



Instructional Coaching Personnel and Program Scalability

David Blazar
University of Maryland
College Park

Doug McNamara
University of Maryland
College Park

Genine Blue
TNTP

While teacher coaching is an attractive alternative to one-size-fits-all professional development, the need for a large number of highly skilled coaches raises potential challenges for scalability and sustainability. Collaborating with a national teacher training organization, our study uses administrative records to estimate the degree of heterogeneity in coach effectiveness at improving teachers' instructional practice, and specific characteristics of coaches that explain these differences. We find substantial variability in effectiveness across individual coaches. The magnitude of the coach-level variation (0.2 to 0.35 standard deviations) is close to the full effect of coaching programs, as identified in other research. We also find that coach-teacher race/ethnicity-matching predicts changes in teacher practice, suggesting that the relational component of coaching is key to success.

VERSION: January 2022

Suggested citation: Blazar, David, Doug McNamara, and Genine Blue. (2022). Instructional Coaching Personnel and Program Scalability. (EdWorkingPaper: 21-499). Retrieved from Annenberg Institute at Brown University: <https://doi.org/10.26300/2des-s681>

Instructional Coaching Personnel and Program Scalability

David Blazar (dblazar@umd.edu)*

Doug McNamara (dmcnama1@umd.edu)

University of Maryland College Park

2311 Benjamin Building, 3942 Campus Drive, College Park, MD 20740

Genine Blue (genine.blue@tntp.org)

TNTP

500 7th Avenue, 8th Floor, New York, NY 10018

Abstract

While teacher coaching is an attractive alternative to one-size-fits-all teacher training and development, the need for a large number of highly skilled coaches raises potential challenges for scalability and sustainability. Collaborating with a national teacher training organization, TNTP, our study uses administrative records to estimate the degree of heterogeneity in coach effectiveness at improving teachers' instructional practice, and characteristics of coaches that explain these differences. We find substantial variability in effectiveness across coaches (0.2 to 0.35 standard deviations), which is close to the full effect of coaching programs as identified in other research. We also find that coach-teacher race/ethnicity-matching predicts changes in teacher practice, suggesting that the relational component of coaching is key to success.

Key words: teacher coaching, teacher training, professional development, race/ethnicity matching

* = corresponding author. We thank our partners and collaborators at TNTP, including Vicky Brady and Bailey Cato, for compiling the data used in this project, and for ongoing brainstorming regarding analyses. We also thank Matthew Kraft for providing valuable feedback on the research questions and on an earlier draft of the manuscript.

Introduction

Instructional coaching has become an attractive alternative to one-size-fits-all teacher training and professional development. Compared to traditional, workshop-based programs that generally are ineffective (Fryer, 2017; Yoon et al., 2007), one-on-one coaching observation and feedback cycles have very large effects on teacher practice (upwards of 0.5 standard deviations [SD]) and on student test scores (upwards of 0.2 SD; Kraft et al., 2018). In fact, after reviewing experimental evidence on an array of educational interventions, Fryer (2017) found that only one-on-one, high-dosage tutoring with students had larger effects on academic outcomes. Because tutoring is more resource intensive per student than coaching, the latter is likely a more cost-effective intervention. Instructional coaching also has gained substantial popularity across the U.S., with coach-student ratios roughly doubling between 2000 and 2010 (Domina et al., 2015), and continued growth of programs since then. In the 2015-16 school year, 66% of public schools nationally had at least one coach (National Center for Education Statistics [NCES], 2016), compared to 57% of schools in 2007-08 (NCES, 2008).

Despite growing consensus on the benefits of coaching as a teacher training and development tool, it is less clear how best to scale programs in a way that also maintains their efficacy. Scalability and sustainability are concerns across the education research space (Slavin & Smith, 2009) but are likely to be particularly pronounced for coach-based teacher training and development that relies primarily on the efficacy of individual coaches. Said another way: coaches *are* the intervention. While coaching programs differ in design features, the intervention model is defined by coaches' engagement with teachers in one-on-one instructional improvement processes that includes time-intensive classroom observation and feedback cycles for each teacher (Joyce & Showers, 1981). The success of these efforts in improving the quality of teachers' classroom instruction is thought to depend on the knowledge, skills, and interpersonal relationship-building that individual coaches bring to their work (Connor, 2017; Denton & Hasbrouck, 2009; Joyce & Showers, 1981; Wong & Nicotera, 2006). For

instructional coaching to be a viable intervention across teacher training organizations, districts, and schools, it is necessary to identify, recruit, and hire very large corps of highly skilled coaches, potentially pulling current, highly effective teachers out of classrooms to serve in these roles (Darling-Hammond, 2017). As such, substantial variability in performance across individual coaches could undermine efforts to make *coaching* a primary—if not the primary—teacher training and development tool.

Empirically, a handful of recent studies provide suggestive evidence on the role of individual coaching personnel in program scale-up and sustainability efforts. Pooling results across all causal evaluations of instructional coaching programs, Kraft et al. (2018) found effects of small-scale programs—enrolling fewer than 100 teachers and led by few coaches—that were roughly twice as large as effects of larger programs—enrolling more than 100 teachers with many more coaches. The authors infer that these differences may be due in part to coaching personnel. However, their work cannot draw definitive conclusions because the number of coaches is conflated with other characteristics of the larger-scale programs and because the data do not include coach-teacher links in order to estimate the effect of individual coaches on desired outcomes. In a separate analysis of one coaching program, Blazar and Kraft (2015, 2019) found very large differences in average treatment effects across multiple cohorts, which they could attribute in part to turnover of coaches. At the same time, the sample size of coaches in this study is quite small ($n = 5$), and so the findings cannot speak to differences in coach efficacy in larger samples that more directly match current educational contexts (NCES, 2016).

In this paper, we build on theory and prior quantitative analyses by estimating the degree of performance heterogeneity across individual coaches in their effectiveness at improving teachers' instructional practice, drawing on secondary data from TNTP (formerly called The New Teacher Project). The collaboration with TNTP is appealing to examine this topic for several reasons. Because TNTP is a national, alternative-route teacher training, development, and certification organization,

our analyses leverage six years of data to examine coach effectiveness across 14 training sites (where sites generally are analogous to school districts), and 317 coaches working with over 3,500 teachers. Thus, in addition to greatly increasing statistical power relative to prior quantitative analyses on this topic, our findings are more generalizable. Relatedly, the context and scope of TNTP's programming speaks directly to the practice and policy question at hand regarding scalability and sustainability. As described by Kraft et al. (2018), many rigorous evaluations of coaching programs have been conducted under best-case scenarios, with relatively small samples of teachers, small numbers of coaches, and where coaches often were the program designers (sometimes also members of the research team). Yet, in real-world settings, teacher training organizations, districts, and schools need to hire much larger corps of coaches and to recruit them from broad labor pools. TNTP's programming closely reflects this context.

We focus on coaching cycles and data collected over the summer prior to individuals' first year as full-time teacher of record. This is TNTP's pre-service training period and, thus, the time during which the organization hires a large number of coaches and where data are collected systematically across sites. During pre-service training, TNTP trainees split their time between teaching in summer-school classrooms and largely coach-based learning activities (Menzes & Maier, 2014). While the literature base on instructional coaching to date often has focused on in-service training and development (Kraft et al., 2018), there is growing attention to coaching in pre-service training settings as well (e.g., Britton & Anderson, 2010; Cohen et al., 2020). We still refer to pre-service trainees as teachers, given that the structure of the alternative-route certification programs means that they are actively teaching students enrolled in summer school. Further, our analyses examine measures of instructional practice during lessons taught by these individuals.

