

Discretion in Hiring

Mitchell Hoffman
Lisa B. Kahn
Danielle Li

Working Paper 16-055



Discretion in Hiring

Mitchell Hoffman

University of Toronto

Lisa B. Kahn

Yale University

Danielle Li

Harvard Business School

Working Paper 16-055

Copyright © 2015 by Mitchell Hoffman, Lisa B. Kahn, and Danielle Li

Working papers are in draft form. This working paper is distributed for purposes of comment and discussion only. It may not be reproduced without permission of the copyright holder. Copies of working papers are available from the author.

Discretion in Hiring

Mitchell Hoffman
University of Toronto

Lisa B. Kahn
Yale University & NBER

Danielle Li
Harvard University

September 25, 2015

PRELIMINARY & INCOMPLETE

Please do not cite or circulate without permission

Abstract

Who should make hiring decisions? We propose an empirical test for assessing whether firms should rely on hard metrics such as job test scores or grant managers discretion in making hiring decisions. We implement our test in the context of the introduction of a valuable job test across 15 firms employing low-skill service sector workers. Our results suggest that firms can improve worker quality by limiting managerial discretion. This is because, when faced with similar applicant pools, managers who exercise more discretion (as measured by their likelihood of overruling job test recommendations) systematically end up with worse hires.

*Correspondence: Mitchell Hoffman, University of Toronto Rotman School of Management, 105 St. George St., Toronto, ON M5S 3E6. Email: mitchell.hoffman@rotman.utoronto.ca. Lisa Kahn, Yale School of Management, 165 Whitney Ave, PO Box 208200, New Haven, CT 06511. Email: lisa.kahn@yale.edu. Danielle Li, Harvard Business School, 211 Rock Center, Boston, MA 02163. Email: dli@hbs.edu. We are grateful to Jason Abaluck, Ricardo Alonso, David Berger, Arthur Campbell, David Deming, Alex Frankel, Jin Li, Liz Lyons, Steve Malliaris, Mike Powell, Kathryn Shaw, Steve Tadelis, and numerous seminar participants. Hoffman acknowledges financial support from the Social Science and Humanities Research Council of Canada. All errors are our own.

1 Introduction

Hiring the right workers is one of the most important and difficult problems that a firm faces. Resumes, interviews, and other screening tools are often limited in their ability to reveal whether a worker has the right skills or will be a good fit. Further, the managers that firms employ to gather and interpret this information may have poor judgement or preferences that are imperfectly aligned with firm objectives.¹ Firms may thus face both information and agency problems when making hiring decisions.

The increasing adoption of “workforce analytics” and job testing has provided firms with new hiring tools.² Job testing has the potential to both improve information about the quality of candidates and to reduce agency problems between firms and human resource (HR) managers. As with interviews, job tests provide an additional signal of a worker’s quality. Yet, unlike interviews and other subjective assessments, job testing provides information about worker quality that is directly verifiable by the firm.

What is the impact of job testing on the quality of hires and how should firms use job tests, if at all? In the absence of agency problems, firms should allow managers discretion to weigh job tests alongside interviews and other private signals when deciding whom to hire. Yet, if managers are biased or if their judgment is otherwise flawed, firms may prefer to limit discretion and place more weight on test results, even if this means ignoring the private information of the manager. Firms may have difficulty evaluating this trade off because they cannot tell whether a manager hires a candidate with poor test scores because he or she has private evidence to the contrary, or because he or she is biased or simply mistaken.

In this paper, we evaluate the introduction of a job test and develop a diagnostic to inform how firms should incorporate it into their hiring decisions. Using a unique personnel dataset on HR managers, job applicants, and hired workers across 15 firms that adopt job testing, we present two key findings. First, job testing substantially improves the match quality of hired workers: those hired with job testing have about 15% longer tenures than

¹For example, a manager could have preferences over demographics or family background that do not maximize productivity. In a case study of elite professional services firms, Riviera (2012) shows that one of the most important determinants of hiring is the presence of shared leisure activities.

²See, for instance, *Forbes*: <http://www.forbes.com/sites/joshbersin/2013/02/17/bigdata-in-human-resources-talent-analytics-comes-of-age/>.

those hired without testing. Second, managers who overrule test recommendations more often hire workers with lower match quality, as measured by job tenure. This second result suggests that managers exercise discretion because they are biased or have poor judgement, not because they are better informed. This implies that firms in our setting can further improve match quality by limiting managerial discretion and placing more weight on the test.

Our paper makes the following contributions. First, we provide new evidence that managers systematically make hiring decisions that are not in the interest of the firm. This generates increased turnover in a setting where workers already spend a substantial fraction of their tenure in paid training rather than in producing output. Second, we show that job testing can improve hiring outcomes not simply by providing more information, but by making information verifiable, and thereby expanding the scope for contractual solutions to agency problems within the firm. Finally, we develop a simple tractable test for assessing the value of discretion in hiring. Our test uses data likely available to any firm with job testing, and is applicable to a wide variety of settings where at least one objective correlate of productivity is available.

We begin with a model in which firms rely on potentially biased HR managers who observe both public and private signals of worker quality. Using this model, we develop a simple empirical diagnostic based on the following intuition: if managers make exceptions to test recommendations because they have superior private information about a worker's quality, then we would expect better informed managers to both be more likely to make exceptions and to hire workers who are a better fit. As such, a positive correlation between exceptions and outcomes suggests that the discretion granted was valuable. If, in contrast, managers who make more exceptions hire workers with worse outcomes, then it is likely that managers are either biased or mistaken, and firms should limit discretion.

We apply this test using data from an anonymous firm that provides online job testing services to client firms. Our sample consists of 15 client firms who employ low-skill service-sector workers. Prior to the introduction of testing, firms employed HR managers involved in hiring new workers. After the introduction of testing, HR managers were also given access to a test score for each applicant: green (high potential candidate), yellow (moderate potential

candidate), or red (lowest rating).³ Managers were encouraged to factor the test into their hiring decisions but were still given discretion to use other signals of quality.

First, we estimate the impact of introducing a job test on the match quality of hired workers. By examining the staggered introduction of job testing across our sample locations, we show that cohorts of workers hired with job testing have about 15% longer tenures than cohorts of workers hired without testing. We provide a number of tests in the paper to ensure that our results are not driven by the endogenous adoption of testing or by other policies that firms may have concurrently implemented.

This finding suggests that job tests contain valuable information about the match quality of candidates. Next, we ask how firms should use this information, in particular, whether firms should limit discretion and follow test recommendations, or allow managers to exercise discretion and make exceptions to those recommendations. A unique feature of our data is that it allows us to measure the exercise of discretion explicitly: we observe when a manager hires a worker with a test score of yellow when an applicant with a score of green goes unhired (or similarly, when a red is hired above a yellow or a green). As explained above, the correlation between a manager's likelihood of making these exceptions and eventual outcomes of hires can inform whether the exercise of discretion is beneficial from the firm's perspective. Across a variety of specifications, we find that the exercise of discretion is strongly correlated with worse outcomes. Even when faced with applicant pools that are identical in terms of test scores, managers that make more exceptions systematically hire workers who are more likely to quit or be fired.

Finally, we show that our results are unlikely to be driven by the possibility that managers sacrifice job tenure in search of workers who have higher quality on other dimensions. If this were the case, limiting discretion may improve worker durations, but at the expense of other quality measures. To assess whether this is a possible explanation for our findings, we examine the relationship between hiring, exceptions, and a direct measure of productivity, daily output per hour, which we observe for a subset of firms in our sample. Based on this supplemental analysis, we no evidence that firms are trading off duration for higher produc-

³Section 2 provides more information on the job test.

tivity. Taken together, our findings suggest that firms could improve both match quality and worker productivity by placing more weight on the recommendations of the job test.

As data analytics becomes more frequently applied to human resource management decisions, it becomes increasingly important to understand how these new technologies impact the organizational structure of the firm and the efficiency of worker-firm matching. While a large theoretical literature has studied how firms should allocate authority, ours is the first paper to provide an empirical test for assessing the value of discretion in hiring.⁴ Our findings provide direct evidence that screening technologies can help resolve agency problems by improving information symmetry, and thereby relaxing contracting constraints. In this spirit, our paper is related to the classic Baker and Hubbard (2004) analysis of the adoption of on board computers in the trucking industry.

We also contribute to a small, but growing literature on the impact of screening technologies on the quality of hires.⁵ Our work is most closely related to Autor and Scarborough (2008), the first paper in economics to provide an estimate of the impact of job testing on worker performance. The authors evaluate the introduction of a job test in retail trade, with a particular focus on whether testing will have a disparate impact on minority hiring. Our paper, by contrast, studies the implications of job testing on the allocation of authority within the firm.

Our work is also relevant to a broader literature on hiring and employer learning.⁶ Oyer and Schaefer (2011) note in their handbook chapter that hiring remains an important open

⁴For theoretical work, see Bolton and Dewatripont (2012) for a survey and Dessein (2002) and Alonso and Matouschek (2008) for particularly relevant instances. There is a small empirical literature on bias, discretion and rule-making in other settings. For example, Paravisini and Schoar (2012) find that credit scoring technology aligns loan offer incentives and improves lending performance. Li (2012) documents an empirical tradeoff between expertise and bias among grant selection committees. Kuziemko (2013) shows that the exercise of discretion in parole boards is efficient, relative to fixed sentences.

⁵Other screening technologies include labor market intermediaries (e.g., Autor (2001), Stanton and Thomas (2014), Horton (2013)), and employee referrals (e.g., Brown et al., (2015), Burks et al. (2015) and Pallais and Sands (2015)).

⁶A central literature in labor economics emphasizes that imperfect information generates substantial problems for allocative efficiency in the labor market. This literature suggests imperfect information is a substantial problem facing those making hiring decisions. See for example Jovanovic (1979), Farber and Gibbons (1996), Altonji and Pierret (2001), and Kahn and Lange (2014).

area of research. We point out that hiring is made even more challenging because firms must often entrust these decisions to managers who may be biased or exhibit poor judgment.⁷

Lastly, our results are broadly aligned with findings in psychology and behavioral economics that emphasize the potential of machine-based algorithms to mitigate errors and biases in human judgement across a variety of domains.⁸

The remainder of this paper proceeds as follows. Section 2 describes the setting and data. Section 3 evaluates the impact of testing on the quality of hires. Section 4 presents a model of hiring with both hard and soft signals of quality. Section 5 evaluates the role of discretion in test adoption. Section 6 concludes.

2 Setting and Data

Firms have increasingly incorporated testing into their hiring practices. One explanation for this shift is that the increasing power of data analytics has made it easier to look for regularities that predict worker performance. We obtain data from an anonymous job testing provider that follows such a model. We hereafter term this firm the “data firm.” In this section we summarize the key features of our dataset. More detail can be found in Appendix A.

The data firm offers a test designed to predict performance for a particular job in the low-skilled service sector. To preserve the confidentiality of the data firm, we are unable to reveal the exact nature of the job, but it is conducted in a non-retail environment and is similar to jobs such as data entry work, standardized test grading and call center work. The data firm sells its services to clients (hereafter, “client firms”) that wish to fill these types of positions. We have 15 such client firms in our dataset.

The job test consists of an online questionnaire comprising a large battery of questions, including those on technical skills, personality, cognitive skills, fit for the job, and various job scenarios. The data firm then matches these responses with subsequent performance in

⁷This notion stems from the canonical principal-agent problem, for instance as in Aghion and Tirole (1997). In addition, many other models of management focus on moral hazard problems generated when a manager is allocated decision rights.

⁸See Kuncel et. al. (2013) for a meta-analysis of this literature and Kahneman (2011) for a behavioral economics perspective.

order to identify the questions or sets of questions that are the most predictive of future workplace success in this setting. These correlations are then aggregated by a proprietary algorithm to deliver a *green, yellow, red* job test score.

In its marketing materials, our data firm emphasizes the ability of its job test to reduce worker turnover, which is a perennial challenge for firms employing low skill service sector workers. To illustrate this concern, Figure 1 shows a histogram of job tenure for completed spells (75% of the spells in our data) among employees in our sample client firms. The median worker (solid red line) stays only 99 days, or just over 3 months. Twenty percent of hired workers leave after only a month. At the same time, our client firms generally report spending the first several weeks training each new hire, during which time the hire is being paid.⁹ Correspondingly, our analysis will also focus on job retention as the primary measure of hiring quality. For a subset of our client firms we also observe a direct measure of worker productivity: output per hour.¹⁰ Because these data are available for a much smaller set of workers (roughly a quarter of hired workers), we report these findings separately when we discuss alternative explanations.

Prior to testing, our client firms gave their managers discretion to make hiring decisions or recommendations based on interviews and resumes.¹¹ After testing, firms made scores available to managers and encouraged them to factor scores into hiring recommendations, but authority over hiring decisions was still typically delegated to managers.¹²

Our data contains information on each hired worker, including hire and termination dates, the reason for the exit, job function, and location for each hire. This information is collected by client firms and shared with the data firm. In addition, once a partnership with the data firm forms, we can also observe applicant test scores, application date, and an identifier for the HR manager responsible for a given applicant.

⁹Each client firm in our sample provides paid training to its workforce. Reported lengths of training vary considerably, from around 1-2 weeks to around a couple months or more.

¹⁰A similar productivity measure was used in Lazear et al., (2015) to evaluate the value of bosses in a comparable setting to ours.

¹¹In addition, the data firm informed us that a number of client firms had some other form of testing before the introduction of the data firm's test.

¹²We do not directly observe authority relations in our data. However, in a survey that the data firm conducted with a number of the client firms, the client firms reported that managers were not required to hire strictly by the test.

