

TAX FARMING REDUX: EXPERIMENTAL EVIDENCE ON PERFORMANCE PAY FOR TAX COLLECTORS*

ADNAN Q. KHAN
ASIM I. KHWAJA
BENJAMIN A. OLKEN

Performance pay for tax collectors has the potential to raise revenues, but might come at a cost if it increases the bargaining power of tax collectors vis-à-vis taxpayers. We report the first large-scale field experiment on these issues, where we experimentally allocated 482 property tax units in Punjab, Pakistan, into one of three performance pay schemes or a control. After two years, incentivized units had 9.4 log points higher revenue than controls, which translates to a 46% higher growth rate. The scheme that rewarded purely on revenue did best, increasing revenue by 12.9 log points (64% higher growth rate), with little penalty for customer satisfaction and assessment accuracy compared to the two other schemes that explicitly also rewarded these dimensions. The revenue gains accrue from a small number of properties becoming taxed at their true value, which is substantially more than they had been taxed at previously. The majority of properties in incentivized areas in fact pay no more taxes, but instead report higher bribes. The results are consistent with a collusive setting in which performance pay increases collectors' bargaining power over taxpayers, who have to either pay higher bribes to avoid being reassessed or pay substantially higher taxes if collusion breaks down. *JEL* Codes: D73, H26, H83, J33.

I. INTRODUCTION

Tax systems in developing countries collect substantially less revenue as a share of GDP than do their counterparts in

*This project is the result of collaboration among many people. We thank the editor, Lawrence Katz, the second editor, Andrei Shleifer, and four anonymous referees for helpful comments. We thank Jon Hill, Donghee Jo, Alyssa Lawther, Kunal Mangal, Wayne Sandholtz, Mahvish Shaukat, Gabriel Tourek, He Yang, and Gabriel Zucker for outstanding research assistance in Cambridge and Zahir Ali, Osman Haq, Turab Hassan, Zahra Mansoor, Obeid Rahman, Shahrukh Raja, Adeel Shafqat, and Sadaqat Shah for outstanding research assistance in Lahore. We thank all the secretaries, director generals, directors, the two project directors from the Punjab Department of Excise and Taxation, the Punjab Finance, Planning and Development departments, and the Chief Secretary and Chief Minister's offices for their support over the many years of this project. Financial support for the evaluation came from 3ie, the IGC, and the NSF (under grant SES-1124134), and financial support for the incentive payments described here came from the Government of the Punjab, Pakistan. The views expressed here are those of the authors and do not necessarily reflect those of the many individuals or organizations acknowledged here.

© The Author(s) 2015. Published by Oxford University Press, on behalf of President and Fellows of Harvard College. All rights reserved. For Permissions, please email: journals.permissions@oup.com

The Quarterly Journal of Economics (2016), 219–271. doi:10.1093/qje/qjv042.

Advance Access publication on November 1, 2015.

developed countries (Gordon and Li 2009; Kleven, Kreiner, and Saez 2014). While there are many differences, an important one is the role played by the tax officials in assessing, enforcing, and auditing taxes. Combined with relatively low wages and limited performance rewards, the temptations for tax inspectors to collude with taxpayers to reduce tax receipts are great.

One possible solution is to tie the compensation of tax staff to the revenue they generate. This is an old idea: historically, states from the Roman empire through the French monarchy (Bartlett 1994; White 2004) sold the rights to collect taxes to “tax farmers,” who then kept a fraction (or in some cases all) of the tax revenue they collected. U.S. states similarly experimented with highly incentivized “tax ferrets” to collect property taxes in the nineteenth century. Though tax officials in these historical regimes were unpopular, and the world has largely moved to salaried tax officials (Parrillo 2013), countries such as Brazil, Peru, Pakistan, and others have begun to reconsider incentives for tax staff (Das-Gupta and Mookherjee 1998; Kahn, Silva, and Ziliak 2001) as a way to improve tax compliance.

The challenge, however, is that by strengthening the bargaining ability of tax collectors, performance incentives may not only lead to taxpayer dissatisfaction, but may also alter the division of rents from collusion without necessarily increasing revenue raised by the government. To see this, consider a simple bargaining setting in which a tax collector colludes with a taxpayer to reduce the tax assessment in exchange for a bribe. If there is no cost to either party from reducing tax liability, then performance pay for tax collectors will simply raise the bribe paid with no impact on revenue, as the taxpayer now has to compensate the collector’s forgone incentive payment with a higher bribe. In more realistic settings, where there is some cost to either party from reduced tax liabilities, there will be two different effects: some taxpayers will continue in the collusive, low tax equilibrium but pay higher bribes, whereas others will end up paying higher taxes and lower bribes as they switch from the collusive, low tax equilibrium to a noncollusive, higher tax equilibrium. Performance pay could thus have heterogeneous effects on tax revenue and bribes among taxpayers. Whether performance pay actually leads to increased revenue—and at what cost in terms of higher bribes and potentially forgone taxpayer satisfaction—is therefore ultimately an empirical question.

In this article, we provide what is to the best of our knowledge the first experimental evidence on these questions. Working with the Punjab, Pakistan provincial government, we randomly allocated tax officials in the entire provincial urban property tax department, which consists of 482 property tax units (known as circles), into one of three versions of performance-based pay schemes or a control group. A total of 218 circles, consisting of about 550 tax personnel, were randomly allocated to one of the three treatment groups for two fiscal years. The incentives were large: the three-person tax team in each treated circle was collectively given an average of 30% of all tax revenues it collected above a historically predicted benchmark. Many personnel in treated areas were able to double their baseline salaries or more through these incentives.

Given concerns about potential negative impacts of high-powered incentives, the three schemes varied in both the extent to which they based performance pay explicitly on taxpayer satisfaction and accuracy of assessment in addition to revenue, and the extent to which they allowed for subjective evaluation on the part of the tax department. The “revenue” scheme provided incentives based solely on revenue collected above a benchmark predicted from historical data. To address multitasking concerns and in an effort to incentivize accurate assessment (and hence also tackle collusion) (Hölmstrom and Milgrom 1991), the “revenue plus” scheme provided incentives exactly as in the revenue scheme, but made adjustments (plus/minus three-fourths of baseline salary) based on whether the circle ranked in the top, middle, or bottom third of circles in terms of taxpayer satisfaction and accuracy of tax assessments, as determined by an independent survey of taxpayers. To allow for more subjective assessments rather than purely formulaic criteria (Baker, Gibbons, and Murphy 1994; MacLeod 2003), the third scheme, “flexible bonus,” took this a step further by rewarding collectors for a much wider set of prespecified criteria set by the tax department and by allowing for subjective adjustments based on period-end overall performance.

We evaluate the impact of the schemes using multiple sources of data. For tax revenues, we obtained administrative data, which we verified by conducting random spot checks against the tax department’s bank records. For outcomes such as perceived corruption and satisfaction with the tax department, we conducted a survey of over 16,000 taxpayers and their

properties throughout the province. For estimating assessment accuracy, the surveyors also directly observed and recorded the property characteristics used in the tax calculation. We then manually matched surveyed properties to the tax rolls to obtain the corresponding tax records for each property. Tax assessment is determined formulaically from these property characteristics, so this allowed us to determine the accuracy of assessments by comparing our survey measurements to those on the official tax rolls.

We find that on average across the three schemes, by the end of the two years performance pay led to an increase in tax revenue of about 9.4 log points based on the administrative data. This translates to a 46% higher growth rate in revenues compared with control areas. We show that this came predominantly through an increase in the reported tax base (i.e., the total assessed value of properties) rather than through increased recovery or changes in exemptions granted. On average, we find little impact of the schemes on taxpayer satisfaction. Specifically, the increased revenue generated as a result of the schemes is not accompanied by a decline in the typical taxpayers' perceptions of the quality of service from the tax office or in their satisfaction with their dealings with the tax office. We also find no overall change in the accuracy of tax assessments. Thus, on average we find that the incentives increase revenue with little obvious downside in terms of overall perception of the tax department in the eyes of the typical taxpayer.

Comparing the three schemes, we find that they differ substantially in terms of their impact on revenue, with relatively small differences on taxpayer satisfaction and perception of the tax department. Specifically, the revenue scheme, which provided incentives purely based on revenue collected, showed 15.2 log points higher current-year revenues relative to controls (57% higher growth rate) by the second year. In comparison, the revenue plus scheme achieved only 8.1 log points, and the flexible bonus scheme only a statistically insignificant 3.5 log point increase in current-year revenue. While the revenue plus scheme did improve perceived customer satisfaction and quality perceptions relative to the revenue and flexible bonus schemes, the differences were small, and the substantially lower revenue collected meant that this scheme had a substantially lower rate of return. The flexible bonus scheme did not do better on any dimension we can measure in our data, and in fact did worse

compared with the control group on perception of the department's quality. Thus, adding multiple dimensions to performance pay substantially diluted the impact on revenue without a substantial corresponding increase in nonmonetary outcomes.

Our survey data suggest that there was indeed a reallocation of rents associated with performance incentives and finds evidence of precisely the sort of heterogeneity suggested by a simple theory of collusion. For most properties in performance pay tax circles, taxpayers were not reassessed and reported no change in tax paid. However, relative to the control group, they reported a Rs. 594 (about US\$6) increase in the going rate for a bribe paid to property tax officers for properties similar to theirs, which represents a roughly 32% increase. Although this does not necessarily imply that every household paid these higher bribes, respondents also indicated that bribe payments were more frequent.

However, for the small number of properties whose tax valuation was formally changed (either newly assessed or reassessed), these taxpayers report paying substantially higher taxes, but do not report the higher bribes that other properties in performance pay circles reported. Moreover, while comparisons between our survey data and corresponding administrative records suggest that typical properties are undertaxed, this does not hold for these reassessed properties, which appear on average to be taxed accurately. There is also an increase in the number of these newly assessed or reassessed properties in performance pay circles. These results are consistent with what one might expect given collusion: performance pay means that inspectors can demand higher bribes to compensate them for their forgone performance pay, but, given the higher bribe now required to maintain collusion, some taxpayers may instead switch from collusion (low tax, high bribe) to noncollusion (high tax, low bribe).

These results suggest that the increase in tax collected under the performance pay schemes is driven by a relatively small number of properties that are (correctly) reassessed and switch from collusion to noncollusion, paying much higher taxes and lower bribes. It is interesting to examine what determines who ends up in this group. In general, we find that these newly reassessed properties have taxable value that is about 67% higher than the typical (nonreassessed) property. In treatment areas, the reassessed properties are even more valuable than reassessed properties elsewhere, by another 33%. Reassessed properties in

general are also more likely to be commercial properties, which are taxed at a higher rate. There is also some suggestive evidence that although property owners with political connections avoid being reassessed in control areas, they lose this degree of protection in treatment areas. On net, the results suggest that tax inspectors focus on a small number of high-value properties to increase revenue, thus potentially raising revenue while minimizing political costs.

From the government's perspective, the relative desirability of the schemes depends on the government's objective function. For a politician who seeks to maximize tax revenues subject to political constraints, the evidence presented here suggests that the revenue scheme is the most effective: it raised the current-year revenue by 15.2 log points (57% higher growth rate), which implies a substantially positive return on investment (35%–51%), and it did not appreciably reduce satisfaction with the tax department compared to controls. While the revenue plus scheme did slightly better on satisfaction than the revenue scheme, it generated a lower (14%–28%) return on investment.

This article builds on several different literatures. First, while there is a substantial tradition of theoretical work on performance pay and compensation for tax officials in the developing world (see, for example, Besley and McLaren 1993; Mookherjee and Png 1995), there is very little empirical evidence on how these types of incentives work in practice.¹ Indeed, although there is a small but growing and exciting empirical literature on tax and development, it has focused to date primarily on how taxpayers respond to different types of enforcement (e.g., Gordon and Li 2009; Pomeranz 2013; Kumler, Verhoogen, and Frasn 2013; Carillo, Pomeranz, and Singhal 2014) and various aspects of the tax code (Best et al., 2013; Kleven and Waseem 2013), rather than on the role of, or how to improve performance of, tax staff. Second, this article is related to several recent papers on improving developing country civil service performance in other contexts and using other tools. Existing work has focused on the role of wages (Dal Bó, Finan, and Rossi 2013), intrinsic motivation (Ashraf, Bandiera, and Jack 2013), and management (Rasul

1. To our knowledge, the best empirical evidence on the impact of performance pay on tax collection is a time-series study of a performance pay reform in Brazil (Kahn, Silva, and Ziliak 2001), which is not able to examine any nonrevenue outcomes such as bribery or taxpayer satisfaction.

and Rogger 2013). The recent work on performance pay has been centered on education and health sectors (Glewwe, Ilias, and Kremer 2010; Muralidharan and Sundararaman 2011; Gertler and Vermeersch 2013), where collusive forces are not as salient. Finally, this article builds on the growing literature on corruption (see Olken and Pande 2012 for a review). For example, a recent paper by Duflo et al. (2013) shows that changing the incentives for third-party auditors to make them more independent increases honesty in their reporting; this article finds similar benefits but also highlights that such incentives have the counter-vailing potential to increase bribes if collusion continues. More generally, it underscores that when there is corruption, output-based incentives for government officials can have very different effects depending on how they affect the downstream bargaining between officials and citizens.

The remainder of this article is structured as follows. Section II describes the relevant features of the property tax administration in Punjab, the setting in which the study takes place. Section III outlines theoretically what impact one might expect from performance pay in a setting with collusion between tax inspectors and taxpayers. Section IV outlines the experimental design, Section V describes the data and empirical approach, and Section VI presents the results. Section VII concludes. All supplementary material is available in the Online Appendix.

II. SETTING

II.A. Property Taxes in Urban Punjab

Punjab is Pakistan's most populous province: its population of over 80 million would rank fifteenth in the world were it a country. Property tax collection in Punjab is roughly a fifth of the level of comparable countries (World Bank 2006) due to a wide variety of problems: not only is the tax base narrow and tax rates low, but also tax evasion and corruption are widespread, distrust in public institutions runs high, and administration is weak (World Bank 2006, 2009; Bahl, Wallace, and Cyan 2008).

