores on a common metric, even with only a partially-overlapping set of test questions, using a set of common items between any two assessments as “anchor” items. This is particularly advantageous when setting tests in samples with possibly large differences in mean achievement (but which have substantial common support in achievement) since it allows for customizing tests to the difficulty level of the particular sample but to still express each individual’s test score on a single continuous metric. This is particularly important in our application in enabling us to compute business-as-usual value-added in the control group.<sup>10</sup>

Third, IRT models also offer a framework to assess the performance of each test item individually which is advantageous for designing tests that include an appropriate mix of items of varying difficulty but high discrimination.

We used the 3-parameter logistic model to score tests. This model posits the relationship between underlying achievement and the probability of correctly answering a given question as a function of three item characteristics: the difficulty of the item, the discrimination of the item, and the pseudo-guessing parameter. This relationship is given by:

$$P_g(\theta_i) = c_g + \frac{1 - c_g}{1 + \exp(-1.7 \cdot a_g \cdot (\theta_i - b_g))} \quad (1)$$

where  $i$  indexes students and  $g$  indexes test questions.  $\theta_i$  is the student’s latent achievement (ability),  $P$  is the probability of answering question  $g$  correctly,  $b_g$  is the difficulty parameter and  $a_g$  is the discrimination parameter (slope of the ICC at  $b$ ).  $c_g$  is the pseudo-guessing parameter which takes into account that, with multiple choice questions, even the lowest ability can answer some questions correctly.

Given this parametric relationship between (latent) ability and items characteristics, this relationship can be formulated as a joint maximum likelihood problem which uses the matrix of  $N \times M$  student responses to estimate  $N + 3M$  unknown parameters. Test scores were generated using the OpenIRT software for Stata written by Tristan Zajonc. We use maximum likelihood estimates of student achievement in the analysis which are unbiased individual measures of ability (results are similar when using Bayesian expected a posteriori scores instead).

---

<sup>10</sup>IRT scores are only identified up to a linear transformation. Without explicitly linking baseline and endline scores, the constant term in our value-added regressions (which we interpret as value-added in the control group) would have conflates the arbitrary linear transformation and value-added in the control group.

## **E.6 Empirical distribution of test scores**

Figure E1 presents the percentage correct responses in both math and Hindi for baseline and endline. It shows that the tests offer a well-distributed measure of achievement with few students unable to answer any question or to answer all questions correctly. This confirms that our achievement measures are informative over the full range of student achievement in this setting.

Figure E2 presents similar graphs for the distribution of IRT test scores. Note that raw percent correct scores in Figure E1 are not comparable over rounds or across booklets because of the different composition of test questions but the IRT scores used in the analysis are.

## **E.7 Item fit**

The parametric relationship between the underlying ability and item characteristics is assumed, in IRT models, to be invariant across individuals (in the psychometrics literature, referred to as no differential item functioning). An intuitive check for the performance of the IRT model is to assess the empirical fit of the data to the estimated item characteristics.

Figure E3 plots the estimated Item Characteristic Curve (ICC) for each individual item in math and Hindi endline assessments along with the empirical fit for treatment and control groups separately. The fit of the items is generally quite good and there are no indications of differential item functioning (DIF) between the treatment and control groups. This indicates that estimated treatment effects do not reflect a (spurious) relationship induced by a differential performance of the measurement model in treatment and control groups.