In the first part of this paper, we examine the impact of testing technology on worker match quality, as measured by tenure. For any given client firm, testing was rolled out gradually at roughly the location level. During the period in which the test is being introduced, not all applicants to the same location received test scores.¹³ We therefore impute a location-specific date of testing adoption. Our preferred metric for the date of testing adoption is the first date in which at least 50% of the workers hired in that month and location have a test score. Once testing is adopted at a location, based on our definition, we impose that testing is thereafter always available.¹⁴ In practice, this choice makes little difference and we are robust to a number of other definitions, for example, whether the individual or any hire in a cohort was tested.

Table 1 provides sample characteristics. Across our whole sample period we have nearly 300,000 hires; two-thirds of these were observed before testing was introduced and one-third were observed after, based on our preferred imputed definition of testing. Once we link applicants to the HR manager responsible for them (only after testing), we have 555 such managers in the data.¹⁵ These managers primarily serve a recruiting role, and are unlikely to manage day-to-day production.¹⁶ Post-testing, when we have information on applicants as well as hires, we have nearly 94,000 hires and a total of 690,000 applicants.

Table 1 also reports worker performance pre- and post-testing, by test score. On average, greens stay 12 days (11%) longer than yellows, who stay 17 days (18%) longer than reds. These differences are statistically significant and hold up to the full range of controls described below. This provides some evidence that test scores are indeed informative about worker performance. Even among the selected sample of hired workers, better test scores predict longer tenures. We might expect these differences to be even larger in the overall applicant population if managers hire red and yellow applicants only when unobserved quality is particularly high. On our productivity measure, output per hour, which averages roughly 8, performance is fairly similar across color.

¹³We are told by the data firm, however, that the intention of clients was generally to bring testing into a location at the same time for workers in that location.

¹⁴This fits patterns in the data, for example, that most locations weakly increase the share of applicants that are tested throughout our sample period.

¹⁵Our data firm links applicants to the manager who was most responsible for handling the applicant's case, though other managers may take part in hiring decisions as well.

¹⁶The HR managers we study are referred to as recruiters by our data provider.

3 The Impact of Testing

3.1 Empirical Strategy

Before examining whether firms should grant managers discretion over how to use job testing information, we first evaluate the impact of introducing testing information itself. To do so, we exploit the gradual roll-out in testing across locations and over time, and examine its impact on worker match quality, as measured by tenure:

$$\text{Outcome}_{lt} = \alpha_0 + \alpha_1 \text{Testing}_{lt} + \delta_l + \gamma_t + \text{Controls} + \epsilon_{lt} \quad (1)$$

Equation (1) compares outcomes for workers hired with and without job testing. We regress a productivity outcome (Outcome_{lt}) for workers hired to a location l , at time t , on an indicator for whether testing was available at that location at that time (Testing_{lt}) and controls. In practice, we define testing availability as whether the median hire at that location-date was tested, though we discuss robustness to other measures. As mentioned above, the location-time-specific measure of testing availability is preferred to using an indicator for whether an individual was tested (though we also report results with this metric) because of concerns that an applicant’s testing status is correlated with his or her perceived quality. We estimate these regressions at the location-time (month-by-year) level, the level of variation underlying our key explanatory variable, and weight by number of hires in a location-date.¹⁷ The outcome measure is the average outcome for workers hired to the same location at the same time.

All regressions include a complete set of location (δ_l) and month by year of hire (γ_t) fixed effects. They control for time-invariant differences across locations within our client firms, as well as for cohort and macroeconomic effects that may impact job duration. We also experiment with a number of additional control variables, described in our results section, below. In all specifications, standard errors are clustered at the location level to account for correlated observations within a location over time.

¹⁷This aggregation affords substantial savings on computation time, and, will produce identical results to those from a worker-level regression, given the regression weights.

Our primary outcome measure, Outcome_{lt} , is the log of the length of completed job spells, averaged across workers hired to firm-location l , at time t . We focus on this, and other outcomes related to the length of job spells, for several reasons. The length of a job spell is a measure that both theory and the firms in our study agree is important. Canonical models of job search (e.g., Jovanovic 1979), predict a positive correlation between match quality and job duration. Moreover, as discussed in Section 2, our client firms employ low-skill service sector workers and face high turnover and training costs: several weeks of paid training in a setting where the median worker stays only 99 days (see Figure 1.) Job duration is also a measure that has been used previously in the literature, for example by Autor and Scarborough (2008), who also focus on a low-skill service sector setting (retail). Finally, job duration is available for all workers in our sample.

3.2 Results

Table 2 reports regression results for the log duration of completed job spells. We later report results for several duration-related outcomes that do not restrict the sample to completed spells. Of the 270,086 hired workers that we observe in our sample, 75%, or 202,728 workers have completed spells (4,401 location-month cohorts), with an average spell lasting 203 days and a median spell of 99 days. The key explanatory variable is whether or not the median hire at this location-date was tested.

In the baseline specification (Panel 1, Column 1 of Table 2) we find that employees hired with the assistance of job testing stay, on average, 0.272 log points, or 31% longer, significant at the 5% level.

Panel 1 Column 2 introduces client firm-by-year fixed effects to control for the implementation of any new strategies and HR policies that firms may have adopted along with testing.¹⁸ In this specification, we compare locations in the same firm in the same year, some of which receive job testing sooner than others. The identifying assumption is that, within a firm, locations that receive testing sooner vs. later were on parallel trends before testing.

¹⁸Our data firm has indicated that it was not aware of other client-specific policy changes, though they acknowledge they would not have had full visibility into whether such changes may have occurred.

Here our estimated coefficient falls by roughly a third in magnitude, and we lose statistical significance.

To account for the possibility that the timing of the introduction of testing is related to trends at the location level, for example, that testing was introduced first to the locations that were on an upward (or downward) trajectory, Column 3 introduces location-specific time trends. These trends also account for broad trends that may impact worker retention, for instance, smooth changes in local labor market conditions. Adding these controls reduces the magnitude of our estimate but also greatly reduces the standard errors. We thus estimate an increased completed job duration of 0.137 log points or 15%, significant at the 5%-level.

Finally, in Column 4, we add controls for the composition of the applicant pool at a location after testing is implemented: fixed effects for the number of green, yellow, and red applicants. Because these variables are defined only after testing, these controls should be thought of as interactions between composition and the post-testing indicator, and are set to zero pre-testing. With these controls, the coefficient α_1 on Testing_{it} is the impact of the introduction of testing, for locations that end up receiving similarly qualified applicants. However, these variables also absorb any impact of testing on the quality of applicants that a location receives. For instance, the introduction of testing may have a screening effect: as candidates gradually learn about testing, the least qualified may be deterred from applying. Our point estimate remains unchanged with the inclusion of this set of controls, but the standard errors do increase substantially. This suggests that match quality improves because testing aids managers in identifying productive workers, rather than by exclusively altering the quality of the applicant pool. Overall, the range of estimates in Table 2 are similar to previous estimates found in Autor and Scarborough (2008).

Panel 2 of Table 2 examines robustness to defining testing at the individual level. For these specifications we regress an individual's job duration (conditional on completion) on whether or not the individual was tested. Because these specifications are at the individual level, our sample size increases from 4,401 location-months to 202,728 individual hiring events. Using these same controls, we find numerically similar estimates. The one exception is Column 4, which is now significant and larger: a 26% increase. From now on, we continue with our preferred metric of testing adoption (whether the median worker was tested).

Figure 2 shows event studies where we estimate the treatment impact of testing by quarter, from 12 quarters before testing to 12 quarters after testing, using our baseline set of controls. The top left panel shows the event study using log length of completed tenure spells as the outcome measure. The figure shows that locations that will obtain testing within the next few months look very similar to those that will not (because they either have already received testing or will receive it later). After testing is introduced, however, we begin to see large differences. The treatment effect of testing appears to grow over time, suggesting either that HR managers and other participants might take some time to learn how to use the test effectively. This alleviates any concerns that any systematic differences across locations drive the timing of testing adoption.

We also explore a range of other duration-related outcomes to examine whether the impact of testing is concentrated at any point in the duration distribution. For each hired worker, we measure whether they stay at least three, six, or twelve months, for the set of workers who are not right-censored.¹⁹ We aggregate this variable to measure the proportion of hires in a location-cohort that meet each duration milestone. Regression results (analogous to those reported in Panel 1 of Table 2) are reported in Appendix Table A1, while event studies are shown in the remaining panels of Figure 2. For each of these measures, we again see that testing improves job durations, and we see no evidence of any pre-trends.

This section thus establishes that the adoption of testing improves outcomes of hired workers. We next ask whether firms should change their hiring practices given they now have access to an apparently valuable signal of applicant quality.

4 Model

We formalize a model in which a firm makes hiring decisions with the help of an HR manager. There are two sources of information about the quality of job candidates. First, interviews generate unverifiable information about a candidate’s quality that is privately observed by the HR manager. Second, the job testing provides verifiable information about

¹⁹That is, a worker will be included in this metric if his or her hire date was at least three, six, or twelve months, respectively, before the end of data collection.

quality that is observed by both the manager and the firm. Managers then make hiring recommendations with the aid of both sources of information.

In this setting, job testing can improve hiring in two ways. First, it can help managers make more informed choices by providing an additional signal of worker quality. Second, because test information is verifiable, it enables the firm to constrain biased managers. Granting managers discretion enables the firm to take advantage of both interview and test signals, but may also leave it vulnerable to managerial biases. Limiting discretion and relying on the test removes scope for bias, but at the cost of ignoring private information. The following model formalizes this tradeoff and outlines an empirical test of whether firms can improve worker quality by eliminating discretion.

4.1 Setup

A mass one of applicants apply for job openings within a firm. The firm’s payoff of hiring worker i is given by the worker’s match quality, a_i . We assume that a_i is drawn from a distribution which depends on a worker’s type, $t_i \in \{G, Y\}$; a share of workers p_G are type G , a share $1 - p_G$ are type Y , and $a|t \sim N(\mu_t, \sigma_a^2)$ with $\mu_G > \mu_Y$ and $\sigma_a^2 \in (0, \infty)$. This match quality distribution enables us to naturally incorporate the discrete test score into the hiring environment. We do so by assuming that the test publicly reveals t .²⁰

The firm’s objective is to hire a proportion, W , of workers that maximizes expected match quality, $E[a|Hire]$.²¹ For simplicity, we also assume $W < p_G$.²²

To hire workers, the firm must employ HR managers whose interests are imperfectly aligned with that of the firm. In particular, a manager’s payoff for hiring worker i is given

²⁰The values of G and Y in the model correspond to test scores green and yellow, respectively, in our data. We assume binary outcomes for simplicity, even though in our data the signal can take three possible values. This is without loss of generality for the mechanics of the model.

²¹In theory, firms should hire all workers whose expected match quality is greater than their cost (wage). In practice, we find that having access to job testing information does not impact the number of workers that a firm hires. One explanation for this is that a threshold rule such as $E[a] > \bar{a}$ is not contractable because a_i is unobservable. Nonetheless, a firm with rational expectations will know the typical share W of applicants that are worth hiring, and W itself is contractable. Assuming a fixed hiring share is also consistent with the previous literature, for example, Autor and Scarborough (2008).

²²This implies that a manager could always fill a hired cohort with type G applicants. In our data, 0.43 of applicants are green and 0.6 of the green or yellow applicants are green, while the hire rate is 19%, so this will be true for the typical pool.

by:

$$U_i = (1 - k)a_i + kb_i.$$

In addition to valuing match quality, managers also receive an idiosyncratic payoff b_i , which they value with a weight k that is assumed to fall between 0 and 1. We assume that $a \perp b$.

The additional quality, b , can be thought of in two ways. First, it may capture idiosyncratic preferences of the manager for workers in certain demographic groups or with similar backgrounds (same alma mater, for example). Second, b can represent manager mistakes that drive them to prefer the wrong candidates.²³

The manager privately observes information about a_i and b_i . First, for simplicity, we assume that b_i is perfectly observed by the HR manager, and is distributed in the population by $N(0, \sigma_b^2)$ with $\sigma_b^2 \in (0, \infty)$. Second, the manager observes a noisy signal of match quality, s_i :

$$s_i = a_i + \epsilon_i$$

where $\epsilon_i \sim N(0, \sigma_\epsilon^2)$ is independent of a_i , t_i , and b_i . The parameter $\sigma_\epsilon^2 \in \mathbb{R}_+ \cup \{\infty\}$ measures the level of the manager's information. A manager with perfect information on a_i has $\sigma_\epsilon^2 = 0$, while a manager with no private information has $\sigma_\epsilon^2 = \infty$.

The parameter k measures the manager's bias, i.e., the degree to which the manager's incentives are misaligned with those of the firm or the degree to which the manager is mistaken. An unbiased manager has $k = 0$, while a manager who makes decisions entirely based on bias or the wrong characteristics corresponds to $k = 1$.

Let M denote the set of managers in a firm. For a given manager, $m \in M$, his or her type is defined by the pair $(k, 1/\sigma_\epsilon^2)$, corresponding to the bias and precision of private information, respectively. These have implied subscripts, m , which we suppress for ease of notation. We assume firms do not observe manager type, nor do they observe s_i or b_i .

Managers form a posterior expectation of worker quality given both their private signal and the test signal. They then maximize their own utility by hiring a worker if and only if the

²³For example, a manager may genuinely have the same preferences as the firm but draw incorrect inferences from his or her interview. Indeed, work in psychology (e.g., Dana et al., 2013) shows that interviewers are often overconfident about their ability to read candidates. Such mistakes fit our assumed form for manager utility because we can always separate the posterior belief over worker ability into a component related to true ability, and an orthogonal component resulting from their error.

expected value of U_i conditional on s_i , b_i , and t_i is at least some threshold. Managers thus wield “discretion” because they choose how to weigh the various signals about an applicant when making hiring decisions. We denote the quality of hires for a given manager under this policy as $E[a|Hire]$ (where an m subscript is implied).