To estimate heterogeneity in coach effectiveness, we take a value-added approach that is similar to the teacher effectiveness literature (Hanushek & Rivkin, 2010). Specifically, we predict

teachers' observed quality of instruction at the end of coaching as a function of baseline observation measures (hence "value-added") and additional covariates that aim to capture the primary avenues through which coaches are matched with teachers (e.g., site, certification area). While a randomized trial—in which coaches are randomly assigned to teachers—would provide stronger evidence of heterogeneity in effectiveness across individual coaches, we find that our value-added approach generally passes falsification tests that estimate the "effect" of coaches on measures that they should not impact (i.e., background teacher demographic characteristics).

Overall, we find substantial variability across coaches in terms of changes in teacher practice. A 1 SD increase in coach effectiveness is associated with a 0.2 to 0.35 SD increase in multiple dimensions of instructional quality measured on TNTP's observation instrument, including the extent to which teachers provide content aligned to appropriate standards for grade and subject and teachers' supports for students to engage meaningfully in classroom activities. A 2 SD increase in coach effectiveness—or the difference between having a coach at the 84th versus the 16th percentile in the performance distribution—is associated with a 0.4 to 0.7 SD increase in teachers' observed quality of instruction. Results of coach-level variation are similar when we nest coaches within sites to test for variation at each level, as well as when we estimate coach-level variation across each of the four largest training sites. These patterns suggest that it is the coach—and not the support, training, and oversight provided by each site—that likely matters most. Our estimates of coach-level variation in changes in teacher practice are close to the full effect of coaching programs, on average, as identified in other research (Kraft et al., 2018).

To further aid teacher-training organizations, schools, and districts looking to implement or scale coaching models through targeted recruitment and coach development, we also examine whether observable characteristics of coaches predict changes in teachers' instructional practice. We find positive associations for coach-teacher race/ethnicity matches, which is consistent with a large

theoretical and empirical literature base on the benefits of race/ethnicity matching between teachers and students (for a review, see Redding, 2019). These patterns also align with theoretical discussion of coaching as a relational activity (Joyce & Showers, 1981), where race/ethnicity matches likely support coach-teacher relationships built on shared understanding, experiences, and perspectives (Bristol & Martin-Fernandez, 2019; Ladson-Billings, 1995). Ultimately, these findings suggest that recruitment of a diverse pool of coaches and screening for coaches' interpersonal skills may be one of the best strategies for scale-up.

Framework and Motivating Literature on Performance Heterogeneity

The motivating hypothesis for our paper is that individual coaches vary substantially in their effectiveness at improving desired educational outcomes, namely the quality of teachers' instruction. If this is the case, then it may be more challenging to scale and sustain coaching programs.

We come to this hypothesis, in part, from recent empirical investigations examining mechanisms of effective instructional coaching programs (summarized above and described in more detail below), as well as from broader lines of theoretical and empirical work that point to substantial heterogeneity in the efficacy of personnel and labor pools. The most immediate link is to the teacher effectiveness literature, where studies consistently show that teachers differ not only in the quality of their classroom instruction (Bell et al., 2012; Hill et al., 2015; Kane & Staiger, 2012) but also in their subsequent impacts on students' test scores and social-emotional development (Blazar & Kraft, 2017; Hanushek & Rivkin, 2010; Nye et al., 2004). Our analyses also align with newer lines of research that find substantively meaningful variation across principals (Grissom et al., 2015) and guidance counselors (Mulhern, 2019) in terms of their effects on student outcomes. Outside of the education sector, examining personnel productivity vis-à-vis performance outcomes has longstanding discussion in the health sector, with doctors linked to patient outcomes (Safran et al., 1998), and in the economics and management literature on firms (Holmstrom & Milgrom, 1991).

One appealing framework derived from this literature is that the effectiveness of individual personnel can be estimated by way of their impacts on key beneficiaries—such as teachers, counselors, and principals linked to student outcomes, and, in this paper, coaches linked to teacher outcomes. A second learning is that we must consider not just whether individuals differ in their performance, but more importantly the magnitude of that variation. Teacher effect estimates of roughly 0.2 to 0.3 SD indicate that, on average, assignment to a teacher at the 84th percentile of effectiveness moves the medium-performing student to roughly the 60th percentile, relative to students’ peers assigned to a teacher at the 50th percentile in the performance distribution. These differences are quite large as benchmarked against students’ average yearly test-score gains, the effect of varied educational interventions, and policy-relevant gaps in achievement between students from different backgrounds (Hill et al., 2008). Findings related to performance heterogeneity across teachers have led to general consensus that teachers are by far the most important within-school resource that we can provide to students.

Applying a framework of performance heterogeneity to instructional coaches, we expect similar if not greater degrees of variability in performance as has been observed for other labor pools such as teachers. After all, at their core coaching programs are meant to be individualized, driven both by the needs of individual teachers with whom they work and one-on-one development work implemented by individual coaches. In their pioneering work describing the theory of action underlying instructional coaching models, Joyce and Showers (1981) note that coaching “represents a continuing problem-solving endeavor between the teacher and the coach...” that relies on “...a collegial approach to the analysis of teaching for the purpose of integrating mastered skills and strategies into: (a) a curriculum, (b) a set of instructional goals, (c) a time span, and (d) a personal teaching style” (p. 170). Aligned to this perspective, additional researchers and practitioners describe coaching as a relational endeavor driven primarily by coaches’ “people skills,” including building relationships and

trust with teachers, and differentiating support for individual teachers' needs (Denton & Hasbrouck, 2009; Wong & Nicotera, 2006).

Despite robust theoretical discussion on this topic, to our knowledge, only Blazar and Kraft (2015, 2019) quantitatively examine differences in effectiveness of individual coaches when linked to teacher outcomes. In that study, the authors found substantial differences in average treatment effects of the coaching program across multiple cohorts of their randomized experiment, with large positive effects in the first cohort but null effects in two subsequent cohorts. After ruling out potential explanations for these patterns related to research design (e.g., characteristics of participating teachers, spillover effects), exploratory analyses suggested that differential treatment effects likely were attributable in part to turnover of coaches and differences in coach effectiveness. On average, the teachers of the most effective coach—who left the program after the first cohort—scored roughly 1.2 SD higher than the teachers of the least effective coach—who entered the program in the second cohort—on instructional quality measures derived from classroom observation, as well as 0.7 SD higher on student-reported measures of classroom experiences. At the same time, these analyses focused on a small sample of five coaches, and so findings cannot speak to an underlying population distribution of coach effectiveness. It may be that very large differences in effectiveness across five coaches are due to sampling idiosyncrasies and potential outliers. A primary goal of our paper and analyses is to examine heterogeneity in coach effectiveness at improving teachers' instructional practice in a much larger sample.

Two personnel and performance management questions related to scaling and sustaining coaching programs also are relevant to the current topic: To the extent that individual coaches differ substantially in their impacts on teacher practice, what are the key domains of coach characteristics that explain these differences? How can these skills be leveraged for recruitment and screening of, and professional learning for coaches? Here, there is a small but growing literature base. By and large,

coaches tend to be expert teachers with a demonstrated track record of success in the classroom, who often enter the role through a career ladder; coaches may come from within a school or district, or from another context (Darling-Hammond, 2017; Wenner & Campbell, 2017). In terms of the specific characteristics and skills of potential coaches to look for, Connor (2017) hypothesizes three areas of effectiveness. First, there must be a strong interpersonal relationship between the coach and teacher. Coaches and teachers who communicate and collaborate more effectively may experience bigger rewards from the coaching relationship. Second, a coach's knowledge of effective teaching and coaching practices may affect teaching outcomes. Similarly, more effective coaches may have content-specific knowledge which they use in the coaching relationship. Knowledge of effective teaching practices plays a direct role in ensuring high-quality observation-feedback cycles. Third, the types of tools (e.g., modeling, providing direct feedback, video observation, etc.) and technologies (e.g. online vs. in-person coaching, bug-in-ear real-time coaching, etc.) a coach uses may matter.