The urban property tax in Punjab is levied on the gross annual rental value (GARV) of the property, which is computed by formula. The GARV is determined by measuring the square footage of land and buildings on the property and then multiplying by standardized values from a table. These valuation tables divide

the province into seven categories (A–G) according to the extent of facilities and infrastructure in the area, with different rates for each category. Rates further vary by residential, commercial, or industrial status, whether the property is owner-occupied or rented, and location (i.e., on or off a main road). Taxes are paid into designated bank branches (through the National Bank of Pakistan). A copy of the receipt of payment is given to the taxpayer at the time of payment, and the bank also provides a copy to the tax collector and a copy to the provincial treasury.

Several distortions place constraints on tax collection and introduce substantial scope for corruption. These distortions include substantially different rates for residential and commercial properties (which can be easily reclassified), as well as granting exemptions to widows, the disabled, owners of plots below 5 marlas (about 125 square meters), retired federal and provincial government employees, and religious charitable institutions (World Bank 2006). The two most notable distortions are between owner-occupied and rented residential properties (the latter are taxed 10 times more) and between residential and commercial properties (the latter are taxed between 3 and 6 times more). Qualitative evidence suggests that these distortions are the main ways tax evasion takes place, both due to the significant impact these margins have on tax assessment and also because it is less easy to verify whether a residential property is being rented or, particularly for mixed usage properties, what fraction of the property is being used for commercial purposes.

For research purposes, a methodological advantage of property taxes is that unlike most taxes, true property tax liability can be independently estimated by the researcher. By comparing official tax payments to an independent assessment by an external survey team, we can determine changes to both the accuracy of tax evasion and the average level of over or under taxation. This approach follows other examples in the corruption literature (e.g., Fisman and Wei 2004; Olken 2007).

II.B. Property Tax Administration

The primary unit of tax collection is the “circle,” a predefined geographical area that covers anywhere from 2,000 to 10,000 unique properties. Within each circle is a team of three tax officers: an “inspector” who leads the team, determines tax assessments, and issues notices that demand payment; a “clerk” in

charge of record keeping; and a “constable” who assists the inspector in the field. Together they maintain a record of all properties and their attributes (size, type of use, etc.), apply the valuation tables to each property, and determine which exemptions apply. The inspector determines each property’s tax liability and sends an annual tax bill to the property owner.

All three officials are part of the provincial career bureaucracy, with wages determined by salary band and length of service. As is common for civil servants in developing economies, tax officials receive fairly low wages that are rarely, if ever, tied to performance. However, since the department has explicit financial targets each year, there is pressure on each circle team to contribute. This occurs typically through each administrative level pressuring lower levels to increase collections. With limited reward mechanisms or vertical mobility, threats of transfers are the primarily tool available to supervisors who want to improve performance. While some inspectors do have strong preferences over their posted circle, these threats have limited effectiveness since transfers are often more politically based than merit based (Piracha and Moore 2015).

The tax administration system leaves considerable opportunities for leakages, collusion, and low collection, especially because there are few independent checks on the actions of the tax circle team and limited audit mechanisms. For example, the property database is manually recorded on physical registers and does not automatically include new properties or property updates. Building permits and rental agreements are not always formally registered, and when they are registered they are not automatically linked to tax rolls, so the tax department learns about new construction or changes in property use only through the efforts of the circle staff. In addition, officials may employ significant discretion in applying valuation tables to individual properties and determining exemptions. For example, properties can be incorrectly designated as owner-occupied when they are being rented out (and as noted, the latter are taxed at a 10 times higher rate), classified as residential when they are in fact commercial, designated as “off road” when they are on a main road, or mismeasured. Finally, the manual system of billing and collection, in which tax bills are hand-written by inspectors and clerks and hand-delivered by tax constables, is prone to errors and/or manipulation in crediting collections.

In this context, performance pay has the potential to induce tax officials to raise collections. Although this could be due to greater effort in tracking new properties and uncovering physical and usage changes that increase tax assessments, anecdotal evidence suggests that collectors likely have substantial private information regarding a property's true tax liability already, which they use for extracting bribes rather than assessing higher taxes. They may choose to only reveal (parts of) this information to the authorities when faced with significant opportunity for rewards. An extreme form of such information disclosure is revealing the existence of (newly constructed) properties and formally adding them to tax registers (recall there is no automatic process though which this happens). In addition, tax collectors may increase valuations by revealing the true, higher tax valuation of a property or denying (incorrectly provided) exemptions. The next section formalizes these incentives to strategically disclose information within a standard model of collusion.

III. THEORY

Consider a simple setting where a taxpayer i faces a true tax liability τ_i^* . The tax inspector knows τ_i^* , but can instead choose to report a lower tax liability to the government, τ_i . The tax inspector receives an incentive payment that is a constant fraction r of actual taxes paid, that is, $r\tau_i$.

Both taxpayer and tax inspector face costs from colluding to report $\tau_i < \tau_i^*$. The taxpayer's cost of accepting a reduced tax liability is $\alpha_i(\tau_i^* - \tau_i)$ and the tax inspector's cost of giving a reduced tax liability is $\beta_i(\tau_i^* - \tau_i)$.² For example, if a taxpayer evades taxation, their property can be sealed and unavailable for use; the utility loss from the expected loss of use of the property for each dollar of tax evasion is captured by α_i . Similarly, a tax official could in theory be jailed for accepting bribes, or he may just experience a utility loss from being dishonest; the utility loss of this

2. While the costs are modeled in terms of deviations from the true tax liability, that is, $\alpha(\tau_i^* - \tau_i)$, an alternative formulation would have costs in terms of bribes, that is, to value bribes $\alpha_i b_i$ instead of b_i , so that bribes are less valuable than cash. This could represent the fact that there is some chance bribes are detected, or that one needs to launder bribe money, which makes it less valuable than legal money. In Online Appendix A, we show similar qualitative results using this alternative specification of costs.

punishment for the tax collector for each dollar of tax evasion is captured by β_i . We allow α_i and β_i to differ for each individual taxpayer-inspector pairing i .

We assume that the taxpayer and the tax inspector engage in Nash bargaining, with the taxpayer potentially paying a bribe b_i as a transfer to tax inspector. If no agreement is reached, the taxpayer receives payoff $-\tau_i^*$ and the tax inspector receives payoff $r\tau_i^*$. If an agreement is reached, the taxpayer receives payoff $-\tau_i - \alpha_i(\tau_i^* - \tau_i) - b_i$ and the tax inspector receives payoff $r\tau_i - \beta_i(\tau_i^* - \tau_i) + b_i$.

To arrive at the solution, note that the joint surplus from agreement is

$$(1) \quad \tau_i^* - \tau_i - \alpha_i(\tau_i^* - \tau_i) + r(\tau_i - \tau_i^*) - \beta_i(\tau_i^* - \tau_i),$$

which can be rewritten as

$$(2) \quad -\tau_i(1 - \alpha_i - \beta_i - r) + (1 - \alpha_i - \beta_i - r)\tau_i^*.$$

This equation shows that if

$$(3) \quad r + \alpha_i + \beta_i < 1,$$

the joint surplus is maximized at $\tau_i = 0$ (full collusion); otherwise the joint surplus is maximized at $\tau_i = \tau_i^*$ (no collusion).

Suppose that γ_i is the bargaining weight of the taxpayer (and $1 - \gamma_i$ is the bargaining weight of the inspector). If collusion takes place, the bribe paid is such that each side receives their outside option plus their share of the surplus. This implies that the bribe the taxpayer pays to the tax inspector is

$$(4) \quad b_i = [(\beta_i + r)\gamma_i + (1 - \gamma_i)(1 - \alpha_i)]\tau_i^*.$$

What are the implications for tax revenue and bribes of moving from no incentive ($r = 0$) to positive incentive payments r ? This simple framework shows that it depends on whether the equilibrium shifts from the collusive equilibrium to the noncollusive equilibrium. So long as $r + \alpha_i + \beta_i < 1$ and $\gamma_i > 0$, increasing the incentive rate increases bribes (since the taxpayer now has to compensate the inspector for the forgone incentive payments). On the other hand, if increasing r means that the threshold is crossed such that $r + \alpha_i + \beta_i > 1$, then collusion disappears, bribes fall to zero, and tax revenue increases from 0 to τ_i^* . The result that bribes increase with incentives to the tax collector, but that collusion may disappear if incentives are

sufficiently great, is closely related to Shleifer and Vishny (1994) and Boycko, Shleifer, and Vishny (1996), who study bribes between politicians and managers of firms as part of their analysis of privatizations.

Since there is heterogeneity in α_i and β_i , the aggregate response of tax revenue will depend on the fraction of households induced to switch from the collusive to noncollusive equilibrium. Denote by $f(\alpha, \beta, \tau^*)$ the joint distribution of α , β , and τ^* in the population. Then the increase in total tax revenue, T , in response to an increase in r is given by

$$\begin{aligned} \frac{dT}{dr} &= \iiint_{r+\alpha+\beta=1, \tau^* \in (0, \infty)} \tau^* f(\alpha, \beta, \tau^*) d\alpha d\beta d\tau^* \\ (5) \quad &= \iint_{\alpha, \tau^*} \tau^* f(\alpha, 1-r-\alpha, \tau^*) d\alpha d\tau^*. \end{aligned}$$

Equation (5) makes clear that the change in tax revenue depends on the density of taxpayers that are just indifferent between colluding and not colluding, and the average tax liability of those households.

The welfare implications depend on the degree to which the costs from avoidance, α and β , are social costs or private costs (Chetty 2009). To the extent that they represent social costs (i.e., they are utility costs from being dishonest or caught cheating, and so represent real costs, as opposed to fines, which would just be transfers), then assuming equal welfare weights on taxpayers and tax inspectors, the social welfare gain from the increase in r is equal to the increase in tax revenue less the cost of the incentive payments.³ To see this, write social welfare as the net surplus of the taxpayer, tax inspector, and government for each taxpayer-tax inspector pair and integrate as follows:

$$(6) \quad W = \iiint_{\alpha, \beta, \tau^*} \left[\underbrace{-\tau - \alpha(\tau^* - \tau) - b}_{\text{taxpayer}} + \underbrace{r\tau - \beta_i(\tau^* - \tau) + b}_{\text{taxinspector}} + \underbrace{\tau - r\tau}_{\text{government}} \right] f(\alpha, \beta, \tau^*) d\alpha d\beta d\tau^*.$$

3. For example, in Pakistan, if one refuses to pay property taxes, the property might be sealed so no one could use it until the taxes are paid; this represents a real social cost, since the property is unusable during this period. If the government has different welfare weights for payments from taxpayers, bribes received by taxpayers, and tax revenues received by the government, then the welfare formula would be more complex and would need to take these differences into account.

Note whether the costs are social costs (like sealing a property) or private costs (like fines) is a policy choice, and the change in social welfare will be different to the extent that these are private costs (see Chetty 2009 for a more detailed discussion).

Assuming equal social welfare weights, so that the change in social welfare is really just the change in deadweight loss, when we change r , the only change in social welfare comes from those taxpayer-inspector pairs i induced to switch from collusion to noncollusion; for pairs that remain collusive, bribes increase, but this is just a transfer and does not affect social welfare. Similarly, for pairs that remain noncollusive, there is no change in overall welfare since overall taxes paid are unchanged (and there are no bribes) and the greater incentive payments made to the tax collector are entirely offset by the cost of these payments to the government. Those taxpayer-inspector pairs induced to change by a marginal increase in r are those that were just indifferent between colluding and not colluding (i.e., had $\alpha_i + \beta_i + r = 1$), so a switch from collusion to noncollusion does not change the sum of taxpayer and tax inspector utility. The government, however, experiences a first-order utility change equal to the tax revenues it collects less the incentive payment it needs to pay out. The fact that social welfare is equal to the net change in the government's fiscal position is related to the classic result by Feldstein (1999), with the exception that in this case, the "fiscal externality" is tax revenue net of incentive payments.

The model presented here was simplified for ease of exposition, in that the costs to reducing tax liability are linear. Linearity is not crucial; as we outline in the model with more general cost functions (which also avoids the corner solutions inherent in the linear case) in Online Appendix A, all we need for the key qualitative patterns is that the marginal costs of collusion to both parties are weakly positively increasing in $\tau_i^* - \tau_i$. The key difference in this generalized model is that the increase in revenue from an increase in r comes not just from households that switch from collusive to non-collusive (as in equation (5)) but also from households that continue to collude, but now collude a bit less than before. Welfare analysis, however, will be similar.

Note that by assuming that the tax inspector knows τ_i^* , we have suppressed both an effort and an overtaxation/extortion margin. The effort margin recognizes that with more effort tax inspectors could discover more properties or learn a property's true tax liability. We suppress it in the model since models of

increased effort under incentives are well understood, and we wish to focus on the bargaining implications, though changes in effort could be possible in our context as well. We give one example in Online Appendix A of how our framework could be extended to include an effort component. As we show in that case, including an effort margin does not yield any qualitatively different insights. Although the overtaxation margin is conceptually interesting, in practice this appears less common (as we show below, our property survey suggests the typical property is in fact undertaxed, not overtaxed) and we therefore do not incorporate it in the model.

IV. DESIGN

This section presents the design of the performance pay mechanisms introduced and the experimental design of the study. Section IV.A describes the performance pay program, and Section IV.B describes the randomization and balance check.

IV.A. Performance Pay Design

Tax circles were randomly allocated into one of three performance pay schemes: the revenue, revenue plus, and flexible bonus schemes. A total of approximately 70 circles were allocated to each scheme (50 in the first year and an additional 20 in the second year). In addition, in the second year, two new treatments were added: a performance pay scheme for supervisory personnel, and an “information-only” scheme that replicated the information, meetings, and perceived salience of the revenue scheme, but without any financial payments. We describe each scheme below and then discuss tax officials’ understanding of the schemes and the schemes’ credibility.