4.2 Model Predictions

Our model focuses on the question of how much firms should rely on their managers, versus relying on hard test information. Firms can follow the set up described above, allowing their managers to weigh both signals and make ultimate hiring decisions (we call this the “Discretion” regime). Alternatively, firms may eliminate discretion and rely solely on test recommendations (“No Discretion”). In this section we generate a diagnostic for when one policy will dominate the other.

Neither retaining nor fully eliminating discretion need be the optimal policy response after the introduction of testing. Firms may, for example, consider hybrid policies such as requiring managers to hire lexicographically by the test score before choosing his or her preferred candidates, and these may generate more benefits. Rather than solving for the optimal hiring policy, we focus on the extreme of eliminating discretion entirely. This is because we can provide a tractable test for whether this counterfactual policy would make our client firms better off, relative to their current practice.²⁴ All proofs are in the Appendix.

Proposition 4.1 *The following results formalize conditions under which the firm will prefer Discretion or No Discretion.*

1. *For any given precision of private information, $1/\sigma_\epsilon^2 > 0$, there exists a $k' \in (0, 1)$ such that if $k < k'$ match quality is higher under Discretion than No Discretion and the opposite if $k > k'$.*
2. *For any given bias, $k > 0$, there exists $\underline{\rho}$ such that when $1/\sigma_\epsilon^2 < \underline{\rho}$, i.e., when precision of private information is low, match quality is higher under No Discretion than Discretion.*

²⁴We also abstract away from other policies the firm could adopt, for example, directly incentivizing managers based on the productivity of their hires or fully replacing managers with the test.

3. For any value of information $\bar{\rho} \in (0, \infty)$, there exists a bias, $k'' \in (0, 1)$, such that if $k < k''$ and $1/\sigma_\epsilon^2 > \bar{\rho}$, i.e., high precision of private information, match quality is higher under Discretion than No Discretion.

Proposition 4.1 illustrates the fundamental tradeoff firms face when allocating authority: managers have private information, but they are also biased. In general, larger bias pushes the firm to prefer No Discretion, while better information pushes it towards Discretion. Specifically, the first finding states that when bias, k is low, firms prefer to grant discretion, and when bias is high, firms prefer No Discretion. Part 2 states that when the precision of a manager’s private information becomes sufficiently small, firms cannot benefit from granting discretion, even if the manager has a low level of bias. Uninformed managers would at best follow test recommendations and, at worst deviate because they are mistaken or biased. Finally, part 3 states that for any fixed information precision threshold, there exists an accompanying bias threshold such that if managerial information is greater and bias is smaller, firms prefer to grant discretion. Put simply, firms benefit from Discretion when a manager has very precise information, but only if the manager is not too biased.

To understand whether No Discretion improves upon Discretion, employers would ideally like to directly observe a manager’s type (bias and information). In practice, this is not possible. Instead, it is easier to observe 1) the choice set of applicants available to managers when they made hiring decisions and 2) the performance outcomes of workers hired from those applicant pools. These are also two pieces of information that we observe in our data.

Specifically, we observe cases in which managers exercise discretion to explicitly contradict test recommendations. We define a hired worker as an “exception” if the worker would not have been hired under No Discretion (i.e., based on the test recommendation alone): any time a Y worker is hired when a G worker is available but not hired.

Denote the probability of an exception for a given manager, $m \in M$, as R_m . Given the assumptions made above, $R_m = E_m[Pr(Hire|Y)]$. That is, the probability of an exception is simply the probability that a Y type is hired, because this is implicitly also equal to the probability that a Y is hired *over* a G .

Proposition 4.2 *Across the set of managers M , the exception rate, R_m , is increasing in both managerial bias, k , and the precision of the manager’s private information, $1/\sigma_\epsilon^2$.*

Intuitively, managers with better information make more exceptions because they then place less weight on the test relative to their own signal of a . More biased managers also make more exceptions because they place more weight on maximizing other qualities, b . Thus, increases in exceptions can be driven by both more information and more bias.

It is therefore difficult to discern whether granting discretion is beneficial to firms simply by examining how often managers make exceptions. Instead, Propositions 4.1 and 4.2 suggest that it is instructive to examine the relationship between how often managers make exceptions and the subsequent match quality of their workers. Specifically, while exceptions (R_m) are increasing in both managerial bias and the value of the manager’s private information, match quality ($E[a|Hire]$) is decreasing in bias. If across managers, $E[a|Hire]$ is negatively correlated with R_m , then it is likely that exceptions are being driven primarily by managerial bias (because bias increases the probability of an exception and decreases the match quality of hires). In this case, eliminating discretion can improve outcomes. If the opposite is true, then exceptions are primarily driven by private information and discretion is valuable. The following proposition formalizes this intuition.

Proposition 4.3 *If the quality of hired workers is decreasing in the exception rate, $\frac{\partial E[a|Hire]}{\partial R_m} < 0$ across M , then firms can improve outcomes by eliminating discretion. If quality is increasing in the exception rate then discretion is better than no discretion.*

The intuition behind the proof is as follows. Consider two managers, one who never makes exceptions, and one who does. If a manager never makes exceptions, it must be that he or she has no additional information and no bias. As such, the match quality of this manager’s hires is equivalent to match quality of workers that would be hired if the firm eliminated discretion by relying only on test information. If increasing the probability of exceptions increases the match quality of hires, then granting discretion improves outcomes relative to no discretion. If match quality declines in the probability that managers make exceptions, then firms can improve outcomes by moving to a regime with no exceptions—that is, by eliminating discretion and using only the test.

5 Managerial Discretion

The model motivates the following empirical question: Is worker tenure increasing or decreasing in the probability of an exception? If decreasing, then No Discretion improves worker outcomes relative to Discretion, and if increasing then Discretion improves upon No Discretion.

In order to implement this test, we must address the empirical challenge that “exceptions” in our data are driven not only by managerial type (bias and information) as in the model, but also by other factors. For example, quality and size of the applicant pools may vary systematically with manager or location quality. We discuss how we apply our theory to the data in the next two subsections. We first define an exception rate to normalize variation across pools that mechanically makes exceptions more or less likely. We then discuss empirical specifications designed to limit remaining concerns.

5.1 Defining Exceptions

Our data provides us with the test scores of applicants post-testing. We use this information to define an “applicant pool” as a group of applicants being considered by the same manager for a job at the same location in the same month.²⁵

We can then measure how often managers overrule the recommendation of the test by either 1) hiring a yellow when a green had applied and is not hired, or 2) hiring a red when a yellow or green had applied and is not hired. We define the exception rate, for a manager m at a location l in a month t , as follows.

$$\text{Exception Rate}_{mlt} = \frac{N_y^h * N_g^{nh} + N_r^h * (N_g^{nh} + N_y^{nh})}{\text{Maximum \# of Exceptions}} \quad (2)$$

N_{color}^h and N_{color}^{nh} are the number of hired and not hire applicants, respectively. These variables are defined at the pool level (m, l, t) though subscripts have been suppressed for notational ease.

²⁵An applicant is under consideration if he or she applied in the last 4 months and had not yet been hired. Over 90% of workers are hired within 4 months of the date they first submitted an application.

The numerator of Exception Rate_{mlt} counts the number of exceptions (or “order violations”) a manager makes when hiring, i.e., the number of times a yellow is hired for each green that goes unhired plus the number of times a red is hired for each yellow and green that goes unhired.

The number of exceptions in a pool depends on both the manager’s choices and on factors related to the applicant pool, such as size and color composition. For example, if a pool has only green applicants, it is impossible to make an exception. Similarly, if the manager hires all available applicants, then there can also be no exceptions. These variations were implicitly held constant in our model, but need to be accounted for in the empirics.

To isolate the portion of variation in exceptions that are driven by managerial decisions, we normalize the number of order violations by the maximum number of violations that could occur, given the applicant pool that the recruiter faces and the number of hires. Importantly, although propositions in Section 4 are derived for the probability of an exception, their proofs hold equally for this definition of an exception rate.²⁶

From Table 1, we have 4,209 applicant pools in our data consisting of, on average 268 applicants.²⁷ On average, 19% of workers in a given pool are hired. Roughly 40% of all applicants in a given pool receive a “green”, while “yellow” and “red” candidates make up roughly 30%, each. The test score is predictive of whether or not an applicant is hired. In the average pool, greens and yellows are hired at a rate of roughly 20%, while only 9% of reds are hired. Still, managers very frequently make exceptions to test recommendations: the exception rate in the average pool (the average applicant is in a pool where an exception rate of) 24% of the maximal number of possible exceptions.

Furthermore, we see substantial variation in the extent to which managers actually follow test recommendations when making hiring decisions.²⁸ Figure 3 shows histograms of the exception rate, at the application pool level, as well as aggregated to the manager and

²⁶Results reported below are qualitatively robust to a variety of different assumptions on functional form for the exception rate.

²⁷This excludes months in which no hires were made.

²⁸According to the survey of client firms done by the data firm (on our behalf) (mentioned in footnote 12), client firms often told their managers that job test recommendations should be used in making hiring decisions but gave managers discretion over how to use the test (though some firms strongly discouraged managers from hiring red candidates).

location levels. The top panels show unweighted distributions, while the bottom panels show distributions weighted by the number of applicants.

In all figures, the median exception rate is about 20% of the maximal number of possible exceptions. At the pool level, the standard deviation is also about 20 percentage points; at the manager and location levels, it is about 11 percentage points. This means that managers very frequently make exceptions and that some managers and locations consistently make more exceptions than others.

5.2 Empirical Specifications

Proposition 4.3 examines the correlation between the exception rate and the realized match quality of hires in the post-testing period:

$$\text{Duration}_{mlt} = a_0 + a_1 \text{Exception Rate}_{mlt} + X_{mlt} \gamma + \delta_l + \delta_t + \epsilon_{mlt} \quad (3)$$

The coefficient of interest is a_1 . A negative coefficient, $a_1 < 0$, indicates that the match quality of hires is decreasing in the exception rate, meaning that firms can improve the match quality of hires by eliminating discretion and relying solely on job test information.

In addition to normalizing exception rates to account for differences in applicant pool composition, we estimate multiple version of Equation (3) that include location and time fixed effects, client-year fixed effects, location-specific linear time trends, and detailed controls for the quality and number of applicants in an application pool.

These controls are important because observed exception rates may be driven by factors other than a manager’s type (bias and information parameters). For example, some locations may be inherently less desirable than others, attracting both lower quality managers and lower quality applicants. In this case, lower quality managers may make more exceptions because they are biased. At the same time, lower quality workers may be more likely to quit or be fired. Both facts would be driven by unobserved location characteristics. Another potential concern is that undesirable locations may have difficulty hiring green workers, even conditional on them having applied. In our data, we cannot distinguish a green worker who

refuses a job offer from one who was never offered the job. As long as these characteristics are fixed or vary smoothly at the location-level, our controls absorb this variation.

A downside of including many fixed effects in Equation (3) is that it increases the extent to which our identifying variation is driven by pool-to-pool variation in the idiosyncratic quality of applicants. To see why this is problematic, imagine an applicant pool with a particularly weak draw of green candidates. In this case, we would expect a manager to make more exceptions, and, that the yellows and reds hired will perform better than the unhired greens from this particular pool. However, they may not perform better than the typical green hired by that manager. In this case, a manager could be using his or her discretion to improve match quality, but exceptions will still be correlated with poor outcomes. That is, when we identify off of pool-to-pool variation in exception rates, we may get the counterfactual wrong because exceptions are correlated with variation in unobserved quality within color.

To deal with the concern that Equation (3) relies too much on pool-to-pool variation in exception rates, we can aggregate exception rates to the manager- or location-level. Aggregating across multiple pools removes the portion of exception rates that are driven by idiosyncratic differences in the quality of workers in a given pool. The remaining variation—differences in the average exception rate across managers or locations—is more likely to represent exceptions made because of managerial type (bias and information). Doing so, however, reduces the amount of within-location variation left in our explanatory variable, making controlling for location fixed effects difficult or impossible.

To accommodate aggregate exception rates, we expand our data to include pre-testing worker observations. Specifically, we estimate whether the *impact* of testing, as described in Section 3, varies with exception rates:

$$\begin{aligned} \text{Duration}_{mlt} = & b_0 + b_1 \text{Testing}_{lt} \times \text{Exception Rate}_{mlt} + b_2 \text{Testing}_{lt} \\ & + X_{mlt} \gamma + \delta_l + \delta_t + \epsilon_{mlt} \end{aligned} \quad (4)$$

Equation (4) estimates how the impact of testing differs when managers make exceptions. The coefficient of interest is b_1 . Finding $b_1 < 0$ indicates that making more exceptions

decreases the improvement that locations see from the implementation of testing, relative to their pre-testing baseline. Because exception rates are not defined in the pre-testing period (there are no test scores in the pre-period), there is no main effect of exceptions in the pre-testing period, beyond that which is absorbed by the location fixed effects δ_l .

This specification allows us to use the pre-testing period to control for location-specific factors that might drive correlations between exception rates and outcomes. It also expands the sample on which we estimate location-specific time trends. This allows us to use exception rates that are aggregated to the manager- or location-level, avoiding small sample variation.²⁹ Aggregating exception rates to the location level also helps remove variation generated by any systematic assignment of managers to applicants within a location that might be correlated with exception rates and applicant quality.³⁰

To summarize, we test Proposition 4.3 with two approaches. First, we estimate the correlation between pool-level exception rates and quality of hires across applicant pools. Second, we estimate the differential impact of testing across pools with different exception rates of hires, where exception rates can be defined at the application pool, manager-, or location-level. In Section 5.4, we describe additional robustness checks.