Figure E1: Distribution of raw percentage correct scores

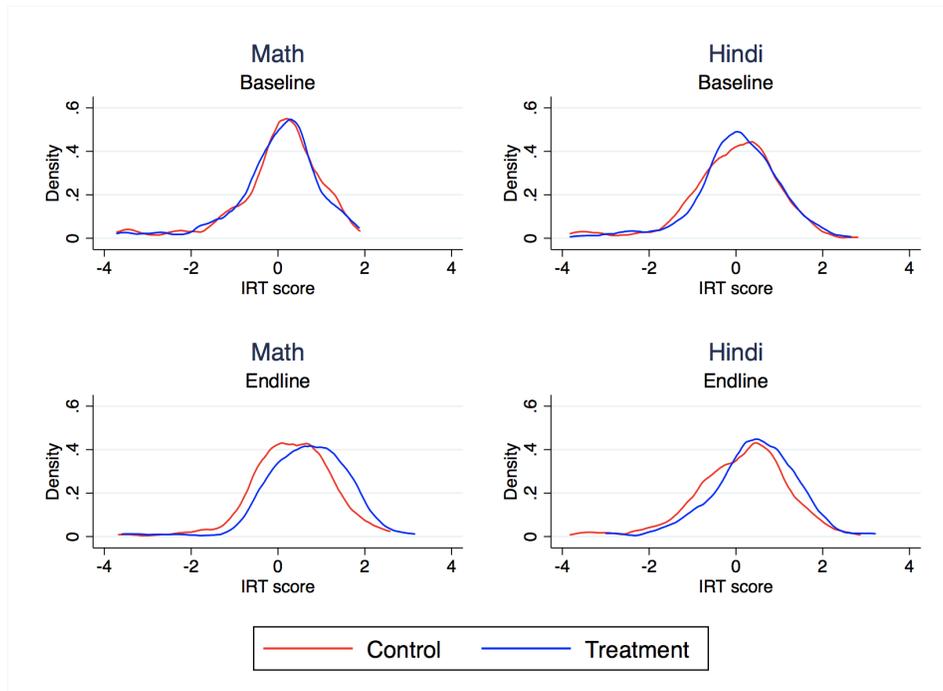