1. Revenue-based. This performance pay group rewarded tax circle staff (inspectors, constables, and clerks) based on the revenue they collected above a predefined benchmark. The benchmark for each circle was generated using historical revenue data for that circle. Specifically, each tax staff member continued to receive his or her current base salary, plus a bonus calculated by the following formula:

$$(7) \quad \text{Bonus}_c = \alpha_c \max(\text{Revenue}_c - \text{Benchmark}_c, 0),$$

where the bonus rate α_c is 40% for circles below the 50th percentile in baseline revenue, 30% for circles between the 50th and 75th percentiles, and 20% for circles above the 75th percentile. The differential bonus rates were put in place for equity considerations, that is, staff in larger circles were compensated at a lower rate than those in smaller circles, where it was perceived to be more difficult to raise a given amount of revenue. It is important to note that this scheme treated increased collections due to expansion of the tax base (new properties) or increased collection on the current base (higher recovery rates) symmetrically. Benchmarks were generated using a three-year average of historical collections, adjusted for the normal rate of increase in collections, and were designed such that most circles would be “in the money” and face linear incentives on the margin.⁴ Since most inspectors are rotated to new circles every two to three years, the use of two to four lags of revenue collection in determining benchmarks means that ratchet effects should not be a first-order concern, since by the time higher revenue collection starts to impact benchmarks substantially, the inspector would likely be in a different circle and not subject to those benchmarks.

As each tax circle staff consists of three members, the bonus was divided 40%-30%-30% among inspector, constable, and clerk, respectively. On net, with a 30% average incentive payment to the group and this division among the three group members, each individual inspector, constable, and clerk faced a roughly 10% individual marginal incentive. Payments for all incentive schemes were restricted to staff who were posted in the circle at

4. Specifically, in the first year (FY11–12), the benchmark was the three-year average of revenues from FY07–08, FY08–09, and FY09–10, plus 10%. Since the rate of increase in collections averaged about 8% a year, the benchmark should be approximately 13% below the average revenue under business as usual, so that almost all circles would be in-the-money and face linear incentives on the margin. The adjustment rate was increased slightly in year 2 in light of the growth rates observed in year 1, so that in the second year (FY12–13), the historical benchmark was the three-year average of revenues from FY08–09, FY09–10, and FY10–11, plus 20%. Note that in the first year, there were separate benchmarks for current-year tax collection and arrears collection, so that the formula was $Incentive_c = \alpha_c \max(CurrentYearRevenue_c - CurrentYearBenchmark_c, 0) + \alpha_c \max(ArrearsRevenue_c - ArrearsBenchmark_c, 0)$. Given that inspectors have some leeway in classifying revenue into current or arrears, but no flexibility in total revenue (since it must match the amount of money deposited into the bank), in the second year, incentives were simplified to be based simply on the total revenue.

the time of randomization, and staff were no longer eligible to receive payments if they were transferred to a nonincentivized circle.

2. *Revenue plus.* The revenue plus scheme was similar to the revenue-based scheme, but included additional incentives (the “plus” component) to help address the multitasking problem inherent in the tax collector’s job (Hölmstrom and Milgrom 1991). Specifically, in addition to maximizing revenue collected, the government also cares about how people feel they are treated by the tax department and whether taxes are assessed accurately.

To address these concerns, in addition to rewarding on revenue using the identical formula as in the revenue scheme, this scheme adjusted pay based on taxpayer satisfaction and accuracy of tax assessments. Circles in the scheme were ranked based on the accuracy and satisfaction measures and divided into three equal-sized groups. Circle staff were paid as in the revenue treatment, but the top group received an additional bonus equal to about 0.75 times their base salary, and the bottom group lost 0.75 times their average base salary.⁵ By design the total payments under the scheme could never be negative (that is, base salary was never at risk; an inspector in the bottom group might receive 0 from the scheme but would not forfeit his base salary); otherwise, (conditional on the same revenue increase) average payments would be identical between the revenue and revenue plus schemes.⁶

The satisfaction and assessment accuracy measures were based on an independent survey of 12,000 randomly sampled properties (described in Section V.A). Taxpayer satisfaction was measured based on two survey questions about the quality and

5. Inspectors in the top group received an extra Rs. 15,000 a month, and constables and clerks received an extra Rs. 11,500 a month; those in the bottom group lost an equivalent amount.

6. To the extent that tax officials are risk-neutral, or they are risk-averse but have CARA utility with cost of effort expressible in monetary terms, the additional variance induced by the plus scheme should not affect their choices other than through the multitasking channel. Relaxing the assumptions of risk-neutrality or CARA utility, however, could allow there to be direct effects from the increased variance due to the plus component on the return to revenue from the effort component.

results of interactions with the tax department.⁷ Accuracy was measured as 1 minus the absolute value of the difference between GARV as measured by the survey and the official GARV as measured from the tax department's administrative records, divided by the average of these two values.⁸ Since this "plus" component relies on third-party surveys and could also lead to losing the performance pay earned due to increased tax collections, it effectively constitutes an audit component (though was not referred to as such so as to maintain better optics).

3. *Flexible bonus.* The third scheme was designed to be analogous to the way bonuses work in the private sector for many complex jobs, such as those in Wall Street firms: managers distributed a fixed bonus pool to talented employees based on all factors (including subjective ones) they observe.⁹ In this treatment, staff were again divided into three groups and pay was determined by group (just as in the revenue plus scheme), but rather than have their pay determined by an ex ante specified formula, they were divided by their performance as ranked by a departmental Performance Evaluation Committee (PEC) composed of senior tax officials. Everyone in the treatment provisionally earned a base salary supplement roughly equal to their average salary.¹⁰ At the end of the year, adjustments were

7. The questions were: "In your opinion, what has been the overall quality of service offered by this department to this property?" and "In your personal dealings with members of this department, how satisfied are you with the outcomes?" Each question was answered on a 1–5 scale.

8. In the first year, this measure was noisier due to survey and measurement logistics that were resolved by the second year. Therefore in the first year we instead calculated accuracy by correlating log GARV in the official register with log GARV according to the survey, which was more robust to being off by a constant.

9. For example, managers might be able to observe effort in addition to outcomes; they also might have information that certain areas were more difficult. While subjective assessments can potentially better match the complexities of real jobs, they can be less effective than formulaic systems if workers do not trust the managers to implement them properly, if managers play favorites, or if managers and workers disagree about the subjective component of performance (Baker, Gibbons, and Murphy 1994; Prendergast and Topel 1996; Prendergast 1999; MacLeod 2003).

10. In the first year, the base salary supplement was Rs. 30,000 for inspectors and Rs. 23,000 for constables and clerks, that is, closer to 1.5 times base pay. This figure was adjusted in the second year to Rs. 22,000 and Rs. 16,500 so that the three schemes generated roughly equal average honorariums.

made just as in the revenue plus scheme: the top third of circles received an additional bonus equal to approximately 0.75 times their base salary, and the bottom group lost approximately 0.75 times their average base salary.

In determining payments under this scheme, the PEC was allowed to use any criteria it chose, as long as it could document a reason behind them, and the committee was provided all of the same information used in the revenue plus treatment (increase in revenue over benchmarks, customer satisfaction, and accuracy of assessments). The main differences between the flexible bonus and revenue plus schemes were that the objective revenue-based formula was replaced by a fixed increase in base salary (with an end of year bonus), and that the grouping was made by the PEC as it saw fit with relatively few restrictions.

Although the official rules allowed the PEC full flexibility in using subjective criteria, they in fact created a (richer) formula for ranking circles, using the following indicators and weights (in parentheses): increase in revenue collected (40%), increase in tax base (25%), accuracy of assessment (15%), subjective director's rating (10%), and customer satisfaction (10%). This was publicized about six months after the intervention began, so by the beginning of year 2, inspectors should have been fully aware of the assessment criteria. The two additional criteria (compared to revenue plus) were tax base increases and the subjective director's assessment. On net, the correlation between the PEC ranking and the ranking of payments that would have been generated under the revenue plus formula was 0.269 in year 2.

4. Additional treatments. Two additional treatments were introduced in the second year of the program (FY12–13). The “information-only treatment” (70 circles) was intended to capture the part of the effect that arises from all aspects of treatment besides the monetary incentives. Staff from these circles went through the same process as in the revenue treatment (including receiving quarterly reports on collections above their historically predicted benchmarks and attending progress quarterly meetings), but with no corresponding incentive payments. While the quarterly reports just repackaged information that staff already had, the reports presented the information in a more systematic format, which may have increased its salience. Furthermore, the act of attending the quarterly meetings may have led circle staff

to believe that they were being monitored more carefully. The information-only scheme nets out these effects from the direct impact of the payments per se in the performance-based incentives.

In addition, a supervisor's performance pay scheme was introduced in the second year. This was identical to the revenue scheme, but applied to assistant excise and taxation officers (AETOs), who supervise the circle staff, and excise and taxation officers (ETOs), who supervise AETOs. Randomization was done by ETO, with 26 treatments and 25 controls. All AETOs working under selected ETOs were included. Payments were calculated based on the average increase in revenue over benchmarks for circles under their supervision. The bonus rate was determined by average circle size, and each supervisor received a 50% share of all imputed bonus payments (recall an inspector's share was 40%). Since this intervention was randomized by ETO (of which there are only 51), whereas the circle-level intervention was randomized by circle (with almost 500 circles), this intervention will have substantially lower statistical power than the main circle-level treatments.

5. Knowledge and credibility. To ensure that collectors understood the specifics of the scheme they were in, we carried out detailed trainings at the start of the year, followed by quizzes and refresher trainings throughout. By seven months after treatments started, quiz results revealed that virtually all inspectors understood the scheme and could accurately calculate payments. A survey of all inspectors (treatment and control) confirmed that inspectors could identify whether they would receive payments, and which scheme they were in. To ensure that inspectors believed that payments would actually be made, the project was officially approved by the Chief Minister (the highest political authority in the province). A small pilot was conducted (and payments made) in 11 circles for an entire year before the main experiment began, and payments were made quarterly throughout the main experiment.

IV.B. Randomization Design and Balance Checks

The randomization was carried out through computerized public lotteries, with representatives from the tax department present. This helped minimize any perceived bias, especially

since the performance pay schemes were popular (most staff wanted to opt in). To reduce any concerns about differential selection across the schemes while maintaining informed consent, the lottery was conducted in two stages. In the first stage, circles were selected to participate in the project and staff consent was sought. Staff were told about the three possible schemes, and were told that a second lottery would determine which scheme they would be assigned to.¹¹ Once consent was obtained, a second lottery assigned consented circles to particular incentive schemes. Over 95% of circle staff selected in the first lottery consented. Given the high consent rates in the first year, both stages were conducted in a single lottery in year 2. The lotteries were held as close as possible after the start of the fiscal year on July 1.

Table I shows the experimental design. In year 1, 160 circles were selected in the first ballot, to be divided equally into one of three treatments. In year 2, an additional 58 were selected and divided into the same three treatments. The circles selected in year 1 remained in the same treatment, and new inspectors who had previously transferred into these circles became eligible for performance pay in year 2.¹² In addition, 70 circles were selected for the information-only treatment. Each lottery was stratified with 19 strata based on the 11 administrative divisions of the province and—for all but the smallest few divisions—circle size.

Online Appendix Table 11 compares the selected circles to controls on baseline characteristics in the administrative data based on the randomization at the end of year 2. Out of the 42 comparisons (7 variables * 6 columns), only 1 is significant at the 5% level (the coefficient on log nonexemption rate for the flexible bonus scheme).¹³

11. Given the crucial role played by the inspector in collecting tax, the circle as a whole could only participate if the inspector consented. Constables or clerks could individually opt out of the scheme, though this rarely happened.

12. Since this was not part of the policy initially (we had made clear that anyone transferring in during the year would not be part of the treatment) there is not much concern that staff were strategically transferring in the hope that they would be eligible in the second year.

13. Looking scheme by scheme, the joint balance test shows statistical significance in one of the schemes (revenue plus) compared with pure controls, even though none of the individual covariates are statistically significantly different. In the Online Appendix tables we show that the main average effects of incentives do not seem to be driven by this one subtreatment (Appendix J), and that controlling for the variables included in the balance table does not meaningfully change the results (Appendix K).

TABLE I
EXPERIMENTAL DESIGN

| | Randomization | | Implementation | |
|----------------|---------------|--------|----------------|--------|
| | Year 1 | Year 2 | Year 1 | Year 2 |
| Revenue | 53 | 72 | 47 | 68 |
| Revenue plus | 54 | 74 | 48 | 68 |
| Flexible bonus | 54 | 73 | 49 | 67 |
| Information | 0 | 70 | 0 | 66 |
| Control | 322 | 194 | 338 | 213 |

Notes. The first two columns (under Randomization) show the number of circles that were assigned to each of the three (or four) treatment types in each year. In cases where staff did not consent to treatment after the first ballot (in year 1), circles were assigned treatment values of 1/3 for each main treatment type (revenue, revenue plus, and flexible bonus). Values are rounded. The second two columns (under Implementation) show the number of circles that were actually implementing the treatment at the end of the fiscal year. Treatment wasn't implemented either because of lack of consent or because the initially selected circle staff were transferred to new posts. See text for more details.

V. DATA AND EMPIRICAL METHODOLOGY

V.A. Data

We use two main sources of data: circle-level administrative data for our main measures of tax performance, and property/taxpayer-level data based on a survey we conducted to obtain measures of accuracy of tax assessment, customer satisfaction, and corruption. Online Appendix B provides further details on both data sets, including the additional verification checks we ran on the administrative data, how we addressed circles with boundary changes, details of the survey, and variable definitions. Here we highlight a few of these aspects.

The administrative data are based on the quarterly reports that each inspector files, which show their overall collections (separately for current year and past years/arrears collections) and the total assessed tax base. We digitized these reports for all tax circles and selected a random sample to be verified each year by aggregating (thousands of) bank-verified receipts of individual payments. We found no statistically or economically significant discrepancy between the administrative data and our independent verifications.