5.3 Results

To gain a sense of the correlation between exception rates and outcome of hires, we first summarize the raw data by plotting both variables at the location level. Figure 4 shows a scatter plot of the location-level average exception rate on the x-axis and the location-level average tenure (log of completed duration) for workers hired post-testing on the y-axis. In the first panel, each location has the same weight; in the second, locations are weighted by the inverse variance of their pre-period mean, which takes into account their size and the confidence of our estimates. In both cases, we see a negative correlation between the extent

²⁹We define a time-invariant exception rate for managers (locations) that equals the average exception rate across all pools the manager (location) hired in (weighted by the number of applicants).

³⁰It also helps us rule out any measurement error generated by the matching of applicants to HR managers. This would be a problem if in some cases hiring decisions are made more collectively, or with scrutiny from multiple managers, and these cases were correlated with applicant quality.

to which managers exercise discretion by hiring exceptions, and the match quality of those hired.

The first two columns of Table 3 present the correlation between exception rates and worker tenure. We use a standardized exception rate with mean 0 and standard deviation 1 and in this panel exception rates are defined at the pool level (based on the set of applicants and hires a manager makes at a particular location in a given month).³¹

Column 1 contains our base specification and indicates that a one standard deviation increase in the exception rate of a pool is associated with a 5% reduction in completed tenure for that group, significant at the 5% level. The coefficient is still sizeable in Column 2 which contains our full set of controls, though it does fall slightly in magnitude and become insignificant.

The remaining columns of Table 3 examine how the impact of testing varies by the extent to which managers make exceptions. Our main explanatory variable is the interaction between the introduction of testing and a post-testing exception rate.

In Columns 3 and 4, we continue to use pool-level exception rates. The coefficient on the main effect of testing represents the impact of testing at the mean exception rate (since the exception rate has been standardized). Including the full set of controls (Column 4), we find that locations with the mean exception rate experience a 0.22 log point increase in duration as a result of the implementation of testing, but that this effect is offset by a quarter (0.05) for each standard deviation increase in the exception rate, significant at the 1% level.

In Columns 5-8, we aggregate exception rates to the manager- and location-level.³² Results are quite consistent, using these aggregations, and the differential effects are even larger in magnitude. Managers and locations that tend to exercise discretion benefit much less from the introduction of testing. A one standard deviation increase in the exception rate reduces the impact of testing by roughly half to two-thirds.³³

³¹We have experimented with different functional forms for the exception rate variable in both panels and obtained qualitatively similar results.

³²We have 555 managers who are observed in an average of 18 pools each (average taken over all managers, unweighted). We have 111 locations with on average 87 pools each (average taken over all locations, unweighted).

³³As we noted above, it is not possible to use these aggregated exceptions rates when examining the post-testing correlation between exceptions and outcomes (as in columns 1 and 2) because they leave little

To better illustrate the variation underlying these results, we plot location-specific treatment effects of testing on the location’s average exception rate. Figure 5 plots these for both an unweighted and a weighted sample.³⁴ The relationship is clearly negative, and does not look to be driven by any particular location.

We therefore find that the match quality of hires is lower for applicant pools, managers, and locations with higher exception rates. It is worth emphasizing that with the controls for the size and quality of the applicant pool, our identification comes from comparing outcomes of hires across managers who make different numbers of exceptions when facing similar applicant pools. Given this, differences in exception rates should be driven by a manager’s own weighting of his or her private preferences and private information. If managers were making these decisions optimally from the firm’s perspective, we should not expect to see (as we do in Table 3) that the workers they hire perform systematically worse. Based on Proposition 4.3, we can infer then that exceptions are largely driven by managerial bias, rather than private information, and these firms could improve outcomes of hires by limiting discretion.

5.4 Additional Robustness Checks

In this section we address several alternative explanations for our findings.

5.4.1 Quality of “Passed Over” Workers

There are several scenarios under which we might find a negative correlation between overall exceptions and outcomes without biased managers. For example, as mentioned above, managers may make more exceptions when green applicants in an applicant pool are idiosyncratically weak. If yellow workers in these pools are weaker than green workers in our sample on average, it will appear that more exceptions are correlated with worse outcomes even though managers are making individual exceptions to maximize match quality. Similarly, or no variation within locations to also identify location fixed effects, which, as we have argued, are quite important.

³⁴In this figure, we weight by the second panel weights by the inverse variance of the error associated with estimating that location’s treatment effect. This upweights locations with more precisely estimated treatment effects.

our results in Table 3 show that locations with more exceptions see fewer benefits from the introduction of testing. An alternative explanation for this finding is that high exception locations are ones in which managers have always had better information about applicants: these locations see fewer benefits from testing because they simply do not need the test.

In these and other similar scenarios, it should still be the case that individual exceptions are correct: a yellow hired as an exception should perform better than a green who is not hired. To examine this, we would like to be able to observe the counterfactual performance of all workers who are not hired. This would allow us to directly assess whether managers make exceptions to reduce false negatives, the possibility that a great worker is left unhired because he or she scored poorly on the test.

While we cannot observe the performance all non-hired greens, we can proxy for this comparison by exploiting the timing of hires. Specifically, we compare the performance of yellow workers hired as exceptions to green workers from the same applicant pool who are not hired that month, but who subsequently begin working in a later month. If it is the case that managers are making exceptions to increase the match quality of workers, then the exception yellows should have longer completed tenures than the “passed over” greens.

Table 4 shows that is not the case. The first panel compares individual durations by restricting our sample to workers who are either exception yellows, or greens who are initially passed over but then subsequently hired, and including an indicator for being in the latter group. Because these workers are hired at different times, all regressions control for hire year-month fixed effects to account for mechanical differences in duration. For the last column, which includes applicant pool fixed effects, the coefficient on being a passed over green compares this group to the specific yellow applicants who were hired before them. The second panel of Table 4 repeats this exercise, comparing red workers hired as exceptions (the omitted group), against passed over yellows and passed over greens.

In both panels, we find that workers hired as exceptions have shorter tenures. Column 3 is our preferred specification because it adds controls for applicant pool fixed effects. This means we compare the green (and yellow) applicants who were passed over one month but eventually hired, to the actual yellow (red) applicants hired first. We find that passed over greens stay about 8% longer than the yellows hired before them in the same pool (top panel

Column 3) and greens and yellows stay almost 19% and 12% longer, respectively, compared to the reds they were passed over for.

The results in Table 4 mean that it is unlikely that exceptions are driven by better information. When workers with better test scores are at first passed over and then later hired, they still outperform the workers chosen first.³⁵

Table 5 provides additional evidence that workers with longer gaps between application date and hire date (which we treat as temporarily passed over applicants) are not simply ones who were delayed because of better outside options. If this were the case, we would expect these workers to have better outcomes once they do begin work. In Table 5, we compare match quality for workers hired immediately (the omitted category), compared to those who waited one, two, or three months before starting, holding constant test score. Because these workers are hired at different times, all regressions again control for hire year-month fixed effects. Across all specifications, we find no significant differences between these groups. If anything we find for greens and yellows that were hired with longer delays have shorter job spells than immediate hires. We thus feel more comfortable interpreting the workers with longer delays as having been initially passed over.

Table 5 also provides insights about how much information managers have, beyond the job test. If managers have useful private information about workers, then we would expect them to be able to distinguish quality within test-color categories: greens hired first should be better than greens who are passed up. Table 5 shows that this does not appear to be the case. We estimate only small and insignificant differences in tenure, within color, across start dates. That is, within color, workers who appear to be a manager's first choice do not perform better than workers who appear to be a manager's last choice. This again suggests the value of managerial private information is small, relative to the test.

³⁵An alternative explanation is that the applicants with higher test scores were not initially passed up, but were instead initially unavailable because of better outside options. However, our data firm thinks it unlikely that applicants would delay starting their positions, both because this is a low-skill job where workers have few outside options, and because client firms want to fill training classes.

5.4.2 Extreme Outcomes

We have thus far assumed that the firm would like managers to maximize the average match quality of hired workers. Firms may instead instruct managers to take a chance on candidates with poor test scores to avoid missing out on an exceptional hire. This could explain why managers take chances on lower test score candidates to the detriment of average match quality; managers may be using discretion to minimize false negatives. Alternatively, firms may want managers to use discretion to minimize the chance of hiring a worker who leaves immediately.

The first piece of evidence that managers do not effectively maximize upside or minimize downside is the fact, already shown, that workers hired as perform worse than the very workers they were originally passed over. To address this more directly, Appendix Table A2 repeats our analysis focusing on performance in the tails, using the 90th and 10th percentiles of log completed durations for a cohort as the dependent variables. These results show that more exceptions imply worse performance even among top hires, suggesting managers who make many exceptions are also unsuccessful at finding star workers. We also show that more exceptions decrease the performance of the 10th percentile of completed durations as well.

5.4.3 Heterogeneity across Locations

Another possible concern is that the usefulness of the test varies across locations and that this drives the negative correlation between exception rates and worker outcomes. Our results on individual exceptions already suggest that this is not the case. However, we explore a couple of specific stories here.

As we have already noted, a location with very good private information pre-testing would have both a high exception rate and a low impact of testing. If exceptions were driven by information, rather than mistakes or bias, we should see that higher exception rate locations were more productive pre-testing. However, Figure 6 plots the relationship between a location's eventual exception rates and the match quality of its hires prior to the introduction of testing, and shows no systematic relationship between the two.³⁶

³⁶We maintain the same weighting as 5 so the two figures are comparable.

Alternatively, in very undesirable locations, green applicants might have better outside options and be more difficult to retain. In these locations, a manager attempting to avoid costly retraining may optimally decide to make exceptions in order to hire workers with lower outside options. Here, a negative correlation between exceptions and performance would not necessarily imply that firms could improve productivity by relying more on testing. However, we see no evidence that the “return” to test score varies across locations. For example, when we split locations by pre-testing worker durations (Appendix Table A3) or by exception rates post-testing (Appendix Table A4) we see no systematic differences in the correlation between test score and job duration of hired workers.

5.4.4 Productivity

Our results show that firms can improve the match quality of their workers, as measured by duration, by relying more on job test recommendations. Firms may not want to pursue this strategy, however, if their HR managers exercise discretion in order to improve worker quality on other metrics. For example, managers may optimally choose to hire workers who are more likely to turn over if their private signals indicate that those workers might be more productive while they are employed.

Our final set of results provides evidence that this is unlikely to be the case. Specifically, for a subset of 62,494 workers (one-quarter of all hires) in 6 client firms, we observe a direct measure of worker productivity: output per hour.³⁷ We are unable to reveal the exact nature of this measure but some examples may include: the number of data items entered per hour, the number of standardized tests graded per hour, and the number of phone calls completed per hour. In all of these examples, output per hour is an important measure of efficiency and worker productivity. Our particular measure has an average of roughly 8 with a standard deviation of roughly 5.

Table 6 repeats our main findings, using output per hour instead of job duration as the dependent variable. We focus on estimates only using our base specification (controlling

³⁷We have repeated our main analyses on the subsample of workers that have output per hour data and obtained similar results.

for date and location fixed effects) because the smaller sample and number of clients makes identifying the other controls difficult.³⁸

Column 1 examines the impact of the introduction of testing, which we find leads to a statistically insignificant increase of 0.7 transactions in an hour, or a roughly 8% increase. The standard errors are such that we can rule out virtually any negative impact of testing on productivity with 90% confidence.

Column 2 documents the post-testing correlation between pool-level exceptions and output per hour, and Columns 3-5 examine how the impact of testing varies by exception rates. In all cases, we find no evidence that managerial exceptions improve output per hour. Instead, we find noisy estimates indicating that worker quality appears to be lower on this dimension as well. For example, in Column 2, we find a tiny, insignificant positive coefficient describing the relationship between exceptions and output. Taking it seriously implies that a 1 standard deviation increase in exception rates is correlated with 0.07 more transactions, or a less than 1% increase. In, Columns 3-5, we continue to find an overall positive effect of testing on output; we find no evidence of a positive correlation between exception rates and the impact of testing. If anything, the results suggest that locations with more exceptions experience slightly smaller impacts of testing. These effects are insignificant.

Taken together, the results in Table 6 provide no evidence that exceptions are positively correlated with productivity. This refutes the hypothesis that, when making exceptions, managers optimally sacrifice job tenure in favor of workers who perform better on other quality dimensions.

6 Conclusion

We evaluate the introduction of a hiring test across a number of firms and locations for a low-skill service sector job. Exploiting variation in the timing of adoption across locations within firms, we show that testing increases the durations of hired workers by about 15%. We then document substantial variation in how managers use job test recommendations. Some managers tend to hire applicants with the best test scores while others make many

³⁸Results are, however, qualitatively similar with additional controls.

more exceptions. Across a range of specifications, we show that the exercise of discretion (hiring against the test recommendation) is associated with worse outcomes.

Our paper contributes a new methodology for evaluating the value of discretion in firms. Our test is intuitive, tractable, and requires only data that would readily be available for firms using workforce analytics. In our setting it provides the stark recommendation that firms would do better to remove discretion of the average HR manager and instead hire based solely on the test. In practice, our test can be used as evidence that the typical manager underweights the job test relative to what the firm would prefer. Based on such evidence, firms may want to explore a range of alternative options, for example, allowing managers some degree of discretion but limiting the frequency with which they can overrule the test, or, adopt other policies to influence manager behavior such as direct pay for performance or selective hiring and firing.

These findings highlight the role new technologies can play in reducing the impact of managerial mistakes or biases by making contractual solutions possible. As workforce analytics becomes an increasingly important part of human resource management, more work needs to be done to understand how such technologies interact with organizational structure and the allocation of decisions rights with the firm. This paper makes an important step towards understanding and quantifying these issues.

