

Call Me Maybe: Experimental Evidence on Frequency and Medium Effects in Microenterprise Surveys*

Robert Garlick[†], Kate Orkin[‡], Simon Quinn[§]

February 9, 2019

PRE-ANALYSIS PLAN · QUESTIONNAIRES

Abstract

We study the effect of differences in survey frequency and medium on microenterprise survey data. We randomly assign enterprises to monthly in-person, weekly in-person, or weekly phone surveys for a 12-week panel. We find few differences across groups in measured means, distributions, or deviations of measured data from an objective data quality standard provided by Benford's Law. However, phone interviews generate higher within-enterprise variation through time in several variables and may be more sensitive to social desirability bias. Higher-frequency interviews do not lead to persistent changes in reporting or increase permanent attrition from the panel but do increase the share of missed interviews. These findings show that collecting high frequency survey data by phone does not substantially data quality. However, researchers who are particularly interested in within-enterprise dynamics should exercise caution when choosing survey medium.

JEL codes: C81, C83, D22, O12, O17

*This project was funded by Exploratory Research Grant 892 from Private Enterprise Development for Low-Income Countries, a joint research initiative of the Centre for Economic Policy Research (CEPR) and the Department for International Development (DFID). The authors thank Bongani Khumalo, Thembele Manyathi, Mbuso Moyo, Mohammed Motala, Egenes Mudzingwa, and fieldwork staff at the Community Agency for Social Enquiry (CASE); Mzi Shabangu and Arul Naidoo at Statistics South Africa; Rose Page and staff at the Centre for Study of African Economies; and Chris Woodruff and the PEDL team. Our thanks to editor David McKenzie, three anonymous reviewers, Markus Eberhardt, Simon Franklin, Markus Goldstein, David Lam, Murray Leibbrandt, Ethan Ligon, Owen Ozier, Duncan Thomas and seminar audiences and conference participants for excellent comments. Our pre-analysis plan can be viewed at <https://www.socialscisearch.org/trials/346>.

[†]Department of Economics, Duke University; robert.garlick@duke.edu.

[‡]Blavatnik School of Government, Centre for the Study of African Economies and Merton College, University of Oxford; kate.orkin@merton.ox.ac.uk

[§]Department of Economics, Centre for the Study of African Economies and St Antony's College, University of Oxford; simon.quinn@economics.ox.ac.uk

1 Introduction

Researchers designing surveys must choose the interview frequency and medium that generate the optimal quality and volume of data given budget constraints. Alternatives to traditional in-person, low-frequency surveys are increasingly widely used. Phone surveys offer cost savings and the ability to reach mobile populations and to collect data during periods of conflict or disease.¹ High-frequency surveys enable better measurement of short-term fluctuations and outcome dynamics.² However, high-frequency or phone surveys may generate systematically different measurements. This might offset the advantages of richer, cheaper data. Experimental comparisons of data collected using different survey methods help researchers to evaluate these tradeoffs.

We run the first randomised controlled trial to compare microenterprise data from surveys of different frequency and medium. We study a representative sample of microenterprises in the city of Soweto in South Africa. We randomly divide them into three groups. The first group is interviewed in person at every fourth week for 12 weeks. This group is most similar to traditional panel surveys. The second group is interviewed in person every week for 12 weeks. We compare the monthly and weekly in-person groups to test the effects of collecting data at higher frequency, holding the interview medium fixed. The third group is interviewed every week by mobile phone for 12 weeks. We compare the weekly phone and in-person groups to test the effects of data collection medium, holding the interview frequency constant. All interviews use an identical questionnaire measuring 14 enterprise outcomes in approximately 20 minutes.

We find three main results. First, for most outcomes there are few frequency or medium effects on the means, on prespecified quantiles of the distribution, or on the frequency of outliers. There

¹For examples, see Bauer et al. (2013), Dillon (2012), Pape (2018), Turay et al. (2015), and van der Windt and Humphreys (2013).

²Researchers can use high-frequency data to study volatility and dynamics in enterprise and household outcomes (Dupas et al., 2018; Collins et al., 2009; McKenzie and Woodruff, 2008), inform models of intertemporal optimisation in response to shocks (Banerjee et al., 2015; Rosenzweig and Wolpin, 1993), illustrate the time path of treatment effects (Jacobson et al., 1993), explore dynamic treatment regimes (Abbring and Heckman, 2007; Robins, 1997), or average over multiple measures to improve power (Frison and Pocock, 1992; McKenzie, 2012). High-frequency surveys also allow researchers to use shorter recall periods without sacrificing comprehensive time-series coverage (Beegle et al., 2012; Das et al., 2012; De Nicola and Giné, 2014; Heath et al., 2017). See Abebe et al. (2016), Beaman et al. (2014), Carranza et al. (2018), Dabalén et al. (2016), Franklin (2017), Leo et al. (2015), and Zwane et al. (2011) for other examples of high-frequency or phone-based surveys.

are no substantial differences for key outcomes like enterprise closure, profit, sales, costs, fixed assets, or numbers of employees. The largest difference is that phone surveys generate lower reported labour supply. This seems to arise because our in-person interviews take place at enterprises and disproportionately miss respondents who work few hours. Phone respondents also report holding less stock and inventory; transferring more money, stock, or services to the household; and using more often written records to answer survey questions. All outcomes with medium effects except the use of written records are “estimating outcomes,” high-valued outcomes where responses are likely to be estimates rather than precise counts (Blair and Burton, 1987; Gibson and Kim, 2007). We see no frequency or medium effects on our smaller list of “counting outcomes,” low-valued outcomes where respondents can feasibly give a precise answer by counting.

Second, objective indicators of data quality do not differ systematically between interview frequencies or media. We measure data quality by comparing the digit distribution in survey responses to Benford’s Law, a statistical regularity often used to test for data manipulation.³ No one method performs consistently better on this metric. We similarly find few differences between groups in a measure of internal consistency between survey answers: the difference between directly and indirectly elicited profit. These comparisons show that there are limited cross-sectional quality differences in data collected monthly or weekly, by phone or in person.

Third, however, we find phone surveys yield more dispersion in within-enterprise data through time than in-person surveys. Phone surveys yield lower one-week autocorrelations and higher within-enterprise standard deviations on roughly half our 14 outcomes, including some flow and some stock outcomes. We conclude that using phone or high-frequency surveys does not systematically raise or lower the quality of data used for cross-sectional or static panels models. However, researchers particularly interested in within-enterprise dynamics should exercise caution when choosing survey medium.

We also document four secondary results that may inform researchers’ choice of survey frequency and medium. First, autocorrelations are higher for outcomes collected at higher frequen-

³We thank an anonymous reviewer for this suggestion. See Judge and Schechter (2009) for a review of multiple survey datasets from developing countries against this benchmark. Schündeln (2018), Mahadevan (2018), and Garlick (2019) use this approach to assess quality of administrative and survey data in development applications.

cies, so the new information generated by additional surveys may be smaller when surveys are closer together in time. Second, phone surveys are cheaper than in-person surveys. Third, respondents miss a higher share of high-frequency interviews, but high-frequency interviews are still more likely to capture respondents in any month. Fourth, the few frequency effects on means and distributions we observe during our panel do not persist in an in-person endline survey we conduct several weeks later. This shows that higher-frequency surveys in this setting do not generate large persistent changes in behaviour or reporting. This may be useful for researchers interested in the data quality implications of conducting high-frequency panels with subsamples of a larger sample.

These findings come with two caveats. First, non-response in our sample is relatively high: we complete slightly more than half the scheduled interviews. This occurs partly because we imposed a maximum of three attempts to contact each respondent for an interview. This avoided a backlog developing across multiple weeks and kept the number of contact attempts consistent across methods. Patterns of non-response might differ in panels that take advantage of low phone costs to make more attempts. Second, our panel lasts for only three months with at most 12 interviews per respondent. More frequent interviews over a longer time period might induce different patterns in reporting, attrition or non-response.

Our results contribute to a literature in development economics exploring effects of survey frequency or data collection mode (Caeyers et al., 2012; Lane et al., 2006; Fafchamps et al., 2012).⁴ To the best of our knowledge, this and concurrent work by Heath et al. (2017) are the first papers experimentally comparing both survey frequency and medium in a developing country context. While we study microenterprises, our results may be relevant for other types of surveys. Some microenterprise outcomes correspond to outcomes in other surveys; for example, item-specific and total costs in a survey of small enterprises may behave similarly to item-specific and total expenditure in a household survey. A comprehensive mapping of our outcomes to outcomes in other types of surveys is outside the scope of this paper. Instead, we report outcome-specific information and allow the reader to decide if and how these map to their outcomes of interest.

⁴A related literature uses experimental variation to test if different questionnaire designs, recall periods, or survey incentives affect reported outcomes, response rates, or data quality (Arthi et al., 2018; Beegle et al., 2012; Beaman and Dillon, 2012; Das et al., 2012; Dillon et al., 2012; Friedman et al., 2017; Gibson and Kim, 2007; Scott and Amenuvegbe, 1991; Stecklov et al., 2017).

Our work also relates to research on survey media in household surveys and opinion polls, mostly from the US. We find that survey medium has limited effects on reported outcomes, consistent with [De Leeuw \(1992\)](#), [Groves \(1990\)](#), [Groves et al. \(2001\)](#), and [Körmendi \(2001\)](#). Unlike these studies, we find that response rates do not differ by medium. This difference may occur because we analyse medium effects in a panel that has been recruited in person, while the US studies often analysed medium effects in “cold-called” cross-sectional samples. Like this work, we find evidence consistent with higher social desirability bias in phone interviews ([Holbrook et al., 2003](#)). Phone respondents, whose actions cannot be seen by enumerators, are more likely to report using written records to help them answer our survey questions.

Our work also relates to the literature on panel conditioning, which shows that being surveyed or being surveyed more frequently sometimes changes behaviour ([Beaman et al., 2014](#); [Crossley et al., 2017](#); [Stango and Zinman, 2014](#); [Zwane et al., 2011](#)). We find little evidence that differences in interview frequency over our three-month panel generate persistent changes in reported outcomes. This may arise because we study enterprise outcomes that are already salient to respondents. [Bach and Eckman \(2018\)](#) and [Franklin \(2017\)](#) similarly find no persistent effects of interview frequency on the already salient outcome of employment.

We describe the experimental design and data collection processes in [Section 2](#) (and [Appendices A and B](#)). In [Sections 3, 4, and 5](#) (and [Appendices C – F](#)), we present the three main results of the paper: frequency and medium effects on respectively outcome means and distributions, data quality, and within-enterprise data patterns through time. [Section 6](#) brings these results together to categorize outcomes based on the pattern of frequency and medium effects. [Section 7](#) (and [Appendices G – I](#)) presents our four secondary results on autocorrelations, costs, non-response, and persistence. [Section 8](#) concludes.

2 Sample, Experimental Design, and Data

2.1 Sampling and Randomisation

We work in Soweto, near Johannesburg, South Africa. This is a city of approximately 1.28 million people in 2011, of whom 99% are Black Africans. 41% of adults 15 and older engage in some

form of economic activity, including occasional work. 19% of households reported receiving no annual income and another 42% reported receiving less than \$10 per day.⁵

We recruited a representative sample of 1,046 households who owned eligible microenterprises and lived in low-income areas of Soweto. We recontacted 895 of these households several months later to complete our baseline survey. We describe the sampling scheme in Appendix A. We use a common definition of microenterprises: enterprises with at most two full-time employees (in addition to the owner) that do not provide a professional service (e.g. medicine). We excluded any enterprise which did not operate at least three days each week, to exclude seasonal or occasional enterprises where there would be limited intertemporal variation in outcomes.

Most of the 895 enterprises operated in food services (43%) or retail (32%). They were relatively well-established (mean age seven years) and had a diversified client base (mean and median client numbers of respectively 34 and 20, varying substantially by sector). However, they were relatively small: 61% had no employees other than the owner and 28% had only one other employee. Very few were formally registered for payroll or value-added tax, but 20% reported keeping written financial records. The sample is similar to five microenterprise samples from the Dominican Republic, Ghana, Nigeria, and Sri Lanka (De Mel et al., 2008; Drexler et al., 2014; Fafchamps et al., 2014; Karlan et al., 2012; McKenzie, 2017), though our enterprises are slightly older and more concentrated in food and retail/trade. Appendix Table A1 shows detailed summary statistics.

The enterprise owners' households had mean monthly income of US\$394 across all sources, falling in the fourth decile for all households across South Africa.⁶ The households had an average of 3.8 other members, with an interdecile range of 1 to 7. In 55% of households, the enterprise accounted for half or more of household income and 63% of owners perceived pressure within their households to share profits. Only 15% had less than some secondary education. All sampled enterprise owners owned mobile phones.⁷

⁵Authors' own calculations, from the 2011 Census public release data. We follow the terminology of Statistics South Africa, which asks census respondents to describe themselves in terms of five population groups: Black African, Coloured, Indian or Asian, Other, and White.

⁶We use an exchange rate of US\$1 = ZAR10.28, the market exchange rate on the first day of data collection in August 2013.

⁷87% of South Africans aged 18 or older own a mobile phone and the rate is higher in cities (Mitullah and Kama, 2013).

After a baseline survey, we divided the 895 enterprises into three data collection groups using stratified random assignment: monthly in-person interviews (298 enterprises), weekly in-person interviews (299 enterprises), and weekly phone interviews (298 enterprises). We describe the randomisation scheme in Appendix A and show in Table A1 that the groups are balanced on 38 of 40 measured baseline characteristics.

2.2 Survey Protocols

We conducted repeated interviews with each enterprise owner between March and July 2014. We attempted to survey microenterprises in the weekly group once per week for 12 weeks, in person or by phone. We attempted to survey enterprises in the monthly group three times every fourth week for 12 weeks, in person. We randomly split the monthly group into four, so 75 enterprises were interviewed each of weeks 1/5/9, 2/6/10, 3/7/11, and 4/8/12, providing a comparison group for each week when the weekly enterprises were interviewed.

We standardised all survey protocols across arms, except for the variations in survey frequency and medium we study. We used the same questionnaire in all rounds, which lasted roughly 20 minutes (see Section 2.4 for a detailed description).⁸ All respondents received similar incentives: a mobile phone airtime voucher worth US\$1.17 transferred to their phone for every fourth interview they completed, as well as after the baseline and endline interviews. The maximum individual payment is 0.3% of mean annual household income, so income effects should be negligible. The incentives were designed to encourage participation, not to precisely equalise compensation for respondent time across arms.⁹ South African mobile phone users are not charged for calls received, so respondents pay no pecuniary cost for completing the surveys.

Enumerators surveyed the same enterprise each week or month to simplify tracking. We randomly assigned enumerators to data collection groups, conditional on languages spoken. We assigned two, eight and four enumerators respectively to the monthly in-person, weekly in-person

⁸We did not use real-time panel data consistency checks to query responses that changed a lot from previous weeks. However, [Fafchamps et al. \(2012\)](#) find that “the overall impact of these consistency checks on the full sample is rather limited.”

⁹We do not test the effect of variation in incentive size. See [Singer and Ye \(2013\)](#) for a review of research into survey incentives, response rates, and data quality.

and weekly phone groups. Enumerator age, gender, experience, and language are balanced across groups. Within each group, enumerators were assigned to enterprises to allow interviews in owners' preferred language (English, seSotho, seTswana, or isiZulu) and to minimise enumerators' travel time between enterprises.

Surveys were conducted at similar times of day, during working hours. Enumerators set up an appointment time to contact their set of respondents before the first repeated interview and tried to use that time each week or month for the remainder of the panel. Enumerators confirmed the time for the next interview at the end of each interview. Enumerator assignments are collinear with treatment groups. We used only 14 enumerators, so readers may be concerned that treatment and enumerator effects are confounded. However, differences in reported outcomes across enumerators appear small. Conditioning on enumerator fixed effects increases the centered R^2 by only 0.004 to 0.078 for our main specifications, except for three variables we discuss below. All our findings are robust to controlling for enumerators' age, gender, experience, and language.

We conducted in-person interviews at the enterprises. If a respondent was scheduled to close their business, enumerators usually moved the interview to another day. All in-person interviews and 86% of phone interviews were conducted at the enterprise location (or respondent home for home-based enterprises). This difference highlights a useful feature of phone surveys – more flexibility in tracking respondents – but might induce differences in selection. We return to this issue in Section 3.

We also conducted an in-person endline interview with each enterprise owner at the enterprise location, 1-4 weeks after the repeated interviews finished. For this interview, we randomly re-assigned enumerators to enterprises. Because all enterprises are interviewed face-to-face for the endline, any differences observed in the endline data must be due to persistent frequency or medium effects from the repeated interviews.

2.3 Tracking Protocols

In this section, we describe our tracking protocols. We describe the patterns of non-response and attrition in Section 7.3 and Appendix H. We show in Section 3 and Appendices C – F that our main

results are robust to accounting for non-response.

We standardised our tracking protocol across groups to ensure that differences between groups reflect frequency and medium effects, rather than tracking effects.¹⁰ Enumerators made three attempts to contact each respondent in each scheduled week/month, as in some Living Standards Measurement Studies and Demographic and Health Surveys (Grosh and Munoz, 1996; McKenzie, 2015). The high frequency of our panel meant that we had to impose a maximum number of contact attempts. Some low-frequency panels continue to attempt to contact respondents for many months (Thomas et al., 2012). Few economic studies report detailed tracking rules so we cannot measure the prevalence of different tracking protocols.

Enumerators were supposed to complete second attempts later on the same day as the first attempt and third attempts 1-2 days after the second attempt. A contact attempt for the in-person groups meant a visit to the enterprise premises. A contact attempt for the phone groups meant talking to the respondent, so a missed call or talking to another person did not count as an attempt. After failing to interview a respondent on the third attempt, enumerators marked them as missing for that week/month. Respondents who missed an interview were always contacted in the next scheduled week or month (except if they asked not to be recontacted).¹¹

2.4 Outcome Measures

We used the same questionnaire for all repeated and endline interviews in all groups.¹² The questionnaire covered both stock variables – replacement costs for stock and inventory and for fixed assets, number of employees, number of paid employees, number of full-time employees – and flow variables – total profit, total sales, nine cost items, hours of enterprise operation, money taken by owner, goods or services by other household members. The questionnaire also asked respon-

¹⁰Stecklov et al. (2017) run a survey experiment adopting a similar strategy on non-response. We could have instead used group-specific tracking protocols that aimed to equate the response rate across groups. However, this would have required strong prior evidence about frequency, medium, and tracking effects on response rates.

¹¹We continued to interview respondents who closed or sold their enterprises using a different questionnaire. We did not track respondents who left the greater Johannesburg region, as we could only interview them by phone and did not want to break comparability between groups.

¹²The questionnaire first asked if the respondent still operated their enterprise. If not, the questionnaire asked what happened to the enterprise and asked about the respondents' current economic activities. Only 2% respondents stopped operating their enterprise during the survey period so we do not analyse data on closed enterprises.

dents if they used written records during the interview and several tracking questions. At the end of the interview, the enumerator assessed whether the respondent answered questions honestly and carefully. We show summary statistics for all outcomes in Appendix B and questionnaires are available in the supplementary materials. All flow measures used a one-week recall period except hours of operation (last day) and sales (both last week and last 4 weeks). The two sales measures allow us to test if frequency or medium effects differ by recall period.

We elicit profits directly, following De Mel et al. (2009), using the question “What was the total income the business earned last week, after paying all expenses (including wages of any employees), but not including any money that you paid yourself? That is, what were the profits of your business for last week?” This measure is more computationally intensive for the respondent. We compare this to sales minus total costs as a measure of consistency in reporting.

Costs are calculated from nine cost subcategories for the previous week: purchase of stock or inventory, wages or salaries, rent and rates for the property where the enterprise is based, repayments on enterprise loans, equipment purchases, fixing and maintaining equipment, transport costs for the enterprise, telephone and internet costs for the enterprise, and all other enterprise expenses.

3 Few Frequency or Medium Effects on Outcome Means or Distributions

In this section, we estimate frequency and medium effects on reported mean outcomes and the distribution of outcomes. We pool observations through time across enterprises and do not yet examine patterns in within-enterprise outcomes through time. Most core enterprise outcomes do not differ by frequency or medium – enterprise closure, sales, costs, and various measures of employment – or differ by small margins in the upper tails – stock/inventory, assets, and profit. There are substantial frequency and medium effects on resources taken from the enterprise by the owner or her family and on hours worked, though the latter effect may be driven by medium-induced sample selection. Only medium effects on stock and inventory, hours worked and household takings survive corrections for multiple testing. Our findings do not differ by recall period and we find little heterogeneity in treatment effects by baseline covariates.

We estimate mean effects of interview frequency and medium using

$$Y_{kit} = \beta_1 \cdot T_{1i} + \beta_2 \cdot T_{2i} + \eta_g + \phi_t + \varepsilon_{kit}, \quad (1)$$

where Y_{ki} is an outcome variable, winsorised at the 95th percentile; i , k and t index respectively enterprises, outcomes, and weeks; T_{1i} and T_{2i} are indicators for respectively the monthly in-person group and the weekly phone group; η_g is a stratification block fixed effect; and ϕ_t is a calendar week fixed effect to capture common shocks.¹³ We cluster standard errors by enterprise and test $\beta_1 = 0$, $\beta_2 = 0$, and $\beta_1 = \beta_2 = 0$. In Appendix C, we adjust our results to account for multiple testing. We estimate sharpened q -values that control the false discovery rate across all outcomes (Benjamini et al., 2006).

We report the mean effects in Table 1, along with two measures of their reliability. First, we report minimum detectable mean differences to show that these comparisons are well-powered. The medians of the minimum detectable mean differences across the binary and continuous outcomes are respectively 8 percentage points and 0.06 standard deviations.¹⁴ Second, we estimate bounds on mean effects that adjust for differences across groups in response rates, following Lee (2009). The median bounds across all binary and continuous measures allow us to rule out differences of, respectively, 11 percentage points and 0.16 standard deviations. These bounds account for differences in response rates across groups but not for the high overall level of non-response and not for any systematic relationship between baseline covariates and non-response. In Appendix C, we show that the results in this section are robust to adjusting for non-response using inverse probability of non-response weights.