Empirically, scholars have started to operationalize domains of coach skill in survey instruments and observation tools to capture the quality of coach-teacher interactions (e.g., Howley et al., 2014), examine variability in how coaches instantiate these practices in their work with teachers (e.g., Shannon et al., 2021), and link coach characteristics and practices to teacher outcomes (e.g., Marsh et al., 2012; Yopp et al., 2019). For example, in the context of a math coaching program in Tennessee, Russell et al. (2020) found that a 1 SD change in the depth and specificity of coaches' conversations with teachers was associated with a 0.2 SD increase in the quality of teachers' instruction. However, much of this work has been conducted in small samples, generally with no more than 30 coaches. Further, because this literature base is quite new, many of the theorized domains of coach effectiveness have not been linked to changes in teacher practice, particularly in samples that can lead to generalizable conclusions. In our study, while we are not able to examine all hypothesized domains of coach effectiveness, we provide suggestive evidence on some of the key skills highlighted

in the theoretical literature. And, we examine heterogeneity in coach effectiveness across a number of U.S. states and school districts.

The TNTP Coaching Model and Pre-Service Training Context

TNTP is an approved alternative-route teacher certification entity, which has trained and certified over 50,000 teachers since opening its doors in 1997. Like other alternative-route teacher certification programs, TNTP partners with school systems to recruit prospective teachers largely from local labor pools, with the goal of filling local teacher vacancies in hard-to-staff subject areas and schools (Walsh & Jacobs, 2007). The nature of the alternative-route certification program means that training is condensed into five to seven weeks prior to becoming a full-time teacher of record, and the practicum component occurs in summer school classrooms.

Aligned to longstanding calls for and trends in teacher education and training reform—within which alternative certification programs have played a key role (Wilson, 2014)—in 2012 TNTP shifted its programming to focus more intentionally on a targeted set of foundational teaching skills, and on providing substantial time for teachers to practice and receive directed feedback on their implementation of these skills in real-world classrooms (Menzes & Maier, 2014). Our study focuses on this post-2012 time period. The prioritized set of instructional skills include: clear delivery of lessons, maintaining high academic and behavioral expectations, and maximizing instructional time. These elements of instructional practice are instantiated in a classroom observation instrument developed at TNTP that guides formative assessment and feedback, as well as summative evaluations to determine whether or not prospective teacher candidates gain certification. In our study, we use this instrument to capture instructional practice outcome measures (see discussion below).

Attention to practice and feedback as key resources for developing teaching skill align closely with the theory of action underlying instructional coaching programs (Joyce & Showers, 1981). On average over the course of TNTP's summer training period, teachers spend at least 32 hours working

with an instructional coach. Training starts with trainers showing teachers examples of what effective classroom environments look like, both through videotapes of exemplar lessons and modeling by coaches. Then, coaching observation and feedback cycles include three core components: active observations, direct and specific feedback, and immediate practice. Coaches typically engage in the process through visits to teachers' classes where they observe instruction. Coaches may also explicitly model a particular teaching skill or guide teachers in more subtle ways, including holding up signs or whispering to the teacher. Following a classroom visit, coaches meet with teachers for debriefing sessions to provide "bite-sized" feedback on one or two observed elements of instruction. These feedback points stem from the classroom observation and are meant to help teachers improve in their next lesson. A goal for the feedback process is to provide teachers with concrete and manageable steps that they can address that day or the next day. Teachers may practice this new technique in front of their coach during the debrief session. For additional details on TNTP's pre-service training and coaching model, see Menzes and Maier (2014).

While TNTP coaching and pre-service training operates under a common organizational model, individual coaches are the program implementers and they do so with guidance from individual site/school district managers. Each summer, sites hire coaches, pulling both from pools of TNTP-trained teachers and local educators. Coaches are expected to have a minimum of two years of successful teaching experience in high-need subject areas, familiarity with the instructional standards associated with the school district in which they are serving, and demonstrated ability to support teacher trainees in developing the teaching techniques emphasized in TNTP's training model. In the spring and early summer, coaches receive up to 40 hours of training from site leads who often coaches themselves in prior years, which includes an overview of the coaching model, practicing coaching, and observing and scoring the quality of classroom instruction. Following training, coaches work

individually with teachers, providing guidance and support aligned to their observations of teachers' instruction in summer-school classrooms and their perceptions of teachers' most immediate needs.

Research Design

In this study, we ask: (1) *To what extent do individual coaches vary in their effectiveness at improving teachers' instructional practice?* (2) *To what extent do observable characteristics of coaches (i.e., years of coaching experience, demographic matches with teachers) explain their effects on teacher practice?* We answer both questions in the context of pre-service training delivered by TNTP.

Empirical Strategy

To answer our research questions, we draw on the teacher effectiveness and value-added literatures to specify a production function of the following form:

$$\text{OBSERVATION}_{ijst} = \beta_0 + \beta_1 \text{OBSERVATION}_{ijsc(t-1)} + \beta_2 I_{j(t-1)} + \delta_{st} + (\mu_j + \epsilon_{ijst}) \quad (1)$$

where the outcome of interest is the end-of-coaching observation score for teacher i working with coach j in site s and year t . The key feature of our model is that we control for a baseline measure of the outcome, $\text{OBSERVATION}_{ijsc(t-1)}$, captured at the beginning of the training period and prior to the start of coaching. Controlling for a baseline measure allows us to estimate changes in teacher practice associated with individual coaches and, most importantly, to account for bias due to non-random sorting of coaches to teachers. To this same end, we further control for baseline teacher characteristics (i.e., gender, race/ethnicity) and certification area, included in the vector, $I_{j(t-1)}$, as well as site-year fixed effects, δ_{st} . According to TNTP, these are the primary avenues and characteristics that drive coach-teacher matches.

Our primary estimate of interest comes from μ_j , which is a coach random effect and can be thought of as the contribution of individual coaches to teacher outcomes above and beyond variables controlled for in the model. The j subscript on μ indicates that the random effect is a random variable, with an effectiveness estimate generated for each coach. We are primarily interested in the underlying

distribution of the coach effects and the degree of dispersion. A large degree of dispersion—as indicated by a large SD of the coach effectiveness distribution—suggests that it makes a large difference for teachers’ instructional practice in terms of the coach with whom they work. Comparatively, a SD of or close to zero would indicate that there is little heterogeneity in effectiveness across coaches. We do not need to calculate the individual coach effects and their distribution, as our random effects model allows us to generate a model-based estimate of the variation in changes in teacher practices associated with individual coaches. Model-based estimation via restricted maximum likelihood produces a consistent estimator for the true variance of coach effects (Guarino et al., 2015; Raudenbush & Bryk, 2002). Calculating the variation across individual coach effect estimates using Ordinary Least Squares regression would bias our variance estimates upward because it would conflate true variation with estimation error. Our random effects models shrink the coach effects back towards the mean based on the precision of those estimates, driven primarily by the number of teachers with whom an individual coach works (mean = 8.2 teachers per coach/year, SD = 2.5).

We specify models that include a coach random effect, μ_j , as well as models that include a coach-year random effect, μ_{jt} . The coach-level random effect considers coach effects as stable across years, while the coach-year random effect allows for variation across years. In our data, 74% of coaches are observed in the data for just one year (see Table 1), suggesting that the coach and coach-year random effects are unlikely to produce vastly different estimates. As shown below, we indeed find that both sets of estimates are quite similar. In some models, we nest the coach-year random effect within a site-year random effect—moving δ_{st} from the fixed to the residual portion of the model—in order to examine whether coaches versus the sites within which they work are a primary driver of changes in teacher outcomes.