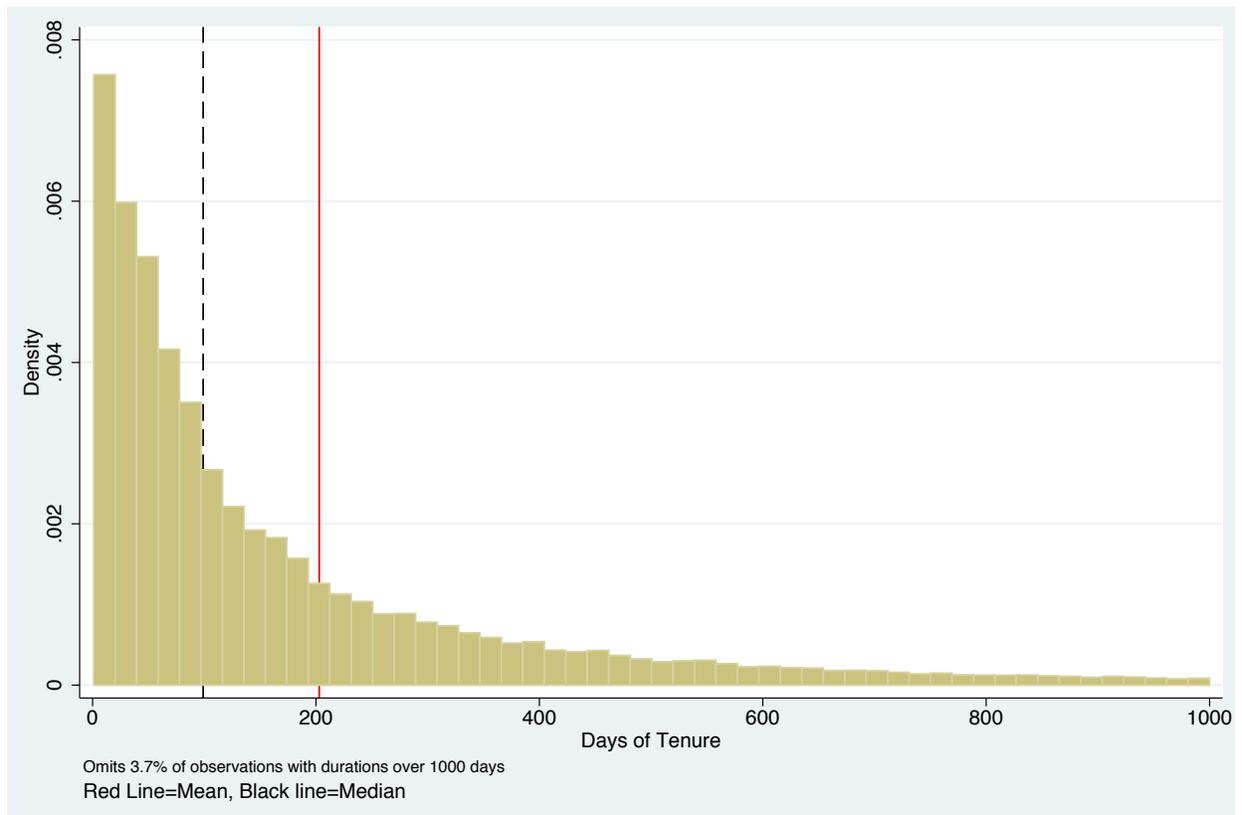
References

- [1] Aghion, Philippe and Jean Tirole (1997), "Formal and Real Authority in Organizations," *The Journal of Political Economy*, 105(1).
- [2] Altonji, Joseph and Charles Pierret (2001), "Employer Learning and Statistical Discrimination," *Quarterly Journal of Economics*, 113: pp. 79-119.
- [3] Alonso, Ricardo and Niko Matouschek (2008), "Optimal Delegation," *The Review of Economic Studies*, 75(1): pp 259-3.
- [4] Autor, David (2001), "Why Do Temporary Help Firms Provide Free General Skills Training?," *Quarterly Journal of Economics*, 116(4): pp. 1409-1448.
- [5] Autor, David and D. Scarborough (2008), "Does Job Testing Harm Minority Workers? Evidence from Retail Establishments," *Quarterly Journal of Economics*, 123(1): pp. 219-277.
- [6] Baker, George and Thomas Hubbard (2004), "Contractibility and Asset Ownership: On-Board Computers and Governance in U.S. Trucking," *Quarterly Journal of Economics*, 119(4): pp. 1443-1479.
- [7] Bolton, Patrick and Mathias Dewatripont (2010) "Authority in Organizations." in Robert Gibbons and John Roberts (eds.), *The Handbook of Organizational Economics*. Princeton, NJ: Princeton University Press.
- [8] Brown, Meta, Elizabeth Setren, and Giorgio Topa (2015), "Do Informal Referrals Lead to Better Matches? Evidence from a Firm's Employee Referral System," *Journal of Labor Economics*, forthcoming.
- [9] Burks, Stephen, Bo Cowgill, Mitchell Hoffman, and Michael Housman (2015), "The Value of Hiring through Employee Referrals," *Quarterly Journal of Economics*, 130(2): pp. 805-839.
- [10] Dana, Jason, Robyn Dawes, and Nathaniel Peterson (2013), "Belief in the Unstructured Interview: The Persistence of an Illusion", *Judgment and Decision Making*, 8(5), pp. 512-520.

- [11] Dessein, Wouter (2002) "Authority and Communication in Organizations," *Review of Economic Studies*. 69, pp. 811-838.
- [12] Farber, Henry. and Robert Gibbons (1996), "Learning and Wage Dynamics," *Quarterly Journal of Economics*, 111: pp. 1007-1047.
- [13] Horton, John (2013), "The Effects of Subsidizing Employer Search," mimeo New York University.
- [14] Jovanovic, Boyan (1979), "Job Matching and the Theory of Turnover," *The Journal of Political Economy*, 87(October), pp. 972-90.
- [15] Kahn, Lisa and Fabian Lange (2014), "Employer Learning, Productivity and the Earnings Distribution: Evidence from Performance Measures," *Review of Economic Studies*, 81(4) pp.1575-1613.
- [16] Kahneman, Daniel (2011). *Thinking Fast and Slow*. New York: Farrar, Strauss, and Giroux.
- [17] Kuncel, Nathan, David Klieger, Brian Connelly, and Deniz Ones (2013), "Mechanical Versus Clinical Data Combination in Selection and Admissions Decisions: A Meta-Analysis," *Journal of Applied Psychology*. Vol. 98, No. 6, 1060–1072.
- [18] Kuziemko, Ilyana (2013), "How Should Inmates Be Released from Prison? an Assessment of Parole Versus Fixed Sentence Regimes," *Quarterly Journal of Economics*. Vol. 128, No. 1, 371–424.
- [19] Lazear, Edward, Kathryn Shaw, and Christopher Stanton (2015), "The Value of Bosses," *Journal of Labor Economics*, forthcoming.
- [20] Li, Danielle. (2012), "Expertise and Bias in Evaluation: Evidence from the NIH" mimeo Harvard University.
- [21] Oyer, Paul and Scott Schaefer (2011), "Personnel Economics: Hiring and Incentives," in the *Handbook of Labor Economics*, 4B, eds. David Card and Orley Ashenfelter, pp. 1769-1823.

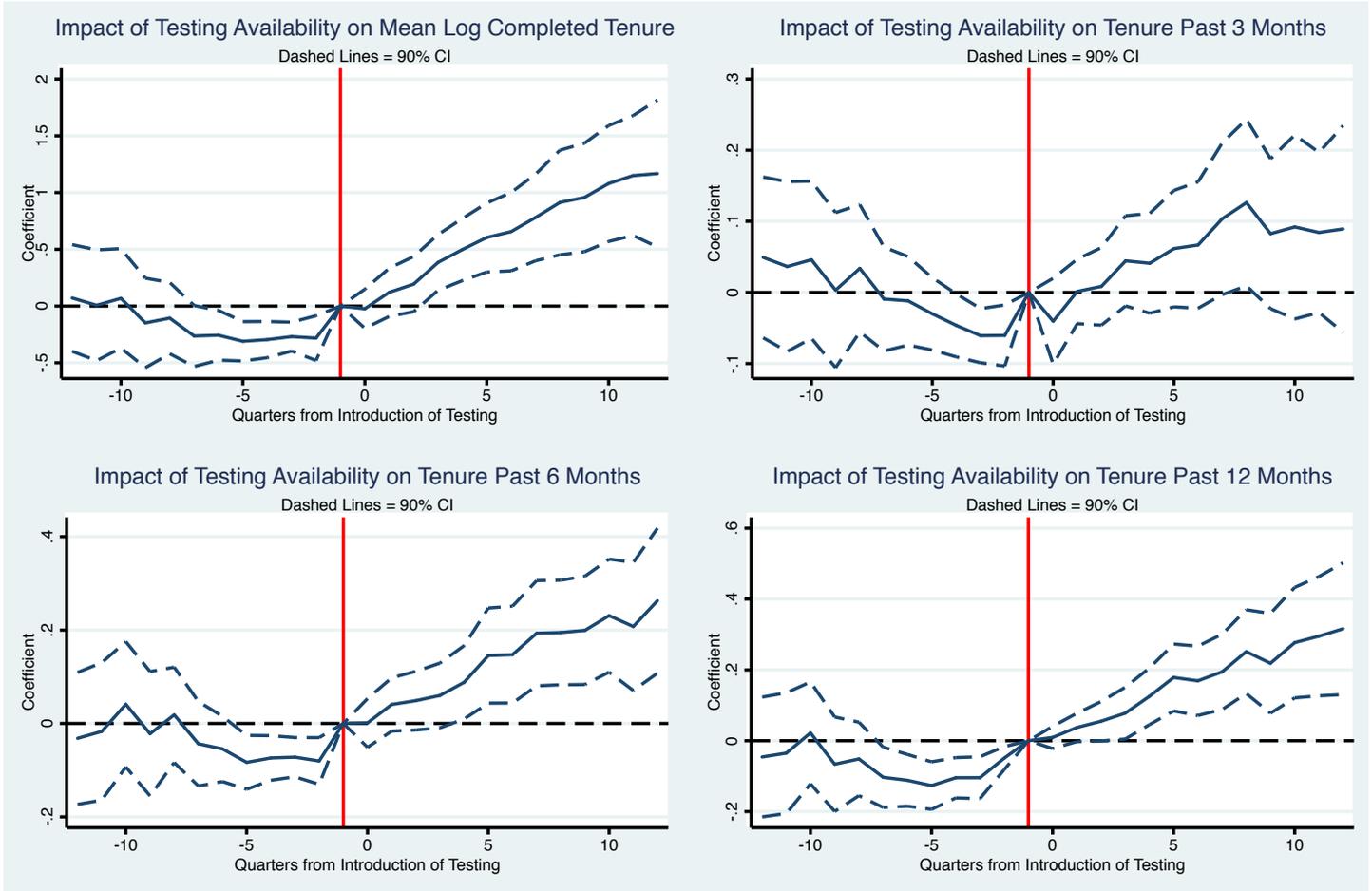
- [22] Pallais, Amanda and Emily Sands (2015), “Why the Referential Treatment? Evidence from Field Experiments on Referrals,” *The Journal of Political Economy*, forthcoming.
- [23] Paravisini, Daniel and Antoinette Schoar (2013) “The Incentive Effect of IT: Randomized Evidence from Credit Committees” NBER Working Paper #19303.
- [24] Riviera, Lauren. (2014) “Hiring as Cultural Matching: The Case of Elite Professional Service Firms.” *American Sociological Review*. 77: 999-1022
- [25] Stanton, Christopher and Catherine Thomas (2014), “Landing The First Job: The Value of Intermediaries in Online Hiring,” mimeo London School of Economics.

FIGURE 1: DISTRIBUTION OF LENGTH OF COMPLETED JOB SPELLS



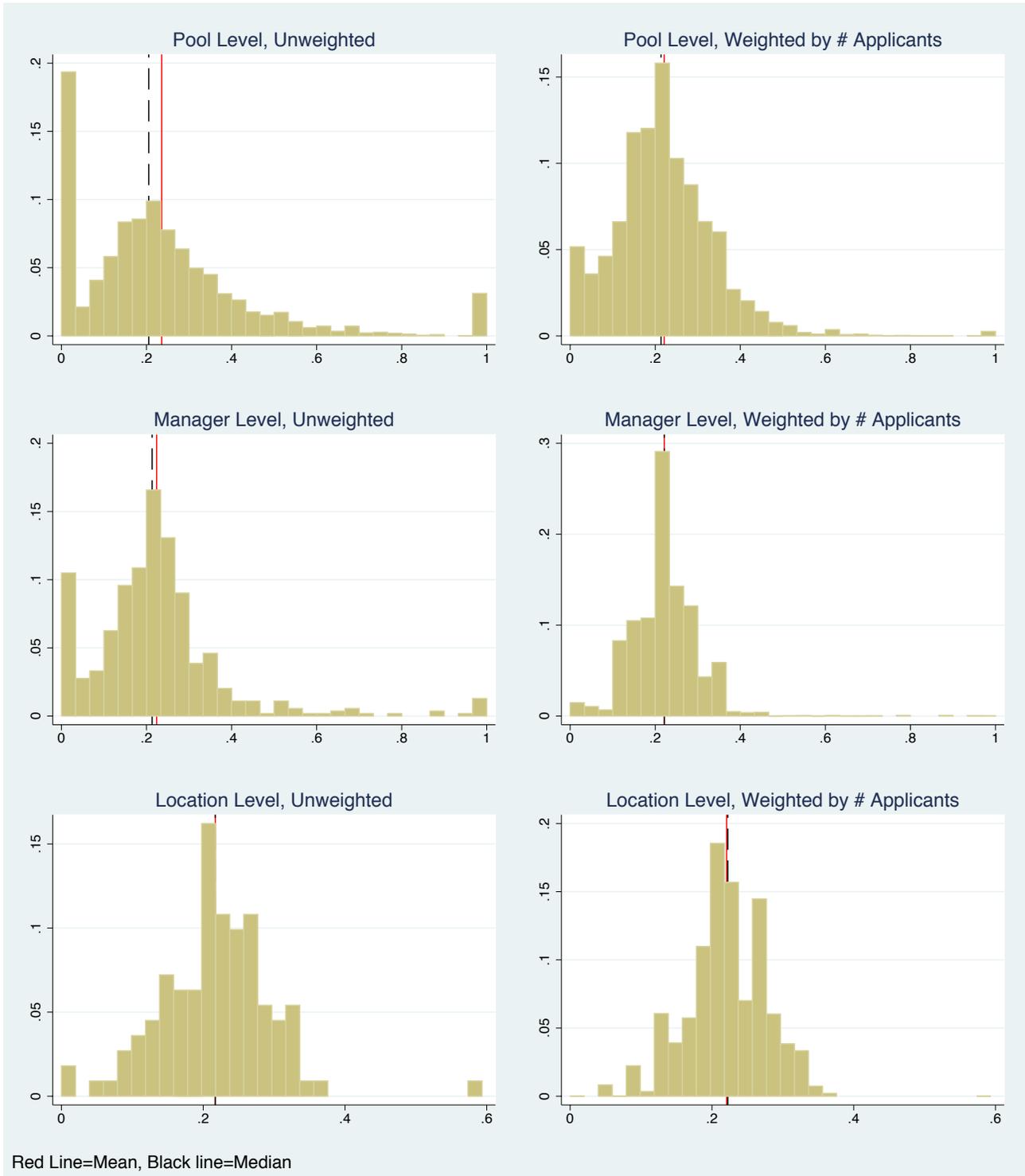
NOTES: Figure 1 plots the distribution of completed job spells at the individual level.