We estimate distributional effects of interview frequency and medium in two ways. First, we estimate the empirical CDFs by group and show the results of quantile regressions testing for differences at prespecified quantiles $\{0.05, 0.1, 0.25, 0.5, 0.75, 0.9, 0.95\}$. We test for frequency

¹³We prespecified trimming outcomes but we subsequently chose to winsorise to reduce the loss of information from real outliers. The trimmed and winsorised results are similar. We standardise all continuous outcomes to have mean zero and standard deviation one in the monthly in-person group. We do not standardise categorical and binary measures. The categorical variables seldom have values greater than one, so we discuss treatment effects on them in percentage point terms.

¹⁴See Appendix D for an explanation of how we calculate minimum detectable effects using the observed experimental data. Note that MDEs calculated using our approach may be smaller than coefficient estimates from the sample data that are not significant at the chosen test size because we aim for 80%, rather than 100% power.

and medium effects at each quantile, using the false discovery rate to control for multiple testing across quantiles (Benjamini et al., 2006). Second, we examine effects on the outcome tails by constructing indicators for observations above the 95th percentile and using these indicators as outcomes in model 1.¹⁵ These distributional measures are important for some research questions, such as analyses of high-performing or fast-growing enterprises.

We report all mean effects in Table 1 and summarize these results in Figure 1. We display CDFs and quantile test results in Figure 2 for outcomes with significant differences at any quantile and Figure A1 for all other outcomes. We report tail effects in Table 2.

There are no frequency or medium effects on means, distributions, or shares of outliers for half our outcomes: enterprise closure; number of total, full-time and paid employees; sales over two recall periods; total costs; and enterprise money kept by the respondent. There are small and marginally statistically significant frequency effects at some higher quantiles of two outcomes: fixed asset value and profit.

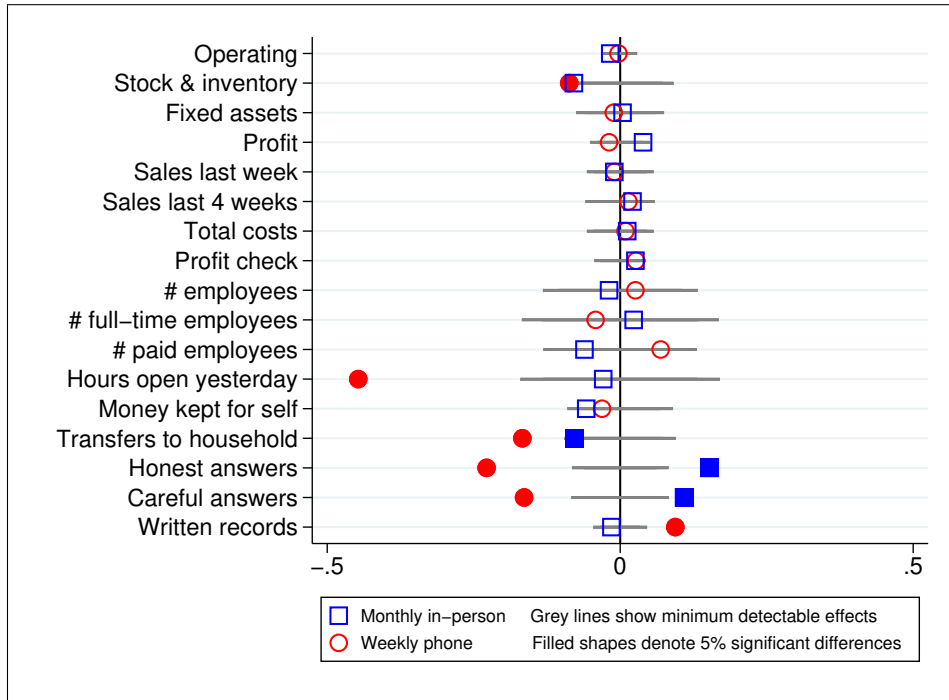
The few substantial differences we find are mostly medium, rather than frequency, effects. There are large medium effects but no frequency effects on two outcomes: phone respondents report working fewer hours and more often using written records to answer our survey. These are robust to corrections for multiple testing. The former effect is driven entirely by a higher probability of working zero hours. This result partly reflects selection induced by the location flexibility of phone surveys.¹⁶ 14% of phone interviews were completed when respondents were away from their enterprises, while in-person interviews all took place at enterprises. 44% of respondents interviewed away from their enterprises reported working zero hours the previous day, more than 20 percentage points more than respondents interviewed at their enterprises. This shows that phone interviews catch respondents who work fewer hours and were more likely to be missed by in-person interviews at enterprise locations. The higher self-reported rate of using written records by phone respondents is surprising given that they are less likely to be interviewed at the enterprise. The effect is also large, an increase of 9 percentage points from an 8 percentage point base, and

¹⁵We focus on the right tails of the outcome distributions because all our measures are truncated below at zero and have substantial numbers of zeros. Results are similar when we focus on the top 10 or 1% of the outcome distributions.

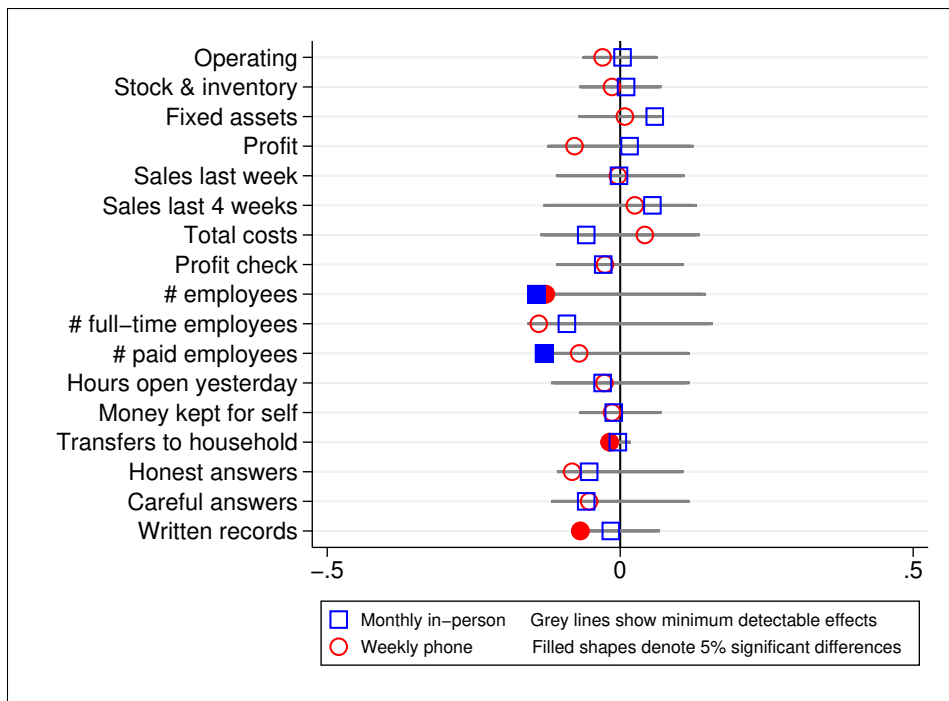
¹⁶We thank the editor for suggesting this explanation.

Figure 1: Frequency and Medium Effects on Mean Outcomes

Panel A: Repeated Interviews



Panel B: Endline Interviews



Coefficients are from regressions of each outcome, winsorised at the 95th percentile, on a vector of data collection group indicators, randomisation stratum fixed effects, and survey week fixed effects (repeated interviews only). Continuous outcomes are standardised to have mean zero and standard deviation one within survey week. Significance tests are based on heteroskedasticity-robust standard errors, clustering by enterprise (repeated interviews only). The lines for each variable show the minimum detectable differences (MDEs) between weekly and monthly in-person interviews in Panel A; the MDEs between weekly in-person and phone interviews are approximately 25% smaller. The MDEs are between weekly in-person and phone interviews in Panel B; the MDEs between weekly and monthly in-person interviews are approximately 5% smaller.

Table 1: Frequency and Medium Effects on Mean Outcomes in Repeated Interviews

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Operating	Stock & inventory	Fixed assets	Profit	Sales last week	Sales last 4 weeks	Total costs	Profit check
Monthly in-person	-0.017 (0.010)	-0.079 (0.040)*	0.004 (0.032)	0.039 (0.022)*	-0.010 (0.025)	0.021 (0.025)	0.012 (0.028)	0.026 (0.023)
Weekly by phone	-0.003 (0.006)	-0.087 (0.029)***	-0.011 (0.027)	-0.019 (0.018)	-0.010 (0.021)	0.014 (0.022)	0.009 (0.022)	0.027 (0.018)
Observations	4070	3989	3987	3986	3985	3987	3987	3984
All treatments equal (<i>p</i>)	0.262	0.011**	0.867	0.032**	0.880	0.677	0.882	0.248
MDE: Monthly in-person	0.029	0.091	0.075	0.051	0.057	0.059	0.057	0.044
MDE: Weekly by phone	0.023	0.073	0.060	0.039	0.045	0.047	0.045	0.034
Lee bound: Monthly in-person (lower)	-0.039	-0.125	-0.045	-0.002	-0.052	-0.023	-0.019	0.000
Lee bound: Monthly in-person (upper)	-0.013	0.175	0.163	0.157	0.125	0.144	0.156	0.148
Lee bound: Weekly by phone (lower)	-0.002	-0.116	-0.040	-0.041	-0.033	-0.001	-0.006	0.020
Lee bound: Weekly by phone (upper)	-0.001	-0.003	0.022	0.011	0.030	0.049	0.024	0.035

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Employees	Full-time	Paid	Hours yesterday	Money kept	Household takings	Honest	Careful	Written records
Monthly in-person	-0.019 (0.057)	0.023 (0.069)	-0.061 (0.058)	-0.029 (0.066)	-0.058 (0.033)*	-0.078 (0.031)**	0.152 (0.028)***	0.109 (0.030)***	-0.015 (0.017)
Weekly by phone	0.026 (0.051)	-0.042 (0.063)	0.069 (0.056)	-0.447 (0.058)***	-0.031 (0.030)	-0.167 (0.029)***	-0.228 (0.031)***	-0.164 (0.030)***	0.094 (0.022)***
Observations	3987	3984	3973	3987	3986	3986	4056	4056	3987
All treatments equal (<i>p</i>)	0.736	0.620	0.096*	0.000***	0.201	0.000***	0.000***	0.000***	0.000***
MDE: Monthly in-person	0.132	0.168	0.131	0.170	0.090	0.095	0.082	0.083	0.046
MDE: Weekly by phone	0.105	0.134	0.105	0.131	0.069	0.073	0.062	0.063	0.034
Lee bound: Monthly in-person (lower)	-0.081	-0.042	-0.126	-0.333	-0.125	-0.131	0.038	0.002	-0.031
Lee bound: Monthly in-person (upper)	0.254	0.359	0.213	0.202	0.155	0.160	0.226	0.200	0.052
Lee bound: Weekly by phone (lower)	0.015	-0.073	0.044	-0.621	-0.063	-0.209	-0.303	-0.229	0.069
Lee bound: Weekly by phone (upper)	0.131	0.064	0.171	-0.378	0.011	-0.026	-0.222	-0.144	0.095

Coefficients are from regressions of each outcome on a vector of data collection group indicators, randomisation stratum fixed effects, and survey week fixed effects. Continuous outcomes are standardised to have mean zero and standard deviation one in the monthly in-person group and winsorised at the 95th percentile. Owners who close their enterprises are included in regressions only for panel A column 1 and panel B columns 7 and 8. Heteroskedasticity-robust standard errors are shown in parentheses. ***, **, * and * denote significance at the 1, 5, and 10% levels.

Figure 2: Frequency and Medium Effects on Outcome Distributions in Repeated Interviews

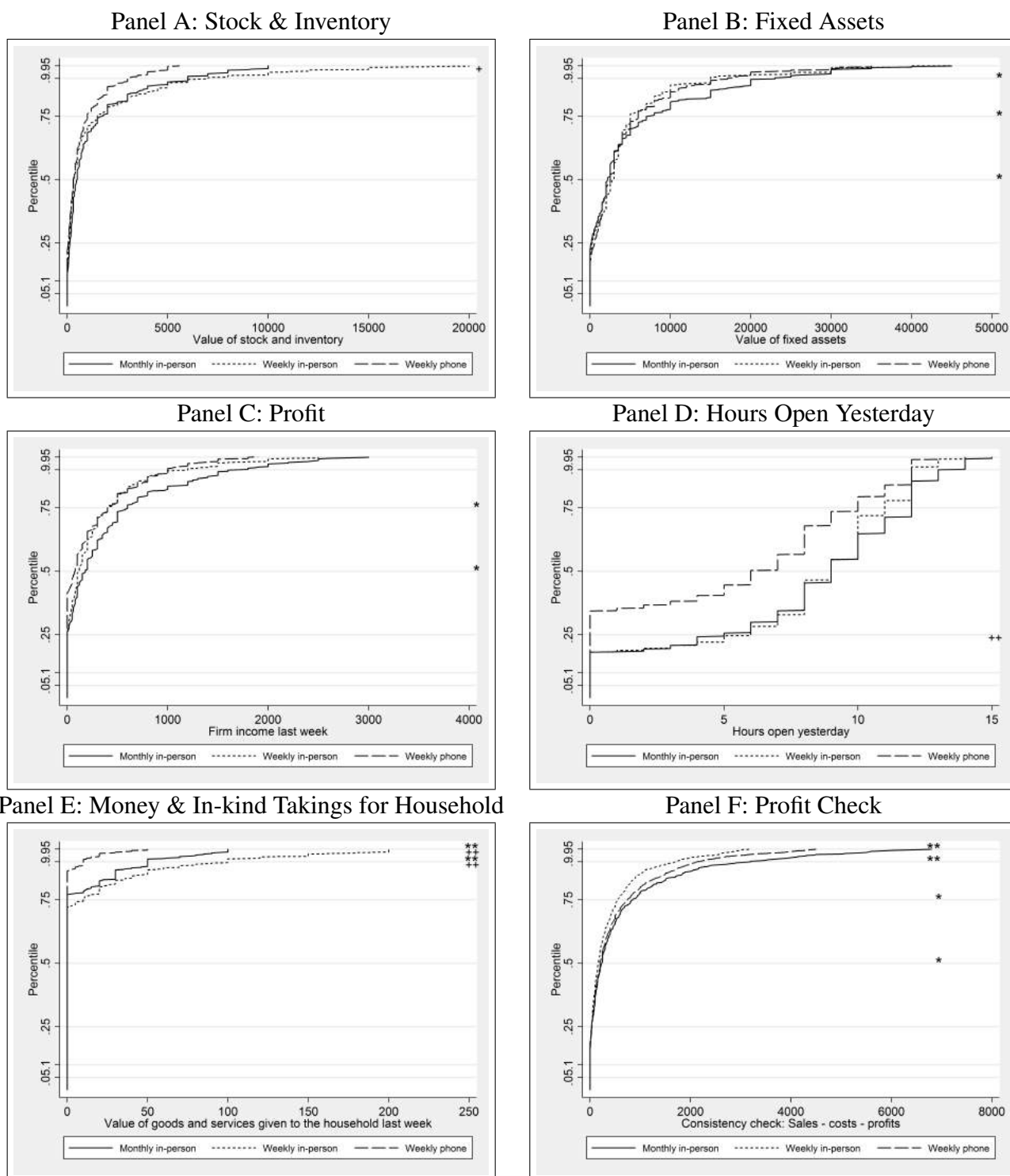


Figure shows empirical CDFs of *outcomes for which there are significant differences across groups at any prespecified quantile*. Empirical CDFs for all other outcomes – sales last week, sales in the last 4 weeks, total costs, money kept by respondent, number of employees, full-time employees and paid employees – are shown in Appendix Figure A1. We use quantile regression to test for differences at each of the quantile shown on the *y*-axis. We cluster by enterprise (Parente and Silva, 2016) and use the false discovery rate (Benjamini et al., 2006) to control for multiple testing across quantiles. + indicates a medium effect: rejection of the null hypothesis that the coefficients for weekly in-person and phone interviews are equal. * indicates a frequency effect: rejection of the null hypothesis that the coefficients for weekly and monthly in-person interviews are equal. +++/**, ++/**, and +/* denote significance at the 1, 5, and 10% levels.

Table 2: Frequency and Medium Effects on Share of Outliers in Repeated Interviews

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Operating	Stock & inventory	Fixed assets	Profit	Sales last week	Sales last 4 weeks	Total costs
Monthly in-person	-	-0.046** (0.020)	-0.005 (0.019)	0.009 (0.019)	0.009 (0.019)	0.026 (0.019)	0.013 (0.019)
Weekly by phone	-	-0.044*** (0.015)	-0.022 (0.018)	-0.014 (0.014)	-0.018 (0.016)	-0.005 (0.016)	0.010 (0.015)
Observations	-	3989	3987	3986	3985	3987	3987
All groups equal (p)	-	0.879	0.429	0.218	0.157	0.124	0.878
	(8)	(9)	(10)	(11)	(12)	(13)	(14)
	Profit check	Employees	Full-time	Paid	Hours yesterday	Money kept	Household takings
Monthly in-person	0.029 (0.018)	-0.011 (0.017)	-0.012 (0.011)	-0.017 (0.015)	0.017 (0.018)	-0.006 (0.014)	-0.039*** (0.014)
Weekly by phone	0.019 (0.013)	-0.012 (0.015)	-0.000 (0.013)	-0.002 (0.015)	-0.008 (0.012)	-0.003 (0.013)	-0.063*** (0.012)
Observations	3984	3987	3984	3973	3987	3986	3986
All groups equal (p)	0.614	0.934	0.295	0.320	0.140	0.776	0.078

Coefficients are from regressing an indicator for being in the top ventile of the distribution on treatment indicators, randomization stratum fixed effects, and survey week fixed effects. Bootstrap standard errors from 1000 iterations are shown in parentheses, resampling by enterprise. ***, **, and * denote significance at the 1, 5, and 10% levels.

predicted by enumerator fixed effects. This result is consistent with social desirability bias and lack of verifiability: respondents may report using written records to please the enumerator and phone interviews make this claim less verifiable. This is consistent with work from the US showing more social desirability bias in phone interviews (Holbrook et al., 2003). This outcome is not completely verifiable even for in-person interviews, as respondents can claim to have used written records to prepare before the interview.

There are substantial frequency and medium effects on the means and distributions of two outcomes: stock/inventory and household takings of money/goods from the enterprise. For both outcomes, weekly in-person interviews yield higher winsorised means and more right-tail outliers than monthly in-person interviews or weekly phone interviews. However, only the medium effects are robust to adjustment for multiple testing (Appendix C). The stock/inventory effect is driven by a longer right tail for weekly in-person interviews. The lower stock/inventory value in the phone group might arise if respondents avoid reporting high values when enumerators cannot visually verify the values. The stock/inventory differences are only 0.08-0.09 standard deviations but this is large in value: roughly US\$22 or 15% of mean winsorised stock/inventory value. The medium effect on household takings is driven by fewer zero values for weekly in-person interviews. This is

consistent with a social desirability bias explanation: enumerators may directly observe household members taking/receiving money/goods from the enterprises during in-person interviews but not during phone interviews. This is also consistent with respondents in phone interviews understating household takings of goods because they have just reported a lower value of stock/inventory, noted above.

There are large frequency and medium effects on two binary variables: enumerators' assessments of respondents' honesty and carefulness. These may show that respondents are most engaged in low-frequency in-person surveys and least engaged in high-frequency phone surveys. But they may also reflect enumerators' subjective impressions of the data collection methods. Consistent with the latter explanation, we find that enumerator effects, conditional on data collection method, strongly predict these two assessments. Enumerator assessments of the quality of a respondent's answers are also weakly related to the objective data quality measures we discuss in Section 4. Hence we place low weight on these two outcomes.

We explore why frequency and medium effects may differ across types of outcomes by aggregating outcomes into two indices based on two strategies for answering questions (Appendix Table A6). Respondents may give an actual count for rare events or outcomes they can easily count ('episodic enumeration') but estimate for higher-frequency events or higher-valued outcomes (Gibson and Kim, 2007). We construct a *counting index* based on number of total, full-time, and permanent employees and an *estimating index* based on values of stock/inventory, fixed assets, profit, sales, costs, money kept for the owner, household takings, and hours worked. Both indices are inverse-covariance weighted averages of the underlying variables, following Anderson (2008).

We find no frequency or medium effects on the counting index. Previous work finds that reported counting measures may be sensitive to factors like the length of the recall period (Blair and Burton, 1987; Gibson and Kim, 2007). Our non-result for the counting index may occur because neither frequency nor medium changes respondents' willingness and ability to count or because our only three counting measures are low-valued stock measures and require little counting. We find a large medium effect on the estimating index, which is 0.3 standard deviations lower for phone respondents. Half of this difference is due to the hours worked measure, discussed above.

Most of the remaining difference is due to stock/inventory and household takings, also discussed above. To the extent that respondents are estimating responses, their estimates are on average lower in phone-based interviews.

We are able to compare frequency and medium effects over different recall periods. Many survey responses are sensitive to recall periods: shorter recall periods can cause undercounting as they miss infrequent events, can cause overcounting as respondents compress events over a longer time period into the recall period ('telescoping'), or can avoid undercounting as respondents forget fewer events in short recall periods (Beegle et al., 2012; Friedman et al., 2017). Theory does not provide a clear guide to how these factors differ by survey frequency and medium. Most of our flow measures use a one-week recall period. We measure one variable, sales, over both one- and four-week recall periods. We find that the relationship between the one- and four-week sales measures does not substantially differ by medium or frequency and neither frequency nor medium effects on the two sales measures are not significantly different. We report a more detailed analysis in Appendix C. For sales at least, our conclusions are not sensitive to the recall period used.