Summary statistics for key variables from the administrative data are shown in Panel A of Table II for the second year of the experiment (FY12–13); summary statistics for additional years and variables can be found in the Online Appendix. Several observations are worth noting. First, current year revenues are

TABLE II
SUMMARY STATISTICS

| | Mean | Std. dev. | N |
|---|--------|-----------|--------|
| Panel A: Administrative data | | | |
| Log revenue (total) | 15.75 | 0.74 | 482 |
| Log revenue (current) | 15.52 | 0.73 | 482 |
| Log revenue (arrears) | 13.91 | 1.17 | 479 |
| Log tax base (total) | 16.14 | 0.80 | 482 |
| Log tax base (current) | 15.86 | 0.73 | 482 |
| Log tax base (arrears) | 14.40 | 1.37 | 479 |
| Log nonexemption rate (total) | -0.23 | 0.20 | 482 |
| Log nonexemption rate (current) | -0.19 | 0.13 | 482 |
| Log nonexemption rate (arrears) | -0.30 | 0.41 | 479 |
| Log recovery rate (total) | -0.16 | 0.18 | 482 |
| Log recovery rate (current) | -0.14 | 0.14 | 482 |
| Log recovery rate (arrears) | -0.19 | 0.29 | 479 |
| Panel B: Survey Data | | | |
| Property successfully found in administrative records (dummy) | 0.84 | 0.37 | 11,971 |
| Quality of tax department (0–1) | 0.53 | 0.22 | 6,050 |
| Satisfaction with tax department (0–1) | 0.55 | 0.23 | 6,050 |
| Inaccuracy | 0.34 | 0.27 | 9,870 |
| Tax gap | -0.099 | 0.42 | 9,870 |
| GARV | 32,302 | 252,426 | 10,787 |
| Self-reported tax payment in FY 2013 | 3,562 | 18,604 | 12,000 |
| Bribe payment | 2,073 | 3,932 | 5,993 |
| Frequency of bribe payment | 0.76 | 0.88 | 4,802 |

Notes. Panel A statistics from administrative data are shown at the end of year 2 of the study (FY 2012–2013). Each observation is one of the 482 circles as defined at the time of randomization. Panel B statistics from the property survey are for properties from the random sample drawn from the field. The inaccuracy and tax gap measures are available for only those properties that could be matched to the administrative records. Subjective variables—quality, satisfaction, bribe payment, and frequency of bribe payment—are reported for circles from the first phase of the survey only (see Online Appendix B for more details).

substantially larger than arrears (i.e., collections against past years' unpaid taxes)—the mean of log current revenues is 15.52 compared with just 13.91 for log arrears, implying that on average, current revenue is about five times as large as arrears. This suggests that the main impacts on total revenue will likely be felt through increases in current year revenue. Second, there is much more variation in arrears—the standard deviation in log arrears is about 1.5 times that of log current revenue—implying that detecting effects on arrears statistically will be more difficult. The log recovery rate (the log of tax revenue divided by the tax base net of exemptions) is -0.14 for current year taxes, which implies

that about 85% of all taxes that are demanded by the government are in fact paid. Thus while nonpayment is a substantial issue (a typical developed country government would not be satisfied with a 15% nonpayment rate of property taxes),¹⁴ it is still the case that the bulk of taxpayers do in fact pay the tax bills they receive. Thus any potential evasion may come from underassessment of properties (as we will see later) rather than just flagrant disregard of issued tax notices.

The second primary data source is the property survey we conducted at the end of the two-year period. This survey provides our main nonrevenue outcomes (taxpayer satisfaction measures and tax assessment accuracy), as well as owner/property characteristics that help examine heterogeneous effects. The survey is based on two distinct samples. The first, which we refer to as the “general population sample,” consists of roughly 12,000 properties selected by randomly sampling five GPS coordinates in each circle and then surveying a total of five properties around that coordinate. These properties therefore represent the picture for the typical property in a tax circle. The second sample, which we refer to as the “reassessed sample,” consists of slightly more than 4,000 properties (roughly 10 per circle) which were sampled from an administrative list of properties that are newly assessed or reassessed. These properties were then located in the field and surveyed. This oversamples the (few) properties that experience such changes each year so we can examine the impacts on such properties separately.

Panel B of Table II presents summary statistics for properties from the general population sample. Several facts are worth noting. First, 84% of properties we randomly sampled in the field were successfully located on the tax registers. Although there are a substantial number of untaxed properties, it is not the case that only a few properties are on the tax rolls. Second, conditional on being on the tax rolls, on average properties appear undertaxed. We focus on the GARV of the property, which is the main measure of a property’s tax value, before exemptions and reductions are

14. For example, on average from 2010 to 2013, the city of Cambridge, Massachusetts, collected almost exactly 100% of all property taxes due (City of Cambridge 2014). The collection rate in Pakistan is more comparable to Detroit, Michigan, which had an 80% average property tax collection rate from 2010 to 2013 and filed for bankruptcy (City of Detroit, Michigan 2013).

applied.¹⁵ To measure under- or overtaxation, we focus on the “tax gap,” defined as

$$(8) \quad TaxGap = \frac{GARV_{Inspector} - GARV_{Survey}}{(GARV_{Inspector} + GARV_{Survey})}.$$

This captures the difference between what the inspector officially reported and what was obtained through our own survey. Our measure of inaccuracy is the absolute value of the tax gap. On average, inaccuracy is 0.34, indicating substantial disagreement between the two measures. The tax gap has a mean value of -0.10 , suggesting that undertaxation is prevalent in our population.¹⁶

Corruption also appears to be prevalent. On average, respondents report that annual bribes paid for a property similar to theirs are around Rs. 2,000 (US\$20)—about half of the amount they report paying in property taxes. Bribes are frequent—when asked how many times a typical property owner would need to bribe the property tax department, the mean is 0.76 bribes paid per year. On the other hand, respondents are not wildly unsatisfied with service from the tax department—on a 0–1 scale, the average response is 0.53 for quality of service and 0.55 for satisfaction.¹⁷ Of course, this could be consistent with corruption: a respondent might be “satisfied” if he was able to reduce his official tax liability by paying a bribe.

In addition to these two primary sources of data, in some of the appendix tables we also make use of a short phone-based

15. We focus on GARV, rather than tax assessed, because nonlinearities in the tax formula mean that there is substantially more measurement error in tax assessed than in GARV. For example, if the land area is less than 5 marla (1,361 square feet), nonrented, residential properties are completely exempt from tax. By contrast, GARV is a continuous function of the underlying property characteristics and hence is much more robust to measurement error.

16. Given the way it is normalized, an average tax gap of -0.10 means that on average, the inspector’s assessment is 19% less than the survey’s estimate.

17. One might be concerned that the quality and satisfaction variables are simply picking up noise. However, Panel A of Online Appendix Table 12 shows that the satisfaction and quality measures are internally consistent: that is, households who report higher satisfaction report higher quality of service, households that report higher quality report lower bribes, and so on. More important, households in a circle tend to agree with each other. Panel B of Online Appendix Table 12 regresses these measures on what other respondents in the same circle report: people report high satisfaction when others in their neighborhood report high satisfaction, report high bribes when others report high bribes, and so on.

survey of inspectors in which we gathered basic information about the self-reported effort and perceived supervisory support and pressure felt by the tax inspectors.

V.B. Empirical Methodology

Since we are evaluating a randomized experiment, the empirical methodology is straightforward. We estimate 2SLS regressions, where the endogenous variable is the treatment status at any point in time and the instruments are the results of the lottery.¹⁸ Our primary specification for assessing circle-level outcomes using the administrative data is

$$(9) \quad \ln Y_{cst} = \alpha_s + \beta \text{Treatment}_{cst} + \gamma \ln Y_{cs0} + \epsilon_{cst},$$

where Y_{cst} is the outcome of interest for circle c in stratum s at time t , and Treatment_{cst} is a continuous variable that takes values from 0 to 1 that represents the fraction of treated circle staff present in circle c in the last quarter of the given fiscal year. Y_{cs0} is the value of the outcome variable at baseline (i.e., in the fiscal year prior to randomization). Treatment is instrumented by a binary variable that represents the circle's randomization status into any one of the three incentive schemes.¹⁹ We include stratum fixed effects (α_s) given the lottery was stratified by these strata. All regressions based on administrative data are run using circle boundaries that existed at the time of randomization. We report robust standard errors clustered at the level of the robust partition of circles, that is, the maximum set of circles that have been involved together in a set of splits and merges since randomization.

18. The reason the treatment status is not exactly equal to the lottery results is that a small number of circles (8 out of 482) did not consent to participate, and because some circle staff lost eligibility to continue in the scheme after they were transferred out to another circle.

19. Reduced-form versions of the main table can be found as Online Appendix Tables 3-G1 and 4-G1. Note also that the information-only scheme is not included as a treatment, but is instead included as part of the control group to maximize statistical power. Online Appendix Tables 3-I through 7-I reestimate the tables in the article, where instead the information treatment is separated out, so performance pay treatments are compared only to pure controls, with qualitatively similar results.

To estimate the impact of the separate subtreatments, we estimate the analogous regression separately by treatment:²⁰

$$(10) \quad \ln Y_{cst} = \alpha_s + \beta_1 \text{Revenue}_{cst} + \beta_2 \text{RevenuePlus}_{cst} + \beta_3 \text{FlexibleBonus}_{cst} + \gamma \ln Y_{cs0} + \epsilon_{cst}.$$

For survey-based outcomes, we run regressions at the property level. When examining the general population sample, we run regressions of the form:

$$(11) \quad Y_{ics} = \alpha_s + \beta \text{Treatment}_{cs} + \epsilon_{ics},$$

where i is an individual property.²¹ As before, we instrument for *Treatment* with the randomization results. We include stratum fixed effects and cluster standard errors at the circle level.²² When available, we include controls for baseline level outcome variables.

For regressions where we are interested in the difference between reassessed and new properties and regular properties, we include both the general sample and the reassessed sample

20. In such regressions, in addition to reporting β_1 , β_2 , and β_3 , we report several other statistics to guide the analysis. In particular, we report the p -values for a test of the joint statistical significance of the incentive schemes (i.e., a test of the null that $\beta_1 = \beta_2 = \beta_3 = 0$) and a test that the three schemes are identical (i.e., a test of the null that $\beta_1 = \beta_2 = \beta_3$). We also report p -values from a test of whether the schemes that dealt with multitasking are identical to those that did not (i.e., a test of the null that $\beta_1 = \frac{\beta_2 + \beta_3}{2}$), and from a test of whether the scheme that used subjective information from the department is identical to the formulaic schemes (i.e., a test of the null that $\beta_3 = \frac{\beta_1 + \beta_2}{2}$).

21. Our sampling strategy was to randomly draw five initial GPS coordinates from within the boundary of a tax circle. We then survey the property closest to that point and then following a left-hand rule (or if that is not possible, a right-hand one) survey an additional four properties. A potential concern is that we may be over-sampling larger properties since a randomly chosen GPS point is more likely to fall inside a larger property. Although this may be true for the first sampled point, we have confirmed that it is not true of subsequent properties, that is, there is very little correlation between the land area of the first property (chosen by GPS point) and the subsequent properties (chosen by moving to the left). As a robustness exercise we therefore redo our estimates after dropping the first sampled point and using only the remaining points, and find that our results are qualitatively similar. See Online Appendix Tables 6-L and 8-L.

22. Regressions based on survey data are run using circles boundaries when the sample of properties was drawn, which happened in the middle of the second fiscal year of the study.

(which includes newly assessed properties and those whose valuations changed), and then estimate:

$$(12) Y_{ics} = \alpha_c + \beta_1 \text{Treatment}_c * \text{ReAssessed}_{ic} + \beta_2 \text{ReAssessed}_{ic} + \epsilon_{ic},$$

where *ReAssessed* is a dummy that is 1 if a property was sampled from the list of properties whose valuation was changed (we do not distinguish in this regression between properties whose tax valuation was changed and newly assessed properties; both are captured by *ReAssessed*). Note that unlike equation (11), we now include circle fixed effects (α_c) to capture fixed differences among circles between properties. We examine the analogue of equation (10) when we examine sub-treatments.

In interpreting equation (12), it is important to note that which properties are reassessed is potentially an outcome of the treatment. As such, the coefficient β_1 includes two margins of treatment effects—an extensive margin effect (i.e., the type/number of properties revalued can be impacted) and an intensive margin effect (a given reassessed property may now be dealt with differently). For example, if Y_{ics} is the amount of bribes paid, the coefficient β_1 in equation (12) shows how the difference in bribes paid between reassessed and nonreassessed properties changes in treatment versus control circles. As outlined in the conceptual framework, this net effect β_1 will include both margins (i.e., (i) the average bribe amount changes as the set/type of people who collude changes and (ii) conditional on collusion, the bribe amount changes). To shed some light on these effects, in Section VI we also examine how the composition of those in the reassessed sample changes by estimating equation (12) on fixed characteristics of reassessed properties.

VI. RESULTS

In Section VI.A, we examine the impacts of the performance pay schemes on the key revenue and nonrevenue outcomes of interest. Section VI.B then probes the mechanisms through which changes in tax base occur in light of the model in Section III. Although we focus on the pay-for-performance aspect of the schemes (i.e., price effects), Section VI.C considers a variety of alternative explanations for the results, such as perceptions of additional monitoring, income effects, and interactions with

supervisors. Section VI.D concludes with a discussion of cost-effectiveness.

VI.A. *Main Impacts*

1. *Impacts on Revenue Outcomes.* Table III considers the impact of the performance pay schemes on (log) revenue at the end of each of the two years of the study. We first consider the impact on total revenue (columns (1) and (4)). The remaining columns break this down into revenue derived from current year taxes and revenue from arrears. Current year revenue is about five times larger than arrears revenue. Arrears revenue is also substantially more variable over time, which is why the standard errors are larger. Panel A reports the impact where we pool all three performance pay schemes, and Panel B shows the impact for the schemes separately.