Figure E2: Distribution of IRT scores, by round and treatment status

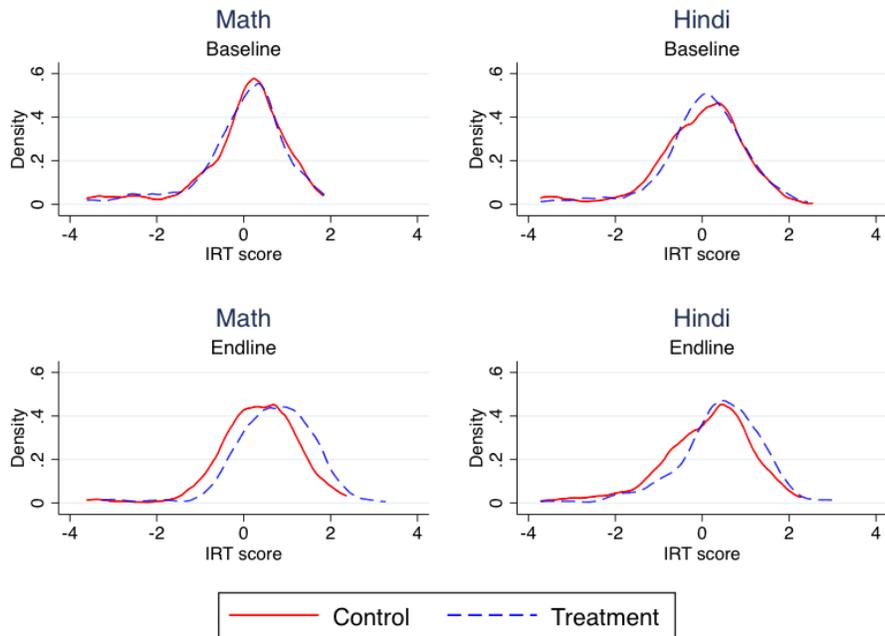