Data and Sample

We fit our models using data collected by TNTP across six years (2014 through 2019) and 14 summer training sites. Our primary sample includes a census of pre-service teachers ($n = 3,526$) and coaches ($n = 317$) with whom TNTP worked during this time period. In Table 1, we show that this sample of teachers is roughly two-thirds female, one-quarter Black, and two-fifths White. (Twenty percent of teachers did not report race/ethnicity information.) These characteristics are more diverse than national characteristics of teachers (NCES, 2020), but are aligned with characteristics of teachers who go through alternative-route teacher certification programs that often operate in urban settings with a goal of decreasing barriers to entry into the profession for historically marginalized groups (NCES, 2016; Shen, 1997). Demographic characteristics of coaches are similar to those of teachers: roughly two-thirds are female, one-quarter are Black, and half are White; three-quarters have one year of experience coaching for TNTP. Coaches may have coaching experience outside of TNTP, which we are not able to capture in the administrative records available for this study.

Trained evaluators rated teachers' instructional practice multiple times over the course of the summer using TNTP's observation rubric (TNTP, 2014). This rubric includes three dimensions of instructional practice, each of which is scored on a scale from 1 (Ineffective) to 3 (Developing): (i) *Culture of Learning* asks whether all students are engaged in the work of the lesson from start to finish, and focuses on the extent to which teachers maximize instructional time and maintain high expectations for student behavior; (ii) *Essential Content* asks whether all students are engaged in content aligned to the appropriate standards of their subject and grade, and focuses on the extent to which teachers plan and deliver content accurately and clearly; and (iii) *Demonstration of Learning* asks whether all students demonstrate that they are learning, and focuses on the extent to which teachers check for student understanding and respond to student misunderstandings. (For additional details on the dimensions of practice, see TNTP, 2013.) We also created a composite measure of effective teaching

practice that is an average of these three dimensions. We standardized observation scores to have a mean of 0 and SD of 1. Observers participated in training during which they rated no fewer than seven full-length instructional videos followed by three to four “check in” points to rate and discuss additional lesson videos or co-observe in classrooms. Overall, observers receive about 40 to 50 hours a year of observation practice.

All domains of teaching practice have been linked to student test score growth in other TNTP-led research projects (TNTP, 2018) and in an external validation study (McEachin et al., 2018). Our own analyses, shown in Table 2, also provide evidence that these scores capture the underlying construct of interest—i.e., the quality of classroom instruction—as opposed to construct-irrelevant sources of variation such as raters. Lesson-level intraclass correlations (ICC) range from 0.36 to 0.49, and are similar to other studies in which trained observers score the quality of teachers’ instruction (Bell et al., 2012; Hill et al., 2012). Our analyses focus on these lesson-level scores as the outcome of interest, though we also note that adjusted teacher-level ICCs that accumulate information across lessons are substantially higher, ranging from 0.55 to 0.69. Measurement error in our dependent variables can limit the precision of our estimates, but will not lead to attenuation bias, as is the case with measurement error in independent variables.

In most instances, teachers’ coaches conducted and scored observations. While this setup closely matches the purpose of coaching models that are organized around observation and feedback cycles led by the coach, it could bias our estimates of variation in coach effectiveness given that the coach is both the key input and the one responsible for measuring outcomes. Changes in teacher practices associated with different coaches may be due to coaches’ underlying effectiveness or to differences in how each coach scores instruction. At the same time, we find that, amongst a set of sites and years in which lessons were observed both by the teachers’ own coach and another observer, interrater agreement rates are comparable to other studies in which trained observers score the quality

of teachers' instruction (Bell et al., 2012; Hill et al., 2012): 70% for *Culture of Learning*, 66% for *Essential Content*, and 51% for *Demonstration of Learning* (see Table 2). Further, in a set of robustness tests that focus only on lessons observed by outside raters, we find that variation in coach effectiveness is larger than in the full sample.

Findings

Heterogeneity in Effectiveness Across Coaches

We begin, in Table 3, by showing the variation in coach effectiveness as measured by changes in each of the four measures of teaching practice (the three individual dimensions and the composite measure), pooling across all sites and years. We find that a 1 SD increase in coach effectiveness is associated with a roughly 0.2 SD increase over the course of the summer in the composite measure of teacher practice (0.19 SD for the coach random effect, and 0.22 SD for the coach-year random effect). Our estimates of variability in coach effectiveness are similar for *Culture of Learning* and *Essential Content*, and slightly larger for *Demonstration of Learning* (0.24 to 0.29 SD).

In Appendix Table 1, we re-estimate coach effects using a subset of site-years in which a rater other than teachers' own coach observed and scored their instruction. We find that the variation in coach effectiveness often is larger than in the full sample: roughly 0.3 SD for the composite measure of teaching practice, roughly 0.22 SD for *Culture of Learning* and *Essential Content*, and 0.33 to 0.35 SD for *Demonstration of Learning*. The latter dimension of practice is where inter-rater agreement rates between a teacher's own coach and another rater were lowest (see Table 2). Therefore, it appears that we are underestimating variation in coach effectiveness by using scores rated by teachers' own coach. That said, as we proceed with our results, we rely on the full sample in order to maximize precision and generalizability. Here and in Table 3, estimates of coach and coach-year variation are quite similar, though the latter often are estimated more precisely. Therefore, we focus primarily on coach-year random effects in the rest of our analyses.

In Table 4, we present additional estimates that examine the extent to which variation in coach effectiveness is driven by specific sites. Even though all sites operate under a common TNTP coaching model and management structure, each site hires its own coaches and provides training, support, and management to them. Given this, one might expect to see variation in changes in teacher practices and coach effectiveness across sites. However, overall, we find that it is the coach and not the site that appears to be primarily responsible for changes in teacher practice. In column 1, we nest coach-years within site-years in our random effects structure, finding negligible and non-significant variation at the site-year level (0.02 SD) and variation at the coach-year level (0.2 SD) that is very similar to our primary analyses that exclude the site-year random effect. In the next four columns, we disaggregate coach effects for the four largest training sites, each of which has a sample of at least 30 coaches when pooling across available years of data. Estimates of the coach-year variation range from 0.17 SD to 0.23 SD.

Coach Characteristics that Predict Changes in Teacher Practice

Knowing that coaches vary substantially in their effects on teacher practice begs the question: What characteristics, knowledge, and skills of coaches explain these differences? TNTP's administrative records include background data on coaches (see Table 1) that align with theory on some key dimensions of coach quality: (i) *years of coaching experience with TNTP* serves as a proxy for the accumulated knowledge and skills coaches build in their work over time (Connor, 2017), while (ii) *coach-teacher demographic matches* may increase the strength of interpersonal relationship between coaches and teachers (Bristol & Martin-Fernandez, 2019; Ladson-Billings, 1995).

In Table 5, we examine whether these characteristics predict changes in teacher outcomes, adding these characteristics to the fixed portion of our model outlined in equation (1) above. Here, we expand our analyses to focus on all four measures of teaching practice—rather than limiting just to the composite measure—given robust theoretical discussion about how race/ethnicity-matching

can be particularly beneficial for building culturally relevant and responsive classroom environments (Ladson-Billings, 1995). Given the composition of our teacher and coach samples (see Tables 1) that are comprised primarily of Black and White individuals, we focus on three race/ethnicity categories: Black, White, and non-Black/non-White. We exclude teachers and coaches who are missing information on race/ethnicity or gender. In the top panel of Table 5, we start with models that include main effects of individual coach characteristics; in the bottom panel, we interact coach demographic characteristics with teacher demographic characteristics to examine the role of matching.