We test for heterogeneous effects by estimating Equation (1) with interactions between the group indicators and six prespecified baseline measures. We find limited evidence of heterogeneous interview frequency or medium effects on six dimensions: respondent education, score on a digit span recall test, score on a numeracy test, keeping written records at baseline, number of employees at baseline, and gender.¹⁷ Owners with better record-keeping capacity (had multiple employees or kept written records at baseline) or better numerical skills (education, digit recall span, and numeracy scores) are no more or less susceptible to interview frequency or medium effects. There are a few scattered differences – male respondents report holding more stock and taking more money from the enterprise for their own use when interviewed weekly, and there are some scattered differences in affect by medium. However, these differences are generally imprecisely estimated and do not follow a clear pattern. Given the number of dimensions of heterogeneity we test and the generally lower power of subgroup analyses, the heterogeneity we observe may simply reflect sampling variation.

¹⁷For education, digit span recall, numeracy and the number of employees, we interact the group indicators with indicator variables equal to one for values above the baseline median.

4 Few Frequency or Medium Effects on Objective Data Quality Measures

Comparing outcome means and distributions by frequency and medium, as in Section 3, does not show which survey methods deliver higher quality data. We therefore examine two measures of data quality in this section. First, we compare the distribution of first digits in our data to a benchmark derived from Benford’s Law. Second, we compare direct and indirect measures of enterprise profit as a measure of internal consistency between survey answers. We find only small differences by frequency and medium in these two measures.

4.1 Comparing Data to Benford’s Law

Benford’s Law is a statistical regularity characterising variables in many datasets. Specifically, Benford’s Law states that the probability that the first significant digit (FSD) of a value, $j > 0$, is approximately $\log_{10}(1 + j^{-1})$. Data seldom exactly follow the distribution, but statisticians routinely use the distance between the actual distribution of FSDs and the distribution under Benford’s Law as a measure of data quality. See [Judge and Schechter \(2009\)](#) and [Schündeln \(2018\)](#) for examples comparing household surveys in developing countries to Benford’s distribution.

We use Benford’s Law to evaluate each continuous variable in our data in two ways. First, we calculate the difference between the observed FSD distribution in each data collection group and the distribution under Benford’s Law. This allows us to rank the ‘quality’ of data produced by each frequency and medium. Following [Cho and Gaines \(2007\)](#) we estimate the Euclidean distance between the distributions, $d = \sqrt{\sum_{j=1}^9 (e_j - \log_{10}(1 + j^{-1}))^2}$ where e_j is the observed share of observations with FSD j and rescale d to have maximum value 1. Second, we test for pairwise equality of the FSD distribution between data collection groups. This allows us to test if differences in data quality are statistically significant between groups. We regress nine indicators for having FSDs 1, 2, . . . , 9 on data collection group indicators using systems estimation, clustering standard errors by firm. We then test if the nine coefficients on each group indicator jointly equal their values implied by Benford’s Law. We exclude categorical measures such as the number of employees as Benford’s Law does not generally hold for low-valued integer measures.

Our data follow Benford’s Law reasonably closely (Table 3, second panel). The 24 d -statistics,

Table 3: Comparing Each Data Collection Group to Benford’s Law

	(1)	(2)	(3)	(4)
	Stock & inventory	Fixed assets	Profit	Sales last week
Panel A: Comparison of first digits across groups				
Monthly = weekly (p)	0.57	0.01	0.89	0.19
Weekly = phone (p)	0.28	0.02	0.34	0.62
Monthly = weekly = phone (p)	0.19	0.00	0.29	0.18
Panel B: Distance of first digit distribution from Benford’s Law				
Monthly in-person	0.07	0.12	0.07	0.06
Weekly in-person	0.05	0.16	0.06	0.03
Weekly phone	0.05	0.08	0.06	0.03
	(5)	(6)	(7)	(8)
	Sales last 4 weeks	Total costs	Money kept	Household takings
Panel A: Comparison of first digits across groups				
Monthly = weekly (p)	0.43	0.60	0.51	0.36
Weekly = phone (p)	0.73	0.43	0.74	0.00
Monthly = weekly = phone (p)	0.22	0.59	0.16	0.00
Panel B: Distance of first digit distribution from Benford’s Law				
Monthly in-person	0.04	0.05	0.13	0.17
Weekly in-person	0.09	0.04	0.10	0.10
Weekly phone	0.05	0.02	0.11	0.17

This table compares distributions of first significant digits (FSDs). The first three rows report p -values from Wald tests that the distributions of FSDs are equal across data collection groups. These statistics are obtained by regressing indicators for each of the nine possible FSDs on group indicators in a system of equations, clustering standard errors by enterprise, and testing if the nine coefficients on each group indicator are jointly equal across groups. The final three rows report Euclidean distances (rescaled to be $\in [0, 1]$) between the observed FSD distribution for each data collection group and the distribution under Benford’s Law, following [Cho and Gaines \(2007\)](#).

one for each continuous variable for each data collection group have interquartile range [0.05,0.10]. To contextualise this range, the statistic is bounded between 0 and 1 by construction and the d -statistics for developing country surveys reviewed by [Judge and Schechter \(2009\)](#) have interquartile range [0.05,0.13].

No single data collection group follows Benford’s Law more closely than the others (Table 3, first panel). We find substantial medium effects on fixed assets and household takings and a smaller frequency effect on fixed assets. The results for household takings should be interpreted with caution as this outcome is zero for most observations and only the positive values are used for the test against Benford’s Law. There are no significant differences across groups for the other six variables. The monthly in-person group is farthest from Benford’s Law for six of the eight variables but is never significantly farther than the weekly in-person group. Taken together, these results show that neither phone nor high-frequency interviewing leads to drops in data quality.

We also use Benford’s Law to show that enumerators’ assessment of respondents’ honesty and carefulness should be treated with caution. We estimate the normalised Euclidean distance between

the observed FSD distribution and the distribution under Benford's Law separately for interviews where the enumerator classified the respondent as honest and not honest. We then test if the FSD distributions differ between interviews classified as honest and not honest. The 'honest' interviews do not generate data whose FSD distribution is closer to Benford's Law. We repeat this exercise for interviews where the enumerator regarded the respondent as careful and not careful. The 'careful' interviews do not generate data whose FSD distribution is closer to Benford's Law. These findings echo Judge and Schetchter's evaluation of enumerators' subjective assessments using Benford's Law. This contributes to our skepticism of these subjective assessments as a data quality measure, first raised in Section 3. See Table A7 for detailed results.

Finally, we use Benford's Law to show that data quality does not decline over the life of the panel. We split the sample into observations from the first and second halves of the panel, test if the FSD distribution differs between the first and second halves, and estimate the deviation of the FSD distribution from Benford's Law in each half of the panel. The FSD distribution in the first half of the panel is not systematically closer to Benford's Law for any of the three data collection groups. This result differs from related work by Schündeln (2018), who finds that data quality in a Ghanaian household survey declines as households are surveyed more often. Schündeln examines even higher-frequency interviews (up to 10 in a single month), so we should be cautious about generalising our result to higher frequencies. See Table A7 for detailed results.

4.2 Consistency Across Multiple Profit Measures

We examine one prespecified measure of reporting consistency within the survey, the difference between two profit measures. We directly elicit profit, sales, and costs, and use these to construct a 'profit check' outcome equal to the absolute value of (sales - costs) - profits.¹⁸ This is not a direct measure of reporting accuracy because we do not observe true profits. However, consistency across two ways of eliciting profits may indicate more accurate reporting, in line with psychometricians' use of consistency across questions to measure construct validity (John and Benet-Martinez, 2014).

¹⁸The correlation between directly measured profit and sales minus costs is 0.29, similar to most studies reviewed in De Mel et al. (2009). This correlation is highest for the weekly in-person interviews but does not significantly differ by interview frequency or medium.

We find limited evidence of frequency and medium effects on reporting consistency. There are no differences across groups in the profit check means (Table 1) or shares of outliers (Table 2), though the right tail of the distribution is higher in the monthly group (Figure 2). The latter result is consistent with the idea that high-frequency surveys raise data quality by allowing respondents to practice calculating or estimating profit from sales and costs and hence avoid large discrepancies. This result does not persist after the repeated interviews end (see Section 7.4), casting some doubt on the practice hypothesis, although comparisons in the endline survey are also less well powered. There are also no differences across groups in the panel structure measures discussed in Section 5 except a slightly higher within-enterprise standard deviation of profit checks in the weekly phone group (Table A10). We conclude that the weekly in-person interviews deliver slightly more consistent measures of profits and (sales - costs), though the differences by frequency and particularly medium are small.

5 Phone Surveys Yield Higher Within-Enterprise Dispersion through Time

In Sections 3 and 4 above, we pool survey outcomes from different enterprises in the same data collection group to estimate group-specific means, distributions, and measures of quality. Researchers may also be interested in the behaviour of outcomes for the same enterprise through time. In this section, we examine the panel structure of outcomes within enterprises, using four measures of panel structure. We focus mainly on medium effects, as the monthly in-person surveys are too widely spaced to estimate frequency effects on some of these measures. We show that, on two of the four measures, phone surveys yield more dispersed data than in-person surveys. This difference is consistent with higher measurement error in phone surveys or better measurement of transient shocks in the phone surveys. We do show in Section 4.1 that this is not driven by greater fatigue-induced decline in data quality during the panel, as there is little evidence of fatigue-induced decline in data quality in either group.

First, we estimate one-week autocorrelations in outcomes and report these in Table 4. The autocorrelations are broadly consistent with the economic expectation that they should be higher for stock than flow measures: between 0.77 and 0.88 for stock measures such as assets and em-

Table 4: Panel Structure of Repeated Interview Data

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Operating	Stock & inventory	Fixed assets	Profit	Sales last week	Sales last 4 weeks	Total costs
Panel A: Autocorrelations							
Weekly in-person	-	0.873	0.829	0.628	0.750	0.764	0.713
	(-)	(0.023)	(0.031)	(0.063)	(0.038)	(0.038)	(0.057)
Weekly by phone	-	0.665	0.861	0.473	0.589	0.737	0.555
	(-)	(0.070)	(0.035)	(0.054)	(0.049)	(0.038)	(0.054)
All groups equal (p)	-	0.004	0.488	0.057	0.008	0.625	0.039
Panel B: Pr(reporting identical value for two weeks)							
Weekly in-person	0.995	0.261	0.666	0.230	0.171	0.149	0.220
	(0.002)	(0.012)	(0.013)	(0.012)	(0.011)	(0.010)	(0.012)
Weekly by phone	0.995	0.215	0.675	0.285	0.179	0.098	0.220
	(0.004)	(0.030)	(0.029)	(0.029)	(0.025)	(0.022)	(0.028)
All groups equal (p)	0.956	0.123	0.747	0.054	0.767	0.020	0.994
Observations	2431	2414	2412	2412	2412	2414	2414
	(8)	(9)	(10)	(11)	(12)	(13)	(14)
	Profit check	Employees	Full-time	Paid	Hours yesterday	Money kept	Household takings
Panel A: Autocorrelations							
Weekly in-person	0.538	0.850	0.834	0.878	0.526	0.513	0.506
	(0.066)	(0.027)	(0.034)	(0.029)	(0.038)	(0.045)	(0.047)
Weekly by phone	0.513	0.771	0.800	0.865	0.515	0.458	0.293
	(0.046)	(0.031)	(0.042)	(0.024)	(0.037)	(0.051)	(0.069)
All groups equal (p)	0.758	0.050	0.523	0.726	0.840	0.420	0.011
Panel B: Pr(reporting identical value for two weeks)							
Weekly in-person	0.108	0.925	0.948	0.950	0.454	0.376	0.688
	(0.009)	(0.007)	(0.006)	(0.006)	(0.014)	(0.014)	(0.013)
Weekly by phone	0.106	0.837	0.930	0.915	0.384	0.445	0.815
	(0.020)	(0.020)	(0.015)	(0.015)	(0.032)	(0.032)	(0.033)
All groups equal (p)	0.907	0.000	0.211	0.021	0.029	0.033	0.000
Observations	2410	2414	2411	2393	2414	2413	2412

Panel A autocorrelations are correlations between week t and $t - 1$ values for each measure, pooling observations across enterprises. Panel A standard errors in parentheses are from 1000 bootstrap iterations, resampling by enterprise. Panel B coefficients are from regressions of an indicator for no change in value between weeks t and $t - 1$ on treatment group indicators, stratification block fixed effects & week fixed effects. Panel B standard errors in parentheses are heteroskedasticity-robust and clustered by enterprise. The still operating outcome is omitted from the autocorrelation analysis because the measure has little variation, with mean = 0.98.

ployment counts; and between 0.29 and 0.76 for flow measures such as profit, sales, costs, hours worked, money kept, and household takings.

Autocorrelations are significantly lower in the phone than in-person group for six outcomes: stock and inventory, profit, sales in the last week, costs, total employees, and household takings. These include both stock and flow outcomes and both estimating and counting outcomes. This may reflect higher measurement error in the phone group or, potentially, more anchoring on past answers in the in-person group. We are not aware of any aspect of the survey administration that would induce more anchoring specifically in the in-person group, so suspect the higher intertemporal variation in the phone group reflects slightly higher measurement error.

Second, we report the group- and outcome-specific probabilities that respondents report identical values for two consecutive weeks in Table 4. This provides a test for differential anchoring on past answers by medium. The probability of reporting identical values two weeks in a row is 0.10-0.29 for flow measures such as profit, sales, and costs where we expect frequent changes. The probability is much higher, 0.67-0.93, for stock measures such as assets and employment counts. Stock/inventory behaves more like the flow than stock variables, consistent with the fact that most enterprises are very small retailers that may restock regularly. The probability is 0.38-0.82 for hours worked, money kept, and household takings because these outcomes have many persistent zeros. These results are in line with economic expectations that stock outcomes should be more persistent than flow outcomes.

The probability of reporting the same value two weeks in a row is significantly lower in the phone group for four variables: sales in the last four weeks, total employees, paid employees, and hours worked. The probability is higher in the phone group for three variables: profit, money kept and household takings. The latter difference is driven by the higher share of zeros for household takings in the phone group. On this measure of panel structure, neither medium generates consistently higher or lower dispersion across outcomes.

Third, we calculate the within-enterprise standard deviation through time for each outcome, estimate treatment effects on the standard deviations, and report these in Table A10. This is the only measure of the panel structure we construct for the monthly in-person group. Phone interviews yield higher standard deviations than in-person surveys on five of thirteen outcomes (four of which are robust to adjustment for multiple testing, as shown in Table A12) and lower standard deviations only for household takings, a measure dominated by zero values. In contrast, we do not find large or robust frequency effects on within-enterprise standard deviations. There are frequency effects on the standard deviations of only three of thirteen variables, these do not have consistent signs, and they are not robust to adjustment for multiple testing.

Fourth, we characterize the panel structure of one flow variable – log profit – and one stock variable – log capital stock – following [Blundell and Bond \(1998\)](#). We estimate dynamic panel models separately for weekly phone and weekly in-person groups, assuming an AR(1) structure

on both error terms and using two lags for profit and four lags for capital stock. We fail to reject equality of the full set of parameters across the two groups. The full results are shown in Table A13. This exercise is driven by both the mean and panel structure of the outcomes, so it is possible that the lack of medium effects on the means offset any medium effects on the panel structure.

6 Classifying Outcomes Based on Interview Frequency and Medium Effects

Combining the results from Sections 3 – 5, we can divide outcomes into five categories. First, we see no frequency effects and at most small medium effects for enterprise closure, assets, profit, our profit consistency check, the number of full-time employees, and money taken from the enterprise by the respondent.¹⁹ Second, there are no frequency or medium effects on the mean or distributions of sales over both recall periods, costs, and the numbers of total employees and paid employees, but phone surveys do generate higher within-enterprise dispersion through time. These eleven outcomes are robust to the frequencies and media we evaluate from the perspective of estimating means, distributions, or average or quantile treatment effects. But the second group of variables are more sensitive to medium choices from the perspective of estimating within-enterprise dynamics.

Third, there is a substantial medium effect, robust to multiple test correction, on the use of written records: phone respondent are more likely to self-report using written records to complete the survey. The same respondents report using written records less often when we interview them in person after the panel ends (Section 7.4). This pattern is most consistent with social desirability bias and lower verifiability in the phone interviews. This suggests researchers asking questions subject to social desirability bias and medium-specific verifiability should be cautious about medium choices. Fourth, there is a large medium effect but no frequency effect on hours worked. As discussed in Section 3, this reflects selection from our in-person interviews disproportionately missing respondents who were seldom working at their enterprises. We view this as an advantage of phone surveys relative to in-person surveys, though allowing in-person interviews at flexible locations may also achieve this.

¹⁹The value of fixed assets is difficult to classify. There are no frequency or medium effects on the mean and only small frequency effects on some quantiles of the distribution. But there are significant differences in digit distributions. The digit distributions for the two in-person groups are quite far from Benford's Law.

Fifth, there are substantial frequency and medium effects on the values of current stock/inventory and money/stock/services given to the household (‘household takings’) although only the medium effects are robust to multiple test adjustment. Their means, distributions, share of outliers, and within-enterprise dynamics are all sensitive to frequency and medium. Household takings also has a digit distribution that differs substantially from Benford’s Law, though stock/inventory does not. It is unclear why these specific variables are the most sensitive. Both are estimating, rather than counting, measures but other estimating measures are less sensitive. Stock/inventory is a stock variable, with high intertemporal persistence, while household takings is a flow variable with most mass at zero and low intertemporal persistence otherwise. It is possible that household takings is subject to the same social desirability bias and differential verifiability as using written records. But we do not observe information about the presence of household members during the interview that might allow us to test this explanation.

7 Other Considerations when Choosing Survey Frequency and Medium

In this section, we discuss four remaining considerations for researchers choosing survey frequency and medium. First, we show that outcome autocorrelations are higher for very closely-spaced surveys — hence the precision gains from averaging multiple survey rounds are decreasing in survey frequency. Second, we document that phone surveys are substantially cheaper than in-person surveys. Third, we show that permanent attrition from the panel does not differ by survey frequency or medium, but that non-response in any given survey round is higher at higher frequencies. Fourth, we show that data collection at different frequencies or using a different medium does not have persistent effects on microenterprises or their owners after the panel has ended. We survey everyone in person at endline, holding the location and survey medium constant and randomly re-assigning enumerators to treatment groups, and find few differences in outcomes between treatment groups.

7.1 Precision Gains from Multiple Measures Are Lower at Higher Frequency

Researchers sometimes collect multiple measures of enterprise performance through time to improve precision by averaging out both transient real shocks and transient measurement error (McKen-

zie, 2012). We show in this section that there are substantial precision gains from averaging high-frequency measures, especially flow measures, but that the precision gains are larger when measures are spaced farther apart.

We focus on two outcomes of particular interest to microenterprise researchers: profit and the value of fixed assets. Both have substantial intertemporal variation: the within-enterprise coefficients of variation through time have interquartile ranges of [0.72,1.41] for standardised profit and [0.17,0.82] for standardised assets. This variation may reflect transient real shocks of interest to researchers or transient measurement error. Our experiment is not designed to separate these explanations but our data quality checks in Section 4 suggest this is not entirely measurement error.

In Appendix Table A9, we show large precision gains from repeated measures, particularly for flow outcomes such as profit.²⁰ Measuring profit and fixed assets two weeks in a row reduces outcome variance by respectively 22 and 8%, while measuring them four weeks in a row reduces outcome variance by respectively 33 and 12%.

These precision gains are slightly larger with longer gaps between measures. Measuring profit and fixed assets two months in a row, rather than two weeks in a row, reduces outcome variance by respectively 27 and 11%, instead of 22 and 8%. Measuring them four months in a row rather than four week in a row reduces outcome variance by respectively 40 and 16%, instead of 33 and 12%. Precision gains are larger with longer gaps because the 4-week autocorrelations are lower than the 1-week autocorrelations for most flow and stock measures. With a fixed budget, researchers gain more power by averaging over three measures a month apart than over three measures a week apart. This is stronger evidence of non-stationarity than in the enterprise datasets reviewed in McKenzie (2012). We may find more evidence of non-stationarity because we evaluate higher-frequency panel data relative to most of the literature.

7.2 Phone Surveys Reduce Costs

Phone interviews reduce our per-interview costs by approximately 25% and larger cost savings should be possible in other settings. We calculate costs by analysing the survey firm's general

²⁰These calculations use 1- and 4-week autocorrelations from pooling the weekly in-person and phone groups, shown in columns (7) and (8) of Appendix Table A9.

ledger entries, which break expenditure down by date and purpose. We exclude the costs of the screening, baseline, and endline interviews (conducted in person for all respondents); fixed costs (e.g. office costs and management salaries); and equipment costs. Each completed phone interview cost US\$4.76 while each completed in-person interview cost US\$7.30 in the monthly group and US\$6.12 in the weekly group.²¹ All costs are per successfully completed interview. More phone than in-person interviews were missed, so this approach overstates the relative cost per attempted phone interview.

Each completed phone interview, relative to a completed interview in the weekly in-person group, saved US\$1.94 on enumerator transport and US\$0.91 on enumerator salaries but cost US\$1.21 more in airtime. The remaining cost differences are due to data capture and respondent incentives, which depend entirely on medium-specific response rates. See Figure A3 for detailed breakdown. Our cost savings are relatively low because we worked in a dense urban area with low transport costs and high airtime costs (roughly US\$1.30 per 15 minute interview). Cost savings from phone interviews will increase as the time and expense of travelling between interviews increase and as the costs of calling mobile phones decrease.