FIGURE 2: EVENT STUDY OF DURATION OUTCOMES



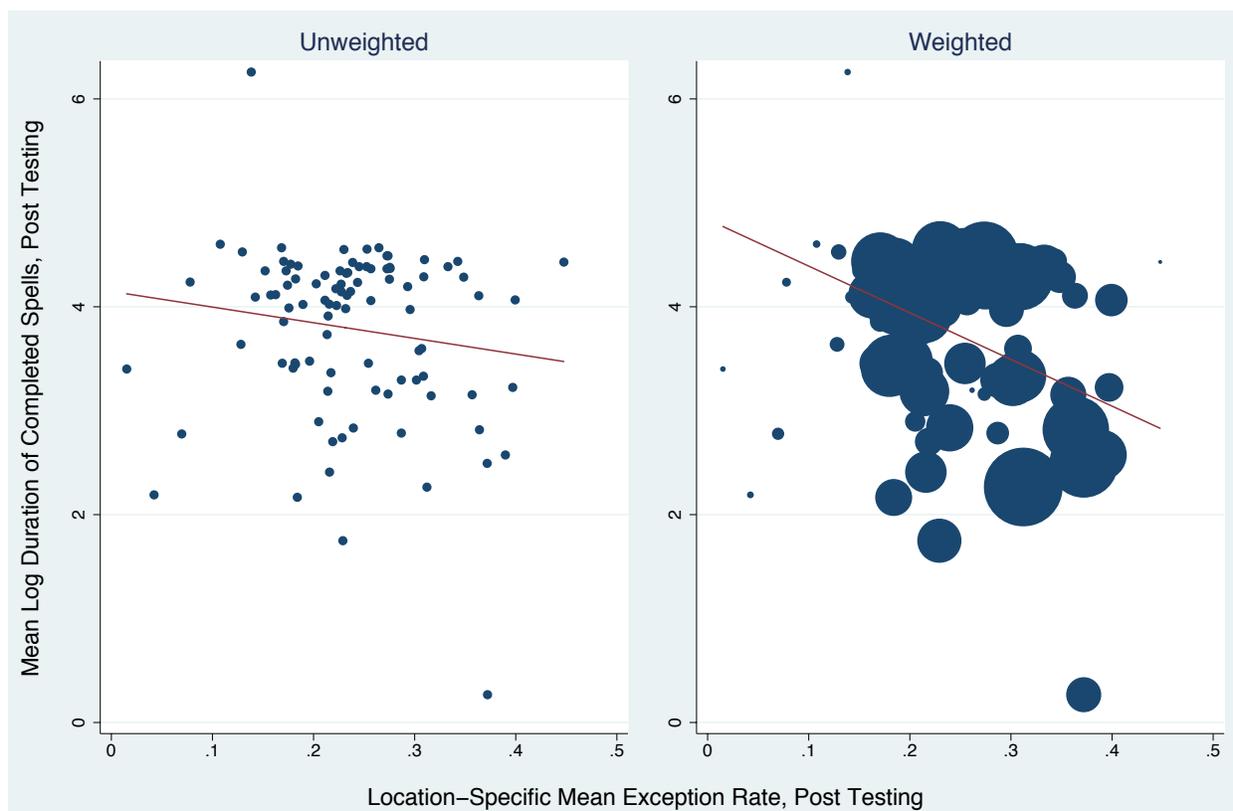
NOTES: These figures plot the average duration outcome for entry cohorts by time (in quarters) until or time after testing is adopted. The underlying estimating equation is given by $\text{Log}(\text{Duration})_{it} = \alpha_0 + I_{it}^{\text{time since testing}} \alpha_1 + \delta_l + \gamma_t + \epsilon_{it}$, where $I_{it}^{\text{time since testing}}$ is a vector of dummies indicating how many quarters until or after testing is adopted, with one quarter before as the omitted category. This regression includes the base set of controls – location (δ_l) and date (γ_t) fixed effects; it does not control for location-specific time trends.

FIGURE 3: DISTRIBUTIONS OF APPLICATION POOL EXCEPTION RATES



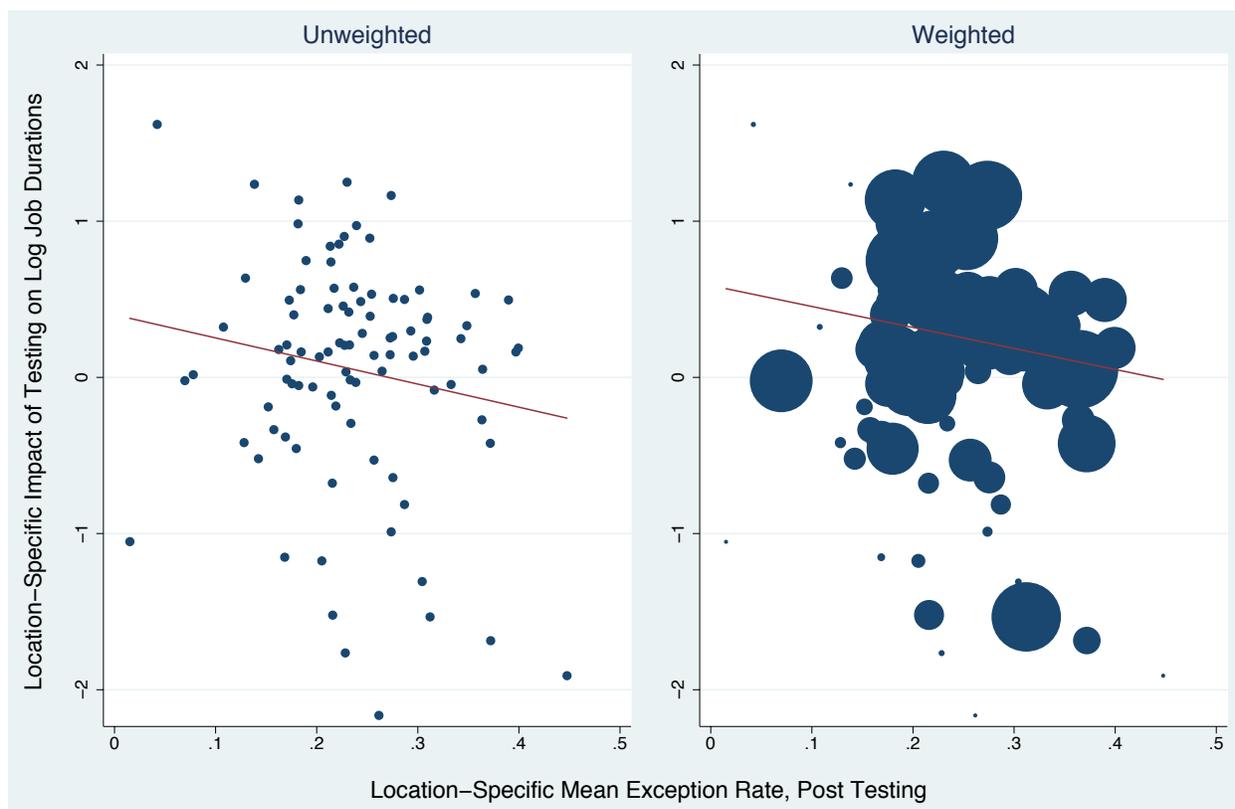
NOTES: These figures plot the distribution of the exception rate, as defined by Equation (2) in Section 5. The leftmost panel presents results at the applicant pool level (defined to be a manager–location–month). The middle panel aggregates these data to the manager level and the rightmost panel aggregates further to the location level. Exception rates are only defined for the post-testing sample.

FIGURE 4: LOCATION-LEVEL EXCEPTION RATES AND POST-TESTING JOB DURATIONS



NOTES: Each dot represents a location. The y-axis is the mean log completed tenure at a given location after the introduction of testing; the x-axis is the average exception rate across all applicant pools associated with a location. The first panel presents unweighed correlations; the second panel weights by the inverse variance of the error associated with estimating that location's treatment effect, to be consistent with Figure 5. The line shows the best linear fit of the scatter plot with corresponding weighting.

FIGURE 5: LOCATION-LEVEL EXCEPTION RATES AND THE IMPACT OF TESTING ON JOB DURATIONS



NOTES: Each dot represents a given location. The y-axis is the coefficient on the location-specific estimate of the introduction of testing on log of the length of completed job spells; the x-axis is the average exception rate across all application pools associated with a location. The first panel presents unweighted correlations; the second panel weights by the inverse variance of the error associated with estimating that location's treatment effect. The line shows the best linear fit of the scatter plot with corresponding weighting.

FIGURE 6: LOCATION-LEVEL EXCEPTION RATES AND PRE-TESTING JOB DURATIONS



NOTES: Each dot represents a given location. The y-axis reports mean of log duration of completed spells at a given location prior to the introduction of testing; the x-axis is the average exception rate across all application pools associated with a location. The first panel presents unweighed correlations; the second panel weights by the inverse variance of the error associated with estimating that location’s treatment effect, to remain consistent with Figure 5. The line shows the best linear fit of the scatter plot with corresponding weighting.

TABLE 1: SUMMARY STATISTICS

Sample Coverage					
	All	Pre-testing	Post-testing		
<i>Sample Coverage</i>					
# Locations	131	116	111		
# Hired Workers	270,086	176,390	93,696		
# Applicants			691,352		
# HR Managers			555		
# Pools			4,209		
# Applicants/Pool			268		
Worker Performance					
			<i>mean (st dev)</i>		
	Pre-testing	Post-testing	Green	Yellow	Red
Duration of Completed Spell (Days) (N=202,728)	247 (314)	116 (116)	122 (143)	110 (131)	93 (122)
# Customers Served/Hr (N=62,494)	8.32 (4.58)	8.41 (5.06)	8.37 (4.91)	8.30 (5.04)	9.09 (5.94)
Applicant Pool Characteristics					
	Post-testing	Green	Yellow	Red	
Share Applicants		0.43	0.29	0.28	
Share Hired		0.19	0.22	0.18	0.09
Exception Rate		0.24			

NOTES: The sample includes all non stock-sampled workers. Post-testing is defined at the location-month level as the first month in which 50% of hires had test scores, and all months thereafter. An applicant pool is defined at the manager-location-month level and includes all applicants that had applied within four months of the current month and not yet hired. Number of applicants reflects the total number in any pool.

TABLE 2: IMPACT OF JOB TESTING ON LENGTH OF COMPLETED JOB SPELLS

	(1)	(2)	(3)	(4)
Panel 1: Location-Cohort Mean Log Duration of Completed Spells				
<i>Post-Testing</i>	0.272** (0.113)	0.178 (0.113)	0.137** (0.0685)	0.142 (0.101)
N	4,401	4,401	4,401	4,401
Panel 2: Individual-Level Log Duration of Completed Spells				
<i>Individual Applicant is Tested</i>	0.195* (0.115)	0.139 (0.124)	0.141** (0.0637)	0.228** (0.0940)
N	202,728	202,728	202,728	202,728
Year-Month FEs	X	X	X	X
Location FEs	X	X	X	X
Client Firm X Year FEs		X	X	X
Location Time Trends			X	X
<u>Size and Composition of Applicant Pool</u>				X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: In Panel 1, an observation is a location-month. The dependent variable is average log duration, conditional on completion, for the cohort hired in that month. Post-testing is defined at the location-month level as the first month in which 50% of hires had test scores, and all months thereafter. Regressions are weighted by the number of hires in that location-month. Standard errors in parentheses are clustered at the location level. In Panel 2, observations are at the individual level. Testing is defined as whether or not an individual worker has a test score. Regressions are unweighted. Size and composition of the applicant pool controls include separate fixed effects for the number of red, yellow, and green applicants at that location-month.