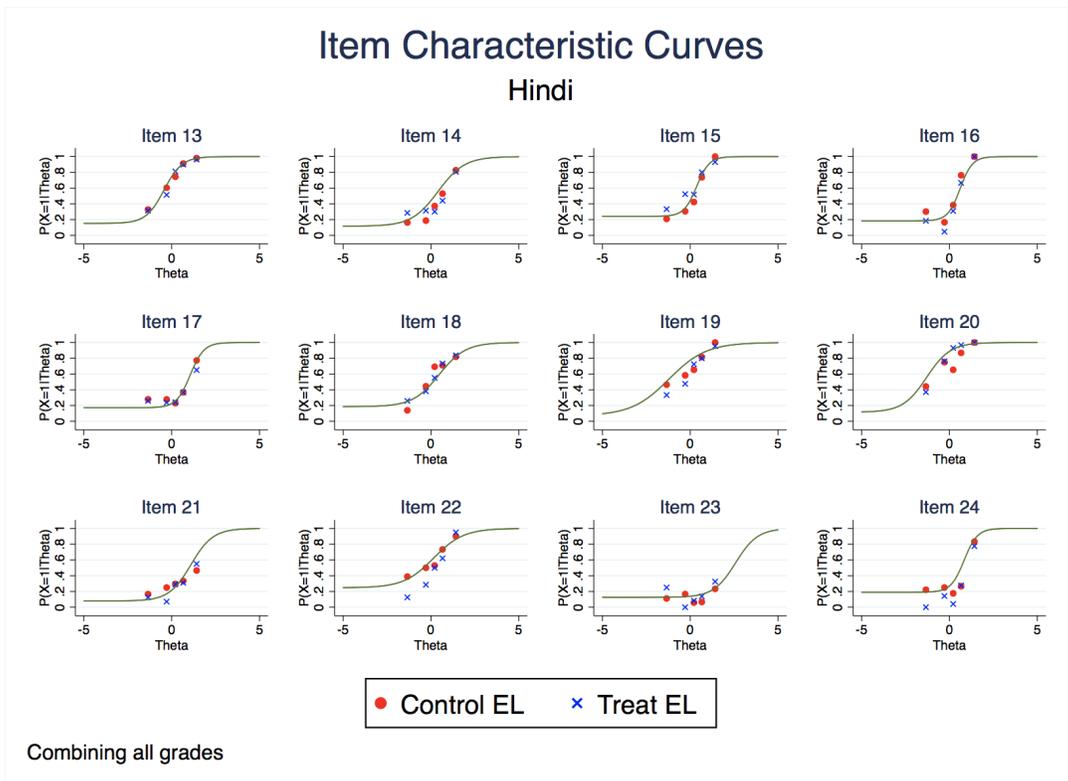
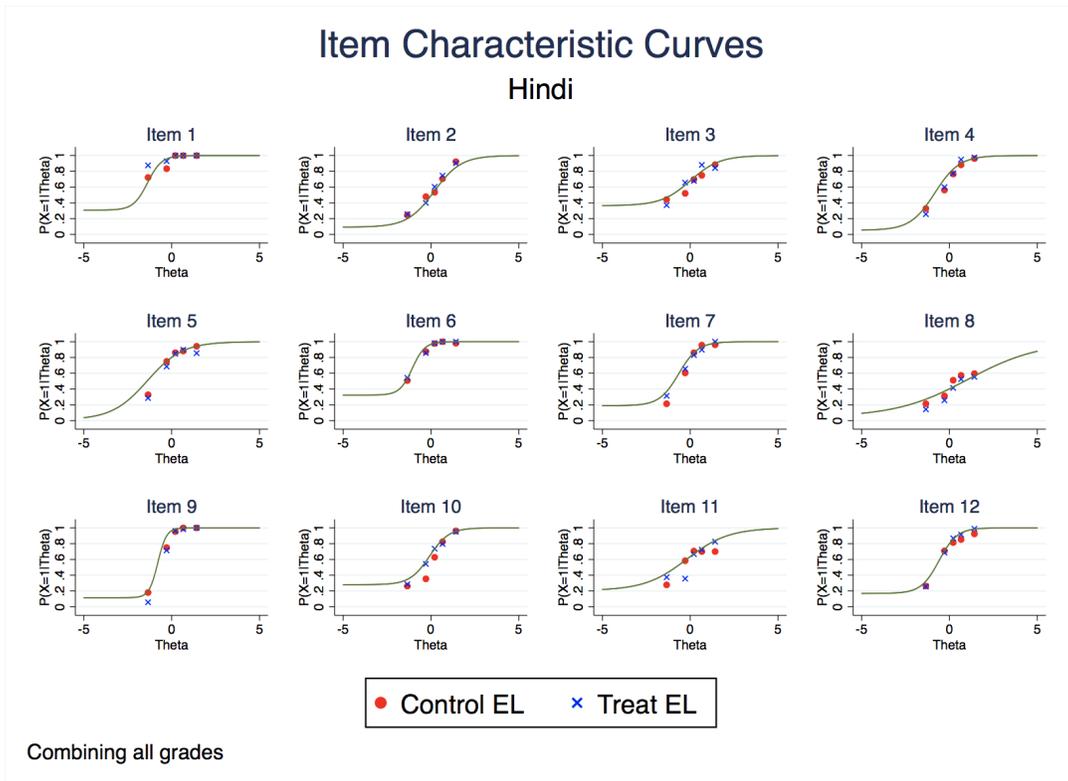
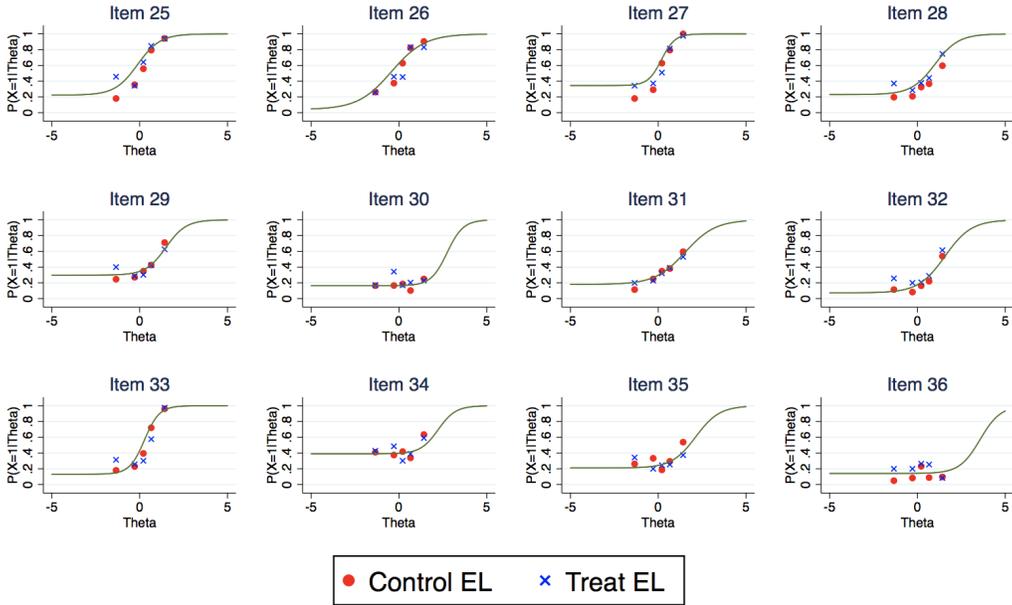