We find substantial impacts of performance pay on total revenue collected. Panel A, column (1) shows that compared to controls, revenue increased by 9.1 log points in treatment circles in the first year, and column (4) shows an increase of 9.4 log points in the second year. To interpret the magnitude, note that on average, control circles experienced an increase in total revenue of about 25 log points between the baseline year and the end of the second year. Exponentiating, this implies that control circles grew by about 28% over the two years, and treatment circles grew by about 41%. Incentives thus led to a 13 percentage point increase in the growth rate, or a 46% higher rate of growth, over the two years of the experiment.

Examining the effects separately by current and arrears revenue, we find that the impact on current year revenue collection is 7.3 log points in year 1 and 9.1 log points in year 2. In contrast, there is a 15.2 log point increase in arrears revenue in year 1, which falls to 11.3 log points (and is no longer statistically significant) in year 2. Although these changes over the years are not statistically distinguishable, the point estimates suggest that inspectors, who exhausted much of the available pools of easily collectable arrears in the first year, switched their focus to increasing current year collection in the second year.

Separating the results by the three compensation schemes (Panel B), we see that as one might expect, schemes that directly reward on revenue collection have a larger impact on revenue collected. Looking at current year revenue (where we have

TABLE III
IMPACTS ON REVENUE COLLECTED

| | (1) | (2) | (3) | (4) | (5) | (6) |
|-----------------------------------|---------------------|---------------------|--------------------|---------------------|---------------------|------------------|
| | Year 1 | | | Year 2 | | |
| | Total | Current | Arrears | Total | Current | Arrears |
| Panel A: Main treatment | | | | | | |
| Any treatment | 0.091*** (0.028) | 0.073*** (0.027) | 0.152** (0.069) | 0.094*** (0.031) | 0.091*** (0.032) | 0.113 (0.083) |
| Panel B: Subtreatments | | | | | | |
| Revenue | 0.118*** (0.035) | 0.109*** (0.034) | 0.134 (0.099) | 0.129*** (0.043) | 0.152*** (0.044) | 0.005 (0.133) |
| Revenue plus | 0.080 (0.053) | 0.086* (0.052) | 0.072 (0.110) | 0.093** (0.045) | 0.081* (0.049) | 0.175 (0.114) |
| Flexible bonus | 0.071* (0.038) | 0.024 (0.035) | 0.243** (0.098) | 0.056 (0.041) | 0.035 (0.042) | 0.148 (0.108) |
| <i>N</i> | 481 | 481 | 481 | 482 | 482 | 479 |
| Mean of control group | 15.671 | 15.379 | 14.030 | 15.745 | 15.518 | 13.915 |
| Rev. vs. multitasking <i>p</i> | 0.323 | 0.193 | 0.830 | 0.233 | 0.049 | 0.262 |
| Objective vs. subjective <i>p</i> | 0.530 | 0.090 | 0.212 | 0.220 | 0.084 | 0.634 |
| Equality of schemes <i>p</i> | 0.562 | 0.143 | 0.433 | 0.359 | 0.086 | 0.527 |
| Joint significance <i>p</i> | 0.004 | 0.010 | 0.073 | 0.012 | 0.005 | 0.305 |

Notes. This table presents results on the impact of the performance pay schemes on revenue-based outcomes. We use instrumental variables regressions, where treatment status is instrumented with randomization results. The unit of observation is a circle, as defined at the time of randomization. Outcome variable is log revenue collection as of the end of the fiscal year, for total revenue (columns (1) and (4)), current year revenue (columns (2) and (5)), and collections against arrears (columns (3) and (6)). Specification follows equation (10) of the main text, and includes stratum fixed effects. "Any treatment" in Panel A includes the three subtreatments in Panel B. The Information treatment is included in the control group. We report *p*-values from tests of equality of coefficients as follows: rev. vs. multitasking tests for equality between revenue and the average of revenue plus and flexible bonus; objective vs. subjective tests for equality of the average of revenue and revenue plus against flexible bonus; equality of schemes tests whether all coefficients are equal; and joint significance tests joint null that all coefficients are equal to 0. Robust standard errors in parentheses. Standard errors are clustered by a robust partition of circles, that is, the group of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < .01$.

much more precise estimates for the aforementioned reasons), Column (5) shows that by the end of year 2, revenue circles collected 15.2 log points more revenue than control circles, compared to an 8.1 log point increase in revenue plus circles and a 3.5 log point increase in flexible bonus circles. We can reject equality of these coefficients at the 10% level. When we test for equality between revenue and an average of the multitasking schemes we are also able to reject equality (p -value $< .05$). The magnitudes for the revenue scheme are large: compared to the 39% average growth in current year revenue in control areas, revenue in revenue circles grew by 62%. This implies that revenue circles had a 58% (23 percentage point) higher growth rate in current revenue

over two years than did controls. The impact on total revenue collection—including arrears—was substantial as well: revenue circles had 62% higher growth than did controls.

Our results show that performance pay schemes did lead to large increases in revenue, with schemes that rewarded more on revenue collected seeing even larger increases. Although our data verification checks gives us confidence that these schemes did in fact bring in real money, one potential concern is that these impacts might be due to temporary (and unreasonable) pressures put on taxpayers that could ultimately be undone through appeals (see, e.g., Das-Gupta and Mookherjee 1998). To investigate this we randomly sampled 22 circles, one incentive and one control in each of the 11 divisions, at the end of the second year of the experiment and investigated all appeals that had been filed to date since the start of the experiment. We find that appeals are much too small (at most 1.5% of annual total revenues) to substantially change the results here, and we find no economically meaningful or statistically significant differences in appeals rates or amounts between treatment and control areas.

2. Impacts on Nonrevenue Outcomes. To the extent that high-powered incentives lead to excessive pressure to collect taxes and/or overtaxation/extortion, one may be concerned that the performance pay schemes—especially the revenue scheme—could adversely impact taxpayer satisfaction and assessment accuracy. Table IV investigates these issues and shows little evidence for such effects.

We examine the impact on measures of taxpayer satisfaction and accuracy of tax assessment using property-level survey data. Columns (1) and (2) in Table IV examine the two measures of taxpayer satisfaction in which we asked the respondent how they rated the “quality of service” of the tax department and how “satisfied” they were with their service. These are the exact measures that were incentivized in the revenue plus scheme, so it is instructive to examine not just whether they worsen in the incentive treatments in general, but whether the revenue plus scheme, and perhaps the flexible bonus scheme, mitigates this effect.

Panel A shows no statistically or economically meaningful treatment effect for either measure. In particular, on a 0–1 scale, the point estimates are -0.006 for quality of service and

TABLE IV
IMPACTS ON NONREVENUE OUTCOMES

| | (1) Quality | (2) Satisfaction | (3) Inaccuracy | (4) Tax gap |
|-----------------------------------|--------------------|---------------------|-------------------|-------------------|
| Panel A: Main treatment | | | | |
| Any treatment | -0.006 (0.022) | -0.011 (0.022) | 0.004 (0.012) | 0.007 (0.022) |
| Panel B: Subtreatments | | | | |
| Revenue | 0.006 (0.036) | -0.006 (0.037) | 0.002 (0.017) | -0.022 (0.029) |
| Revenue plus | 0.040 (0.026) | 0.029 (0.027) | 0.028* (0.016) | 0.015 (0.032) |
| Flexible bonus | -0.060* (0.031) | -0.053* (0.032) | -0.016 (0.018) | 0.029 (0.031) |
| <i>N</i> | 6050 | 6050 | 9870 | 9870 |
| Sample | Phase 1 | Phase 1 | Full | Full |
| Mean of control group | 0.538 | 0.555 | 0.339 | -0.103 |
| Rev. vs. multitasking <i>p</i> | 0.683 | 0.876 | 0.813 | 0.159 |
| Objective vs. subjective <i>p</i> | 0.015 | 0.064 | 0.099 | 0.315 |
| Equality of schemes <i>p</i> | 0.014 | 0.059 | 0.090 | 0.344 |
| Joint significance <i>p</i> | 0.035 | 0.129 | 0.160 | 0.533 |

Notes. This table presents results on the impact of the performance pay schemes on nonrevenue outcomes. We use instrumental variables regressions, where treatment status is instrumented with randomization results. Unit of observation is a property. Specification follows equation (12) of the main text and includes stratum fixed effects. Quality and Satisfaction were measured on a five-point Likert scale and rescaled to a [0,1] interval. Tax gap is the difference in the official gross annual rental value (GARV) minus our estimated GARV, divided by the sum of these. Tax gap measures over-/undertaxation, with positive coefficients indicating overtaxation. Inaccuracy is the absolute value of tax gap. Sample is restricted to phase 1 of the survey for subjective outcomes (quality and satisfaction). The Information treatment is included in the control group. We report *p*-values from tests of equality of coefficients as follows: rev. vs. multitasking tests for equality between revenue and the average of revenue plus and flexible bonus; objective vs. subjective tests for equality of the average of revenue and revenue plus against flexible bonus; equality of schemes tests whether all coefficients are equal; and joint significance tests joint null that all coefficients are equal to 0. Standard errors are clustered by robust partition of circles, that is, the group of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < .01$.

-0.011 for satisfaction, and we can reject a change in either measure of about 0.04 or larger.

Panel B examines the impacts separately for each scheme and finds the estimates for the flexible bonus are negative (-0.060 and -0.053 for quality and satisfaction, respectively), whereas the point estimates for revenue plus are positive (0.040 and 0.029, respectively). Although the results for each scheme are generally not statistically significant, one can reject the null hypothesis of equality of the three schemes, or the null that that the flexible bonus scheme is equal to the other schemes. The

estimates suggest that the revenue plus treatment, which explicitly incentivized quality and satisfaction, may have in fact led to higher levels of both compared to the revenue and flexible bonus incentive schemes, though the magnitude of this impact is relatively small. The flexible bonus not only had the lowest performance in terms of revenue raised for the government but also had worse outcomes on these other dimensions.

The zero average results on quality and satisfaction are quite robust. In particular, we show in Online Appendix Table 4-G2 that the results are qualitatively unchanged if we use ordered probit models instead of the linearized variable with OLS or control for observable property characteristics (area, usage, etc.).

In addition to these satisfaction measures, we examined other metrics that may reflect general attitudes toward the government, such as quality and satisfaction with other departments and stated preference for the incumbent party (based on self-reported voting behavior). These are shown in Online Appendix Table 13. In general, none of these metrics show meaningful differences between treatment and control. The only notable difference is that the pattern that revenue plus areas show higher satisfaction and quality of service appears generalized to other departments beyond just tax, suggesting that there may be positive spillovers, which is consistent with citizens attributing a positive interaction in one government service to other related services.

Columns (3) and (4) in Table IV examine the second main nonrevenue dimension, the inaccuracy of tax assessment of the property. The results show no changes in inaccuracy or the tax gap overall (Panel A). When we explore the subtreatments (Panel B), we do get some indication that revenue plus may have increased overall inaccuracy, although this does not seem to have an impact on the tax gap, which suggests that it may have raised both under- and overtaxation for the full sample of properties. It is important to note, however, that this is the average effect for all properties. One potential reason we may not detect changes in this metric is that the number of properties affected may be small; we explore this in more detail when we focus on reassessed properties later.

On net, there are two key conclusions from the results thus far. First, compared to control circles, we find that the incentives overall have a substantial, positive effect on revenue, with little detectable downside in terms of taxpayer satisfaction and the

accuracy of tax perceptions for the typical property. Second, performance pay schemes with clearly defined objective criteria and with fixed proportional incentives tend to do better than more subjective, potentially uncertain, and multidimensional schemes. Comparing the revenue and revenue plus scheme, we find that by year 2 the revenue scheme had increased revenue by about 13 log points, whereas the revenue plus scheme increased current revenue by only about 9 log points; on the other hand, customer satisfaction appears slightly higher in the revenue plus scheme. The flexible bonus scheme did worse than either revenue or revenue plus on all dimensions measured here. This provides suggestive evidence against subjective, potentially uncertain, and more multidimensional assessments and in favor of clearer, predictable, formula-based assessments that consider fewer metrics. This may be especially so in contexts where there may be concerns about credibility and how the more complex, subjective, and flexible assessments may be applied (see Baker, Gibbons, and Murphy 1994; Prendergast and Topel 1996; Prendergast 1999; MacLeod 2003 for related theoretical work on subjective bonuses).

VI.B. Changes in Tax Assessments and Rent Sharing

The model in Section III illustrates how taxpayers and tax collectors may collude to not pay taxes. The model shows how performance pay can make collusion harder and lead to higher tax collection and a switch from the collusive (high bribe, low tax) equilibrium to a noncollusive (low bribe, high tax) one, which could explain the increase in tax revenues. But the model also suggests that other taxpayers, who remain in the collusive equilibrium, would instead have to pay higher bribes to compensate tax inspectors for their forgone incentive pay. This section explores these issues in more detail.

1. How Many Properties Have Valuation Changed? If bargaining breaks down, the theory suggests this should result in a change in the official tax valuation recorded by the government, τ . To explore changes in τ in the data, we examine impact on the number and composition of properties whose official tax valuation was changed. We refer to these properties as “reassessed,” which includes both properties added to the official tax rolls for the first

time as well as previously taxed properties whose tax valuation is updated.