In both the top and bottom panels, we do not find evidence that increased experience as a TNTP coach is associated with larger changes in teacher practice, relative to teachers whose coach has less experience. Estimates linking a dummy indicator for having a coach in their third year of experience with TNTP or higher (compared to having a first- or second-year coach) for the composite measure of instructional practice and *Culture of Learning* both are positive but not statistically significantly different from zero. Our sample is composed primarily of coaches with limited experience in this role at TNTP (see Table 1), and so we cannot make claims about whether having a coach who has many years of experience may make a difference for teachers' instructional practice. We also cannot observe when coaches had experience in a similar position but outside of TNTP. Limited variation in the coach experience variable can also limit statistical power. Nonetheless, the point estimates are small.

We find some evidence that having a male coach is related to changes in *Culture of Learning* (top panel), though the estimate (0.09 SD) only is statistically significant at the $p = 0.1$ threshold. In turn, we also examine male coach-teacher matches (bottom panel), finding positive point estimates when predicting all four teaching practice measures; however, none of these estimates is statistically significantly different from zero. We observe similar patterns for the main effect of having a Black coach: all four point estimates are positive but none are statistically significantly different from zero.

Comparatively, we find that assignment of a Black coach to a Black teacher is associated with a 0.18 SD increase in the composite measure of effective instruction, and a 0.22 SD increase in *Culture of Learning*. These estimates compare Black teachers with a Black coach to their Black peers with a White coach, White teachers with a non-White coach, and non-Black/non-White teachers with a White coach or a non-Black/non-White coach. Results are almost identical when we change the reference category. In these models, we also control for the main effect of having a Black coach. Though not shown in Table 5, none of these estimates are statistically significantly different from zero, which is consistent with patterns shown in the top panel. We also find that Black teachers assigned to a non-Black/non-White coach outperform their peers (0.26 SD for the composite measure and 0.25 SD for *Culture of Learning*). We do not find any statistically significant relationships of race-matching for White teachers working with a White coach.

Identification Check

The internal validity of our findings rests on the assumption that coach-teacher assignments are random, conditional on covariates included in the model (i.e., baseline measure of teaching practice, teacher demographics, and site-year and certification area fixed effects). We assess this assumption in Appendix Table 2 by conducting a falsification test that estimates the “effect” of coaches on observable background teacher characteristics (i.e., gender, race/ethnicity), still controlling for a baseline measure of the outcome, and site-year and certification area fixed effects. Positive and statistically significant estimates here do not invalidate our value-added methodology, but rather point to potential sorting bias that is not fully accounted for with the set of available covariates (Goldhaber & Chaplin, 2015). We find that the coach-level variation is zero or very close to zero when predicting each of the race/ethnicity dummy variables. Random effects models have known challenges when estimates are close to zero (Harville, 1977). For example, when the estimated variance approaches zero, the standard error is undefined (i.e., estimates in Appendix Table 2 predicting dummy indicators

for Black teacher and White teacher). To confirm that our estimates are true zeros, we estimated results to 10 decimal places, finding similar results.

When predicting teacher gender, we observe non-zero variation at the coach or coach-year level, but the estimate is roughly a third as large as when predicting teacher practices. These patterns suggest that our covariates likely have accounted for potential sorting bias, of coaches to teachers generally (relevant for analyses of individual coaches effects) and of coaches to teachers of different races or ethnicities (relevant for analyses of race/ethnicity matches).

Discussion and Conclusion

Using a value-added approach similar to the teacher effectiveness literature, we present evidence that individual coaches are the key ingredient for success of coaching programs. Across a range of models and specifications, we observe substantial variation across coaches in how teachers improve in their instructional practice. The magnitude of the coach-level variation as measured by changes in teacher practice is particularly large when compared to the full effect of coaching programs. We find that a 1 SD increase in coach effectiveness is associated with a 0.2 to 0.35 SD increase in multiple dimensions of teaching practice, and a 2 SD increase in coach effectiveness is associated with a 0.4 to 0.7 SD increase in teachers' instructional quality. Comparatively, meta-analytic estimates indicate that coaching programs, on average, improve teacher practice by roughly 0.5 SD (Kraft et al., 2018). In other words, variation in effectiveness across individual coaches explains almost the full effect of coaching programs.

One on hand, our study's focus on several hundred coaches working in school districts across the U.S. increases generalizability relative to other similar studies conducted with a small number of coaches or in a single setting (Blazar & Kraft, 2015, 2019; Russell et al., 2020). On the other hand, we focus only on the pre-service component of teacher training, and so we cannot make claims regarding variation in coach effects during in-service professional development. While pre-service teacher

coaching has less coverage in the empirical literature base compared to in-service programs, recent experimental evidence of pre-service coaching identifies effects on teacher practice that are on par with or larger than effects of in-service coaching (Cohen et al., 2020; Kraft et al., 2018).

Further, aligned to the work of other scholars (Connor, 2017; Denton & Hasbrouck, 2009; Joyce & Showers, 1981; Wong & Nicotera, 2006), we theorize that there are multiple potential mechanisms that might explain differences in coach effectiveness: the knowledge and skill that coaches bring to their work with teachers, coaches' interpersonal relationships with a given teacher, and the types of tools the coaches use in their interactions with teachers (e.g., providing direct feedback, modeling instruction). While use of administrative records means that we have a limited set of variables to capture these varied skills, we find initial evidence that the second avenue related to interpersonal relationships may be key to coach effectiveness and coaching program success. We find that Black teachers assigned to a Black or to a non-Black/non-White coach outperformed their peers in terms of changes in instructional practice; these differences are driven primarily by changes in classroom climate and cultural components of high-quality teaching. Drawing from the theoretical literature on teacher-student racial/ethnic matches (Bristol & Martin-Fernandez, 2019; Ladson-Billings, 1995; Redding, 2019), we argue that these patterns may be driven by the unique interpersonal relationships that teachers and coaches can develop when they have similar shared experiences and understandings. Comparatively, additional years of coaching experience—a proxy for the background knowledge and skill that coaches bring to their work—is not associated with increased teaching quality. The latter finding may be driven, in part, by our sample that is comprised largely of first-year TNTTP coaches.

To confirm and extend these findings, future research might estimate coach effects under experimental conditions, where coaches are randomly assigned to teachers. This design then could be paired with more extensive data collection on the various theorized dimensions of coach quality and

skill, with each dimension then linked to teacher outcomes. Identifying specific coach practices and skills that improve teachers' delivery of rigorous content and teachers' work with students around that content would help build on our findings. While our estimates of coach-level variation are largest when predicting *Demonstration of Learning*, we did not find that observable coach characteristics available in TNTP's administrative records predicted changes in this measure. Future research might also link individual coaches and their skills to student-level outcomes, in addition to teacher-level ones. Estimates of coach effects on student outcomes almost certainly will be smaller than coach effects on teacher-level outcomes, given that the former are more distal than the latter in the instructional improvement process. That said, the magnitude of variability in coach effectiveness associated with changes in teaching practices (upwards of 0.35 SD) suggests that relationships may further translate into changes in student outcomes. These lines of inquiry could be conducted both during pre-service training and in-service development provided to more veteran teachers.

Ultimately our findings have broader implications for teacher training organizations, schools, and districts interested in developing or expanding their coaching programs. Currently, school districts spend approximately \$18 billion on teacher development programs each year (Education Next, 2018) for the 3.5 million full-time teachers in the United States (NCES, 2020). The costs of teacher education and pre-service training also are substantial, with some calculations of alternative-route teacher certification programs costing upwards of \$40,000 per candidate (Kaufman et al., 2020). However, these dollars generally are found to have very little, if any, return on investment (Fryer, 2017; Harris & Sass, 2011; Yoon et al., 2007). Coaching provides an attractive alternative, achieving some of the largest impacts on teacher and student outcomes across all of the education intervention literature (Kraft et al., 2018).