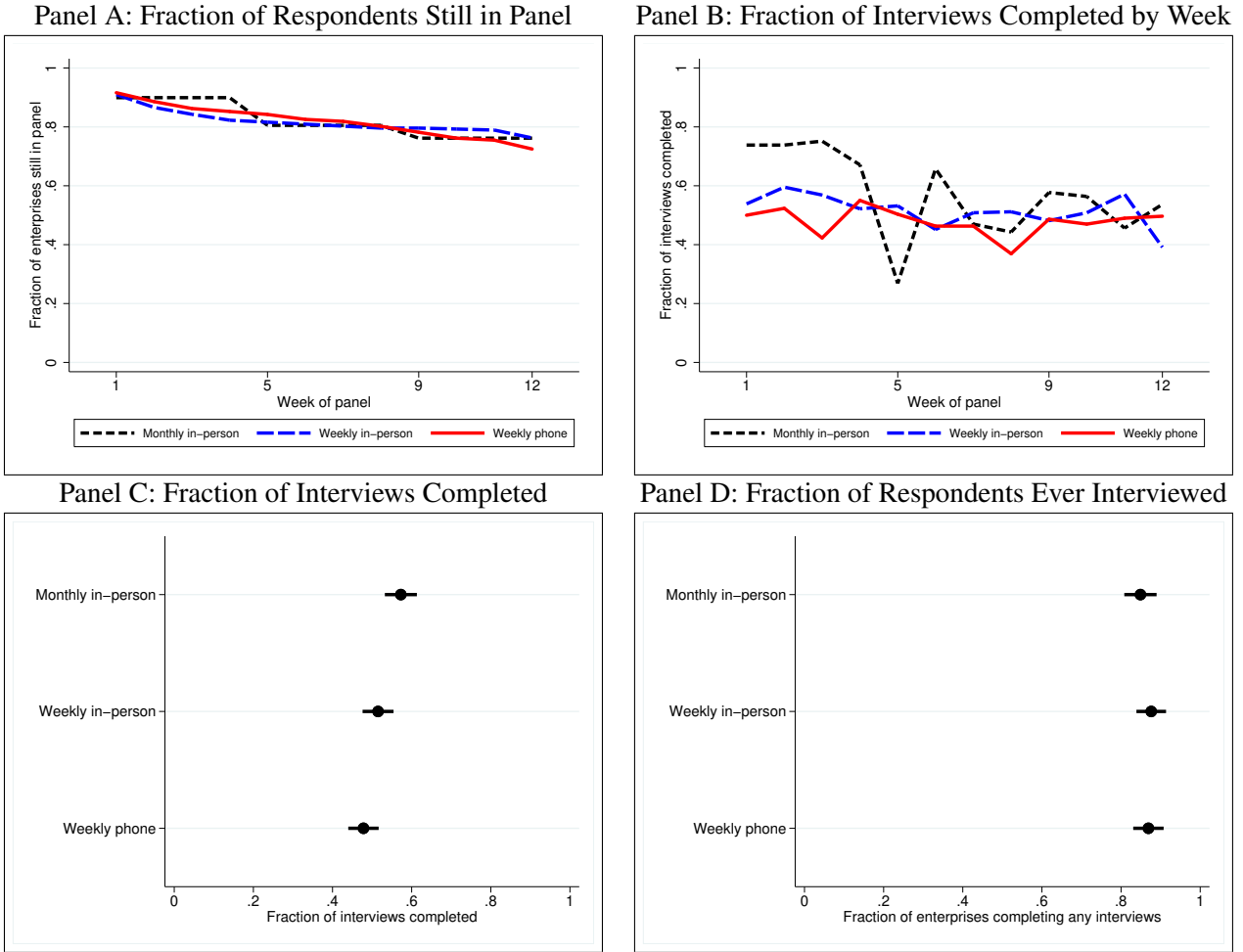
7.3 High-Frequency Measures Risk Higher Non-Response but Not Higher Attrition

Both interview frequency and medium may in principle change respondents' participation in interviews. We briefly outline differences in participation in this section and report more detailed results in Appendix H. We distinguish between two types of non-participation: *permanent attrition* and *non-response*. We define a respondent as a permanent attriter from round $t + 1$ if she is interviewed in round t but not in any round $s \geq t$, including the endline interview. Roughly 20% of respondents attrit by week 12 of the panel. This rate does not differ by frequency or medium (Figure 3, Panel A).

We define non-response as missing an interview in a specific round. Non-response is fairly high: we completed 4070 of 8058 scheduled repeated interviews (51%). There are no medium effects on non-response or on the probability of completing any interviews (Figure 3, Panels C and D). There

²¹This is similar to the per-interview cost range of US\$4.10 – US\$7.10 for mobile phone interviews in a Dar es Salaam panel study (Croke et al., 2014).

Figure 3: Response Rates and Attrition by Data Collection Group



Panel A shows the fraction of respondents in each data collection group in each week $t \in \{1, \dots, 12\}$ who are interviewed in at least one week $s \geq t$. This equals one minus the rate of permanent attrition in the panel. Panel B shows the fraction of respondents in each data collection group who are interviewed in each week. Note that the set of respondents in the monthly in-person group is different in weeks 1/5/9, 2/6/10, 3/7/11, and 4/8/12 due to the staggered start dates. Panel C shows the fraction of interviews completed by each respondent, separately by treatment group. The p -values for testing equality of this measure across groups are 0.045 for the monthly and weekly in-person groups and 0.187 for the weekly in-person and phone groups. Panel D shows the fraction of respondents that complete at least one interview, separately by treatment group. The p -values for testing equality of this measure across groups are 0.334 for the monthly and weekly in-person groups and 0.794 for the weekly in-person and phone groups.

is also little difference in the panel structure of responses: the autocorrelations in non-response in the weekly groups are -0.017 and 0.039 for, respectively, the in-person and phone groups (p -value of difference = 0.078), after conditioning on respondent-specific response rates. There is a substantial frequency effect on non-response. Respondents complete 6 percentage points more

interviews in the monthly in-person group than the weekly in-person group (Figure 3, Panel C). This difference occurs only in the first four weeks of repeated interviews (Figure 3, Panel A). This timing, and the fact that permanent attrition does not differ by frequency, shows that the frequency effect on non-response is not driven by survey fatigue or exhaustion.

Although respondents miss more weekly interviews, weekly interviews are more likely to find all respondents at least once in a given period. The fraction of respondents that are interviewed at least once in each x -week period is higher in both weekly groups than in the monthly group for all values of x (Table A15).²² This presents a trade-off: weekly interviews deliver a higher volume of information, but this information may be less representative in some weeks. If the non-response in any one period is close to random, then the greater volume of information will more than offset the lower response rate in each week. We show in Appendix H that differences in non-response by frequency are weakly related to baseline characteristics and that the marginal respondents who are captured only by higher-frequency surveys are not systematically different to the inframarginal respondents who are captured by high and low-frequency surveys.

Non-response in our weekly panel is comparable to other high-frequency surveys with representative samples (e.g. Croke et al. 2014; Gallup 2012). However, non-response in our weekly panel is higher than in surveys of samples with revealed willingness to persist in panel surveys (e.g. Arthi et al. 2018; Beaman et al. 2014; Heath et al. 2017). This reflects a potential trade-off between high response rates in the panel and gathering a representative sample, though lower-frequency panel surveys are able to achieve both goals (e.g. Thomas et al. 2012). See Appendix H for more detail on these benchmarks.

7.4 Frequency and Medium Effects on Means Do Not Persist

We test if interview frequency or medium during the 12-week panel has persistent effects on data collected several weeks after the panel has ended. Persistent effects may occur for two reasons: survey methods may persistently change how respondents report real outcomes or may actually change real outcomes, potentially by changing behaviour through reminder or salience effects. We

²²The coverage rate for monthly interviews is mechanically lower for $x < 4$. But the lower coverage rate in the monthly group over longer time periods is not mechanical and is informative.

expect that frequency effects are more likely than medium effects but we test for both. This issue is particularly important for researchers using high-frequency panels for a subsample of a broader survey sample who want to preserve comparability between the subsample and the full sample (e.g. [Franklin 2017](#)).

We conduct an endline survey one to four weeks after the panel ends, surveying respondents from all three groups in person and randomly re-assigning enumerators to respondents. Any differences in outcomes measured at this stage must reflect persistent effects of prior interview methods. We regress respondent-level outcomes on indicators for the monthly in-person and weekly phone groups, conditional on stratification block fixed effects and using heteroskedasticity-robust standard errors. We plot the estimates in [Figure 1](#) and show the detailed results in [Table A20](#).

We find few frequency or medium effects. We are powered to detect moderate differences: the median MDEs for binary and continuous outcomes are respectively 11 percentage points and 0.11 standard deviations. Monthly in-person and weekly phone respondents both report fewer employees than weekly in-person respondents, and weekly phone respondents report very slightly lower household takings and using written records less often. These differences are no longer statistically significant after adjusting for multiple testing ([Table A22](#)). These findings are not consistent with the most obvious prediction of a persistence model: frequency and medium effects during the panel should also be visible in the endline. Only the household takings result has the same sign during the panel and the magnitude in the panel is much higher. More generally, the estimated mean effects during the panel and in the endline are not similar. We have 34 mean estimates in total: phone and monthly estimates for each of 17 outcomes. The correlation between mean effects during the panel and in the endline across the 34 estimates is 0.004.

The largest persistent frequency and medium effects are on reported employment. This is driven by the share of respondents reporting zero versus one employees ([Figure A4](#)). This is a puzzling finding. It is unlikely that different survey methods induce large enough behavioural changes to shift real employment. It is possible that prior interaction with enumerators changes respondents' understanding of the definition of employee, inducing a persistent change in how they respond to this question. But the employment differences are driven by full-time and paid employees, which

are easier to define than part-time or unpaid employees. Given that these effects are not statistically significant after adjusting for multiple testing, they may simply reflect noise.

How do we reconcile these results with prior research, discussed in Section 1, which shows that participation in panel interviews can change respondents' behaviour, even over relatively short panels? A likely explanation is that behaviour change has been documented particularly in domains that are not already salient to respondents or where the surveys provide information about previously unknown options: small change management for enterprise owners (Beaman et al., 2014), savings and borrowing (Crossley et al., 2017; Stango and Zinman, 2014), water chlorination (Zwane et al., 2011), or participation in active labour market programmes (Bach and Eckman, 2018). When outcomes are already salient, such as whether a respondent has a job, being surveyed more frequently does not change reporting (Bach and Eckman, 2018; Franklin, 2017).

8 Conclusion

This paper reports the first randomised controlled trial to compare microenterprise data from surveys of different frequency and medium. To study a representative sample of microenterprises in Soweto, South Africa, randomly assigning enterprises to be interviewed in person each month, in person each week, or by phone each week.

We find three main results. First, we find few effects of frequency or medium on the means or distributions of reported outcomes. In particular, we find no substantial differences for enterprise closure, profit, sales, costs, fixed assets, or employment. We do find substantial medium effects on stock/inventory, money/goods/services given to the household, hours worked, and self-reported use of written records. Second, we use Benford's Law to show that data quality does not differ systematically between survey frequencies and media. Third, we find that phone interviews generate higher within-enterprise dispersion through time for some flow and some stock measures.

We conclude that using phone or high-frequency surveys does not systematically raise or lower the quality of microenterprise data used for cross-sectional or static panels models. However, researchers particularly interested in within-enterprise dynamics should exercise caution when choosing survey medium. These results can help researchers choose the interview frequency and

medium that generate the optimal quality and volume of data given budget constraints. In particular, our findings suggest researchers can use phone surveys to reduce costs and high-frequency surveys to collect richer panel data that captures transient shocks and informs models of intertemporal optimisation without substantially reducing data quality.

References

- ABBRING, J. AND J. HECKMAN (2007): “Econometric Evaluation of Social Programs, Part III: Distributional Treatment Effects, Dynamic Treatment Effects, Dynamic Discrete Choices, and General Equilibrium Policy Evaluation,” in *Handbook of Econometrics Volume 6B*, ed. by J. Heckman and E. Leamer, Elsevier, 5145–5303.
- ABEBE, G., S. CARIA, M. FAFCHAMPS, P. FALCO, S. FRANKLIN, AND S. QUINN (2016): “Curse of Anonymity or Tyranny of Distance? The Impacts of Job-Search Support in Urban Ethiopia,” *NBER Working Paper No. 22409*.
- ANDERSON, M. (2008): “Multiple Inference and Gender Differences in the Effects of Early Intervention: A Re-evaluation of the Abecedarian, Perry Preschool, and Early Training Projects,” *Journal of the American Statistical Association*, 103, 1481–1495.
- ARTHI, V., K. BEEGLE, J. DE WEERDT, AND A. PALACIOS-LOPEZ (2018): “Not Your Average Job: Measuring Farm Labor in Tanzania,” *Journal of Development Economics*, 130, 160–172.
- BACH, R. AND S. ECKMAN (2018): “Participating in a Panel Survey Changes Respondents’ Labour Market Behaviour,” *Journal of the Royal Statistical Society Series A*, forthcoming.
- BANERJEE, A., E. DUFLO, R. GLENNERSTER, AND C. KINNAN (2015): “The Miracle of Microfinance? Evidence from a Randomized Evaluation,” *American Economic Journal: Applied Economics*, 7, 22–53.
- BAUER, J.-M., K. AKAKPO, M. ENLUND, AND S. PASSERI (2013): “A New Tool in the Toolbox: Using Mobile Text for Food Security Surveys in a Conflict Setting,” *Humanitarian Practice Network Online Exchange*, 1–2.
- BEAMAN, L. AND A. DILLON (2012): “Do Household Definitions Matter in Survey Design? Results from a Randomized Survey Experiment in Mali,” *Journal of Development Economics*, 98, 124–135.
- BEAMAN, L., J. MAGRUDER, AND J. ROBINSON (2014): “Minding Small Change: Limited Attention among Small Firms in Kenya,” *Journal of Development Economics*, 108, 69–86.
- BEEGLE, K., J. DE WEERDT, J. FRIEDMAN, AND J. GIBSON (2012): “Methods of Household Consumption Measurement Through Surveys: Experimental Results from Tanzania,” *Journal of Development Economics*, 98, 3–18.

- BENJAMINI, Y., A. M. KRIEGER, AND D. YEKUTIELI (2006): “Adaptive Linear Step-Up Procedures that Control the False Discovery Rate,” *Biometrika*, 93, 491–507.
- BLAIR, E. AND S. BURTON (1987): “Cognitive Processes Used by Survey Respondents to Answer Behavioral Frequency Questions,” *Journal of Consumer Research*, 14, 280–288.
- BLUNDELL, R. AND S. BOND (1998): “Initial Conditions and Moment Restrictions in Dynamic Panel Data Models,” *Journal of Econometrics*, 87, 115–143.
- CAEYERS, B., N. CHALMERS, AND J. DE WEERDT (2012): “Improving Consumption Measurement and Other Survey Data through CAPI: Evidence from a Randomized Experiment,” *Journal of Development Economics*, 98, 19–33.
- CARRANZA, E., R. GARLICK, K. ORKIN, AND N. RANKIN (2018): “Job Search and Hiring with Two-sided Limited Information about Workseekers’ Skills,” Working paper.
- CHO, W. AND B. GAINES (2007): “Breaking the (Benford) Law: Statistical Fraud Detection in Campaign Finance,” *The American Statistician*, 61, 1–6.
- COLLINS, D., J. MORDUCH, S. RUTHERFORD, AND O. RUTHVEN (2009): *Portfolios of the Poor: How the World’s Poor Live on \$2 a Day*, Princeton: Princeton University Press.
- CROKE, K., A. DABALEN, G. DEMOMBYNES, M. GIUGALE, AND J. HOOGEVEEN (2014): “Collecting High Frequency Panel Data in Africa using Mobile Phone Interviews,” *Canadian Journal of Development Studies*, 35, 186–207.
- CROSSLEY, T., J. DE BRESSER, L. DELANEY, AND J. WINTER (2017): “Can Survey Participation Alter Household Saving Behaviour?” *Economic Journal*, 127, 2332–2357.
- DABALEN, A., A. ETANG, J. HOOGEVEEN, E. MUSHI, Y. SCHIPPER, AND J. VON ENGELHARDT (2016): *Mobile Phone Panel Surveys in Developing Countries: A Practical Guide for Microdata Collection*, World Bank Directions in Development.
- DAS, J., J. HAMMER, AND C. SÁNCHEZ-PARAMO (2012): “The Impact of Recall Periods on Reported Morbidity and Health Seeking Behavior,” *Journal of Development Economics*, 98, 76–88.
- DE LEEUW, E. (1992): *Data Quality in Mail, Telephone and Face to Face Surveys*, Amsterdam: TT Publikaties.
- DE MEL, S., D. MCKENZIE, AND C. WOODRUFF (2008): “Returns to Capital in Microenterprises: Evidence from a Field Experiment,” *Quarterly Journal of Economics*, 123, 1329–1372.
- DE MEL, S., D. MCKENZIE, AND C. WOODRUFF (2009): “Measuring Microenterprise Profits: Must We Ask How the Sausage is Made?” *Journal of Development Economics*, 88, 19–31.
- DE NICOLA, F. AND X. GINÉ (2014): “How Accurate are Recall Data? Evidence from Coastal India,” *Journal of Development Economics*, 106, 52–65.

- DILLON, A., E. BARDASI, K. BEEGLE, AND P. SERNEELS (2012): “Explaining Variation in Child Labor Statistics,” *Journal of Development Economics*, 98, 136–147.
- DILLON, B. (2012): “Using Mobile Phones to Collect Panel Data in Developing Countries,” *Journal of International Development*, 24, 518–27.
- DREXLER, A., G. FISCHER, AND A. SCHOAR (2014): “Keeping it Simple: Financial Literacy and Rules of Thumb,” *American Economic Journal: Applied Economics*, 6, 1–31.
- DUPAS, P., J. ROBINSON, AND S. SAAVEDRA (2018): “The Daily Grind: Cash Needs and Labor Supply,” Working paper.
- FAFCHAMPS, M., D. MCKENZIE, S. QUINN, AND C. WOODRUFF (2012): “Using PDA Consistency Checks to Increase the Precision of Profits and Sales Measurement in Panels,” *Journal of Development Economics*, 98, 51–57.
- (2014): “Microenterprise Growth and the Flypaper Effect: Evidence from a Randomized Experiment in Ghana,” *Journal of Development Economics*, 106, 211–226.
- FRANKLIN, S. (2017): “Location, Search Costs and Youth Unemployment: Experimental Evidence from Transport Subsidies,” *Economic Journal*, 128, 2353–2379.
- FRIEDMAN, J., K. BEEGLE, J. DE WEERDT, AND J. GIBSON (2017): “Decomposing Response Errors in Food Consumption Measurement: Implications for Survey Design from a Randomized Survey Experiment in Tanzania,” *Food Policy*, 72, 94–111.
- FRISON, L. AND S. POCOCK (1992): “Repeated Measures in Clinical Trials Analysis using Mean Summary Statistics and its Implications for Design,” *Statistics in Medicine*, 11, Statistics in Medicine.
- GALLUP (2012): “The World Bank Listening to LAC (L2L) Pilot: Final Report,” *Gallup Report*.
- GARLICK, R. (2019): “The Effects of Nationwide Tuition Fee Elimination on Enrollment and Attainment,” Working paper.
- GIBSON, J. AND B. KIM (2007): “Measurement Error in Recall Surveys and the Relationship Between Household Size and Food Demand,” *American Journal of Agricultural Economics*, 89, 473–489.
- GROSH, M. AND J. MUNOZ (1996): *A Manual for Planning and Implementing the Living Standards Measurement Study Survey*, The World Bank.
- GROVES, R. (1990): “Theories and Methods of Telephone Surveys,” *Annual Review of Sociology*, 16, 221–240.
- GROVES, R., P. BIEMER, L. LYBERG, J. MASSEY, W. NICHOLLS, AND J. WAKSBERG (2001): *Telephone Survey Methodology*, Wiley.
- HEATH, R., G. MANSURI, D. SHARMA, B. RIJKERS, AND W. SEITZ (2017): “Measuring Employment: Experimental Evidence from Ghana,” Working paper.

- HOLBROOK, A. L., M. C. GREEN, AND J. A. KROSNICK (2003): “Telephone vs. Face-to-Face Interviewing of National Probability Samples With Long Questionnaires: Comparisons of Respondent Satisficing and Social Desirability Response Bias,” *Public Opinion Quarterly*, 67, 79–125.
- JACOBSON, L., R. LALONDE, AND D. SULLIVAN (1993): “Earnings Losses of Displaced Workers,” *American Economic Review*, 83, 685–709.
- JOHN, O. AND V. BENET-MARTINEZ (2014): “Measurement,” in *Handbook of Research Methods in Social and Personality Psychology*, ed. by H. Reis and C. Judd, Cambridge University Press, 473–503.
- JUDGE, G. AND L. SCHECHTER (2009): “Detecting Problems in Survey Data Using Benford’s Law,” *Journal of Human Resources*, 44, 1–24.
- KARLAN, D., R. KNIGHT, AND C. UDRY (2012): “Hoping to Win, Expected to Lose: Theory and Lessons on Micro Enterprise Development,” *National Bureau of Economic Research Working Papers*, 1–54.
- KÖRMENDI, E. (2001): “The Quality of Income Information in Telephone and Face-to-Face Surveys,” in *Telephone Survey Methodology*, ed. by R. M. Groves, P. P. Biemer, L. E. Lyberg, J. T. Massey, W. L. Nicholls, and J. Waksberg, New York: John Wiley and Sons.
- LANE, S. J., N. M. HEDDLE, E. ARNOLD, AND I. WALKER (2006): “A Review of Randomized Controlled Trials Comparing the Effectiveness of Hand Held Computers with Paper Methods for Data Collection,” *BMC Medical Informatics and Decision Making*, 6, 1–10.
- LEE, D. S. (2009): “Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects,” *The Review of Economic Studies*, 76, 1071–1102.
- LEO, B., R. MORELLO, J. MELLON, T. PEIXOTO, AND S. DAVENPORT (2015): “Do Mobile Phone Surveys Work in Poor Countries?” *Centre for Global Development Working Paper Series*, 398, 1–65.
- MAHADEVAN, M. (2018): “The Price of Power: Costs of Political Corruption in Indian Electricity,” Working paper.
- MCKENZIE, D. (2012): “Beyond Baseline and Follow-up: The Case for More T in Experiments,” *Journal of Development Economics*, 99, 210–221.
- (2015): “Three Strikes and They Are Out? Persistence and Reducing Panel Attrition among Firms,” [Http://blogs.worldbank.org/impactevaluations/three-strikes-and-they-are-out-persistence-and-reducing-panel-attrition-among-firms](http://blogs.worldbank.org/impactevaluations/three-strikes-and-they-are-out-persistence-and-reducing-panel-attrition-among-firms).
- (2017): “Identifying and Spurring High-Growth Entrepreneurship: Experimental Evidence from a Business Plan Competition,” *American Economic Review*, 107, 2278–2307.
- MCKENZIE, D. AND C. WOODRUFF (2008): “Experimental Evidence on Returns to Capital and Access to Finance in Mexico,” *World Bank Economic Review*, 22, 457–82.