Table V shows the total number of reassessed properties, broken down by properties reported as assessed for the first time and those who had previously been on the tax rolls but whose valuation was updated. We obtained these data directly from the underlying tax registers. We control for the number of new and reassessed properties added in the baseline year (i.e., FY10–11) to capture heterogeneity across circles in their underlying rate of change of properties.²³

The results show a substantial increase in the number of properties whose valuations were changed in response to the treatment. On average (over the two year treatment period), there are 83 more properties per circle with new or updated valuations in treatment tax circles compared to controls, about an 86% increase over the control group. Most of this increase comes from properties that are newly reported. Column (2) shows that treatment circles add about 74 more newly valued properties to the tax rolls than did controls (202% increase over the control group), whereas an additional 9 properties see their valuations updated. Note, however, that most of these properties are not actually new—53% of these newly assessed properties were built before 2011 and a third constructed prior to 2006 (i.e., more than five years prior to the start of our experiment). In our field visits accompanying tax collectors, it was clear that they made visits to their tax circles frequently and were aware of where properties were located and their status (around two-thirds of the supposedly new properties were within 500 meters of a property that reported having been visited by the tax collector). It therefore seems more likely that the tax inspector was aware of these properties and they were strategically added to the rolls once performance pay incentives were introduced.

Although these numbers document a substantial increase in activity in treatment circles compared to control circles in percentage terms, it is worth noting that the absolute numbers are still relatively small compared to the total number of properties in

23. Note that since obtaining this data required a separate, detailed count of a different set of administrative records, we have this data only for a randomly sampled set (approximately 50%) of circles.

TABLE V
IMPACTS ON NUMBER OF REASSESSED PROPERTIES

| | (1) | (2) | (3) |
|-----------------------|--|--|---|
| | Total number of section 9 properties added to tax rolls in treatment period | Number of new properties added to tax rolls in treatment period | Number of reassessed properties added to tax rolls in treatment period |
| Treatment | 83.0* (45.27) | 74.0** (34.39) | 9.0 (22.35) |
| <i>N</i> | 234 | 234 | 234 |
| Mean of control group | 96.7 | 36.7 | 60.0 |

Notes. This table presents results on the impact of performance pay schemes on the number of properties that experience a change in tax status. Column (1) presents treatment effects on the total number of such properties added. The next two columns disaggregate this effect by whether the property is reported to have been previously registered on the tax rolls (column (3)) or not (column (2)). The sample consists of circles surveyed in phase 2 (see text for details). Specification includes stratum fixed effects and controls for number of new and reassessed properties added in the pretreatment (FY 2011) fiscal year. Standard errors are clustered by robust partition, the partition of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < .01$.

the circle: 74 new properties represents about 3% of the average number of taxable properties in the circle.²⁴ Nevertheless, as we explain shortly, these changing tax valuations are sufficient to explain essentially all of the change in revenue collected from the experiment. Increases in tax revenue can come through several margins: increasing the tax base by adding new properties to the tax rolls or updating their valuations, reducing exemptions granted to move tax bills closer to the gross tax base, or increasing the recovery rate of issued tax bills. Using the administrative data, we can decompose the increases in tax revenue into these three components. Doing so reveals that the vast majority of the increase in tax revenue is attributable to an increase in the tax base (see Online Appendix D), which is the type of reassessment documented in this section.

24. Note that these reassessments are not necessarily exhausting a fixed supply of new or modified properties. Our property survey indicates that 1%–2% of properties are built new each year, and an additional 2% have been either renovated or changed use in the past year, and since these properties tend to be larger and more valuable than the average property, there is substantial scope for additional ongoing increases in revenue collection.

2. *Changes in Collusion?* The model of collusion in Section III suggests that the treatment effects on taxes and bribes paid should be heterogeneous among properties. For properties that switch from collusive to noncollusive equilibrium, we would expect to see an increase in taxes paid and a reduction in bribes. For properties that remain in the collusive equilibrium, we have more ambiguous predictions: the sum of bribes plus taxes paid should go up, but whether this comes from an increase in bribes, taxes, or some combination is less theoretically clear.²⁵ For properties that were in the noncollusive equilibrium before, and remain there, we would expect no changes.

To investigate these effects, in Table VI, Panel A, we first estimate equation (11) on the general population of properties to capture how typical properties in treatment areas differ on tax and bribe payments compared to equivalent properties in control areas. For the typical property, we find that tax payments are essentially unchanged (column (1)). Note, though, that since the change in official revenues observed in the administrative data comes from a very small number of properties (as shown already), we would not necessarily detect it by looking across all properties, and indeed, we cannot reject the null of an average increase in taxes paid of the magnitude found in the administrative data.

In our model of collusive corruption, a low reassessment rate is consistent with many properties rebargaining bribes as a result of the incentive treatment. Indeed, columns (2) and (3) show that bribe rates—measured as the typical amount a property owner would pay in unofficial payments to the tax department over the course of the year for a similar property—increase substantially, by Rs. 594 (US\$6, or about 32% higher compared with the average control area property).²⁶ The frequency of bribe payments also increases substantially. The one metric of corruption that does not change is the overall perception of corruption in the tax department.²⁷

25. Note that in the simple linear framework in Section III, bribes unambiguously increase for properties that remain in the collusive equilibrium, but in the extension in Online Appendix A with convex costs, the prediction on bribes becomes ambiguous.

26. Note that the increase in average bribe payments comes entirely from the intensive margin, as we would expect from a shift in the collusive equilibrium. See Online Appendix Table 15 for more details.

27. Note that we experimented in a pilot survey with asking directly whether the respondent had paid bribes. We experienced low response rates

TABLE VI
 IMPACTS ON TAX PAYMENTS AND CORRUPTION, BY REASSESSED STATUS

| | (1) | (2) | (3) | (4) |
|---|------------------------------|-------------------|----------------------------------|-----------------------------|
| | Self-reported tax payment | Bribe payment | Frequency of bribe payment | Perception of corruption |
| Panel A: General population sample only | | | | |
| Treatment | -62.81 (264.7) | 594.1* (341.7) | 0.2021** (0.0951) | 0.0113 (0.0254) |
| <i>N</i> | 11,586 | 5,993 | 4,802 | 6,050 |
| Mean of control group | 4,069.425 | 1,874.542 | 0.683 | 0.644 |
| Panel B: Reassessed and general population sample | | | | |
| Reassessed * treatment | 1,884* (1,083) | -557.4 (380.1) | -0.1592* (0.0942) | -0.0031 (0.0221) |
| Reassessed | 2,763*** (572.9) | -66.38 (177.5) | 0.0137 (0.0403) | -0.0191* (0.0107) |
| <i>N</i> | 16,353 | 8,207 | 6,993 | 8,268 |
| Sample | Full | Phase 1 | Phase 1 | Phase 1 |
| Mean of control group in gen. pop. sample | 3928.252 | 1874.542 | 0.683 | 0.644 |

Notes. This table considers how the average property in treatment areas differs in terms of the tax payments and bribes it reports (Panel A) as well as asking whether these outcomes differ for reassessed properties (Panel B). In both cases we present instrumental variables regressions, where treatment status is instrumented with randomization results. Unit of observation is a property. Bribe payment is the respondent's response to how much bribe they think others would pay for a similar property. Frequency of bribe payment and perception of corruption are graded on a five-point rubric and scaled to the interval [0,1]. Panel A uses only properties from the random sample drawn from the field, and Panel B includes properties that were selected from the official register of reassessments. The reassessed dummy in Panel B denotes such (reassessed) properties. The specifications in Panel A follow equation (12) of the main text, with the exception of column (1), which controls for self-reported baseline (FY 2011) tax payment. Specifications in Panel B follow equation (12) of the main text. For columns (2)-(4), sample is restricted to circles from the first phase of the survey (see text for details). In both Panels A and B, specifications include a control for whether the response came from the short version of the survey, and the phase of the survey (if applicable). The information treatment is included in the control group. Robust standard errors are in parentheses. Standard errors are clustered by robust partition of circles, that is, the group of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < .01$.

to this question, and found that respondents were much more forthcoming when we asked the question indirectly, that is, what the going bribe rate was for a property that was "similar" to theirs. Note that this phrasing does not necessarily yield a precise average bribe paid, since respondents may answer the question either conditional or unconditional on paying a bribe and the wording of the question is not precise enough to reliably distinguish between the two. Since the frequency of bribes paid also goes up, however, this implies that even though we may not be able to estimate the precise magnitude, average bribe payments do in general increase.

Given limited statistical power in being able to detect the increase in overall taxes paid in treatment circles due to the low frequency of reassessed properties, we now turn to specifically focusing our analysis on these reassessed properties. In Panel B of Table VI, we estimate equation (12), which examines the differential impact between typical properties and reassessed properties (i.e., those whose tax bill changed, who may be disproportionately those that switch from one equilibrium to another). The coefficient β_2 from equation (12), that is, the coefficient on *ReAssessed*, captures how properties that are reassessed (or newly entered on the tax rolls) differ from the general population of properties in control circles, and β_1 , the coefficient on *Reassessed * Treatment*, captures any additional difference in treatment circles (the treatment dummy is absorbed by the circle fixed effect).

There are several key results. First, compared with nonreassessed properties, properties in control circles whose valuations were changed pay substantially higher taxes—Rs. 2,763, or about 70% higher than the control group mean for random properties. This is even more true in treatment areas, where reassessed properties pay an additional Rs. 1,884 more than nonreassessed properties. The results here are consistent with the treatment effect on revenue we see in the administrative data.²⁸ On the other hand, the increase seen in bribes in treatment areas is not seen for reassessed properties, that is, the coefficient on *ReAssessed * Treatment* is negative, and completely

28. To see this, note that average tax in a circle is a weighted average of tax paid by reassessed and non-reassessed properties, that is,

$$E[\text{TaxPayment}] = E[\text{TaxPayment}|\text{Reassessed}]P(\text{Reassessed}) \\ + E[\text{TaxPayment}|\text{NonReassessed}]P(\text{NonReassessed}).$$

Based on our estimates here and data on reassessment rates (9% of taxable properties were reassessed in the cumulative two-year treatment period in control circles, and for simplicity we treat our general population sample as composed only of non-reassessed properties), this average in control areas is $(0.09)(3928 + 2763) + (0.91)(3928) = 4177$. This gives an average tax per property of Rs. 4,177 in control areas. Using our treatment effect estimates (i.e., increases in the number of reassessed properties and the greater payments received from such properties), the analogous average tax in treatment circles is given by $(0.128)(3928 + 2763 + 1884) + (0.872)(3928) = 4523$. An increase in the average tax per property from Rs. 4,177 to Rs. 4,523 represents a 8.3% increase in tax collection, which is quite close to the observed effect from our administrative data of about 9% (9.3 log points).

offsets the treatment effect for bribes on random properties shown in Panel A.²⁹

Thus these results show, as suggested in the model, that performance pay for tax collectors leads to heterogeneous effects: increases in bribes for the majority of properties, but no increases in bribes with substantial increases in tax revenue for a small number of properties that switch from collusion to noncollusion. Although bribes do not fall to zero for reassessed properties, as would be predicted by the linear model if collusion were avoided entirely, the qualitative pattern from the model emerges. This also underscores that the increased revenue as a result of the performance pay schemes is due to a small number of properties moving from a collusive to a noncollusive equilibrium and the corresponding substantial increase in taxes paid by such properties.

Table VII examines whether there is an analogous differential response on nonrevenue outcomes, that is, satisfaction, inaccuracy, and the tax gap. The key results are for inaccuracy and the tax gap. Column (3) shows that reassessed properties are more accurately (i.e., less inaccurately) assessed compared to nonreassessed properties. That is, there is a closer match between the tax liability computed by our independent surveyors and that computed by the tax department. Moreover, column (4) shows that while the typical (i.e., randomly selected) property in the control group is undertaxed, this is eliminated in reassessed properties (i.e., adding the coefficient of 0.122 on reassessment to the mean of -0.103 yields a net result of 0.019, which is not statistically significant from zero; p -value of .191); that is, reassessed properties are, on average, taxed at the amount our independent survey team would predict. Although these effects are similar in both treatment and control areas, they confirm the view of reassessment as a bargaining breakdown: unlike typical randomly selected properties, which in general are undertaxed, reassessed properties are assessed more accurately and are neither over- nor undertaxed on average.

Reassessed properties are not, broadly speaking, unsatisfied with the tax department. In fact, Table VII shows that reassessed properties generally appear more satisfied with the tax

29. Online Appendix Table 6-H repeats analysis of Table VI broken down by the three subtreatments. The results do not show substantial differences in these dimensions among the three subtreatments.

TABLE VII
 IMPACTS ON SATISFACTION AND ACCURACY, BY REASSESSED STATUS

| | (1) | (2) | (3) | (4) |
|--|---------------------|---------------------|----------------------|---------------------|
| | Quality | Satisfaction | Inaccuracy | Tax gap |
| Reassessed * treatment | 0.009 (0.024) | 0.005 (0.024) | 0.001 (0.017) | -0.005 (0.028) |
| Reassessed | 0.049*** (0.013) | 0.044*** (0.013) | -0.061*** (0.009) | 0.122*** (0.015) |
| <i>N</i> | 8,268 | 8,268 | 14,173 | 14,173 |
| Sample | Phase 1 | Phase 1 | Full | Full |
| Mean of control group in gen. pop. sample | 0.538 | 0.555 | 0.339 | -0.103 |

Notes. This table examines whether nonrevenue-based outcomes differ for reassessed properties. The unit of observation is a property. Specification follows equation (12) of the main text, and controls for whether the response came from the short version of the survey. Columns (1) and (2) restrict the sample circles from the first phase of the survey (see Online Appendix B for details). The information treatment is included in the control group. Robust standard errors are in parentheses. Standard errors are clustered by robust partition, that is, the group of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < .01$.

department, and this is not different between treatment and control. One reason that there may be no change in satisfaction for these properties between treatment and control—even though they pay fewer bribes but much higher taxes in treatment areas—is that the theory predicts that those who are reassessed and switch between the collusive and noncollusive equilibrium in response to the treatment are those who are closest to being indifferent between the two regimes. The switch from collusive to noncollusive equilibrium may therefore represent a second-order utility change for these property owners, even though it yields a first-order change in revenue for the government.

All told, the results here paint a picture consistent with the theoretical framework: in pay for performance regimes, most properties pay no more taxes but do pay somewhat higher bribes; but some properties switch from the collusive to noncollusive equilibrium. Those properties that are reassessed do not experience the increase in bribes, but instead pay substantially higher taxes, are assessed more accurately, and are no longer underassessed relative to what our independent survey reveals.