Further, the overall costs of coaching programs are comparable to other training and development offerings. Knight and Skrtic (2021) find that the primary ingredients of coaching

programs are the coach salary and teacher time, with average costs ranging from \$5,300-\$10,500 per teacher per year. Examining coaching in an alternative-route teacher certification context, Kaufman et al. (2020) estimate that coaching comprises roughly a third of total per-teacher costs, at roughly \$13,000. The literature on costs of more traditional teacher development and training is older, but suggests that expenditures are similar, at \$3,100 to \$11,700 per teacher per year (Miles et al., 2004). (All cost estimates are adjusted to 2021 dollars.) In other words, coaching is likely to be substantially more cost effective than more traditional programs. Further, because coaching purposefully is individualized and differentiated, it likely makes sense to provide coaching only to some teachers who need it most and only in some school years. This approach would further decrease the overall coaching program costs from the district perspective.

At the same time, adopting and scaling instructional coaching in lieu of traditional teacher training and development is a risky proposition without knowing how to identify effective coaches—whose salary is the key cost driver of coaching programs (Kaufman et al., 2020; Knight & Skrtic, 2021)—and how to recruit, train, and support more of them. Based on findings from our study, we offer several recommendations for policy and practice. First, our value-added methodology offers one way to identify effective coaches. Like in the teacher effectiveness realm, these measures could be used to make ongoing personnel decisions related to retention and salary. Second, positive relationships between coach-teacher demographic matches and changes in teaching practice suggest that recruitment efforts may focus on building a diverse corps of coaches whose characteristics match demographics of local teacher workforces. We recognize that efforts to diversify coach workforces may work against simultaneous efforts to diversify the teacher workforce, given that coaches often are current or former teachers in the same or a nearby district (Darling-Hammond, 2017; Wenner & Campbell, 2017). That said, large effects of virtual coaching programs (e.g., Allen et al., 2011) suggest that hiring could occur outside of a local area. Further, aligned to the literature on teacher-student

race/ethnicity matching (Bristol & Martin-Fernandez, 2019; Ladson-Billings, 1995), we hypothesize that mechanisms underlying coach-teacher demographic matches likely are related to interpersonal relationships. Thus, school districts—and researchers—may focus on designing instruments to screen and train this skill set, particularly in instances where matching coach and teacher demographics may not be possible.

Rigorous empirical evidence indicates that coaching should be at the forefront of instructional improvement efforts. Scaling and sustaining these programs is doable (Kraft et al., 2018), but will require strategic planning that focuses primarily on building a corps of highly skilled coaches.

References

- Allen, J. P., Pianta, R. C., Gregory, A., Mikami, A. Y., & Lun, J. (2011). An interaction-based approach to enhancing secondary school instruction and student achievement. *Science*, *333*(6045), 1034-1037.
- Bell, C. A., Gitomer, D. H., McCaffrey, D. F., Hamre, B. K., Pianta, R. C., & Qi, Y. (2012). An argument approach to observation protocol validity. *Educational Assessment*, *17*(2-3), 62-87.
- Blazar, D., & Kraft, M. A. (2017). Teacher and teaching effects on students' attitudes and behaviors. *Educational Evaluation and Policy Analysis*, *39*(1), 146-170.
- Blazar, D., and Kraft, M. A. (2019). Balancing Rigor, Replication, and Relevance: A Case for Multiple-Cohort, Longitudinal Experiments. *AERA Open*, *5*(3).
- Blazar, D., & Kraft, M. A. (2015). Exploring mechanisms of effective teacher coaching: A tale of two cohorts from a randomized experiment. *Educational Evaluation and Policy Analysis*, *37*(4), 542-566.
- Bristol, T. J., & Martin-Fernandez, J. (2019). The added value of Latinx and Black teachers for Latinx and Black students: Implications for policy. *Policy Insights from the Behavioral and Brain Sciences*, *6*(2), 147-153.
- Britton, L. R., & Anderson, K. A. (2010). Peer coaching and pre-service teachers: Examining an underutilised concept. *Teaching and Teacher Education*, *26*(2), 306-314.
- Cohen, J., Wong, V., Krishnamachari, A., & Berlin, R. (2020). Teacher coaching in a simulated environment. *Educational Evaluation and Policy Analysis*, *42*(2), 208-231.
- Connor, C. M. (2017). Commentary on the special issue on instructional coaching models: Common elements of effective coaching models. *Theory into Practice*, *56*(1), 78-83.
- Darling-Hammond, L. (2017). Teacher education around the world: What can we learn from international practice?. *European Journal of Teacher Education*, *40*(3), 291-309.

- Denton, C. A., & Hasbrouck, J. A. N. (2009). A description of instructional coaching and its relationship to consultation. *Journal of Educational and Psychological Consultation*, 19(2), 150-175.
- Domina, T., Lewis, R., Agarwal, P., & Hanselman, P. (2015). Professional sense-makers: Instructional specialists in contemporary schooling. *Educational Researcher*, 44(6), 359-364.
- Education Next. (2018, June 12). EdStat: \$18 Billion a Year is Spent on Professional Development for U.S. Teachers. *Education Next*. Retrieved from: <http://www.educationnext.org/edstat-18-billion-year-spent-professional-development-u-s-teachers/>
- Fryer, J., Roland G. (2017). The production of human capital in developed countries: Evidence from 196 randomized field experiments. In *Handbook of economic field experiments* (Vol. 2, pp. 95-322). North-Holland.
- Goldhaber, D., & Chaplin, D. D. (2015). Assessing the “Rothstein Falsification Test”: Does it really show teacher value-added models are biased?. *Journal of Research on Educational Effectiveness*, 8(1), 8-34.
- Grissom, J. A., Kalogrides, D., & Loeb, S. (2015). Using student test scores to measure principal performance. *Educational Evaluation and Policy Analysis*, 37(1), 3-28.
- Guarino, C. M., Maxfield, M., Reckase, M. D., Thompson, P. N., & Wooldridge, J. M. (2015). An evaluation of empirical Bayes’s estimation of value-added teacher performance measures. *Journal of Educational and Behavioral Statistics*, 40(2), 190-222.
- Hanushek, E. A., & Rivkin, S. G. (2010). Generalizations about using value-added measures of teacher quality. *American Economic Review*, 100(2), 267-71.
- Harris, D. N., & Sass, T. R. (2011). Teacher training, teacher quality and student achievement. *Journal of Public Economics*, 95(7-8), 798-812.
- Harville, D. A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *Journal of the American Statistical Association*, 72(358), 320-338.