- MITULLAH, W. AND P. KAMA (2013): *The Partnership of Free Speech and Good Governance in Africa*, vol. 3, Cape Town: Afrobarometer, University of Cape Town.
- PAPE, U. (2018): “Informing Rapid Emergency Response by Phone Surveys,” [Http://blogs.worldbank.org/developmenttalk/informing-rapid-emergency-response-phone-surveys](http://blogs.worldbank.org/developmenttalk/informing-rapid-emergency-response-phone-surveys).
- PARENTE, P. M. AND J. M. S. SILVA (2016): “Quantile Regression with Clustered Data,” *Journal of Econometric Methods*, 5, 1–15.
- ROBINS, J. (1997): “Causal Inference from Complex Longitudinal Data,” in *Latent Variable Modeling and Applications to Causality*, ed. by M. Berkane, Springer-Verlag, 69–117.
- ROSENZWEIG, M. AND K. WOLPIN (1993): “Credit Market Constraints, Consumption Smoothing and the Accumulation of Durable Production Assets in Low-Income Countries: Investments in Bullocks in India,” *Journal of Political Economy*, 101, 223–244.
- SCHÜNDELN, M. (2018): “Multiple Visits and Data Quality in Household Surveys,” *Oxford Bulletin of Economics and Statistics*, 80, 380–405.
- SCOTT, C. AND B. AMENUVEGBE (1991): “Recall Loss and Recall Duration: An Experimental Study in Ghana,” *Inter-Stat*, 4, 31–55.
- SINGER, E. AND C. YE (2013): “The Use and Effects of Incentives in Surveys,” *Annals of The American Academy of Political and Social Science*, 645, 112–141.
- STANGO, V. AND J. ZINMAN (2014): “Limited and Varying Consumer Attention: Evidence from Shocks to the Salience of Bank Overdraft Fees,” *Review of Financial Studies*, 27.
- STECKLOV, G., A. WEINREB, AND C. CARLETTO (2017): “Can Incentives Improve Survey Data Quality in Developing Countries? Results from a Field Experiment in India,” *Journal of the Royal Statistical Society: Series A (Statistics and Society)*.
- THOMAS, D., F. WITOELAR, E. FRANKENBERG, B. SIKOKI, J. STRAUSS, C. SUMANTRI, AND W. SURIASTINI (2012): “Cutting the Costs of Attrition: Results from the Indonesia Family Life Survey,” *Journal of Development Economics*, 98, 108–123.
- TURAY, A., S. TURAY, R. GLENNESTER, K. HIMELEIN, N. ROSAS, T. SURI, AND N. FU (2015): “The Socio-Economic Impacts of Ebola in Sierra Leone: Results from a High Frequency Cell Phone Survey,” Note prepared by Statistics Sierra Leone, the World Bank, and Innovations for Poverty Action.
- VAN DER WINDT, P. AND M. HUMPHREYS (2013): “Crowdseeding Conflict Data,” *Working paper: Columbia University*.
- ZWANE, A. P., J. ZINMAN, E. VAN DUSEN, W. PARIENTE, C. NULL, E. MIGUEL, M. KREMER, D. KARLAN, R. HORNBECK, X. GINÉ, E. DUFLO, F. DEVOTO, B. CREPON, AND A. BANERJEE (2011): “Being Surveyed can Change Later Behavior and Related Parameter Estimates,” *Proceedings of the National Academy of Sciences*, 108, 1821–1826.

Call Me Maybe: Experimental Evidence on Frequency and Medium Effects in Microenterprise Surveys

Appendices for Online Publication

- Appendix [A](#) describes the sampling scheme in more detail.
- Appendix [B](#) reports summary statistics for all baseline, repeated, and endline measures.
- Appendix [C](#) reports additional analyses of interview frequency and medium effects on mean and distributions, building on Section [3](#) of the main paper.
- Appendix [D](#) discusses and derives the formulae used to estimate the *ex post* minimum detectable effect sizes reported in Section [3](#) of the main paper.
- Appendix [E](#) reports additional data quality checks, building on Section [4](#) of the main paper.
- Appendix [F](#) discusses additional features of the panel structure of outcomes, building on Sections [5](#) and [7.1](#) of the main paper.
- Appendix [G](#) reports detailed cost breakdowns for each interview frequency and medium, building on Section [7.2](#) of the main paper.
- Appendix [H](#) reports additional analyses of attrition and non-response, building on Section [7.3](#) of the main paper.
- Appendix [I](#) reports interview frequency and medium effects on mean and distributions in the endline interviews, building on Section [7.4](#) of the main paper.

As pre-specified, we analysed heterogeneity by six pre-specified dimensions: respondent education, score on a digit span recall test, score on a numeracy test, use of written records at baseline, number of employees at baseline, and gender. We find no economically meaningful patterns on any of the dimensions of heterogeneity. These results are available on request.

A Sampling and Randomisation

We constructed the sample using a three-stage clustered sampling scheme. The scheme was designed to sample 8,000 households in order to identify 1,000 eligible microenterprises, under the assumption that approximately one household in eight ran an eligible microenterprise. We used Statistics South Africa's October 2011 population census as a sampling frame. We began with a list of the 94 census 'subplaces' and 1,391 'small area layers' in Soweto. Subplaces contain 13,500 people and occupy 2.13 km² on average. Small area layers (SALs) are nested within subplaces and contain 914 people and occupy 0.14km² on average; each contains a minimum of 500 people by construction. The first stage of the sampling scheme randomly selected subplaces without replacement and with probability proportional to population size (as of October 2011) until 16,000 households were included in the sample. This realised a sample of six subplaces containing 142 SALs.

The second stage of the sampling scheme randomly ordered SALs within these subplaces without replacement and with probability proportional to population size. The random ordering defined the order in which SALs were sampled and allowed us to continue sampling until we identified 1,000 eligible microenterprises without departing from the original sampling design. We originally planned to sample SALs containing 50% of the population within each subplace (8,000 of 16,000 households) but ultimately sampled almost all SALs due to lower than expected rates of microenterprise ownership. In this second stage we excluded SALs that exceeded a pre-defined income threshold to obtain a sample of 'low-income' areas of Soweto. Specifically, we excluded SALs where more than 20% of households were in Johannesburg's top income quartile or fewer than 20% of households were in Johannesburg's bottom income quartile. The 25th and 75th percentiles of monthly household income in Johannesburg were ZAR800 and ZAR12,817 respectively or US\$101 and US\$1612 at the time of the census. These specific cutoffs were chosen in part because Statistics South Africa had only released SAL-level income bands, not microdata, at the time of the study. The bottom and top quartile rules excluded respectively 38 and 0 of the 142 SALs in the sampled subplaces.

The third stage of the sampling scheme was a census of all households within each SAL. The field team identified dwellings by visual inspection and using aerial photographs of SALs obtained from Statistics South Africa. The field team interviewed at least one person from each dwelling to identify all households living in the dwelling (or split across multiple nearby dwellings) and surveyed an available member of each household.

The probability proportional to population size sampling ensures that the sample is self-weighting. The first stage of the design (sampling subplaces) was a purely pragmatic consideration to reduce data collection costs. The sample is designed to be representative of microenterprise owners living in households in low-income areas of Soweto.

This yielded a sample of 1,046 households owning 1,081 microenterprises. In households which owned multiple eligible enterprises, we randomly selected one for the baseline. We then conducted a baseline at the enterprise premises to verify that the enterprises existed between December 2013 and February 2014, locating 895 enterprises. Of those enterprises we could not find, 67% could not be contacted using phone calls or home visits, 18% closed their enterprise between screening and baseline, 8% moved outside Soweto, 6% refused re-interview, and 1% did not answer key questions in the baseline interview.

We assigned enterprises to data collection groups using stratified random assignment. We first created strata based on gender of owner, number of employees, enterprise sector and enterprise location (Bruhn and McKenzie, 2009), yielding 149 strata with one to 51 enterprises each. We then split each stratum randomly between the three data collection groups. Residual enterprises, from strata whose size was not divisible by three, were randomly assigned to data collection groups with a restriction that a pair of residual enterprises from the same stratum would always go into separate groups. We used the census subplace in which the enterprise was located as the location block. This generally differed from the census subplace in which the household was located, which we used for the initial sampling scheme. Table A1 shows that the treatment groups are balanced on 38 of 40 measured baseline characteristics. We cannot reject joint equality of the means of all characteristics across all groups. Differences between groups are also small in magnitude. For each variable, we calculate the maximum pairwise difference between any two group means and

divide this by the standard deviation of the variable, following [Imbens \(2015\)](#). This measure is 0.08 on average and exceeds 0.2 for only 2 of the 40 variables.

B Summary Statistics

Table A1: Sample Description and Balance Test Results

	(1)	(2)	(3)	(4)	(5)	(6)
	Full Sample		Monthly	Weekly	Weekly	p-value for
	Mean	Std Dev.	In-person	In-person	Phone	balance test
Panel A: Variables Used in Stratification						
Owner age	44.8	12.7	44.5	44.7	45.2	0.805
% owners female	0.617	0.486	0.601	0.629	0.621	0.769
# employees at enterprise	0.498	0.685	0.510	0.492	0.493	0.937
% enterprises in trade	0.318	0.466	0.312	0.311	0.332	0.824
% enterprises in food	0.426	0.495	0.423	0.438	0.416	0.857
% enterprises in light manufacturing	0.103	0.304	0.104	0.100	0.104	0.985
% enterprises in services	0.088	0.284	0.094	0.084	0.087	0.904
% enterprises in agriculture/other sector	0.065	0.246	0.067	0.067	0.060	0.929
Panel B: Other Owner Demographic Variables						
% owners Black African	0.993	0.082	0.990	0.997	0.993	0.576
% owners another race	0.007	0.082	0.010	0.003	0.007	0.576
% owners from South Africa	0.923	0.267	0.916	0.936	0.916	0.533
% owners from Mozambique	0.046	0.209	0.047	0.037	0.054	0.597
% owners from another country	0.031	0.174	0.037	0.027	0.030	0.778
% owners who speak English	0.065	0.246	0.064	0.087	0.044	0.096
% owners who speak Sotho	0.213	0.410	0.211	0.217	0.211	0.979
% owners who speak Tswana	0.084	0.277	0.077	0.087	0.087	0.876
% owners who speak Zulu	0.482	0.500	0.493	0.482	0.470	0.849
% owners who speak another language	0.156	0.363	0.154	0.127	0.188	0.124
# years lived in Gauteng	40.2	16.7	39.9	40.2	40.3	0.956
# years lived in Soweto	39.2	17.2	39.3	39.3	39.1	0.990
Panel C: Other Owner Education & Experience Variables						
% with at most primary education	0.152	0.359	0.124	0.181	0.151	0.157
% with some secondary education	0.469	0.499	0.487	0.482	0.440	0.450
% with completed secondary education	0.304	0.460	0.322	0.244	0.346	0.015
% with some tertiary education	0.075	0.263	0.067	0.094	0.064	0.353
% financial numeracy questions correct	0.511	0.264	0.513	0.508	0.512	0.970
Digit recall test score	6.271	1.489	6.333	6.220	6.260	0.632
% owners with previous wage employment	0.760	0.427	0.785	0.773	0.721	0.169
Panel D: Other Owner Household Variables						
Owner's HH size	4.785	2.683	4.745	4.756	4.856	0.852
# HH members with jobs	0.720	0.979	0.728	0.716	0.715	0.984
Owner's total HH income (ZAR)	4049	4285	3994	3957	4191	0.799
% owners whose enterprise supplies $\leq 1/2$ of HH income	0.554	0.497	0.581	0.515	0.567	0.238
% owners with primary care responsible for children	0.544	0.498	0.493	0.542	0.597	0.038
% owners perceive pressure within HH to share profits	0.634	0.482	0.607	0.635	0.658	0.444
% owners perceive pressure outside HH to share profits	0.565	0.496	0.581	0.605	0.510	0.053
Panel E: Other Enterprise Variables						
Enterprise age	7.187	7.511	7.302	7.278	6.980	0.842
% enterprises registered for payroll tax or VAT	0.079	0.270	0.081	0.060	0.097	0.232
% owners who keep written financial records for enterprise	0.196	0.397	0.195	0.167	0.225	0.207
% owners who want to grow enterprise in next five years	0.762	0.426	0.752	0.766	0.768	0.876
% owners who do business by phone at least weekly	0.568	0.496	0.554	0.579	0.570	0.823
# clients for the enterprise	33.7	71.4	28.9	40.8	31.3	0.189
Sample size		895	298	299	298	
Joint balance test statistic over groups				70.9 (0.380)		
Joint balance test statistic over enumerators				793.1 (0.000)		

This table shows summary statistics for 40 variables collected in the screening and baseline interviews in columns 1 and 2. Columns 3 – 5 show the mean values of the variables for each of the three data collection groups. Column 6 shows the p -value for the test that all three groups have equal means. The first eight variables are used in the stratified random assignment algorithm and so are balanced by construction.

Table A2: **Summary Statistics for Repeated Interview Measures**

Variable	# obs	Mean	Std deviation	10 th pctile	90 th pctile
Enterprise operating	4070	0.981	0.136	.	.
Value of stock/inventory	3989	3420	16443	0	5000
Value of fixed assets	3988	24376	950625	0	17150
Profit last week	3986	539	1704	0	1200
Sales last week	3985	1282	3188	0	3000
Sales last 4 weeks	3987	3440	8184	99	7690
Costs last week	3948	1226	3398	0	2880
Sales - costs - profits last week	3945	-479	2837	-1480	500
Abs. value of sales - costs - profits last week	3945	923	2725	0	2046
# employees	3987	0.558	1.011	0	2
# full-time employees	3984	0.329	0.750	0	1
# paid employees	3973	0.446	0.954	0	2
Hours enterprise opened yesterday	3987	7.059	5.251	0	12
Money taken from enterprise for self	3986	214	588	0	500
Goods/services taken from enterprise for self/household	3986	30	178	0	50
Respondent answered honestly (interviewer assessment)	4056	0.550	0.498	.	.
Respondent answered carefully (interviewer assessment)	4056	0.899	0.301	.	.
Respondent consulted written records during survey	3987	0.119	0.323	.	.

Notes: This table shows unwinsorised summary statistics for measures collected during repeated interviews. The total cost measure is the total of nine specific cost items. The ‘honest’ and ‘careful’ answer measures are binary variables equal to one if and only if the response to on a five-point Likert scale is at or above the sample median. All other measures are taken directly from survey items. All measures except ‘operating’, ‘honest,’ and ‘careful’ are asked only if the respondent is still operating their microenterprise (98% of the sample).

Table A3: Summary Statistics for Endline Interview Measures

Variable	# obs	Mean	Std deviation	10 th pctile	90 th pctile
Enterprise operating	591	0.919	0.273	.	.
Value of stock/inventory	543	3007	17420	0	5000
Value of fixed assets	543	11368	44990	0	20000
Profit last week	543	680	3131	0	1800
Sales last week	543	1938	9519	0	4000
Sales last 4 weeks	543	3660	8122	97	8000
Costs last week	543	961	2271	0	2590
Sales - costs - profits last week	543	297	6478	-1126	755
Abs. value of sales - costs - profits last week	543	1159	6381	0	2202
# employees	543	0.519	1.095	0	2
# full-time employees	543	0.326	0.770	0	1
# paid employees	543	0.448	1.899	0	2
Hours enterprise opened yesterday	543	8.109	11.419	0	13
Money taken from enterprise for self	543	330	2243	0	600
Goods/services taken from enterprise for self/household	543	56	575	0	60
Respondent answered honestly (interviewer assessment)	590	0.647	0.478	.	.
Respondent answered carefully (interviewer assessment)	590	0.561	0.497	.	.
Respondent consulted written records during survey	543	0.101	0.302	.	.

Notes: This table shows unwinsorised summary statistics for measures collected during endline interviews. The total cost measure is the total of nine specific cost items. The ‘honest’ and ‘careful’ answer measures are binary variables equal to one if and only if the response to on a five-point Likert scale is at or above the sample median. All other measures are taken directly from survey items. All measures except ‘operating’, ‘honest,’ and ‘careful’ are asked only if the respondent is still operating their microenterprise (92% of the sample).

C Interview Frequency and Medium Effects on Means and Distributions

This appendix discusses five additional issues about interview frequency and medium effects on means and distributions, building on Section 3 in the main paper.

First, we provide more detail on interview effects on outcomes measured with different recall periods. Most of our flow measures use a one-week recall period. However, we measure sales over one- and four-week recall periods and use this to compare frequency and medium effects over different recall periods. We also construct an indirect measure of four-week sales by aggregating one-week sales reports for the 288 respondents who complete four weekly interviews in succession. Comparing the direct one-week sales measure, direct four-week sales measure, and indirect four-week sales measure shows two important patterns.

The first pattern is that reported sales are higher over shorter time periods: weekly sales are on average 35% of four-weekly sales, rather than 25%. Conditional on this level difference, the two recall periods yield relatively similar information. The direct one- and four-week sales measures have correlation 0.739. The direct and indirect four-week sales have correlation 0.787 and the interdecile range of the absolute difference is [0.011,0.901] standard deviations.

The second pattern is that the relationship between the one- and four-week sales measures does not substantially differ by medium or frequency. The correlation between the two measures is highest for the monthly in-person group but none of the pairwise differences between correlations is statistically significant ($p > 0.353$). We also explore the relationship between direct and indirect four-week sales measures by estimating model (1) with the absolute difference between direct and indirect four-week sales measures as an outcome. We find a very small medium effect of 0.018 (standard error 0.051). We cannot estimate a frequency effect on this outcome because the indirect 4 week sales measure is not observed for monthly respondents. We thus also estimate our main specification using as an outcome the direct four-week sales measure for monthly respondents and the indirect four-week sales measure for weekly respondents. We find a small frequency effect (0.102 with standard error 0.103) and a negligible medium effect (0.015 with standard error 0.104). We conclude that our data do not show differences in frequency or medium effects by recall period.

Second, we show estimates of mean frequency and medium effects using inverse probability of response weights. If response and non-response depend entirely on baseline characteristics included in the weighting model, this approach fully corrects for attrition. We construct the weights by regressing response indicators on the full list of baseline variables in Table A1 and predicting the probability of response for each respondent. We then estimate interview frequency and medium effects on mean outcomes, weighting by the inverse probability of response. Table A4 shows that estimates are largely unchanged by weighting. There are still large medium effects on hours worked and enumerator evaluations; smaller medium effects on stock/inventory, household takings, and use of written records; large frequency effects on enumerator evaluations; and minor frequency effects on stock/inventory and profit.

Third, we adjust inferences on mean frequency and medium effects to account for multiple testing. Table A5 reports the same point estimates as Table 1 in the main paper with sharpened q -values that control the false discovery rate across all seventeen outcomes. With this adjustment, there are still statistically significant medium effects on stock/inventory, hours worked, household takings, and enumerator evaluations and frequency effects on enumerator assessments, and a marginally significant effect on household takings ($q = 0.07$). Given the problems with enumerator assessments discussed in Appendix E, we can conclude that frequency effects are largely eliminated by this multiple test adjustment.

Fourth, we show outcome distributions by interview medium and frequency in Figure A1 for outcomes where there are no statistically significant frequency or medium effects. This is the companion to Figure 2 in the main paper, which shows distributions for outcomes where there are statistically significant frequency or medium effects.

Fifth, Table A6 reports mean frequency and medium effects on two indices constructed from the individual outcomes: an index of ‘counting’ measures that respondents can plausibly count directly and an index of ‘estimating’ measures that respondents are more likely to approximate. As discussed in the main paper, we see no survey method effects on the estimating index and a negative medium effect on the estimating index, driven by stock/inventory and particularly hours worked. The latter effect likely reflects the ability of phone interviews to catch respondents who

Table A4: Frequency and Medium Effects on Mean Outcomes in Repeated Interviews with Inverse Probability Weights

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Operating	Stock & inventory	Fixed assets	Profit	Sales last week	Sales last 4 weeks	Total costs	Profit check
Monthly in-person	-0.018 (0.011)*	-0.067 (0.040)*	0.015 (0.033)	0.045 (0.022)**	-0.005 (0.024)	0.028 (0.025)	0.020 (0.027)	0.034 (0.023)
Weekly by phone	-0.005 (0.006)	-0.084 (0.029)***	-0.010 (0.028)	-0.020 (0.018)	-0.011 (0.021)	0.012 (0.022)	0.009 (0.022)	0.026 (0.018)
Observations	4070	3989	3987	3986	3985	3987	3987	3984
All treatments equal (<i>p</i>)	0.220	0.014**	0.746	0.013**	0.871	0.547	0.761	0.201

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Employees	Full-time	Paid	Hours yesterday	Money kept	Household takings	Honest	Careful	Written records
Monthly in-person	-0.012 (0.058)	0.046 (0.070)	-0.046 (0.059)	-0.028 (0.067)	-0.050 (0.033)	-0.071 (0.032)**	0.151 (0.028)**	0.110 (0.030)***	-0.015 (0.017)
Weekly by phone	0.017 (0.052)	-0.048 (0.064)	0.060 (0.057)	-0.441 (0.058)***	-0.030 (0.030)	-0.171 (0.029)***	-0.224 (0.031)**	-0.159 (0.030)***	0.099 (0.022)***
Observations	3987	3984	3973	3987	3986	3986	4056	4056	3987
All treatments equal (<i>p</i>)	0.885	0.406	0.220	0.000***	0.286	0.000***	0.000***	0.000***	0.000***

Coefficients are from regressions of each outcome on a vector of data collection group indicators, randomisation stratum fixed effects, and survey week fixed effects. Regressions are weighted by the inverse probability of response. Continuous outcomes are standardised to have mean zero and standard deviation one in the monthly in-person group and winsorised at the 95th percentile. Owners who close their enterprises are included in regressions only for panel A column 1 and panel B columns 7 and 8. Heteroskedasticity-robust standard errors are shown in parentheses, clustering by enterprise. ***, **, and * denote significance at the 1, 5, and 10% levels.