3. *Who Gets Reassessed?* If these reassessments represent bargaining breakdowns, an interesting question is which property-inspector pairs are affected. In the model, equation (5) shows

that the increase in taxes comes from those properties on the margin of switching—that is, those properties with taxpayer and tax collector values of the disutility of evasion parameters α_i and β_i such that they are close to indifferent between the high bribe, low tax equilibrium and the low bribe, high tax equilibrium.

To examine who these marginal properties are in the data, we consider how reassessed properties differ from typical properties, and how this differs in treatment versus control areas. The results, estimated using equation (12), are presented in Tables VIII and IX, where Table VIII examines characteristics of the property and Table IX examines characteristics of the owner. Table VIII shows that reassessed properties are generally those (in both treatment and control areas) that are subject to higher tax rates than typical property. For example, according to the data we obtain from our independent survey, they have a GARV (i.e., tax base, before exemptions are applied) that is 67% higher than the mean property in control areas. They also have more floors and are more likely to have been recently renovated, to belong to a more expensive tax bracket (tax category), to be commercial (which is taxed at a higher rate), and to be rented (which is also taxed at a higher rate).

Examining whether any of these margins change further in treatment circles, the point estimates suggest that on net, reassessed properties in treatment areas have a GARV that is an additional 33% larger (p -value of .21) than the average reassessed property in control areas. Therefore, reassessed properties in treatment areas have a 122% higher GARV than the typical property in control areas from the general population sample.³⁰ Incentivized staff also seem to focus more on commercial rented properties, which have the highest assessments per square foot of area. One interpretation is that commercial properties have a higher disutility from paying bribes (i.e., higher α_i) than residential properties and hence are more marginal.

Panel B considers differences in owner characteristics. One interesting finding is that those owners who report a close personal (family/friend) relationship with a politician are 1.3

30. Another way to look at this is to plot, nonparametrically, the relationship between the probability of being reassessed and tax density, which is the tax valuation per unit of covered area. Online Appendix Figure F.1 shows that high tax density properties appear more likely to be reassessed in treatment areas.

TABLE VIII
SELECTION EFFECTS ON REASSESSMENTS

| | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) | (9) | (10) |
|--|---------------|---------------------|---|----------------------------|--|--------------|-----------------|--|--|------------------|
| GARV | | Number of floors | Last renovation was ≤ 2 years ago | Land area (sq. feet) | Total covered area (sq. feet) | Main Road | Tax Category | Percent of property commercial and rented | Percent of property commercial and rented | Tax Liability |
| Reassess * | 20,137.796 | 0.002 | -0.005 | -32.599 | 852.092 | -0.002 | -0.226** | 0.018 | 0.075** | 3,897.980 |
| treatment | (16,187.550) | (0.050) | (0.020) | (82.473) | (771.516) | (0.048) | (0.088) | (0.037) | (0.029) | (3,539.474) |
| Reassess | 24,683.609*** | 0.078*** | 0.094*** | 37.396 | -156.619 | 0.064*** | 0.212*** | 0.217*** | 0.176*** | 5,503.481*** |
| | (7,944.915) | (0.026) | (0.011) | (57.199) | (379.299) | (0.024) | (0.044) | (0.019) | (0.015) | (1,754.013) |
| N | 15,090 | 16,352 | 16,354 | 16,352 | 16,352 | 16,352 | 15,090 | 16,226 | 16,227 | 15,090 |
| Mean of control group in gen. pop. sample | 36,808.77 | 1.57 | 0.02 | 301.13 | 2,779.82 | 0.46 | 3.78 | 0.35 | 0.17 | 6,642.00 |

Notes. Property-level 2SLS regressions. Specifications follow equation (12) of the main text and includes a control for whether the response came from the short version of the questionnaire. This table looks at selection effects on property characteristics. The characteristics labeled components of GARV are those that directly enter into the formula used to calculate GARV. Tax category (column (7)) is seven-tiered categorical variable with 7 being the most expensive tax bracket and 1 being the cheapest. Standard errors are clustered by robust partition, the partition of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < 0.01$.

TABLE IX
SELECTION EFFECTS ON REASSESSMENTS

| | (1) | (2) | (3) | (4) | (5) | (6) |
|---|--------------------------------|----------------------------------|-------------------------|---|----------------------------|--|
| | Approximate age of owner | Owner's level of education | Per capita wages | Predicted expenditure given assets | Connected to politician | Connected to politician, government, police |
| Reassess * treatment | -0.348 (0.799) | -0.523* (0.317) | -821.749 (1,078.191) | 110.798 (213.234) | 0.021* (0.012) | 0.005 (0.027) |
| Reassess | -0.656* (0.398) | 0.303* (0.157) | 13.126 (510.006) | -94.529 (122.380) | -0.013** (0.006) | 0.005 (0.014) |
| <i>N</i> | 13,406 | 16,254 | 13,765 | 13,954 | 16,354 | 16,354 |
| Mean of control group in gen. pop. sample | 50.70 | 9.19 | 16,281.55 | 6,292.58 | 0.05 | 0.36 |

Notes. Property-level 2SLS regressions. Specifications follow equation (12) of the main text and includes a control for whether the response came from the short version of the questionnaire. This table looks at selection effects on owner/tenant characteristics. Per capita wages (column (3)) are self-reported household expenditures divided by the total number of working household members. Predicted expenditure given assets (column (4)) is the predicted value of a regression of household expenditure on series of dummy variables indicating various household assets. Standard errors are clustered by robust partition, the partition of circles such that all circles that merged or split with each other are included within the same partition. * $p < .10$, ** $p < .05$, *** $p < .01$.

percentage points (over a baseline value of 5.3% in control circles) less likely to be reassessed than typical properties. However, this effect is undone in treatment areas—so that while connected owners seem to enjoy an advantage in general, this no longer the case in treatment areas. We should state the caveat that this particular result must be interpreted with caution, given that it is only one out of many coefficients examined. However, it is interesting that a similar pattern holds for education: educated owners are generally more likely to be reassessed, but this effect is undone in treatment areas. On net, the results in this section do paint a consistent picture: the performance incentives led inspectors to concentrate on a relatively small number of high-value properties.

VI.C. *Alternative Channels*

1. Changes in Collusion versus Greater Inspector Search Effort? We have interpreted our results so far in the context of changes in collusive behavior as a result of introducing performance pay. This is not meant to imply that changing effort on the part of tax staff might not play a role as well. Instead, we posit that we are unlikely to observe the above pattern of results

without a change in collusion along the lines of our model. That is, although a tax collector may respond to treatment by working harder at uncovering the true tax liability of a taxpayer or in getting recovery against that liability (as in standard moral hazard models such as Hölmstrom 1979), these forces alone are not consistent with the results we find on bribes.

Furthermore, evidence from self-reported behavior by inspectors does not seem to indicate such effort was important: in Online Appendix Table 16 we find little observable change in effort (total hours spent working per day, etc.) reported by inspectors in treatment areas. The only result is that inspectors seem to be spending more time in the office and less in the field. While it is possible that time in the office is correlated with higher effort (e.g., filling out paperwork), it is not a priori what one would have expected in terms of effort, especially to the extent that the relevant margin was uncovering recent property changes. However, changes in collusion could quite plausibly imply more time in the office to change corresponding paperwork.

In addition, Online Appendix Table 23 presents further evidence to suggest that the tax collector likely had prior knowledge about properties that were newly added to the tax registers. We compare attributes of properties newly added to the register to attributes of properties that were verified as new based on the third-party property survey. We find that the former are more likely to be rented, commercial, and larger in area (all of which imply a higher tax liability), and have been occupied by the current tenant for longer. In fact, properties that are newly assessed are more likely to be rented (which is harder to observe without prior interactions) even conditional on other property characteristics. These results suggest properties newly added to the tax rolls are being systematically targeted, something that would not be feasible unless the tax collector had prior (private) information that is now being revealed.

2. Mechanisms beyond Price Effects. We have thus far interpreted our results as due to the increased marginal incentives (i.e., price effects) provided to collect more taxes. However, the schemes have other aspects that could also enhance performance. In this section we briefly discuss several of these alternate channels (additional details can be found in Online Appendix D). The objective of this section is not to definitely rule out these

channels—it is likely they do contribute to some extent—but rather to see how significant they may be. We conclude that although some of these channels may partly contribute, the price effects of the incentives still seem to be the primary way in which the incentives had an impact.

With any performance-based payment scheme that pays over baseline salary, two potential confounds are perceived monitoring effects and income effects. Perceived monitoring refers to the fact that the incentive scheme may affect agents' beliefs about how they will be monitored even outside of the explicit financial incentives. Income effects refer to the fact that a top-up incentive scheme increases the overall income of an agent, in addition to changing the price of the incentivized action. In this case, if honesty is a normal good (i.e., inspectors take bribes because they have a high marginal utility of income) or there are efficiency wage effects as in Becker and Stigler (1974), one could imagine that our effects are also due to income (and not just price) effects.

For both potential alternate channels, an important piece of evidence of the primary role of price effects is the difference between the three incentive schemes. In particular, all three schemes were equally salient to inspectors and all three performance-pay schemes were designed to generate approximately similar expected income (and indeed did so),³¹ yet we saw already that they generate very different impacts on revenue: the increase in tax revenue in revenue was almost double what was in revenue plus, and the flexible bonus scheme produced no detectable tax impacts. These simple facts suggest *prima facie* that the different prices implicit in the different schemes are primarily what are driving the results, not the income transfer per se.

Several additional tests also confirm that price effects seem to be the primary explanation for our results. Starting in year 2, we introduced an “information-only” scheme that provided identical information to the revenue scheme, but without financial payments. As discussed in Online Appendix E, this scheme produced no statistically significant effects on total or current revenue, and suggests that perceived monitoring explained no more

31. Average payments to inspectors in year 1 were: Rs. 255,608 in revenue, Rs. 247,283 in revenue plus, and Rs. 297,370 in flexible bonus. Average payments in year 2 were: Rs. 255,773 in revenue, Rs. 282,490 in revenue plus, and Rs. 255,977 in flexible bonus.

than one-third the total impact. To investigate income effects, Online Appendix E tests for any income effects using the fact that benchmarks in the revenue and revenue plus schemes we determined based on the second, third, and fourth lags of revenue, but not the first lag, so the first lag of revenue generates a shock to income but not price effects. This approach also finds no strong evidence of income effects in this context.

Supervisory effort could also respond to the incentive scheme. Although this is not a confound *per se*, it entails a slightly different interpretation of the results. However, we show in Online Appendix E that a separate treatment scheme that explicitly rewarded supervisors (in an analogous manner to the revenue scheme) had no effects on average, nor did it have substantial interactions with the direct incentives. We also show that inspectors do not report being more extensively pressured by their supervisors to work harder in areas with inspector incentives, so on net this does not appear to be an important part of the story.

Given that the randomization was conducted publicly (to ensure a perception of fairness), this meant that both control and treatment circles knew their respective identities. One potential concern is that control group inspectors may have become discouraged and performed worse, leading us to overestimate treatment effects. Though any such spillovers are less of a concern when comparing among the three treatment schemes, it is also worth noting that in general, the rate of growth of revenue in control circles during the experimental years was greater than during the previous years, suggesting that overall discouragement was not a first-order concern. In addition, to test for spillovers more directly, Online Appendix Table 10 examines the impact of the treatment on nearby, neighboring control circles, where the treatment would be particularly salient, compared with control circles further away with whom inspectors interacted less often. If spillovers were to have occurred, they would likely have occurred locally, as inspectors in nearby circles share the same physical office space. We cannot reject the null of no spillovers.

Finally, one may also be concerned that the performance pay scheme could have increased inspectors' security of tenure within their circles. However, tax staff were told explicitly that normal transfer policy would be in effect during the study, and we show in Online Appendix Table 17 that there are no statistically

significant differences in transfer rates among treatment and control staff.

VI.D. Cost-Effectiveness

From the government's and broader policy perspective, a natural question is whether these schemes were cost-effective, that is, whether the additional revenue received in taxes exceeded the amount paid as incentives. As shown in Section III, under the assumption that bribes represent a transfer rather than an efficiency loss, the change in net revenue for the government—revenue received in taxes less the amount paid in incentives—is a measure of the change in social welfare from the program.

For the revenue and revenue plus scheme, which pay out to staff a percentage of revenue collected over a fixed benchmark, one would expect the net revenue to be positive as long as the benchmark was set sufficiently high that one is not paying for inframarginal collections. Of course, benchmarks cannot be set too high or else staff would not be in the money and would not be receiving incentives on the margin, so setting the benchmark is nontrivial. For the flexible bonus scheme, the payments were fixed in advance, so it is less clear *ex ante* whether net revenue for the government would be positive.

We focus on cost-effectiveness in the second year of the program, when it was at scale. For each circle, we predict the revenue at the end of year 2 using our estimated treatment effects for each scheme.³² We use the estimates to calculate the predicted additional revenue in treatment circles due to the treatment, and then sum this across treatment circles to obtain total additional revenue. The total costs are the actual performance-based payments paid out under each of the schemes. Net revenue is the difference between predicted additional revenue and the incentive costs.

The results are shown in Table X. Since the point estimates are slightly different depending on whether the information treatment is included as part of the control group (as in Table III) or not (as in Online Appendix Table 3-I), we report the results both ways (Panels A and B, respectively). Taken

32. The only change from our main specification is that we estimate reduced form treatment effects, where we weight each circle by the circle's revenue in the baseline year to account for any heterogeneity in treatment effects across circles of different sizes, which matters substantially for the impact on total revenue raised.