Figure E3: Item Characteristic Curves: Hindi



# Item Characteristic Curves

Hindi

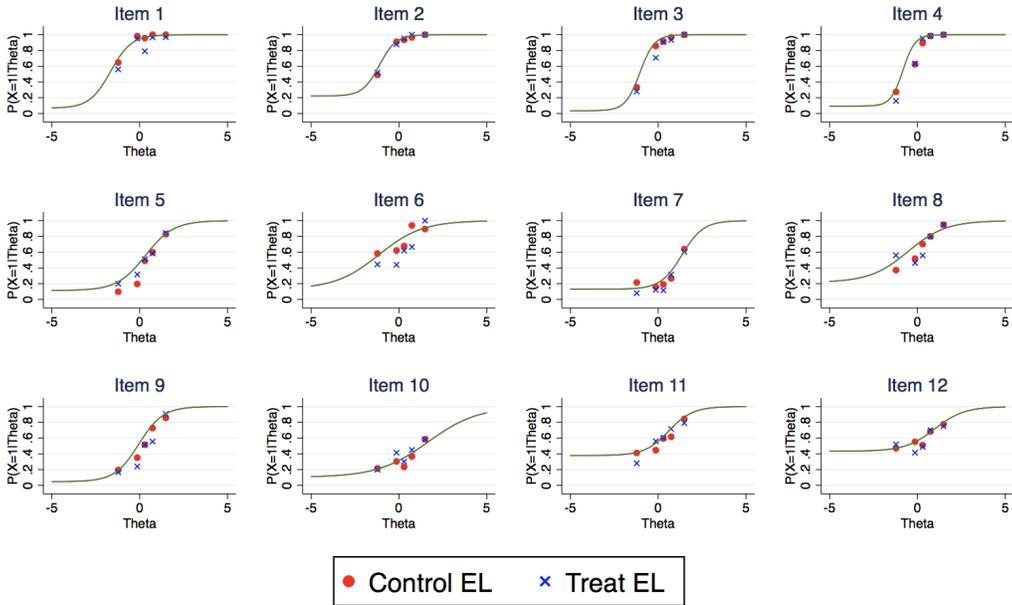


Combining all grades

Figure E4: Item Characteristic Curves: Math

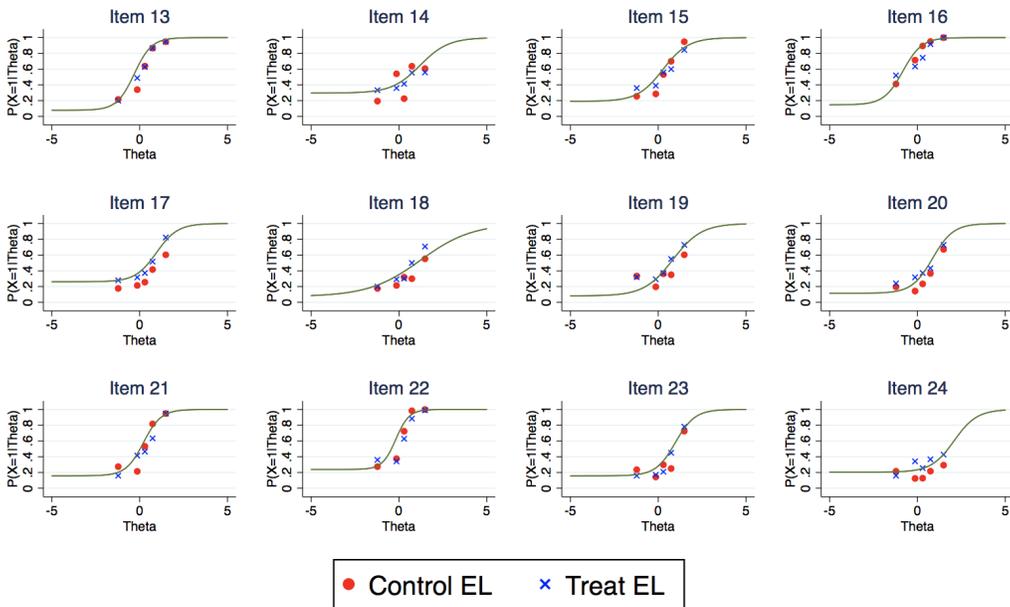
# Item Characteristic Curves

Mathematics



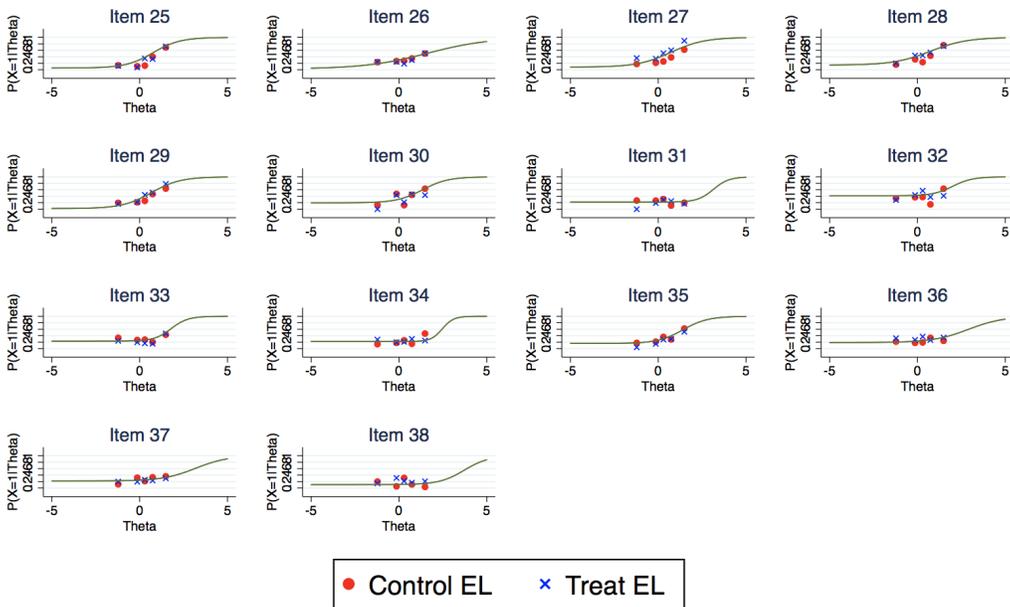
Combining all grades

## Item Characteristic Curves Mathematics



Combining all grades

## Item Characteristic Curves Mathematics



Combining all grades

Table E1: Distribution of questions by grade-level difficulty across test booklets

		Booklets				
		Baseline			Endline	
		Math				
		G4-5	G6-7	G8-9	G4-7	G8-9
Number of questions at each grade level	G2	2	0	0	2	0
	G3	14	6	4	6	6
	G4	13	7	4	9	8
	G5	4	10	3	10	10
	G6	1	10	10	5	6
	G7	1	2	11	2	3
	G8	0	0	3	0	2
		Hindi				
		G4-5	G6-7	G8-9	G4-7	G8-9
Number of questions at each grade level	G2	5	2	1	1	0
	G3	3	4	2	1	1
	G4	7	3	3	8	8
	G5	8	7	2	5	6
	G6	0	2	3	11	11
	G7	0	5	9	0	4
	G8	7	7	7	4	0
	G9	0	0	3	0	0

*Note:* Each cell presents the number of questions by grade-level of content across test booklets. The tests were designed to capture a wide range of student achievement and thus were not restricted to grade-appropriate items only. The grade-level of test questions was established ex-post with the help of a curriculum expert.

Table E2: Distribution of common questions across test booklets

Math				
	BL G6-7	BL G8-9	EL G4-7	EL G8-9
BL G4-5	16	10	14	14
BL G6-7		15	10	10
BL G8-9			7	7
EL G4-7				31

Hindi				
	BL G6-7	BL G8-9	EL G4-7	EL G8-9
BL G4-5	18	10	11	9
BL G6-7		17	13	13
BL G8-9			9	8
EL G4-7				24

*Note:* Each cell presents the number of questions in common across test booklets. Common items across booklets are used to anchor IRT estimates of student achievement on to a common metric.