Table A5: Frequency and Medium Effects on Mean Outcomes in Repeated Interviews with Multiple Testing Adjustments

Outcome	Monthly in-person			Weekly phone			All group equal	
	β	p	q	β	p	q	p	q
Operating	-0.017	0.103	0.260	-0.003	0.566	0.674	0.262	0.274
Stock & inventory	-0.079	0.051	0.219	-0.087	0.003	0.006	0.011	0.022
Fixed assets	0.004	0.891	1.000	-0.011	0.682	0.674	0.867	0.683
Profit	0.039	0.076	0.224	-0.019	0.268	0.434	0.032	0.054
Sales last week	-0.010	0.697	0.877	-0.010	0.644	0.674	0.880	0.683
Sales last 4 weeks	0.021	0.412	0.599	0.014	0.525	0.674	0.677	0.683
Total costs	0.012	0.661	0.877	0.009	0.684	0.674	0.882	0.683
Profit check	0.026	0.252	0.460	0.027	0.125	0.244	0.248	0.274
Employees	-0.019	0.747	0.877	0.026	0.615	0.674	0.736	0.683
Full-time	0.023	0.741	0.877	-0.042	0.503	0.674	0.620	0.683
Paid	-0.061	0.294	0.485	0.069	0.220	0.433	0.096	0.138
Hours yesterday	-0.029	0.663	0.877	-0.447	0.000	0.001	0.000	0.001
Money kept	-0.058	0.078	0.224	-0.031	0.297	0.434	0.201	0.258
Household takings	-0.078	0.013	0.069	-0.167	0.000	0.001	0.000	0.001
Honest	0.152	0.000	0.001	-0.228	0.000	0.001	0.000	0.001
Careful	0.109	0.000	0.003	-0.164	0.000	0.001	0.000	0.001
Written records	-0.015	0.369	0.585	0.094	0.000	0.001	0.000	0.001

This table shows, for each outcome, means in the repeated interviews in the listed groups relative to the weekly in-person group (β 's). The table also shows a p -value for the null hypothesis that each difference equals zero and a sharpened q -value that controls the probability of incorrectly rejecting all null hypotheses (Benjamini, Krieger, and Yekutieli, 2006). The q -values adjust across outcomes for a specific null hypothesis – zero frequency effect, zero medium effect, zero frequency and medium effects – but do not adjust across null hypotheses.

Figure A1: Frequency and Medium Effects on Outcome Distributions

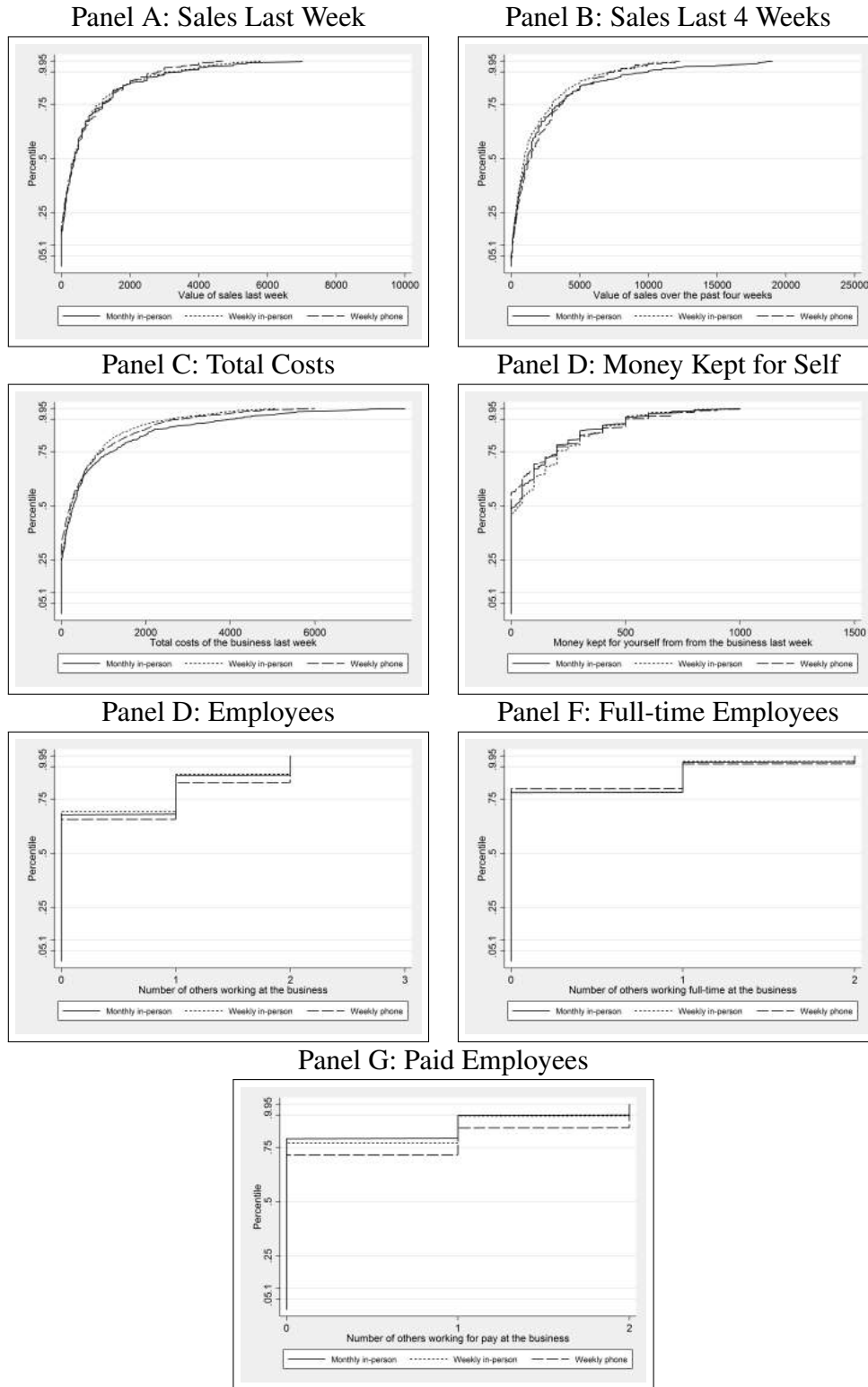


Figure shows empirical CDFs of *outcomes for which there are no significant differences across groups at any pre-specified quantile*. Empirical CDFs for all other outcomes – stock/inventory, fixed assets, profit, profit check, household takings, and hours worked – are shown in Figure 2. We use quantile regression to test for differences at each of the quantile shown on the *y*-axis. We cluster by enterprise (Parente and Silva, 2016) and use the false discovery rate (Benjamini, Krieger, and Yekutieli, 2006) to control for multiple testing across quantiles. + indicates a medium effect: rejection of the null hypothesis that the coefficients for weekly in-person and phone interviews are equal. * indicates a frequency effect: rejection of the null hypothesis that the coefficients for weekly and monthly in-person interviews are equal. +++/***, ++/**, and +/* denote significance at the 1, 5, and 10% levels.

work fewer hours and are hence harder to interview in person at their enterprise.

Table A6: Indices of Counting and Estimating Outcomes by Data Collection Group

	(1) Operating	(2) Counting measures	(3) Estimating measures	(4) Estimating excl. hours worked
Monthly in-person	-0.017 (0.010)	0.001 (0.074)	-0.019 (0.073)	-0.022 (0.077)
Weekly by phone	-0.003 (0.006)	-0.008 (0.067)	-0.311*** (0.057)	-0.170*** (0.060)
Observations	4070	3987	3989	3989
All treatments equal (p)	0.262	0.990	0.000	0.010

Coefficients are from regressing indices on treatment indicators, randomization stratum fixed effects, and survey week fixed effects. Heteroskedasticity-robust standard errors are shown in parentheses, clustering by enterprise. ***, **, and * denote significance at the 1, 5, and 10% levels. Counting index includes number of total, full-time and permanent employees. Estimating index includes values of stock/inventory, fixed assets, sales in the last week and last four weeks, costs, money kept for self, household takings, and hours worked. p-values for equal monthly effects are 0.810 for cols 1 vs 2, 0.981 for cols 1 vs 3, 0.945 for cols 1 vs 4, 0.850 for cols 2 vs 3, and 0.827 for cols 2 vs 4. p-values for equal phone effects are 0.944 for cols 1 vs 2, 0.000 for cols 1 vs 3, 0.006 for cols 1 vs 4, 0.001 for cols 2 vs 3, and 0.072 for cols 2 vs 4.

D Minimum Detectable Effect Sizes

We calculated minimum detectable sizes of the frequency and medium effects on mean outcomes before starting the experiment, as is standard in the literature. These calculations used hypothesised values of sample objects such as the response rate and outcome variance. We update the power calculations using realised values of these quantities and report the updated minimum detectable effects in tables 1 and A20.

Repeated interviews: We first discuss the minimum detectable size of β_1 , the mean outcome difference between enterprises in the weekly in-person and phone groups. We estimate the MDE for β_1 using the sample of enterprises assigned to in-person interviews and the formula $MDE_{1k} = (\tau_{1-\kappa} + \tau_{\alpha/2}) \cdot \sqrt{\frac{\sigma_{Y_k}^2}{\sigma_T^2} \cdot \frac{1}{N} \cdot \left(\rho_{Y_k} + \frac{1-\rho_{Y_k}}{N_W} \right)}$. σ_T^2 is the variance of the weekly phone group indicator T_1 , $\sigma_{Y_k}^2$ is the variance of outcome Y_k conditional on fixed effects, ρ_{Y_k} is the intra-enterprise correlation in outcome Y_k conditional on the fixed effects, N is the total number of enterprises assigned to the in-person groups and N_W is the mean number of completed interviews per enterprise. We estimate the MDE for β_2 , the mean outcome difference between enterprises in the monthly and weekly in-person groups, using the sample of enterprises assigned to weekly interviews and an analogous formula.

This approach simply updates ex ante power calculations using the realised values of σ_T^2 , $\sigma_{Y_k}^2$, ρ_{Y_k} , and N_W . We prefer this approach to calculating power for the observed coefficient estimates. The latter approach is uninformative, as the retrospective power of a test is a one-to-one function of the p -value (see [Scheiner and Gurevich \(2001\)](#) for a detailed discussion on this issue). Note that MDEs calculated using our approach may be smaller than coefficient estimates from the sample data that are not significant at the chosen test size. This occurs because we set power at 80%, rather than 100%, and because the MDE formula does not account for heterogeneous treatment effects that change $\sigma_{Y_k}^2$.

Endline interviews: We first discuss the minimum detectable size of β_1 , the mean outcome difference between enterprises in the weekly in-person and phone groups. We estimate the MDE for β_1 for outcome k using the sample of enterprises from the two in-person groups who completed

the endline and the formula $MDE_{1k} = (\tau_{1-\kappa} + \tau_{\alpha/2}) \cdot \sqrt{\sigma_{Y_k}^2 / (\sigma_T^2 \cdot N)}$. Here $\sigma_{Y_k}^2$ is the variance of Y_k conditional on the stratification fixed effects, σ_T^2 is the variance of T_{1i} , N is the number of enterprises in the endline, and we set $\tau_{1-\kappa} + \tau_{\alpha/2} = 2.8$ for a test with 5% size and 80% power. We estimate the MDE for β_2 for each outcome k using the sample of enterprises assigned to weekly interviews and an analogous formula.

E Objective Data Quality Measures

This section examines whether enumerator assessments of respondent honesty and carefulness correspond to an objective measure of data quality. The objective measure of data quality is the closeness of the observed distribution of first significant digits of outcomes to the distribution predicted under Benford's Law. The first panel of Table A7 shows that the data from interviews where the enumerator assessed the respondent as honest are not closer to the Benford's Law benchmark. The data from interviews where the enumerator assessed the respondent as careful are also not closer to the Benford's Law benchmark. We conclude that enumerator assessments provide a weak measure of data quality.

This section also examines whether data quality declines over time, potentially due to respondent fatigue and frustration. The second panel of Table A7 shows that data from interviews in the first half of the panel are not closer to the Benford's Law benchmark relative to interviews from the second half of the panel. We conclude that data quality does not appear to decline over the life of the panel.

Table A7: Evaluating Enumerator Assessments and Respondent Fatigue/Conditioning Using Benford's Law

	(1) Stock & inventory	(2) Fixed assets	(3) Profit	(4) Sales last week	(5) Sales last 4 weeks	(6) Total costs	(7) Money kept	(8) Household takings
Honest = not honest (<i>p</i>)	0.29	0.00	0.03	0.68	0.23	0.38	0.71	0.50
Honest: Benford distance (<i>d</i>)	0.06	0.12	0.04	0.03	0.07	0.03	0.10	0.08
Not honest: Benford distance (<i>d</i>)	0.05	0.10	0.09	0.02	0.06	0.05	0.12	0.10
Careful = not careful (<i>p</i>)	0.55	0.00	0.00	0.15	0.48	0.55	0.78	0.80
Careful: Benford distance (<i>d</i>)	0.06	0.13	0.04	0.03	0.07	0.03	0.10	0.08
Not careful: Benford distance (<i>d</i>)	0.05	0.09	0.08	0.03	0.07	0.04	0.11	0.09
<i>Monthly in-person:</i>								
First half = Second half (<i>p</i>)	0.31	0.38	0.63	0.72	0.75	0.30	0.00	0.00
First half: Benford distance (<i>d</i>)	0.08	0.15	0.07	0.07	0.04	0.07	0.15	0.21
Second half: Benford distance (<i>d</i>)	0.07	0.09	0.08	0.08	0.08	0.08	0.11	0.20
<i>Weekly in-person:</i>								
First half = Second half (<i>p</i>)	0.05	0.26	0.02	0.84	0.87	0.54	0.69	0.05
First half: Benford distance (<i>d</i>)	0.07	0.17	0.07	0.03	0.09	0.05	0.09	0.11
Second half: Benford distance (<i>d</i>)	0.04	0.15	0.07	0.04	0.09	0.04	0.12	0.10
<i>Weekly phone:</i>								
First half = Second half (<i>p</i>)	0.62	0.55	0.00	0.63	0.28	0.61	0.92	0.00
First half: Benford distance (<i>d</i>)	0.07	0.06	0.07	0.05	0.07	0.03	0.12	0.14
Second half: Benford distance (<i>d</i>)	0.05	0.09	0.10	0.03	0.05	0.05	0.10	0.23

This table reports analyses of the distributions of first significant digits (FSDs). The first six rows compare enumerators' assessments of data quality to data quality measures based on Benford's Law. The first row reports *p*-values from Wald tests that the distribution is equal between interviews in which the enumerator assessed that the respondent was 'honest' and surveys in which the enumerator did not. These statistics are obtained by regressing indicators for each of the nine possible FSDs on 'honest' and 'not honest' indicators using a system of equations, clustering standard errors by enterprise, and testing if the nine coefficients are jointly equal across indicators. The second and third rows report Euclidean distances (rescaled to be $\in [0, 1]$) between the observed FSD distributions for 'honest' and 'not honest' interviews and the distribution under Benford's Law, following [Cho and Gaines \(2007\)](#). The fourth, fifth, and sixth rows report the same results for interviews in which the enumerator assessed that the respondent was 'careful' and surveys in which the enumerator did not. The next nine rows compare data quality in the first and second halves of the panel against Benford's Law to test for fatigue or panel conditioning effects.

F Panel Structure of Outcomes at Different Interview Frequencies

This appendix presents four additional analyses of interview frequency effects on the panel structure of outcomes.

First, Tables A8 and A9 inform the discussion of frequency effects on panel structure in Section 5 and precision gains from combining multiple surveys in Section 7.1. The former table reports 4-week autocorrelations in outcomes and tests if these are different across data collection groups. The latter table tests if, within each data collection group, the 1- and 4-week autocorrelations are different to each other.

Second, Table A10 shows frequency and medium effects on within-enterprise standard deviations of outcomes through time, discussed in Section 5. Table A11 shows estimates with inverse probability of response weights. The weights make little difference to the estimated frequency or medium effects, with only the frequency effect on employees losing significance. Table A12 shows the same estimates with sharpened q -values that control the false discovery rate across all seventeen outcomes. The multiple testing adjustment makes little difference to the frequency effects but renders almost all medium effects statistically insignificant at conventional levels.

Third, Figure A2 shows the distribution of the within-enterprise interquartile range during the repeated interviews for log profit and log capital stock. This depicts the substantial variation through time that is better captured by high- than low-frequency interviews.

Fourth, Table A13 reports parameter estimates of dynamic panel models following Blundell and Bond (1998), as discussed in Section 5. These models allow us to characterise the panel structure of one flow variable – log profit – and one stock variable – log capital stock. We report parameter estimates separately for weekly phone and weekly in-person groups, assuming an AR(1) structure on both error terms and using two lags for profit and four lags for capital stock. We fail to reject equality of the full set of parameters across the two groups.

Table A8: 4-week Autocorrelations by Data Collection Group

	(1)	(2)	(3)	(4)	(5)
	Monthly in-person	Weekly in-person	Weekly phone	p: monthly = weekly	p: in-person = phone
Operating	-	-	-	-	
	(-)	(-)	(-)		
Stock & inventory	0.799	0.810	0.506	0.873	0.001
	(0.059)	(0.059)	(0.091)		
Fixed assets	0.717	0.784	0.780	0.490	0.961
	(0.088)	(0.088)	(0.052)		
Profit	0.550	0.504	0.439	0.713	0.461
	(0.107)	(0.107)	(0.059)		
Sales last week	0.562	0.656	0.580	0.410	0.310
	(0.101)	(0.101)	(0.053)		
Sales last 4 weeks	0.733	0.663	0.687	0.502	0.721
	(0.089)	(0.089)	(0.038)		
Total costs	0.637	0.640	0.580	0.981	0.474
	(0.112)	(0.112)	(0.056)		
Profit check	0.574	0.521	0.512	0.622	0.914
	(0.091)	(0.091)	(0.054)		
Employees	0.657	0.814	0.662	0.033	0.002
	(0.068)	(0.068)	(0.040)		
Full-time	0.675	0.775	0.744	0.287	0.589
	(0.085)	(0.085)	(0.044)		
Paid	0.706	0.847	0.773	0.079	0.112
	(0.076)	(0.076)	(0.036)		
Hours yesterday	0.372	0.356	0.383	0.854	0.672
	(0.076)	(0.076)	(0.044)		
Money kept	0.084	0.463	0.393	0.000	0.365
	(0.090)	(0.090)	(0.055)		
Household takings	0.193	0.427	0.205	0.027	0.002
	(0.095)	(0.095)	(0.050)		

This table shows 4-week autocorrelations by group and tests if these are equal across groups. This table informs the in-text discussion of stability of panel structure across groups. Autocorrelations are calculated as correlations between week t and $t - 4$ values for each measure, pooling observations across enterprises. We account for missing values in week $t - 4$ by using the value from week $t - 3$ or $t - 5$ or their average if both are observed. Standard errors in parentheses are from 1000 bootstrap iterations, resampling by enterprise. The still operating outcome is omitted from the autocorrelation analysis because the measure has little variation, with mean = 0.98.

Table A9: Differences between 1-week and 4-week Autocorrelations by Data Collection Group

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Weekly in-person group			Weekly phone group			Both weekly groups		
	1-week	4-week	p:equality	1-week	4-week	p:equality	1-week	4-week	p:equality
Operating	-	-	-	-	-	-	-	-	-
	(-)	(-)		(-)	(-)		(-)	(-)	
Stock & inventory	0.873	0.810	0.003	0.665	0.506	0.012	0.821	0.730	0.000
	(0.023)	(0.032)		(0.070)	(0.094)		(0.026)	(0.036)	
Fixed assets	0.829	0.784	0.114	0.861	0.780	0.034	0.844	0.781	0.007
	(0.032)	(0.043)		(0.035)	(0.055)		(0.024)	(0.034)	
Profit	0.628	0.504	0.002	0.473	0.439	0.482	0.561	0.469	0.005
	(0.061)	(0.066)		(0.053)	(0.059)		(0.043)	(0.045)	
Sales last week	0.750	0.656	0.011	0.589	0.580	0.849	0.679	0.622	0.049
	(0.037)	(0.051)		(0.048)	(0.055)		(0.031)	(0.038)	
Sales last 4 weeks	0.764	0.663	0.004	0.737	0.687	0.072	0.752	0.675	0.000
	(0.036)	(0.051)		(0.038)	(0.040)		(0.026)	(0.032)	
Total costs	0.713	0.640	0.035	0.555	0.580	0.623	0.637	0.609	0.369
	(0.054)	(0.063)		(0.057)	(0.059)		(0.041)	(0.043)	
Profit check	0.538	0.521	0.655	0.513	0.512	0.983	0.526	0.517	0.773
	(0.061)	(0.058)		(0.047)	(0.055)		(0.038)	(0.040)	
Employees	0.850	0.814	0.160	0.771	0.662	0.000	0.809	0.733	0.000
	(0.027)	(0.032)		(0.033)	(0.038)		(0.021)	(0.025)	
Full-time	0.834	0.775	0.028	0.800	0.744	0.134	0.816	0.759	0.012
	(0.033)	(0.037)		(0.041)	(0.044)		(0.026)	(0.028)	
Paid	0.878	0.847	0.142	0.865	0.773	0.000	0.872	0.807	0.000
	(0.028)	(0.031)		(0.026)	(0.034)		(0.018)	(0.024)	
Hours yesterday	0.526	0.356	0.000	0.515	0.383	0.000	0.547	0.403	0.000
	(0.036)	(0.044)		(0.039)	(0.044)		(0.026)	(0.030)	
Money kept	0.513	0.463	0.299	0.458	0.393	0.166	0.485	0.423	0.058
	(0.044)	(0.054)		(0.052)	(0.055)		(0.034)	(0.039)	
Household takings	0.506	0.427	0.057	0.293	0.205	0.209	0.483	0.391	0.006
	(0.050)	(0.053)		(0.073)	(0.052)		(0.045)	(0.046)	

This table shows 1- and 4-week autocorrelations by group and tests if these are equal, which holds if outcomes are covariance-stationary. This table informs the in-text calculation of variance reductions from collecting multiple measures of enterprise outcomes at different points in time. Autocorrelations are calculated as correlations between week t and either $t - 1$ or $t - 4$ values for each measure, pooling observations across enterprises. We account for missing values in week $t - 4$ by using the value from week $t - 3$ or $t - 5$ or their average if both are observed. Columns (7)-(9) pool observations from weekly in-person and weekly phone groups. Standard errors in parentheses are from 1000 bootstrap iterations, resampling by enterprise. The still operating outcome is omitted from the autocorrelation analysis because the measure has little variation, with mean = 0.98.