TABLE X
COST-EFFECTIVENESS OF INCENTIVES

| | (1) Additional revenue | (2) Cost of incentives | (3) ROI |
|--------------------------------------|------------------------------|------------------------------|------------|
| Panel A: Information in controls | | | |
| Any treatment | 124,961,461 | 108,387,160 | 15.29 |
| Revenue | 50,578,024 | 37,349,784 | 35.42 |
| Revenue plus | 40,671,290 | 35,549,342 | 14.41 |
| Flexible bonus | 30,555,313 | 35,488,035 | -13.90 |
| Panel B: Information out of controls | | | |
| Any treatment | 140,973,016 | 108,387,160 | 30.06 |
| Revenue | 56,269,064 | 37,349,784 | 50.65 |
| Revenue plus | 45,539,845 | 35,549,342 | 28.10 |
| Flexible bonus | 35,571,720 | 35,488,035 | 0.24 |

Notes. This table estimates the economic return generated by the performance pay schemes. Column (1) estimates the additional revenue due to treatment, calculated with a reduced-form regression of log total revenue on log total baseline revenue, weighting observations by baseline revenue (in levels). For each treated observation, we generate a prediction of revenue collection under treatment and a prediction of revenue collection in absence of treatment and subtract to calculate the additional revenue due to treatment. The total additional revenue collection due to treatment is the sum of additional revenue collection across treated observations. Panels A and B show how the calculation changes depending on whether the Information treatment is included in the controls (Panel A) or dummied out (Panel B). Column (2) gives the actual costs of the incentive payments paid to circle staff under each scheme. Column (3) then presents return on investment (ROI), which is simply the percent increase in additional revenue above costs.

together the results show that net revenue is positive, so the schemes are cost-effective. Dividing the net gain (revenue less costs) by costs to calculate a return on investment for the government shows a return of 15% (Panel A) to 30% (Panel B). For the revenue scheme, which raised the most revenue, the return at the end of year 2 ranges from 35% (Panel A) to 51% (Panel B). The revenue plus scheme earns 14% to 28% ROI, and the flexible bonus scheme loses money for the government.

Note that since a main channel seems to be an increase in net demand (i.e., new properties added to the tax rolls), to the extent these changes are permanent and last even after the treatments are discontinued, the long-run cost-effectiveness from a time-limited/temporary introduction of performance-based pay could be substantially higher than the numbers reported here.

VII. CONCLUSION

Our article examines the impact of introducing performance pay schemes in taxation. Taxation is interesting not only because

it is feasible to design outcome-based pay mechanisms, but also because it presents interesting challenges in considering incentive pay mechanisms. Given the potential for collusion between the civil servant and the citizen, high-powered incentives are not simply about increasing worker effort to achieve the desired (i.e., incentivized) outcome. Instead, in such contexts incentives can increase the bargaining power of the civil servant with respect to the taxpayer, leading to potentially less desirable outcomes.

Our results suggest that while these effects on bargaining are present, on net performance pay mechanisms can be quite effective in raising additional taxes, and they can do so without generating too much animosity toward the tax department, as was often associated with tax farming historically. Although it is possible that such costs may show up over a time frame longer than two years, it is nevertheless instructive to examine why such costs might not be as high in our performance-pay schemes. In standard contract theory, a principal has to better incentivize an agent to the extent that the agent's objective function differs from the principal's. In taxation, to the extent that there is collusion—and our results suggest that this is an important margin—there is a clear wedge in such objectives in terms of raising taxes. Performance pay can therefore reduce this wedge by directly making the tax collector a (partly) residual claimant on taxes collected.

But what about divergences in political objectives between the politician/government and tax collectors? The historical tax farming literature suggests that collectors may have been less sensitive to political costs they imposed when raising taxes. However, tax collectors in our context may not be as free to raise taxes—they are not so locally powerful that they are unaffected by the displeasure of the population they tax. In fact, more often than not they may have weaker socioeconomic and political influence compared to those they are meant to tax, and so may also be quite concerned about the potential costs that raising excessive taxes may induce. Qualitatively, conversations with tax collectors suggested that this was a concern, that is, tax collectors would justify lower collections by noting that the taxpayers could get them transferred or otherwise sanctioned both because the individual taxpayer may be quite influential and/or because they may collectively be powerful (e.g., shop-keepers' local associations). In fact, quite often (perhaps as a tacit means of justifying collusion), tax collectors would express sympathy to taxpayers'

unwillingness to pay taxes, particularly in poorer localities, given the general level of dissatisfaction taxpayers would have about how their taxes are used (locally) by the state.

So how might tax collectors balance their increased incentives to raise more taxes due to performance pay schemes with a need to not increase taxpayer dissatisfaction? One could imagine two different types of potential responses. One response is to tax a large number of (poorer) property owners, who may have less influence or ability to push back, and to spare the more connected, wealthier owners of larger properties. Alternatively, inspectors could focus their efforts on a small number of high-value owners. This would generate the largest return per property, and avoid alienating a large number of people, but could be risky if it alienates influential people. In a sense, this is a trade-off between two types of influence: since each person gets one vote, smallholders have more votes per dollar, and hence more influence democratically, but largeholders may have more influence. The results here suggest that inspectors took the latter approach: focusing on a small number of high-value property owners.

In terms of how collusion mediates the impact of performance pay, we find evidence that it indeed strengthens the bargaining power of the tax collector. For the majority of taxpayers, tax payments remain unaffected, although they end up paying higher rents to the tax collector as they rebargain. Although some taxpayers do end up paying more taxes and collusion breaks down, generating more revenue for the government, these results offer a word of caution that the effects of incentives are more complex than they would be in a world where the only margin is effort and there is no collusion. If the goal is to both increase performance/collections and reduce rent-seeking, one may need to accompany a performance pay mechanism with stricter monitoring and direct penalties for rent-seeking.

Taken together, the results suggest that, notwithstanding historical concerns regarding tax farming and the relative absence of such high-powered incentives in developed economies, performance pay schemes in taxation may be a promising avenue to explore for developing economies seeking to raise revenues. The remaining question for governments is whether they can mitigate the potentially undesirable effects of the increased bargaining power tax staff have over taxpayers by more direct audit-based processes that can effectively detect and penalize collusion. The fact that our results show impacts on the tax base

suggest that a promising direction may be to introduce high-powered incentives for short durations and at times when revealing information to the government is particularly important (such as when a major revaluation of properties or similar such reform is under way), and such schemes may need to be accompanied by complementary efforts at reducing corruption and better third-party data verification processes. To the extent these concerns can be addressed, our results demonstrate that such schemes can be an important and financially and politically feasible way for emerging economies to undertake the essential and necessary task of raising tax revenue and enlarging their tax base.

SUPPLEMENTARY MATERIAL

An Online Appendix for this article can be found at QJE online (qje.oxfordjournals.org).

LONDON SCHOOL OF ECONOMICS
HARVARD KENNEDY SCHOOL, HARVARD UNIVERSITY
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

REFERENCES

- Ashraf, N., O. Bandiera, and K. Jack, "No Margin, No Mission? A Field Experiment on Incentives for Public Service Delivery," Unpublished manuscript, London School of Economics, 2013.
- Bahl, R., S. Wallace, and M. Cyan, "Pakistan: Provincial Government Taxation," Technical Report, International Center for Public Policy, Andrew Young School of Policy Studies, Georgia State University, 2008.
- Baker, G., R. Gibbons, and K. J. Murphy, "Subjective Performance Measures in Optimal Incentive Contracts," *Quarterly Journal of Economics*, 109 (1994), 1125–1156.
- Bartlett, B., "How Excessive Government Killed Ancient Rome," *Cato Journal*, 14 (1994), 287–303.
- Becker, G. S., and G. J. Stigler, "Law Enforcement, Malfeasance, and Compensation of Enforcers," *Journal of Legal Studies*, 3 (1974), 1–18.
- Besley, T., and J. McLaren, "Taxes and Bribery: The Role of Wage Incentives," *Economic Journal*, 103 (1993), 119–141.
- Best, M. C., A. Brockmeyer, H. J. Kleven, J. Spinnewijn, and M. Waseem, "Production vs Revenue Efficiency with Limited Tax Capacity: Theory and Evidence from Pakistan," Mimeo, 2013.
- Boycko, M., A. Shleifer, and R. W. Vishny, "A Theory of Privatisation," *Economic Journal*, 106 (1996), 309–319.
- Carillo, P., D. Pomeranz, and M. Singhal, "Tax Me if You Can: Evidence on Firm Misreporting Behavior and Evasion Substitution," Technical Report, Harvard Kennedy School, 2014.
- Chetty, R., "Is the Taxable Income Elasticity Sufficient to Calculate Deadweight Loss? The Implications of Evasion and Avoidance," *American Economic Journal: Economic Policy*, 1 (2009), 31–52.
- City of Cambridge, *Annual Budget 2014–2015*, 2014, available at <http://goo.gl/rMTDjB>.

- City of Detroit, Michigan, "Comprehensive Annual Financial Report for the Fiscal Year Ended June 30, 2013," Technical Report, 2013, available at <http://goo.gl/5NjnIs>.
- Dal Bó, E., F. Finan, and M. A. Rossi, "Strengthening State Capabilities: The Role of Financial Incentives in the Call to Public Service," *Quarterly Journal of Economics*, 128 (2013), 1169–1218.
- Das-Gupta, A., and D. Mookherjee, *Incentives and Institutional Reform in Tax Enforcement: An Analysis of Developing Country Experience* (New York: Oxford University Press, 1998).
- Duflo, E., M. Greenstone, R. Pande, and N. Ryan, "Truth-Telling by Third-Party Auditors and the Response of Polluting Firms: Experimental Evidence from India," *Quarterly Journal of Economics*, 128 (2013), 1499–1545.
- Feldstein, M., "Tax Avoidance and the Deadweight Loss of the Income Tax," *Review of Economics and Statistics*, 81 (1999), 674–680.
- Fisman, R., and S.-J. Wei, "Tax Rates and Tax Evasion: Evidence from 'Missing Imports' in China," *Journal of Political Economy*, 112 (2004), 471–500.
- Gertler, P., and C. Vermeersch, "Using Performance Incentives to Improve Medical Care Productivity and Health Outcomes," NBER Working Paper 19046, 2013, available at <http://www.nber.org/papers/w19046>.
- Glewwe, P., N. Ilias, and M. Kremer, "Teacher Incentives," *American Economic Journal: Applied Economics*, 2 (2010), 205–227, available at <http://www.jstor.org/stable/25760225>.
- Gordon, R., and W. Li, "Tax Structures in Developing Countries: Many Puzzles and a Possible Explanation," *Journal of Public Economics*, 93 (2009), 855–866.
- Hölmstrom, B., "Moral Hazard and Observability," *Bell Journal of Economics*, 10 (1979), 74–91.
- Hölmstrom, B., and P. Milgrom, "Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership, and Job Design," *Journal of Law, Economics, and Organization*, 7 (1991), 24–52.
- Kahn, C. M., E. C. Silva, and J. P. Ziliak, "Performance-Based Wages in Tax Collection: The Brazilian Tax Collection Reform and its Effects," *Economic Journal*, 111 (2001), 188–205.
- Kleven, H. J., C. T. Kreiner, and E. Saez, "Why Can Modern Governments Tax So Much? An Agency Model of Firms as Fiscal Intermediaries," Technical Report, 2014.
- Kleven, H. J., and M. Waseem, "Using Notches to Uncover Optimization Frictions and Structural Elasticities: Theory and Evidence from Pakistan," *Quarterly Journal of Economics*, 128 (2013), 669–723.
- Kumler, T., E. Verhoogen, and J. A. Fras, "Enlisting Employees in Improving Payroll-Tax Compliance: Evidence from Mexico," NBER Working Paper 19385, 2013, available at <http://www.nber.org/papers/w19385>.
- MacLeod, B. W., "Optimal Contracting with Subjective Evaluation," *American Economic Review*, 93 (2003), 216–240.
- Mookherjee, D., and I. P.-L. Png, "Corruptible Law Enforcers: How Should they Be Compensated?," *Economic Journal*, 105 (1995), 145–159.
- Muralidharan, K., and V. Sundararaman, "Teacher Performance Pay: Experimental Evidence from India," *Journal of Political Economy*, 119 (2011), 39–77.
- Olken, B. A., "Monitoring Corruption: Evidence from a Field Experiment in Indonesia," *Journal of Political Economy*, 115 (2007), 200–249.
- Olken, B. A., and R. Pande, "Corruption in Developing Countries," *Annual Review of Economics*, 4 (2012), 479–509.
- Parrillo, N. R., *Against the Profit Motive: The Salary Revolution in American Government, 1780–1940* (New Haven, CT: Yale University Press, 2013).
- Piracha, M. M., and M. Moore, "Understanding Low-Level State Capacity: Property Tax Collection in Pakistan," Technical Report 33, ICTD, 2015.
- Pomeranz, D., "No Taxation without Information: Deterrence and Self-Enforcement in the Value Added Tax," NBER Working Paper 19199, 2013, available at <http://www.nber.org/papers/w19199>.
- Prendergast, C., "The Provision of Incentives in Firms," *Journal of Economic Literature*, 37 (1999), 7–63.

- Prendergast, C., and R. H. Topel, "Favoritism in Organizations," *Journal of Political Economy*, 104 (1996), 958–978.
- Rasul, I., and D. Rogger, "Management of Bureaucrats and Public Service Delivery: Evidence from the Nigerian Civil Service," Working Paper, University College London, 2013.
- Shleifer, A., and R. W. Vishny, "Politicians and Firms," *Quarterly Journal of Economics*, 109 (1994), 995–1025.
- White, E. N., "From Privatized to Government-Administered Tax Collection: Tax Farming in Eighteenth-Century France," *Economic History Review*, 57 (2004), 636–663.
- World Bank, "Property Taxes in the Punjab, Pakistan," Technical Report, 2006, available at <https://openknowledge.worldbank.org/handle/10986/8277>.
- , "Government of the Punjab Property Tax Decentralisation Program: Scope Evaluation Report," Technical Report, 2009, available at <https://openknowledge.worldbank.org/handle/10986/12378>.

This page intentionally left blank