TABLE 3: RELATIONSHIP BETWEEN JOB DURATIONS AND EXCEPTION RATES

	<i>Log Duration of Completed Spells</i>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Post-Testing Sample				Full Sample			
					Level of Aggregation for Exception Rate			
					<i>Pool</i>	<i>Manager</i>	<i>Location</i>	
<i>Post-Testing</i>	--	--	0.277** (0.112)	0.217** (0.0876)	0.281** (0.114)	0.243*** (0.0869)	0.306** (0.123)	0.252*** (0.0880)
<i>Exception Rate*Post-Testing</i>	-0.0491** (0.0223)	-0.0310 (0.0213)	-0.101** (0.0446)	-0.0477*** (0.0182)	-0.211** (0.0907)	-0.131*** (0.0335)	-0.290** (0.112)	-0.170** (0.0744)
N	3,839	3,839	6,869	6,869	6,942	6,942	6,956	6,956
Year-Month FEs	X	X	X	X	X	X	X	X
Location FEs	X	X	X	X	X	X	X	X
Client Firm X Year FEs		X		X		X		X
Location Time Trends		X		X		X		X
Size and Composition of Applicant Pool		X		X		X		X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: Each observation is a manager-location-month, observations are weighted by number of hires, and standard errors are clustered by location. Post-testing is defined at the location-month level as the first month in which 50% of hires had test scores, and all months thereafter. Columns 1 and 2 restrict to the post-testing sample only. The pool-level exception rate is the number of times a yellow is hired above a green or a red is hired above a yellow or green in a given applicant pool, divided by the maximum number of such violations. It is standardized to be mean zero and standard deviation one. Columns 3-4 and 5-6 aggregate the pool-level exception rate to the manager- and location-level, respectively. Exception rates are only defined post testing and are set to 0 pre testing. Size and composition of the applicant pool controls include separate fixed effects for the number of red, yellow, and green applicants at that location-manager-month. See text for additional details.

TABLE 4: MATCH QUALITY OF EXCEPTIONS VS. PASSED OVER APPLICANTS

	<i>Log Duration of Completed Spells</i>		
	(1)	(2)	(3)
Panel 1: Quality of Yellow Exceptions vs. Passed over Greens			
<i>Passed Over Greens</i>	0.0436*** (0.0140)	0.0470*** (0.0140)	0.0743*** (0.0246)
N	59,462	59,462	59,462
Panel 2: Quality of Red Exceptions vs. Passed over Greens and Yellows			
<i>Passed Over Greens</i>	0.131*** (0.0267)	0.145*** (0.0268)	0.173*** (0.0332)
<i>Passed Over Yellows</i>	0.0732*** (0.0265)	0.0948*** (0.0267)	0.114*** (0.0326)
N	44,456	44,456	44,456
Hire Month FEs	X	X	X
Location FEs	X	X	X
Client Firm X Year FEs		X	X
Application Pool FEs			X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: Regressions are at the individual level on the post-testing sample. Standard errors are clustered by location. Panel 1 includes yellow exceptions (who were hired when a green applicant was available but not hired in that month) and passed over green applicants who were later hired. The omitted category are yellow exceptions. The second panel includes red exceptions (who were hired when a green or yellow applicant was available but not hired in that month) and passed over greens and yellows only. Red exceptions are the omitted category. Application pool fixed effects are defined for a given location-manager-month.

TABLE 5: JOB DURATION OF WORKERS, BY LENGTH OF TIME IN APPLICANT POOL

	<i>Log Duration of Completed Spells</i>		
	(1)	(2)	(3)
Green Workers		Green Workers	
<i>Waited 1 Month</i>	-0.00908 (0.0262)	-0.00312 (0.0247)	0.00627 (0.0204)
<i>Waited 2 Months</i>	-0.0822 (0.0630)	-0.0713 (0.0596)	-0.0446 (0.0385)
<i>Waited 3 Months</i>	-0.000460 (0.0652)	0.000469 (0.0638)	-0.0402 (0.0639)
N	41,020	41,020	41,020
		Yellow Workers	
<i>Waited 1 Month</i>	-0.00412 (0.0199)	0.00197 (0.0184)	0.00773 (0.0243)
<i>Waited 2 Months</i>	-0.0100 (0.0448)	-0.00381 (0.0440)	-0.0474 (0.0509)
<i>Waited 3 Months</i>	0.103 (0.0767)	0.110 (0.0764)	0.114 (0.0979)
N	22,077	22,077	22,077
		Red Workers	
<i>Waited 1 Month</i>	0.0712 (0.0520)	0.0741 (0.0519)	0.0531 (0.0617)
<i>Waited 2 Months</i>	0.0501 (0.0944)	0.0509 (0.0941)	0.0769 (0.145)
<i>Waited 3 Months</i>	0.103 (0.121)	0.117 (0.125)	0.149 (0.168)
N	4,919	4,919	4,919
Year-Month FEs	X	X	X
Location FEs	X	X	X
Client Firm X Year FEs		X	X
Application Pool FEs			X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: Each observation is an individual hired worker for the post-testing sample. The first panel restricts to green workers only, with green workers who are hired immediately serving as the omitted group. The other panels are defined analogously for yellow and red. Standard errors are clustered by location. Application pool fixed effects are defined for a given location-manager-month.

TABLE 6: TESTING, EXCEPTION RATES, AND OUTPUT PER HOUR

	Impact of testing	Exceptions and outcomes, post testing	Impact of testing, by exception rate		
			Level of aggregation for exception rate		
			<i>Pool</i>	<i>Manager</i>	<i>Location</i>
	(1)	(2)	(3)	(4)	(5)
Dependent Variable: Output per hour					
<i>Post-Testing</i>	0.669 (0.428)		0.696 (0.432)	0.713 (0.438)	0.733* (0.431)
<i>Exception Rate*Post-Testing</i>		0.073 (0.085)	0.032 (0.084)	-0.042 (0.120)	-0.107 (0.247)
N	2,699	1,527	2,699	2,699	2,699
Year-Month FEs	X	X	X	X	X
Location FEs	X	X	X	X	X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: This table replicates the baseline specifications in Tables 2 and 3, using the number of transactions per hour (mean 8.38, std. dev. 3.21) as the dependent variable. All regressions are weighted by number of hires and standard errors are clustered by location. Column 1 studies the impact of the introduction of testing. Column 2 examines the correlation between output per hour and pool-level exception rates in the post-testing sample only. Columns 3 through 5 study the differential impact of the introduction of testing, with the exception rate aggregated at the pool, manager, and location levels, respectively.

FOR ONLINE PUBLICATION ONLY

A Data Appendix

Firms in the Data. The data were assembled for us by the data firm from records of the individual client firms. The client firms in our sample employ workers in the same occupation, but the specific duties vary from firm to firm. For example, at one firm, workers engage in a relatively high-skilled version of the job we study. Our primary estimates of the relationship between exceptions and duration are similar when individual firms are excluded one by one.³⁹ Among the non-tested workers, the data include a small share of workers outside of entry-level positions, but we checked our results our robust when we repeat our analyses controlling for employee position type.

Pre-testing Data. One downside of the pre-testing data is that the completeness of the data is idiosyncratic across the different client firms. For example, at some client firms, we may have a stock-sampling problem: the pre-testing data contain information about both new workers and incumbent workers. This may generate a survivor bias for incumbent workers, relative to new workers. For example, consider a firm that provided pre-testing data on new hires going back to Jan. 2010. For this firm, we would observe the full set of workers hired at each date after Jan. 2010, but for those hired before, we would only observe the subset who survived to Jan. 2010. We do not explicitly observe the date at which the firm began providing information on all new hires; instead, we conservatively proxy this date using the date of first recorded termination. We label all workers hired before this date as “stock sampled” because we cannot be sure that we observe their full entry cohort. We drop these workers from our primary sample, but have experimented with including them along with flexible controls for being stock sampled in our regressions.⁴⁰

Productivity. In addition to hire and termination dates, with which we calculate our primary outcome measure, some client firms provide data on output per hour. This is available for 31% of hired workers in our sample, and is mentioned by our data firm in its advertising, alongside duration. We trim instances where average transaction time in a given day is less than 1 minute.⁴¹

³⁹Specifically, we estimated Column 1 of Table 3 excluding each firm one by one.

⁴⁰In addition to the issue of stock-sampling, the quantity of pre-hire data varies across firms. However, as noted above, our main result on exceptions and duration is robust to excluding individual firms from our sample.

⁴¹This is about one percent of transactions. Our results are stronger if we do not trim. Some other productivity variables are also shared with our data provider, but each variable is only available for an even smaller share of workers than is output per hour. Such variables would likely face significant statistical power issues if subjected to the analyses in the paper (which involve clustering standard errors at the location level).

Test Scores. As describe in the text, applicants are scored as Red, Yellow, or Green. Applicants may receive multiple scores (e.g., if they are being considered for multiple roles). In these cases, we assign applicants to the maximum of their scores.

HR Manager. We do not have data on the characteristics of HR managers (we only see an individual identifier). When applicants interact with more than one HR manager during the recruitment process, they are assigned to the manager with whom they have the most interactions.⁴²

In our sample, information on HR manager is missing for about a third of the job-tested applicants. Individuals with missing information on HR manager are assigned to a separate category for a given location. We have also repeated our main analyses while excluding job-tested individuals where HR manager is missing and obtained qualitatively similar results.

Race, Gender, Age. Data on race, sex, and age are not available for this project. However, Autor and Scarborough (2008) show that job testing does not seem to affect worker race, suggesting that changes in worker demographics such as race are not the mechanism by which job testing improves durations.⁴³

Location Identifiers. In our dataset, we do not have a common identifier for workplace location for workers hired in the pre-testing period and applicants applying post-testing. Consequently, we develop a crosswalk between anonymized location names (used for workers in the pre-testing period) and the location IDs in the post-testing period. We drop workers from our sample where the merge did not yield a clean location variable.⁴⁴

Hiring Practice Survey. Our data firm also conducted an informal survey of 6 client firms on our behalf to better understand hiring practices once testing was adopted. The survey indicated that firms encouraged managers to hire workers with higher scores (and some firms had policies on not hiring low-scored candidates), but left substantial leeway for managers to overrule testing recommendations.

Job Offers. As discussed in the main text, our data for this project do not include information on the receipt of job offers, only on realized job matches. The data firm has a small amount of information on offers received, but is only available for a few firms and a small share of the total applicants in our sample, and would be too sparse to be of use for this project.

⁴²This is performed excluding interactions where information on the HR manager is missing. If there is a tie for most interactions, applicants are assigned the last manager among those interacted with.

⁴³Autor and Scarborough (2008) do not look at impact of testing on gender. However, they show there is little differential impact of testing by gender (or race).

⁴⁴This includes locations in the pre-testing data where testing is never later introduced.

B Proofs

B.1 Preliminaries

We first provide more detail on the firm's problem, to help with the proofs.

Under Discretion, the manager hires all workers for whom $U_i = (1-k)E[a|s_i, t_i] + kb_i > \underline{u}$ where \underline{u} is chosen so that the total hire rate is fixed at W .

We assume b_i is perfectly observable, that $a|t \sim N(\mu_t, \sigma_a^2)$, and that $s_i = a_i + \epsilon_i$ where $\epsilon \sim N(0, \sigma_\epsilon^2)$ and is independent of a and b .

Thus $E[a|s, t]$ is normally distributed with known parameters. Also, since $s|t$ is normally distributed and the assessment of a conditional on s and t is normally distributed, the assessment of a unconditional on s (but still conditional on t) is also normally distributed with a mean μ_t and variance $\sigma = \frac{(\sigma_a^2)^2}{\sigma_\epsilon^2 + \sigma_a^2}$. Finally, define U_t as the manager's utility for a given applicant, conditional on t . The distribution of U_t unconditional on the signals and b , follows a normal distribution with mean $(1-k)\mu_t$ and variance $(1-k)^2\sigma + k^2\sigma_b^2$.

Thus, the probability of being hired is as follows, where $\tilde{z}_t = \frac{\underline{u} - (1-k)\mu_t}{\sqrt{(1-k)^2\sigma + k^2\sigma_b^2}}$.

$$W = p_G(1 - \Phi(\tilde{z}_G)) + (1 - p_G)(1 - \Phi(\tilde{z}_Y)) \quad (5)$$

The firm is interested in expected match quality conditional on being hired under Discretion. This can be expressed as follows, where $\lambda(\cdot)$ is the inverse Mill's ratio of the standard normal and $z_t(b_i) = \frac{\underline{u} - kb_i - \mu_t}{\sigma}$, i.e., the standard-normalized cutpoint for expected match quality, above which, all applicants with b_i will be hired.

$$E[a|Hire] = E_b[p_G(\mu_G + \lambda(z_G(b_i))\sigma) + (1 - p_G)(\mu_Y + \lambda(z_Y(b_i))\sigma)] \quad (6)$$

Inside the expectation, $E_b[\cdot]$, we have the expected value of a among all workers hired for a given b_i . We then take expectations over b .

Under No Discretion, the firm hires based solely on the test. Since we assume there are plenty of type G applicants, the firm will hire among type G applicants at random. Thus the expected match quality of hires equals μ_G .

B.2 Proof of Proposition 3.1

The following results formalize conditions under which the firm will prefer Discretion or No Discretion.

1. For any given precision of private information, $1/\sigma_\epsilon^2 > 0$, there exists a $k' \in (0, 1)$ such that if $k < k'$ match quality is higher under Discretion than No Discretion and the opposite if $k > k'$.
2. For any given bias, $k > 0$, there exists $\underline{\rho}$ such that when $1/\sigma_\epsilon^2 < \underline{\rho}$, i.e., when precision of private information is low, match quality is higher under No Discretion than Discretion.
3. For any value of information $\bar{\rho} \in (0, \infty)$, there exists a bias, $k'' \in (0, 1)$, such that if $k < k''$ and $1/\sigma_\epsilon^2 > \bar{\rho}$, i.e., high precision of private information, match quality is higher under Discretion than No Discretion.