Table A10: Frequency and Medium Effects on Within-Respondent Standard Deviations in Repeated Interviews

	(1) Stock & inventory	(2) Fixed assets	(3) Profit	(4) Sales last week	(5) Sales last 4 weeks	(6) Total costs	(7) Profit check
Monthly in-person	-0.049 (0.020)**	-0.005 (0.015)	0.001 (0.017)	-0.020 (0.014)	-0.015 (0.013)	-0.004 (0.016)	-0.002 (0.016)
Weekly by phone	-0.009 (0.016)	0.005 (0.012)	0.009 (0.013)	0.005 (0.012)	0.008 (0.011)	0.029 (0.014)**	0.025 (0.012)**
Observations	610	609	608	608	608	608	608
All treatments equal (<i>p</i>)	0.049**	0.787	0.768	0.200	0.186	0.059*	0.091*

	(1) Employees	(2) Full-time	(3) Paid	(4) Hours yesterday	(5) Money kept	(6) Household takings
Monthly in-person	0.078 (0.037)**	0.043 (0.047)	0.008 (0.033)	0.051 (0.051)	-0.007 (0.030)	-0.104 (0.030)***
Weekly by phone	0.146 (0.029)***	0.027 (0.035)	0.080 (0.029)***	0.117 (0.035)***	0.014 (0.025)	-0.120 (0.027)***
Observations	608	608	608	608	608	608
All treatments equal (<i>p</i>)	0.000***	0.569	0.018**	0.003***	0.774	0.000***

Coefficients are from regressions of the within-enterprise standard deviation of each outcome through time on a vector of data collection group indicators and randomisation stratum fixed effects. Continuous outcomes are standardised to have mean zero and standard deviation one in the monthly phone group and winsorised at the 95th percentile before calculating the standard deviations. Heteroskedasticity-robust standard errors are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels.

Table A11: Frequency and Medium Effects on Within-Respondent Standard Deviations in Repeated Interviews with Inverse Probability Weights

	(1) Stock & inventory	(2) Fixed assets	(3) Profit	(4) Sales last week	(5) Sales last 4 weeks	(6) Total costs	(7) Profit check
Monthly in-person	-0.054 (0.021)**	-0.014 (0.016)	-0.004 (0.017)	-0.020 (0.014)	-0.018 (0.014)	-0.008 (0.017)	-0.007 (0.016)
Weekly by phone	-0.009 (0.016)	0.006 (0.013)	0.009 (0.013)	0.003 (0.012)	0.007 (0.011)	0.027 (0.014)**	0.023 (0.012)*
Observations	610	609	608	608	608	608	608
All treatments equal (<i>p</i>)	0.031**	0.471	0.682	0.257	0.167	0.063*	0.096*

	(1) Employees	(2) Full-time	(3) Paid	(4) Hours yesterday	(5) Money kept	(6) Household takings
Monthly in-person	0.039 (0.037)	0.037 (0.050)	0.008 (0.035)	0.047 (0.051)	-0.005 (0.030)	-0.114 (0.030)***
Weekly by phone	0.141 (0.028)***	0.024 (0.036)	0.077 (0.029)***	0.122 (0.034)***	0.013 (0.026)	-0.125 (0.026)***
Observations	608	608	608	608	608	608
All treatments equal (<i>p</i>)	0.000***	0.669	0.023**	0.002***	0.819	0.000***

Coefficients are from regressions of the within-enterprise standard deviation of each outcome through time on a vector of data collection group indicators and randomisation stratum fixed effects, weighted by the inverse probability of response. Continuous outcomes are standardised to have mean zero and standard deviation one in the monthly phone group and winsorised at the 95th percentile before calculating the standard deviations. Heteroskedasticity-robust standard errors are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels.

Table A12: Frequency and Medium Effects on Within-Respondent Standard Deviations in Repeated Interviews with Multiple Testing Adjustments

Outcome	Monthly in-person			Weekly phone			All group equal	
	β	p	q	β	p	q	p	q
Stock & inventory	-0.049	0.017	0.114	-0.009	0.570	0.583	0.049	0.097
Fixed assets	-0.005	0.741	1.000	0.005	0.664	0.583	0.787	0.405
Profit	0.001	0.968	1.000	0.009	0.502	0.583	0.768	0.405
Sales last week	-0.020	0.151	0.610	0.005	0.684	0.583	0.200	0.184
Sales last 4 weeks	-0.015	0.258	0.970	0.008	0.446	0.583	0.186	0.184
Total costs	-0.004	0.786	1.000	0.029	0.035	0.067	0.059	0.097
Profit check	-0.002	0.892	1.000	0.025	0.045	0.073	0.091	0.134
Employees	0.078	0.035	0.150	0.146	0.000	0.001	0.000	0.001
Full-time	0.043	0.359	0.970	0.027	0.437	0.583	0.569	0.405
Paid	0.008	0.815	1.000	0.080	0.006	0.015	0.018	0.046
Hours yesterday	0.051	0.318	0.970	0.117	0.001	0.003	0.003	0.013
Money kept	-0.007	0.820	1.000	0.014	0.583	0.583	0.774	0.405
Household takings	-0.104	0.001	0.008	-0.120	0.000	0.001	0.000	0.001

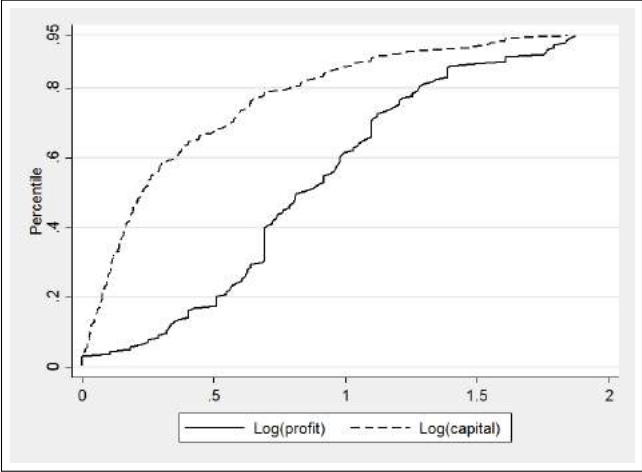
This table shows, for each outcome, standard deviations in the repeated interviews in the listed groups relative to the weekly in-person group (β 's). The table also shows a p -value for the null hypothesis that each difference equals zero and a sharpened q -value that controls the probability of incorrectly rejecting all null hypotheses (Benjamini, Krieger, and Yekutieli, 2006). The q -values adjust across outcomes for a specific null hypothesis – zero frequency effect, zero medium effect, zero frequency and medium effects – but do not adjust across null hypotheses.

Table A13: Estimates of Dynamic Panel Structure, Following Blundell-Bond

	(1) Profit <i>In-person</i>	(2) Profit <i>Phone</i>	(3) Capital <i>In-person</i>	(4) Capital <i>Phone</i>
Lag 1	0.477*** (0.070)	0.363*** (0.079)	0.459*** (0.115)	0.380*** (0.103)
Lag 2	0.305*** (0.114)	0.261*** (0.076)	0.194*** (0.075)	0.195* (0.115)
Lag 3			0.158*** (0.048)	0.137*** (0.046)
Lag 4			0.175*** (0.065)	0.104*** (0.036)
Time dummies	✓	✓	✓	✓
Observations	520	289	420	418
Enterprises	132	87	116	89
Arellano-Bond: AR(1) (<i>p</i> -value)	0.000***	0.000***	0.024**	0.003***
Arellano-Bond: AR(2) (<i>p</i> -value)	0.302	0.168	0.325	0.140
Arellano-Bond: AR(3) (<i>p</i> -value)	0.340	0.319	0.783	0.967
Arellano-Bond: AR(4) (<i>p</i> -value)	0.242	0.392	0.323	0.497
Hansen test (<i>p</i> -value)	0.199	0.988	0.288	0.596
H_0 : Equal parameter estimates (<i>p</i> -value)		0.548		0.907

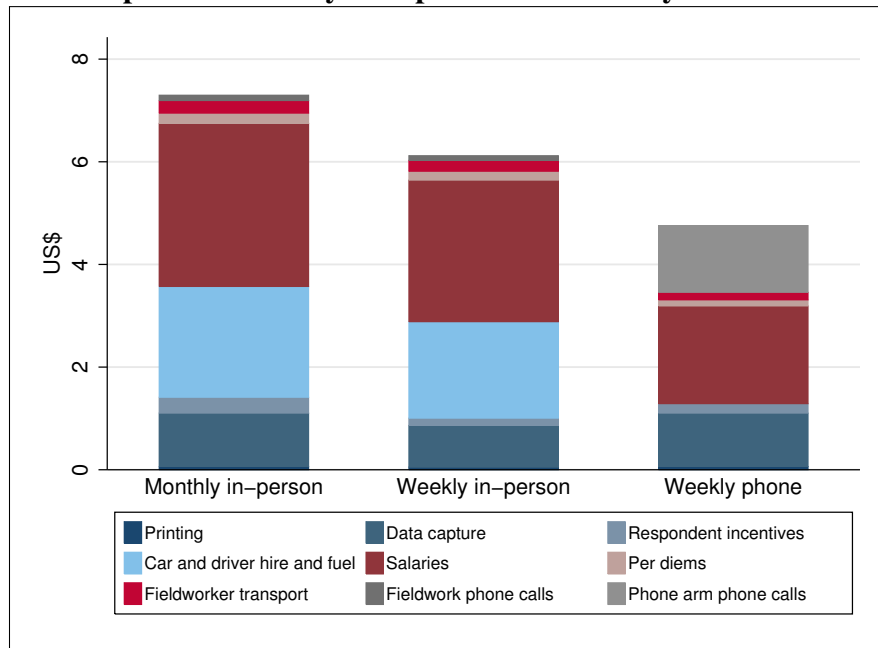
This table shows the parameter estimates for a model of the panel structure of log profit and log capital stock, estimated separately for the weekly in-person and weekly phone groups. Heteroskedasticity-robust standard errors are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels.

Figure A2: **Within-enterprise Dispersion of Repeated Interview Measures**



Notes: We calculate the interquartile range for each of log profit and log capital stock over the repeated interview phase for each enterprise. We show the CDF of this measure here. Capital stock includes both stock/inventory and fixed assets. The range for profit exceeds one log point for 40% of firms and the range for capital exceeds one log point for 15% of firms.

Figure A3: **Cost per Successfully Completed Interview by Data Collection Group**



This figure shows the average cost per completed interview by data collection group. The average cost is constructed by summing the nine cost categories shown above, separately by data collection group, then dividing this by the number of successfully completed interviews in each group. This calculation excludes the costs of the in-person screening, baseline, and endline interview; and excludes fixed costs such as office rental and management salaries. US\$ values are calculated using at the South African Reserve Bank exchange rate on 31 August 2013: 1 US\$D = 10.27 ZAR.

G More Detail on Cost Analysis

Phone interviews reduce our per-interview costs by approximately 25% and larger cost savings should be possible in other settings. As outlined in Section 7.2, we calculate costs by analysing the survey firm’s general ledger entries, which break expenditure down by date and purpose. We can easily allocate most ledger entries to treatment arms: enumerator salaries, per diems and transport allowances are reported for each enumerator and enumerators worked in only one treatment group; one car was hired for each of the weekly in-person and the monthly in-person interview teams; and phone calls to arrange in-person interviews and phone interviews were conducted on different phone lines. We allocate the total cost for printing, data capture, and respondent incentives across the three arms by dividing the total cost by the number of successfully completed interviews in each arm. More phone than in-person interviews were missed, so this approach overstates the relative cost per *attempted* phone interview.

Figure A3 shows that our largest cost saving came from transport costs. Enumerators doing phone interviews received an allowance for transport to the office, at a cost of US\$0.15 per phone interview, while enumerators doing in-person interviews met in a central place close to their houses and were transported to the enterprises, costing US\$2.09 per in-person interview. Phone surveys also saved costs on enumerator salaries and per diems (US\$2.02 for phone and US\$2.93 for in-person). Enumerators did not need to travel to respondents, so fewer enumerators were needed. Enumerators were paid the same daily rate and per diem for phone and in-person interviews to avoid differences in motivation and incentives.

Our cost savings are relatively low because we worked in a dense urban area with low transport costs and high airtime costs (roughly US\$1.30 per 15 minute interview). Cost savings from phone interviews will increase as the time and expense of travelling between interviews increase and as the costs of calling mobile phones decrease. For example, a Tanzanian high-frequency household survey with farmers in remote rural areas spent US\$97 per in-person baseline interview and US\$7 per phone follow-up interview (Dillon, 2012). Phone interviews will also yield larger cost savings from renting, losing, and breaking fewer mobile devices for data collection.

H Interview Frequency and Medium Effects on Attrition and Non-response

Both interview frequency and medium may change survey response rates. More frequent interviews might frustrate respondents, leading to lower response rates, or build closer rapport between respondents and enumerators, leading to higher response rates. Phone interviews make it easier to contact respondents, leading to higher response rates, or change the nature of respondent-enumerator interactions, with ambiguous effects on response rates.

We distinguish between *permanent attrition* and *non-response*. We define a respondent as a permanent attriter from round $t + 1$ if she is interviewed in round t but not in any round $s \geq t$, including the endline interview. Roughly 20% of respondents attrit by week 12 of the panel. This rate does not differ by frequency or medium (Figure 3, Panel A). We include the endline interview in the set of possible s but the differences between the data collection groups are robust to excluding the endline.

We define non-response as missing an interview in a specific round. A respondent who attrits in round t is also a non-responder in all subsequent rounds. Non-response is fairly high: we completed 4070 of 8058 scheduled repeated interviews (51%). There are no medium effects on non-response. Respondents in both weekly groups complete roughly 50% of scheduled repeated interviews, 87% of respondents complete at least one interview, and 7% of respondents complete all interviews (Table A14). The response rate differs between the two groups in 4 of the 12 weeks but there is no clear time pattern to these differences (Figure 3, Panel B). There is also little difference in the panel structure of responses: the autocorrelations in non-response in the weekly groups are -0.017 and 0.039 for respectively the in-person and phone groups (p -value of difference = 0.078), after conditioning on respondent-specific response rates.

Although response rates do not differ by medium, our in-person interviews are only conducted at enterprises while 14% of phone interviews were conducted away from enterprises. This means that in-person interviews are disproportionately likely to miss respondents who are travelling or working short hours at their enterprises, consistent with the labour supply effects discussed in Section 3. Combining these patterns suggests that the phone interviews may have lower response

Table A14: **Response Rates by Data Collection Group**

	(1)	(2)	(3)	(4)
Interview completed	% of repeated	Any repeated	All repeated	Endline
Monthly in-person	0.573 (0.021)	0.849 (0.021)	0.315 (0.027)	0.664 (0.027)
Weekly in-person	0.515 (0.020)	0.876 (0.019)	0.060 (0.014)	0.726 (0.026)
Weekly phone	0.478 (0.020)	0.869 (0.020)	0.081 (0.016)	0.591 (0.029)
# enterprises	895	895	895	895
p-value for monthly in-person = weekly in-person	0.045	0.334	0.000	0.104
p-value for weekly in-person = weekly phone	0.187	0.794	0.332	0.000

This table shows coefficients from respondent-level linear regressions on data collection group indicators with heteroscedasticity-robust standard errors. The dependent variables are calculated at respondent level.

rates than in-person interviews with a more flexible (and expensive) tracking protocol than the approach we used.

We do find a substantial frequency effect on non-response. Respondents complete 6 percentage point more interviews in the monthly in-person group than the weekly in-person group (Table A14). This difference is concentrated in the first four weeks of repeated interviews (Figure 3, Panel B). But there is no frequency effect on the probability of completing at least one interview (Table A14).

Although respondents miss more weekly interviews, weekly interviews are more likely to find all respondents at least once in a given period. The fraction of respondents that are interviewed at least once in each x -week period is higher in both weekly groups than in the monthly group for all values of x (Table A15). The coverage rate for monthly interviews is mechanically lower lower for $x < 4$. But the lower coverage rate in the monthly group over longer time periods is not mechanical and is informative.

This presents a trade-off: weekly interviews deliver a higher volume of information but this information may be less representative in some weeks. If the non-response in any one period is close to random, then the greater volume of information will more than offset the lower response rate in each week. Non-response in our sample is correlated with only 2 of 34 baseline measures (Table A16). But the correlations are not consistent with simple economic models of non-response. In particular, prior experience conducting business on the phone does not predict non-response

Table A15: Coverage Rates by Data Collection Group

	(1)	(2)	(3)	(4)	(5)
	2 weeks	4 weeks	6 weeks	8 weeks	10 weeks
Monthly in-person	0.283 (0.010)	0.543 (0.021)	0.628 (0.022)	0.723 (0.022)	0.792 (0.022)
Weekly in-person	0.648 (0.023)	0.713 (0.023)	0.746 (0.022)	0.778 (0.022)	0.817 (0.020)
Weekly phone	0.604 (0.021)	0.710 (0.022)	0.765 (0.022)	0.799 (0.021)	0.830 (0.020)
# enterprises	895	895	895	895	895
p-value: monthly in-person = weekly in-person	0.000	0.000	0.000	0.080	0.397
p-value: weekly in-person = weekly phone	0.155	0.901	0.557	0.481	0.658

This table shows coefficients from respondent-level linear regressions on data collection group indicators with heteroscedasticity-robust standard errors in parentheses. The dependent variable equals the fraction of periods of x weeks in which the respondent completes at least one interview, with different x in each column. For example, the dependent variable in column 1 equals one if the enterprise completes at least one interview in each of the following periods: weeks 1-2, weeks 2-3, weeks 3-4, etc. The design of the data collection means that the dependent variable is less than one for enterprises in the monthly in-person interview group when $x \leq 3$.

in the full sample or in the phone group. The response rate does not vary with proxies for the opportunity cost of time: income, childcare responsibilities, or enterprise sector. Nor does the response rate vary with proxies for the cognitive costs of interview participation: keeping written records, registration for tax, financial literacy test results, or digit span recall test results. The relationship between non-response and education is statistically significant but non-monotonic.

There are a few differences between data collection methods in the relationship between non-response to targeted interviews and baseline characteristics (Table A17). Phone interviews capture slightly fewer older respondents while monthly interviews capture more respondents with high digit recall scores and fewer respondents who have held permanent employment. Despite these patterns within each interview round, higher-frequency interviews still capture more respondents in any given month. The marginal respondents who are captured only by higher-frequency surveys are not systematically different to the inframarginal respondents who are captured by high and low-frequency surveys. Specifically, we generate a respondent-month dataset of 2685 observations and restrict the sample to respondents who are interviewed at least once in the relevant month (511 respondent-months from the monthly group and 1300 respondent-months from the weekly groups). We then compare the baseline characteristics of respondents in this sample between weekly and monthly groups and report the results in Table A18. The two groups differ on 3 of 34 characteristics

Table A16: Response Predictors in Repeated and Endline Interviews

Characteristic	Repeated		Endline	
Respondent's age	0.003	(0.002)	-0.003	(0.002)
Respondent female	0.026	(0.028)	0.035	(0.039)
Respondent was born in Mozambique	0.056	(0.066)	0.105	(0.089)
Respondent was born in another country	-0.083	(0.088)	-0.061	(0.117)
Respondent speaks Sotho	-0.019	(0.052)	-0.136	(0.088)
Respondent speaks Tswana	0.015	(0.061)	-0.072	(0.099)
Respondent speaks Zulu	-0.003	(0.049)	-0.180**	(0.084)
Respondent speaks another language	-0.055	(0.057)	-0.053	(0.093)
Respondent has some secondary education	0.105***	(0.038)	0.080	(0.050)
Respondent has finished secondary education	0.052	(0.044)	0.034	(0.059)
Respondent has some tertiary education	0.007	(0.062)	0.098	(0.088)
Years respondent has lived in Gauteng	-0.002	(0.003)	-0.004	(0.004)
Years respondent has lived in Soweto	0.001	(0.003)	0.008**	(0.003)
% of financial literacy questions respondent correctly answers	0.021	(0.015)	-0.017	(0.021)
Respondent's digit recall test score	0.002	(0.009)	0.006	(0.013)
Respondent has ever held regular paid employment	0.007	(0.030)	0.026	(0.043)
Respondent's household size	-0.002	(0.005)	0.008	(0.007)
Respondent's household's total income	0.000	(0.000)	0.000	(0.000)
Missing value for respondent's household's total income	-0.127***	(0.041)	-0.082	(0.055)
Enterprise provides at most half of household income	-0.026	(0.026)	-0.015	(0.036)
Respondent has primary responsibility for childcare	0.007	(0.026)	0.018	(0.036)
Respondent perceives pressure within HH to share profits	-0.019	(0.027)	-0.052	(0.038)
Respondent perceives pressure outside HH to share profits	0.024	(0.027)	-0.011	(0.037)
Food sector	-0.031	(0.028)	-0.002	(0.040)
Light manufacturing sector	-0.035	(0.047)	-0.081	(0.060)
Services sector	0.005	(0.046)	-0.028	(0.063)
Agriculture/other sector	0.018	(0.056)	-0.010	(0.072)
# employees	-0.019	(0.018)	-0.065***	(0.025)
Enterprise age	-0.002	(0.002)	0.001	(0.002)
Respondent keeps written financial records	-0.004	(0.032)	-0.037	(0.044)
Enterprise is registered for payroll tax or VAT	-0.065	(0.046)	-0.025	(0.066)
Respondent plans to grow enterprise in next five years	0.003	(0.029)	0.088**	(0.039)
Respondent conducts business by phone at least weekly	-0.009	(0.025)	-0.081**	(0.035)
# clients	0.000	(0.000)	0.000	(0.000)
All coefficients are zero: χ^2 test statistic	61.051		74.609	
All coefficients are zero: p -value	0.003		0.000	
# enterprises	895		895	
Mean value of outcome	52.2		66.0	

Notes: This table shows response rates for repeated and endline interviews. The first column shows marginal effects from a fractional logit regression of the respondent-level response rate for repeated interviews (Papke and Wooldridge, 1996). The third column show marginal effects from a logit regression of an endline response indicator. All marginal effects for continuous variables are evaluated at their sample means. The second and fourth columns show heteroskedasticity-robust standard errors in parentheses calculated using the delta method. Omitted categories are South Africa for country of birth, English for home language, incomplete primary for education and trade/retail for enterprise type. ***, **, and * denote significance at the 1, 5, and 10% levels.

and these differences are not jointly statistically significant.