- Hill, H. C., Blazar, D., & Lynch, K. (2015). Resources for teaching: Examining personal and institutional predictors of high-quality instruction. *AERA Open*, 1(4).
- Hill, C. J., Bloom, H. S., Black, A. R., & Lipsey, M. W. (2008). Empirical benchmarks for interpreting effect sizes in research. *Child Development Perspectives*, 2(3), 172-177.
- Hill, H. C., Charalambous, C. Y., Blazar, D., McGinn, D., Kraft, M. A., Beisiegel, M., ... & Lynch, K. (2012). Validating arguments for observational instruments: Attending to multiple sources of variation. *Educational Assessment*, 17(2-3), 88-106.
- Holmstrom, B., & Milgrom, P. (1991). Multitask principal-agent analyses: Incentive contracts, asset ownership, and job design. *JL Econ. & Org.*, 7, 24.
- Howley, A. A., Dudek, M. H., Rittenberg, R., & Larson, W. (2014). The development of a valid and reliable instrument for measuring instructional coaching skills. *Professional Development in Education*, 40(5), 779-801.
- Joyce, B. R., & Showers, B. (1981). Transfer of training: The contribution of “coaching”. *Journal of Education*, 163(2), 163-172.
- Kane, T. J., & Staiger, D. O. (2012). *Gathering feedback for teaching: Combining high-quality observations with student surveys and achievement gains*. Research Paper. MET Project. Bill & Melinda Gates Foundation.
- Kaufman, J. H., Master, B. K., Huguet, A., Yoo, P. Y., Faxon-Mills, S., Schulker, D., & Grimm, G. E. (2020). *Growing teachers from within: Implementation, impact, and cost of an alternative teacher preparation program in three urban school districts*. Research Report. RR-A256-1. RAND Corporation. Retrieved from: <https://eric.ed.gov/?id=ED609341>
- Knight, D. S., & Skrtic, T. M. (2021). Cost-effectiveness of instructional coaching: Implementing a design-based, continuous improvement model to advance teacher professional development. *Journal of School Leadership*, 31(4), 318-342.

- Kraft, M. A., Blazar, D., and Hogan, D. (2018). The effect of teacher coaching on instruction and achievement: A Meta-Analysis of the Causal Evidence: *Review of Educational Research*, 88(4) 547-588.
- Ladson-Billings, G. (1995). But that's just good teaching! The case for culturally relevant pedagogy. *Theory into practice*, 34(3), 159-165.
- Marsh, J. A., McCombs, J. S., & Martorell, F. (2012). Reading coach quality: Findings from Florida middle schools. *Literacy Research and Instruction*, 51(1), 1-26.
- McEachin, A., Schweig, J. D., Perera, R., & Opper, I. M. (2018). *Validation study of the TNTP Core Teaching Rubric*. RAND. Retrieved from: https://www.rand.org/content/dam/rand/pubs/research_reports/RR2600/RR2623/RAND_RR2623.pdf
- Menzes, A., & Maier, A. (2014). Fast Start: Training Better Teachers Faster, with Focus, Practice and Feedback. *TNTP*. Retrieved from: <https://files.eric.ed.gov/fulltext/ED559704.pdf>
- Mulhern, C. (2019). Beyond teachers: Estimating individual guidance counselors' effects on educational attainment. *Cambridge, MA: Harvard University*. Retrieved January, 26, 2020.
- National Center for Education Statistics. (2008). *School and Staffing Survey*. Retrieved from: https://nces.ed.gov/pubs2009/2009321/tables/sass0708_2009321_s12n_06.asp
- National Center for Education Statistics. (2016). *National Teacher and Principal Survey*. Retrieved from: https://nces.ed.gov/surveys/ntps/tables/Table_5_042617_fl_school.asp
- National Center for Education Statistics. (2020). *Characteristics of Public School Teachers*. https://nces.ed.gov/programs/coe/indicator_clr.asp
- Nye, B., Konstantopoulos, S., & Hedges, L. V. (2004). How large are teacher effects? *Educational Evaluation and Policy Analysis*, 26(3), 237-257.

- Miles, K. H., Odden, A., Fermanich, M., & Archibald, S. (2004). Inside the black box of school district spending on professional development: Lessons from five urban districts. *Journal of Education Finance, 30*(1), 1-26.
- Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical linear models: Applications and data analysis methods* (Vol. 1). Sage.
- Redding, C. (2019). A teacher like me: A review of the effect of student–teacher racial/ethnic matching on teacher perceptions of students and student academic and behavioral outcomes. *Review of Educational Research, 89*(4), 499-535.
- Russell, J. L., Correnti, R., Stein, M. K., Thomas, A., Bill, V., & Speranzo, L. (2020). Mathematics coaching for conceptual understanding: Promising evidence regarding the Tennessee math coaching model. *Educational Evaluation and Policy Analysis, 42*(3), 439–466.
- Safran, D. G., Taira, D. A., Rogers, W. H., Kosinski, M., Ware, J. E., & Tarlov, A. R. (1998). Linking primary care performance to outcomes of care. *Journal of Family Practice, 47*, 213-220.
- Shannon, D. K., Snyder, P. A., Hemmeter, M. L., & McLean, M. (2021). Exploring Coach–Teacher Interactions Within a Practice-Based Coaching Partnership. *Topics in Early Childhood Special Education, 40*(4), 229-240.
- Shen, J. (1997). Has the alternative certification policy materialized its promise? A comparison between traditionally and alternatively certified teachers in public schools. *Educational Evaluation and Policy Analysis, 19*(3), 276-283
- Slavin, R., & Smith, D. (2009). The relationship between sample sizes and effect sizes in systematic reviews in education. *Educational Evaluation and Policy Analysis, 31*(4), 500-506.
- TNTP. (2013). *Leap year: Assessing and supporting effective first-year teachers*. Retrieved from: https://tntp.org/assets/documents/TNTP_Leap_Year_2013.pdf

- TNTP. (2014). *TNTP Core Teaching Rubric: A tool for conducting Common Core-aligned classroom observations*. Retrieved from: <https://tntp.org/publications/view/tntp-core-teaching-rubric-a-tool-for-conducting-classroom-observations>
- TNTP. (2018). *The opportunity myth: Technical appendix*. Retrieved from: <https://files.eric.ed.gov/fulltext/ED590222.pdf>
- Walsh, K., & Jacobs, S. (2007). *Alternative certification isn't alternative*. Thomas B. Fordham Institute. Retrieved from: https://www.nctq.org/nctq/images/Alternative_Certification_Isnt_Alternative.pdf
- Wenner, J. A., & Campbell, T. (2017). The theoretical and empirical basis of teacher leadership: A review of the literature. *Review of Educational Research*, 87(1), 134-171.
- Wilson, S. M. (2014). Innovation and the evolving system of US teacher preparation. *Theory into Practice*, 53(3), 183-195.
- Wong, K., & Nicotera, A. (2006). Peer coaching as a strategy to build instructional capacity in low performing schools. In K. Wong and S. Rutledge (Eds.), *System-wide efforts to improve student achievement*. Greenwich, CT: Information Age Publishing.
- Yoon, K. S., Duncan, T., Lee, S. W. Y., Scarloss, B., & Shapley, K. L. (2007). Reviewing the Evidence on How Teacher Professional Development Affects Student Achievement. Issues & Answers. REL 2007-No. 033. *Regional Educational Laboratory Southwest (NJ1)*.
- Yopp, D. A., Burroughs, E. A., Sutton, J. T., & Greenwood, M. C. (2019). Variations in coaching knowledge and practice that explain elementary and middle school mathematics teacher change. *Journal of Mathematics Teacher Education*, 22(1), 5-36.

Tables

Table 1. Characteristics of Teachers and Coaches

	Teachers	Coaches
<u>Demographics</u>		
Female	0.66	0.67
Male	0.30	0.21
Missing Gender	0.03	0.12
Asian	0.03	0.03
Black	0.26	0.25
Hispanic	0.04	0.04
White	0.40	0.52
Multiple Races/Ethnicities	0.06	0.04
Missing Race/Ethnicity	0.20	0.12
<u>Certification Area</u>		
Early Childhood Education	0.07	NA
Elementary School	0.24	NA
English Language Arts (ELA)	0.11	NA
Math	0.08	NA
Science	0.09	NA
Social Studies	0.01	NA
English as a Second Language	0.04	NA
Special Education	0.15	NA
Foreign Language	0.01	NA
Missing Certification Area	0.20	NA
<u>Coaching Experience with TNTP</u>		
Total yrs.	NA	1.36
1 yr. Experience	NA	0.74
2 yrs. Experience	NA	0.19
3 or more yrs. Experience	NA	0.07
Persons (<i>n</i>)	3,526	317

Table 2. Descriptive Statistics for Observation Scores

Observation Scores (1 to 3 Scale)	Univariate Statistics				Reliability		
	Last Score		First Score		Lesson- Level ICC	Teacher- Level Adjusted ICC	Inter-Rater Agreement
	Mean	SD	Mean	SD			
Composite	2.51	0.50	2.25	0.53	0.49	0.69	NA
Culture of Learning	2.51	0.63	2.28	0.68	0.47	0.68	70%
Essential Content	2.72	0.52	2.50	0.63	0.31	0.55	66%
Demonstration of Learning	2.31	0.70	1.97	0.71	0.36	0.61	51%

Note: ICC = intraclass correlation. Following a generalizability framework, teacher-level ICCs are adjusted for the median number of lessons per teacher. Inter-rater agreement is not calculated for the composite, as researchers (not observers) calculated the composite as an average of the other three dimensions of teaching practice.