For this proof we make use of the following lemma:

Lemma B.1 *The expected match quality of hires for a given manager, $E[a|Hire]$, is decreasing in managerial bias, k .*

Proof A manager will hire all workers for whom $(1 - k)E[a|s_i, t_i] + kb_i > \underline{u}$, i.e., if $b_i > \frac{\underline{u} - (1-k)E[a|s_i, t_i]}{k}$. Managers trade off b for a with slope $-\frac{1-k}{k}$. Consider two managers, Manager 1 and Manager 2, where $k_1 > k_2$, i.e., Manager 1 is more biased than Manager 2. Manager 2 will have a steeper (more negative) slope ($\frac{1-k_2}{k_2} > \frac{1-k_1}{k_1}$) than Manager 1. There will thus be some cutoff \hat{a} such that for $E[a|s_i, t_i] > \hat{a}$ Manager 2 has a lower cutoff for b and for $E[a|s_i, t_i] < \hat{a}$, Manager 1 has a lower cutoff for b .

That is, some candidates will be hired by both managers, but for $E[a|s_i, t_i] > \hat{a}$, Manager 2 (less bias) will hire some candidates that Manager 1 would not, and for $E[a|s_i, t_i] < \hat{a}$ Manager 1 (more bias) will hire some candidates that Manager 2 would not. The candidates that Manager 2 would hire when Manager 1 would not, have high expected values of a , while the candidates that Manager 1 would hire where Manager 2 would not have low expected values of a . Therefore the average a value for workers hired by Manager 2, the less biased manager, must be higher than that for those hired by Manager 1. $E[a|Hire]$ is decreasing in k .

We next prove each item of Proposition 3.1

1. For any given precision of private information, $1/\sigma_\epsilon^2 > 0$, there exists a $k' \in (0, 1)$ such that if $k < k'$ match quality is higher under Discretion than No Discretion and the opposite if $k > k'$.

Proof When $k = 1$, the manager hires based only on b , which is independent of a . So $E[a|Hire] = p_G\mu_G + (1 - p_G)\mu_Y$. The firm would do better under No Discretion (where

match quality of hires equals μ_G). When $k = 0$, the manager hires only applicants whose expected match quality, a , is above the threshold. In this case, the firm will at least weakly prefer Discretion. Since the manager's preferences are perfectly aligned, he or she will always do at least as well as hiring only type G .

Thus, Discretion is better than No Discretion for $k = 0$ and the opposite is true for $k = 1$. Lemma B.1 shows that the firm's payoff is decreasing in k . There must therefore be a single cutpoint, k' , where, below that point, the firm's payoff for Discretion is large than that for No Discretion, and above that point, the opposite is true.

2. *For any given bias, $k > 0$, there exists $\underline{\rho}$ such that when $1/\sigma_\epsilon^2 < \underline{\rho}$, i.e., when precision of private information is low, match quality is higher under No Discretion than Discretion.*

Proof When $1/\sigma_\epsilon^2 = 0$, i.e., the manager has no information, and $k = 0$, he or she will hire based on the test, resulting in an equal payoff to the firm as No Discretion. For all $k > 0$, the payoff to the firm will be worse than No Discretion, thanks to lemma B.1. Thus when the manager has no information the firm prefers No Discretion to Discretion.

We also point out that the firm's payoff under Discretion, expressed above in equation (6), is clearly continuous in σ (which is continuous in $1/\sigma_\epsilon^2 = 0$).

Thus, when the manager has no information, the firm prefers No Discretion and the firm's payoff under Discretion is continuous in the manager's information. Therefore there must be a point $\underline{\rho}$ such that, for precision of manager information below that point, the firm prefers No Discretion to Discretion.

3. *For any value of information $\bar{\rho} \in (0, \infty)$, there exists a bias, $k'' \in (0, 1)$, such that if $k < k''$ and $1/\sigma_\epsilon^2 > \bar{\rho}$, i.e., high precision of private information, match quality is higher under Discretion than No Discretion.*

Proof First, we point out that when $k = 0$, the firm's payoff under Discretion is increasing in $1/\sigma_\epsilon^2$. An unbiased manager will always do better (from the firm's perspective) with more information than less. Second, we have already shown that for $k = 0$, Discretion is always preferable to No Discretion, regardless of the manager's information, and when σ_ϵ^2 approached ∞ , there is no difference between Discretion and No Discretion from the firm's perspective.

Define $\Delta(\sigma_\epsilon^2, k)$ as the difference in match quality of hires under Discretion, compared to no Discretion, for fixed manager type (σ_ϵ^2, k) . We know that $\Delta(\sigma_\epsilon^2, 0)$ is positive and

decreasing in σ_ϵ^2 , and approaches 0 as σ_ϵ^2 approaches ∞ . Also, since the firm's payoff under discretion is continuous in both k and $1/\sigma_\epsilon^2$ (see equation 6 above), $\Delta(\cdot)$ must also be continuous in these variables.

Fix any $\bar{\rho}$ and let $\bar{\sigma}_\epsilon^2 = 1/\bar{\rho}$. Let $y = \Delta(\bar{\sigma}_\epsilon^2, 0)$. We know that $\Delta(\sigma_\epsilon^2, 0) > y$ for all $\sigma_\epsilon^2 < \bar{\sigma}_\epsilon^2$.

Let $d(k) = \max_{\sigma_\epsilon^2 \in [0, \bar{\sigma}_\epsilon^2]} \Delta(\sigma_\epsilon^2, k) - \Delta(\sigma_\epsilon^2, 0)$. We know $d(k)$ exists because $\Delta(\cdot)$ is continuous wrt σ_ϵ^2 and the interval over which we take the maximum is compact. We also know that $d(0) = 0$, i.e., for an unbiased manager, the return to discretion is maximized when managers have full information. Finally, $d(k)$ is continuous in k because $\Delta(\cdot)$ is.

Therefore, we can find $k'' > 0$ such that $d(k) = d(k) - d(0) < y$ whenever $k < k''$. This means that $\Delta(\sigma_\epsilon^2, k) > 0$ for $\sigma_\epsilon^2 < \bar{\sigma}_\epsilon^2$. In other words, at bias k and $\rho > \underline{\rho}$, Discretion is better than No Discretion.

B.3 Proof of Proposition 3.2

Across M , the exception rate, R_m , is increasing in both managerial bias, k , and the precision of the manager's private information, $1/\sigma_\epsilon^2$

Proof Because the hiring rate is fixed at W , $E[\text{Hire}|Y]$ is a sufficient statistic for the probability that an applicant with $t = Y$ is hired *over* an applicant with $t = G$, i.e., an exception is made.

Above, we defined U_t , a manager's utility of a candidate conditional on t , and showed that it is normally distributed with mean $(1 - k)\mu_t$ and variance $\Sigma = (1 - k)^2\sigma + k^2\sigma_b^2$. A manager will hire all applicants for whom U_t is above \underline{u} where the latter is chosen to keep the hire rate fixed at W .

Consider the difference in expected utility across G and Y types. If $\mu_G - \mu_Y$ were smaller, more Y types would be hired, while fewer G types would be hired. This is because, at any given quantile of U_G , there would be more Y types above that threshold.

Let us now define $\tilde{U}_t = \frac{U_t}{\sqrt{\Sigma}}$. This transformation is still normally distributed but now has mean $\frac{(1-k)\mu_t}{\sqrt{\Sigma}}$ and variance 1. This rescaling of course does nothing to the cutoff \underline{u} , and it will still be the case that the probability of an exception is decreasing in the difference in expected utilities across \tilde{U}_G and \tilde{U}_Y : $\Delta_U = \frac{(1-k)(\mu_G - \mu_Y)}{\sqrt{\Sigma}}$.

It is easy to show (with some algebra) that $\frac{\partial \Delta_U}{\partial k} = \frac{-(\mu_G - \mu_Y)\sigma_b^2}{\Sigma^{3/2}}$, which is clearly negative. When k is larger, the expected gap in utility between a G and a Y narrows so the probability of hiring a Y increases.

Similarly, it is each to show that $\frac{\partial \Delta_U}{\partial \sigma_\epsilon^2} = \frac{(1-k)^3(\mu_G - \mu_Y)(\sigma_a^2)^2}{2\Sigma^{3/2}(\sigma_\epsilon^2 + \sigma_a^2)^2}$, which is clearly positive. The gap in expected utility between G and Y widens when managers have less information. It

thus narrows when managers have better private information, as does the probability of an exception.

B.4 Proof of Proposition 3.3

If the quality of hired workers is decreasing in the exception rate, $\frac{\partial E[a|Hire]}{\partial R_m} < 0$ across M , then firms can improve outcomes by eliminating discretion. If quality is increasing in the exception rate then discretion is better than no discretion.

Proof Consider a manager who makes no exceptions even when given discretion: Across a large number of applicants, this only occurs if this manager has no information and no bias. Thus the quality of hires by this manager is the same as that of hires under a no discretion regime, i.e., hiring decisions made solely on the basis of the test. Compare outcomes for this manager to one who makes exceptions. If $\frac{\partial E[a|Hire]}{\partial R_m} < 0$, then the quality of hired workers for the latter manager will be worse than for the former. Since the former is equivalent to hires under no discretion, it then follows that the quality of hires under discretion will be lower than under no discretion. If the opposite is true and the manager who made exceptions, thereby wielding discretion, has better outcomes, then discretion improves upon no discretion.

APPENDIX TABLE A1: THE IMPACT OF JOB TESTING FOR COMPLETED JOB SPELLS
ADDITIONAL OUTCOMES

	Mean Completed Duration (Days, Mean=211; SD=232)	>3 Months (Mean=0.62; SD=0.21)	>6 Months (Mean=0.46; SD=0.24)	>12 Months (Mean=0.32; SD=0.32)				
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
<i>Post-Testing</i>								
	88.89** (35.91)	29.66* (15.40)	0.0404*** (0.00818)	0.0895** (0.0379)	0.0906*** (0.00912)	0.107** (0.0463)	0.107*** (0.00976)	0.0652* (0.0363)
N	4,401	4,401	4,505	4,505	4,324	4,324	3,882	3,882
Year-Month FEs	X	X	X	X	X	X	X	X
Location FEs	X	X	X	X	X	X	X	X
Client Firm X Year FEs		X		X		X		X
Location Time Trends		X		X		X		X
Size and Composition of Applicant Pool		X		X		X		X

*** p<0.1, ** p<0.05, *p<0.1

NOTES: See notes to Table 2. The dependent variables are the mean length of completed job spells in days and the share of workers in a location-cohort who survive 3, 6, or 12 months, among those who are not right-censored.

APPENDIX TABLE A2: TESTING, EXCEPTIONS, AND TAIL OUTCOMES

	Impact of testing	Exceptions and outcomes, post testing	Impact of testing, by exception rate		
			Level of aggregation for exception rate		
			<i>Pool</i>	<i>Manager</i>	<i>Location</i>
(1)	(2)	(3)	(4)	(5)	
Dependent Variable: 90th Percentile of Log Duration of Completed Spells					
<i>Post-Testing</i>	0.126* (0.069)		0.150** (0.0703)	0.141** (0.0697)	0.131* (0.0716)
<i>Exception Rate*Post-Testing</i>		-0.0426* (0.0242)	-0.0461** (0.0215)	-0.0764*** (0.0270)	-0.0311 (0.0563)
Dependent Variable: 10th Percentile of Log Duration of Completed Spells					
<i>Post-Testing</i>	0.239* (0.144)		0.186 (0.145)	0.260* (0.144)	0.314** (0.141)
<i>Exception Rate*Post-Testing</i>		-0.0322 (0.0336)	-0.0540* (0.0322)	-0.204*** (0.0658)	-0.456*** (0.170)
N	6,956	3,839	6,869	6,942	6,956
Year-Month FEs	X	X	X	X	X
Location FEs	X	X	X	X	X
Client Firm X Year FEs	X	X	X	X	X
Location Time Trends	X	X	X	X	X
Size and Composition of Applicant Pool	X	X	X	X	X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: See notes to Tables 2 and 3.

APPENDIX TABLE A3: IMPACT OF COLOR SCORE ON JOB DURATION BY PRE-TESTING LOCATION DURATION

	<i>Log Duration of Completed Spells</i>					
	High Duration	Low Duration	High Duration	Low Duration	High Duration	Low Duration
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Green</i>	0.165*** (0.0417)	0.162*** (0.0525)	0.161*** (0.0406)	0.172*** (0.0541)	0.175*** (0.0535)	0.170*** (0.0528)
<i>Yellow</i>	0.0930** (0.0411)	0.119** (0.0463)	0.0886** (0.0403)	0.130*** (0.0481)	0.108** (0.0513)	0.112** (0.0495)
N	23,596	32,284	23,596	32,284	23,596	32,284
Year-Month FEs	X	X	X	X	X	X
Location FEs	X	X	X	X	X	X
Client Firm X Year FEs			X	X	X	X
Application Pool FEs					X	X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: Each observation is an individual, for hired workers post testing only. The omitted category is red workers. Locations are classified as high duration if their mean duration pre-testing was above median for the pre-testing sample.

APPENDIX TABLE A4: IMPACT OF COLOR SCORE ON JOB DURATION BY
LOCATION-SPECIFIC EXCEPTION RATES

	<i>Log Duration of Completed Spells</i>					
	High Exception Rate	Low Exception Rate	High Exception Rate	Low Exception Rate	High Exception Rate	Low Exception Rate
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Green</i>	0.173*** (0.0317)	0.215*** (0.0689)	0.172*** (0.0307)	0.205*** (0.0711)	0.181*** (0.0331)	0.171*** (0.0642)
<i>Yellow</i>	0.112*** (0.0287)	0.182** (0.0737)	0.112*** (0.0280)	0.174** (0.0760)	0.117*** (0.0296)	0.128* (0.0711)
N	36,088	31,928	36,088	31,928	36,088	31,928
Year-Month FEs	X	X	X	X	X	X
Location FEs	X	X	X	X	X	X
Client Firm X Year FEs			X	X	X	X
Application Pool FEs					X	X

*** p<0.1, ** p<0.05, * p<0.1

NOTES: Each observation is an individual, for hired workers post testing only. The omitted category is red workers. Locations are classified as high exception rate if their mean exception rate post-testing was above median for the post-testing sample.