Table 3. Standard Deviation of Coach-Level Variation, Pooling Across Sites

	Composite	Culture of Learning	Essential Content	Demonstration of Learning
Coach-Year Random Effect	0.219*** (0.023)	0.217*** (0.024)	0.202*** (0.026)	0.288*** (0.023)
Coach Random Effect	0.191*** (0.024)	0.198*** (0.025)	0.170*** (0.028)	0.241*** (0.025)
Teachers (<i>n</i>)	3,526	3,526	3,526	3,526
Coach-Years (<i>n</i>)	430	430	430	430
Coaches (<i>n</i>)	317	317	317	317

Notes: Each estimate comes from a separate multilevel model of teachers' end-of-summer observation score on baseline scores for all three dimensions of practice, teacher gender and race/ethnicity, certification area fixed effects, and site-year fixed effects. *** $z > 3.29$, where z equals the ratio of a given random effects parameter estimate to its standard error. These z -scores do not correspond precisely to p -values as in a traditional linear regression framework, as the confidence interval for a random effect estimate is not symmetric around the estimate (random effect estimates have a lower bound of zero). These traditional markers of significance are included as an aid to the reader, but they should be interpreted with caution.

Table 4. Standard Deviation of Coach-Level Variation on Composite Measure of Instructional Practice, by Site

	All Sites	Site 1	Site 2	Site 3	Site 4
Site-Year Random Effect	0.023 (0.085)	NA	NA	NA	NA
Coach-Year Random Effect	0.196*** (0.024)	0.219*** (0.042)	0.218*** (0.050)	0.171~ (0.089)	0.225*** (0.074)
Teachers (<i>n</i>)	3,526	873	719	326	399
Coach-Years (<i>n</i>)	430	96	90	45	46
Coaches (<i>n</i>)	317	59	47	36	32

Notes: Estimates in each column come from separate multilevel models of teachers' end-of-summer observation score on baseline scores for all three dimensions of practice, teacher gender and race/ethnicity, certification area fixed effects, and site-year or year fixed effects. ~ $z > 1.64$, *** $z > 3.29$, where z equals the ratio of a given random effects parameter estimate to its standard error. These z -scores do not correspond precisely to p -values as in a traditional linear regression framework, as the confidence interval for a random effect estimate is not symmetric around the estimate (random effect estimates have a lower bound of zero). These traditional markers of significance are included as an aid to the reader, but they should be interpreted with caution.

Table 5. Predictive Power of Coach Characteristics

	Composite	Culture of Learning	Essential Content	Demonstration of Learning
<u>Main Effects</u>				
3 or more yrs. Experience	0.046 (0.104)	0.073 (0.108)	-0.010 (0.112)	-0.010 (0.116)
Black Coach	0.081 (0.052)	0.082 (0.054)	0.065 (0.056)	0.047 (0.059)
Non-Black/Non-White Coach	0.050 (0.066)	0.062 (0.069)	0.038 (0.071)	-0.010 (0.074)
Male Coach	0.061 (0.050)	0.091~ (0.052)	-0.02 (0.054)	0.015 (0.057)
<u>Demographic Matching</u>				
3 or more yrs. Experience	0.036 (0.104)	0.064 (0.108)	-0.017 (0.112)	-0.012 (0.116)
Black Teacher*Black Coach	0.181~ (0.107)	0.222* (0.111)	0.073 (0.119)	-0.016 (0.116)
Black Teacher*Non-Black/Non-White Coach	0.261* (0.131)	0.246~ (0.136)	0.168 (0.146)	0.088 (0.141)
White Teacher*White Coach	-0.135 (0.097)	-0.158 (0.101)	-0.086 (0.108)	0.007 (0.104)
Male Teacher*Male Coach	0.058 (0.079)	0.024 (0.082)	0.083 (0.088)	0.026 (0.085)
Coaches (<i>n</i>)	265	265	265	265
Teachers (<i>n</i>)	2,591	2,591	2,591	2,591

Notes: Estimates in each panel and column come from separate multilevel models that include coach-year random effects. All models control for baseline scores for all three dimensions of practice, teacher gender and race/ethnicity, certification area fixed effects, and site-year fixed effects. In models with coach-teacher demographic match indicators, main effects of coach and teacher demographics also included as controls. *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$, ~ $p < 0.1$.

Appendix Tables

Appendix Table 1. Standard Deviation of Coach-Level Variation in Sample where Raters are not Teachers' Coach

	Composite	Culture of Learning	Essential Content	Demonstration of Learning
Coach-Year Random Effect	0.301*** (0.051)	0.224*** (0.060)	0.211*** (0.058)	0.352*** (0.052)
Coach Random Effect	0.288*** (0.056)	0.220*** (0.059)	0.222*** (0.058)	0.327*** (0.059)
Teachers (<i>n</i>)	749	749	749	749
Coach-Years (<i>n</i>)	92	92	92	92
Coaches (<i>n</i>)	81	81	81	81

Notes: Each estimate comes from a separate multilevel model of teachers' end-of-summer observation score on baseline scores for all three dimensions of practice, teacher gender and race/ethnicity, certification area fixed effects, and site-year fixed effects. *** $z > 3.29$, where z equals the ratio of a given random effects parameter estimate to its standard error. These z -scores do not correspond precisely to p -values as in a traditional linear regression framework, as the confidence interval for a random effect estimate is not symmetric around the estimate (random effect estimates have a lower bound of zero). These traditional markers of significance are included as an aid to the reader, but they should be interpreted with caution.

Appendix Table 2. Falsification Tests

	Female	Asian	Black	Hispanic	White
Coach-Year Random Effect	0.074*** (0.015)	0.020* (0.008)	0.000 --	0.014 (0.012)	0.000 --
Coach Random Effect	0.068*** (0.014)	0.011 (0.013)	0.000 --	0.020* (0.008)	0.000 --
Coaches (<i>n</i>)	317	317	317	317	317
Teachers (<i>n</i>)	3,526	3,526	3,526	3,526	3,526

Notes: Each estimate comes from a separate multilevel model of teachers' end-of-summer observation score on baseline scores for all three dimensions of practice, certification area fixed effects, and site-year fixed effects. When female is the outcome, a missing gender dummy and race/ethnicity dummies also are included as controls; when race/ethnicity dummies are the outcomes, a missing race/ethnicity dummy and gender dummies are included as controls. "--" indicates that the relevant parameter could not be estimated. * $z > 1.96$, *** $z > 3.29$, where z equals the ratio of a given random effects parameter estimate to its standard error. These z -scores do not correspond precisely to p -values as in a traditional linear regression framework, as the confidence interval for a random effect estimate is not symmetric around the estimate (random effect estimates have a lower bound of zero). These traditional markers of significance are included as an aid to the reader, but they should be interpreted with caution.