Non-response in our weekly panel is comparable to other high-frequency surveys with representative samples: [Croke, Dabalén, Demombynes, Giugale, and Hoogeveen \(2014\)](#) complete 55% of weekly mobile phone interviews with a random sample of Dar es Salaam households and [Gallup \(2012\)](#) complete approximately 42% of scheduled phone/text interviews with nationally representative households in Honduras and Peru. However, non-response in our weekly panel is lower than in surveys of samples selected (directly or indirectly) for their willingness to participate in interviews. For example, [Heath, Mansuri, Sharma, Rijkers, and Seitz \(2017\)](#) obtain weekly response rates above 77% in a high-frequency panel, but work with respondents who have already completed up to eight annual waves of a household panel survey in Ghana. [Arthi, Beegle, de Weerd, and Palacios-Lopez \(2018\)](#) conduct a high-frequency panel of rural households but replace the 8% of households that miss any interview in the first five weeks of the panel. [Beaman, Magruder, and Robinson \(2014\)](#) begin with a baseline of 1195 microenterprises but, after a screening exercise that partly reflects ease of contacting respondents, continue a high-frequency panel with only 508 of these enterprises. If we had dropped respondents who missed the first one or first two scheduled interviews from our panel, our response rate would have risen by respectively 14 or 29 percentage points.

We complete endline interviews with 66% of the 895 respondents (Table [A14](#)). The rise in response rate from repeated to endline interviews may occur because respondents were more willing to complete one final interview or because the interviews were spread over several weeks, allowing enumerators more flexibility. The response rate is lowest for the weekly phone group, 59%, partly because these respondents' physical locations had not been tracked during the repeated interview phase and enumerators struggled to find them for in-person endline interviews. This points to a potential cost of switching survey medium during a panel. Endline non-response is correlated with 5 of 34 baseline measures but these correlations are again not consistent with simple economic models of non-response (Table [A16](#)).

Table A17: Response Predictors in Repeated Interviews by Data Collection Group

Characteristic	Monthly in-person	Weekly in-person	Weekly phone	Test statistic for coefficient equality across groups
Respondent's age	0.006* (0.003)	0.008*** (0.003)	-0.003 (0.003)	10.341 [0.016]
Respondent female	0.038 (0.048)	0.011 (0.048)	0.024 (0.049)	0.903 [0.825]
Respondent was born in Mozambique	0.131 (0.111)	-0.136 (0.136)	0.202* (0.103)	6.228 [0.101]
Respondent was born in another country	0.046 (0.141)	-0.231 (0.186)	-0.049 (0.150)	1.758 [0.624]
Respondent speaks Sotho	0.038 (0.097)	-0.091 (0.088)	0.030 (0.101)	1.303 [0.729]
Respondent speaks Tswana	0.006 (0.114)	-0.063 (0.101)	0.130 (0.121)	1.531 [0.675]
Respondent speaks Zulu	0.033 (0.094)	-0.075 (0.080)	0.062 (0.098)	1.391 [0.708]
Respondent speaks another language	-0.008 (0.106)	-0.319*** (0.102)	0.082 (0.108)	10.281 [0.016]
Respondent has some secondary education	0.234*** (0.065)	0.128* (0.067)	0.041 (0.067)	17.052 [0.001]
Respondent has finished secondary education	0.077 (0.080)	0.133* (0.075)	-0.019 (0.078)	4.133 [0.247]
Respondent has some tertiary education	-0.015 (0.110)	0.109 (0.113)	-0.001 (0.111)	0.947 [0.814]
Years respondent has lived in Gauteng	0.007 (0.007)	-0.009 (0.006)	0.001 (0.005)	3.477 [0.324]
Years respondent has lived in Soweto	-0.007 (0.007)	0.006 (0.005)	-0.002 (0.005)	2.640 [0.450]
% of financial literacy questions respondent correctly answers	-0.001 (0.024)	0.003 (0.029)	0.056* (0.029)	3.719 [0.293]
Respondent's digit recall test score	0.036** (0.016)	-0.028* (0.016)	-0.007 (0.015)	8.456 [0.037]
Respondent has ever held regular paid employment	-0.141** (0.059)	0.051 (0.055)	0.055 (0.048)	7.994 [0.046]
Respondent's household size	0.003 (0.008)	-0.002 (0.008)	-0.006 (0.009)	0.772 [0.856]
Respondent's household's total income	-0.000 (0.000)	-0.000 (0.000)	-0.000 (0.000)	2.105 [0.551]
Missing value for respondent's household's total income	-0.104 (0.078)	-0.202*** (0.068)	-0.076 (0.077)	11.676 [0.009]
Enterprise provides at most half of household income	-0.057 (0.047)	0.050 (0.043)	-0.083* (0.046)	5.995 [0.112]
Respondent has primary responsibility for childcare	0.008 (0.045)	0.035 (0.045)	-0.031 (0.045)	1.121 [0.772]
Respondent perceives pressure within HH to share profits	0.004 (0.050)	-0.026 (0.050)	0.029 (0.045)	0.672 [0.880]
Respondent perceives pressure outside HH to share profits	0.086* (0.049)	0.002 (0.047)	-0.066 (0.044)	5.413 [0.144]
Food sector	-0.033 (0.049)	-0.006 (0.046)	-0.024 (0.049)	0.714 [0.870]
Light manufacturing sector	-0.075 (0.080)	-0.147* (0.083)	0.071 (0.079)	4.782 [0.188]
Services sector	-0.007 (0.083)	-0.105 (0.081)	0.157** (0.080)	5.554 [0.135]
Agriculture/other sector	-0.009 (0.092)	0.042 (0.103)	-0.045 (0.091)	0.422 [0.936]
# employees	-0.048 (0.032)	-0.019 (0.033)	0.031 (0.031)	3.590 [0.309]
Enterprise age	-0.007** (0.003)	0.001 (0.004)	-0.000 (0.003)	6.367 [0.095]
Respondent keeps written financial records	0.060 (0.061)	-0.012 (0.064)	-0.036 (0.053)	1.463 [0.691]
Enterprise is registered for payroll tax or VAT	-0.180** (0.080)	0.063 (0.102)	-0.072 (0.076)	6.338 [0.096]
Respondent plans to grow enterprise in next five years	0.045 (0.052)	-0.023 (0.053)	-0.008 (0.051)	0.975 [0.807]
Respondent conducts business by phone at least weekly	0.026 (0.048)	-0.072 (0.044)	-0.018 (0.043)	3.143 [0.370]
# clients	-0.000 (0.000)	-0.000 (0.000)	0.000 (0.000)	1.095 [0.778]
Joint equality across rows				102.821 [0.004]
Joint zero down column	87.104 [0.000]	56.348 [0.009]	36.311 [0.361]	
Sample size	298	299	298	

Notes: This table shows response rates for repeated interviews separately for each data collection group. The outcome is the fraction of scheduled interviewed successfully completed. All coefficients are marginal effects from fractional logit regressions evaluated at the sample mean of the relevant variable. Heteroscedasticity-robust standard errors calculated using the Delta method are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels. The final column shows χ^2 test statistics and p -values in brackets for the null hypothesis that the marginal effects are equal across data collection groups. The final row shows χ^2 test statistics and p -values in brackets for the null hypothesis that the group-specific marginal effects are jointly equal to zero. Omitted categories are Black African for race, South Africa for country of birth, English for home language, incomplete primary for education and trade/retail for business type.

Table A18: Difference between Baseline Characteristics of Monthly and Weekly Respondents Interviewed at Least Monthly

Characteristic	Coefficient
Respondent's age	0.001 (0.003)
Respondent female	-0.007 (0.041)
Respondent was born in Mozambique	-0.046 (0.097)
Respondent was born in another country	0.128 (0.143)
Respondent speaks Sotho	-0.000 (0.072)
Respondent speaks Tswana	-0.056 (0.085)
Respondent speaks Zulu	0.019 (0.068)
Respondent speaks another language	0.014 (0.082)
Respondent has some secondary education	0.111** (0.049)
Respondent has finished secondary education	0.087 (0.058)
Respondent has some tertiary education	-0.011 (0.080)
Years respondent has lived in Gauteng	-0.004 (0.004)
Years respondent has lived in Soweto	0.004 (0.003)
% of financial literacy questions respondent correctly answers	-0.015 (0.022)
Respondent's digit recall test score	0.023* (0.012)
Respondent has ever held regular paid employment	-0.001 (0.043)
Respondent's household size	0.004 (0.007)
Respondent's household's total income	-0.000 (0.000)
Missing value for respondent's household's total income	-0.068 (0.057)
Enterprise provides at most half of household income	0.050 (0.037)
Respondent has primary responsibility for childcare	-0.058 (0.038)
Respondent perceives pressure within HH to share profits	-0.064 (0.039)
Respondent perceives pressure outside HH to share profits	0.096** (0.037)
Food sector	-0.006 (0.040)
Light manufacturing sector	-0.024 (0.063)
Services sector	0.031 (0.069)
Agriculture/other sector	0.014 (0.076)
# employees	-0.001 (0.025)
Enterprise age	0.000 (0.003)
Respondent keeps written financial records	0.024 (0.046)
Enterprise is registered for payroll tax or VAT	-0.023 (0.066)
Respondent plans to grow enterprise in next five years	-0.018 (0.040)
Respondent conducts business by phone at least weekly	0.002 (0.036)
# clients	-0.000 (0.000)
Joint significance test	36.807 [0.345]
# respondent-months	1811
# respondents	774

Notes: This table shows the difference in observed baseline characteristics between respondents in the monthly and weekly groups, conditional on being interviewed at least once in a month. The outcome is an indicator equal to one for respondents in the monthly group, so a positive coefficient on a characteristic indicates that characteristic is more common amongst interviewed respondents in the monthly group. The unit of observation is the respondent \times month. All coefficients are marginal effects from a logit regression evaluated at the sample mean of the relevant variable. Heteroscedasticity-robust standard errors calculated using the Delta method are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels. The final row shows χ^2 test statistics and p -values in brackets for the null hypothesis that the coefficients are jointly equal to zero. Omitted categories are Black African for race, South Africa for country of birth, English for home language, incomplete primary for education and trade/retail for business type.

Table A19: Response Predictors in Endline Interviews by Data Collection Group

Characteristic	Monthly in-person	Weekly in-person	Weekly phone	Test statistic for coefficient equality across groups
Respondent's age	0.004 (0.005)	0.004 (0.004)	-0.018*** (0.005)	13.395 [0.004]
Respondent female	-0.004 (0.069)	-0.019 (0.066)	0.169** (0.078)	4.857 [0.183]
Respondent was born in Mozambique	0.276* (0.145)	0.002 (0.143)	0.171 (0.165)	4.672 [0.197]
Respondent was born in another country	0.060 (0.212)	-0.147 (0.191)	-0.169 (0.237)	1.179 [0.758]
Respondent speaks Sotho	-0.243 (0.177)	-0.039 (0.125)	-0.055 (0.163)	2.107 [0.551]
Respondent speaks Tswana	-0.177 (0.182)	0.009 (0.144)	0.162 (0.189)	1.672 [0.643]
Respondent speaks Zulu	-0.311* (0.163)	-0.056 (0.118)	-0.083 (0.157)	4.157 [0.245]
Respondent speaks another language	0.115 (0.189)	-0.132 (0.136)	0.039 (0.171)	1.377 [0.711]
Respondent has some secondary education	0.185** (0.092)	-0.014 (0.077)	0.120 (0.099)	5.533 [0.137]
Respondent has finished secondary education	0.025 (0.111)	0.002 (0.094)	0.120 (0.121)	1.034 [0.793]
Respondent has some tertiary education	-0.074 (0.150)	0.092 (0.141)	0.295 (0.192)	3.021 [0.388]
Years respondent has lived in Gauteng	-0.007 (0.010)	-0.005 (0.005)	0.004 (0.006)	1.940 [0.585]
Years respondent has lived in Soweto	0.007 (0.010)	0.005 (0.005)	0.006 (0.005)	2.970 [0.396]
% of financial literacy questions respondent correctly answers	-0.031 (0.037)	-0.010 (0.036)	-0.070 (0.044)	3.298 [0.348]
Respondent's digit recall test score	0.032 (0.023)	-0.010 (0.023)	-0.003 (0.023)	2.150 [0.542]
Respondent has ever held regular paid employment	-0.043 (0.091)	0.097 (0.071)	-0.010 (0.075)	2.104 [0.551]
Respondent's household size	0.018 (0.012)	0.018 (0.012)	-0.008 (0.014)	4.616 [0.202]
Respondent's household's total income	-0.000 (0.000)	-0.000 (0.000)	-0.000 (0.000)	4.275 [0.233]
Missing value for respondent's household's total income	-0.067 (0.101)	-0.144* (0.085)	-0.038 (0.114)	3.404 [0.333]
Enterprise provides at most half of household income	0.035 (0.064)	-0.000 (0.059)	-0.012 (0.075)	0.314 [0.957]
Respondent has primary responsibility for childcare	0.073 (0.063)	-0.001 (0.057)	0.023 (0.074)	1.427 [0.699]
Respondent perceives pressure within HH to share profits	-0.089 (0.070)	0.068 (0.059)	-0.118 (0.072)	5.601 [0.133]
Respondent perceives pressure outside HH to share profits	-0.004 (0.069)	-0.064 (0.061)	-0.002 (0.071)	1.101 [0.777]
Food sector	-0.017 (0.071)	0.049 (0.067)	-0.015 (0.077)	0.623 [0.891]
Light manufacturing sector	-0.024 (0.113)	-0.166* (0.095)	-0.059 (0.125)	3.337 [0.343]
Services sector	-0.016 (0.106)	-0.118 (0.103)	0.128 (0.131)	2.307 [0.511]
Agriculture/other sector	-0.116 (0.129)	0.015 (0.116)	0.151 (0.151)	1.823 [0.610]
# employees	-0.073 (0.044)	-0.063 (0.041)	-0.071 (0.051)	6.986 [0.072]
Enterprise age	-0.002 (0.004)	0.003 (0.004)	0.002 (0.005)	0.925 [0.819]
Respondent keeps written financial records	0.029 (0.086)	-0.030 (0.072)	-0.038 (0.085)	0.485 [0.922]
Enterprise is registered for payroll tax or VAT	0.049 (0.115)	0.133 (0.113)	-0.258** (0.129)	5.590 [0.133]
Respondent plans to grow enterprise in next five years	0.174** (0.072)	0.121* (0.063)	-0.026 (0.079)	9.695 [0.021]
Respondent conducts business by phone at least weekly	-0.035 (0.068)	-0.116** (0.057)	-0.104 (0.071)	6.547 [0.088]
# clients	0.003** (0.001)	-0.000 (0.000)	0.001 (0.001)	7.549 [0.056]
Joint equality across rows				80.433 [0.144]
Joint zero down column	55.422 [0.012]	46.292 [0.078]	50.000 [0.038]	
Sample size	298	299	298	

Notes: This table shows response rates for endline interviews separately for each data collection group. The outcome is an indicator equal to one if the respondent successfully completed the endline interview. All coefficients are marginal effects from logit regressions evaluated at the sample mean of the relevant variable. Heteroscedasticity-robust standard errors calculated using the Delta method are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels. The final column shows χ^2 test statistics and p -values in brackets for the null hypothesis that the marginal effects are equal across data collection groups. The final row shows χ^2 test statistics and p -values in brackets for the null hypothesis that the group-specific marginal effects are jointly equal to zero. Omitted categories are Black African for race, South Africa for country of birth, English for home language, incomplete primary for education and trade/retail for business type.

I Interview Frequency and Medium Effects in Endline Interviews

This section reports mean frequency and medium effects on outcomes in endline interviews, discussed in Section 7.4.

Table A20 shows mean effects of each interview method discussed in Section 7.4. Table A21 shows the same results from regressions weighted by the inverse probability of non-response. The weights do not shift any coefficient estimate by more than 0.02 standard deviations or 1 percentage point. Table A22 shows the same estimates with sharpened q -values that control the false discovery rate across all seventeen outcomes. None of the few statistically significant frequency and medium effects remain statistically significant at conventional levels after adjusting for multiple testing.

Figures A4 and A5 show the empirical distribution of continuous outcomes by frequency and medium. These show that the differences in employment are explained entirely by the difference between zero and one workers, in addition to the owner. The difference in household takings is driven top quartile of the distribution, particularly for some very large values in the monthly in-person group.

Table A20: Frequency and Medium Effects in Endline Surveys

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Operating	Stock & inventory	Fixed assets	Profit	Sales last week	Sales last 4 weeks	Total costs	Profit check
Monthly in-person	0.004 (0.026)	0.010 (0.031)	0.059 (0.039)	0.016 (0.059)	-0.002 (0.052)	0.055 (0.060)	-0.058 (0.067)	-0.029 (0.051)
Weekly by phone	-0.030 (0.028)	-0.014 (0.030)	0.008 (0.034)	-0.078 (0.058)	-0.004 (0.056)	0.025 (0.060)	0.042 (0.078)	-0.026 (0.054)
Observations	594	546	546	546	546	546	546	546
All treatments equal (<i>p</i>)	0.437	0.747	0.270	0.244	0.998	0.659	0.424	0.817
MDE: Monthly in-person	0.062	0.067	0.068	0.120	0.106	0.125	0.130	0.105
MDE: Weekly by phone	0.064	0.070	0.071	0.125	0.110	0.131	0.136	0.109
Lee bound: Monthly in-person (lower)	0.001	0.010	0.021	0.004	-0.006	0.055	-0.106	-0.032
Lee bound: Monthly in-person (upper)	0.014	0.023	0.044	0.050	0.048	0.068	0.020	0.019
Lee bound: Weekly by phone (lower)	-0.097	-0.054	-0.053	-0.153	-0.077	-0.053	-0.052	-0.088
Lee bound: Weekly by phone (upper)	-0.014	0.121	0.140	0.199	0.249	0.286	0.364	0.210

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Employees	Full-time	Paid	Money kept	Household takings	Hours yesterday	Honest	Careful	Written records
Monthly in-person	-0.143 (0.059)**	-0.091 (0.072)	-0.130 (0.059)**	-0.011 (0.034)	-0.004 (0.008)	-0.030 (0.047)	-0.053 (0.045)	-0.058 (0.048)	-0.016 (0.030)
Weekly by phone	-0.127 (0.064)**	-0.139 (0.074)*	-0.070 (0.068)	-0.014 (0.035)	-0.018 (0.007)**	-0.027 (0.049)	-0.082 (0.046)*	-0.053 (0.049)	-0.068 (0.029)**
Observations	546	546	546	546	546	546	593	593	546
All treatments equal (<i>p</i>)	0.036**	0.157	0.091*	0.914	0.031**	0.785	0.191	0.397	0.049**
MDE: Monthly in-person	0.140	0.152	0.114	0.067	0.017	0.113	0.104	0.114	0.065
MDE: Weekly by phone	0.146	0.158	0.119	0.070	0.018	0.118	0.108	0.118	0.068
Lee bound: Monthly in-person (lower)	-0.195	-0.129	-0.173	-0.024	-0.008	-0.093	-0.108	-0.106	-0.019
Lee bound: Monthly in-person (upper)	-0.065	0.012	-0.068	-0.007	-0.004	-0.010	-0.041	-0.025	-0.001
Lee bound: Weekly by phone (lower)	-0.226	-0.223	-0.132	-0.058	-0.027	-0.199	-0.213	-0.175	-0.098
Lee bound: Weekly by phone (upper)	0.167	0.246	0.227	0.145	0.017	0.106	-0.016	0.036	0.063

Coefficients are from regressions of each outcome, winsorised at the 95th percentile, on a vector of data collection group indicators and randomisation stratum fixed effects. Continuous outcomes are standardised to have mean zero and standard deviation one. Owners who close their enterprises are included in regressions only for panel A column 1 and panel B columns 7 and 8. Heteroskedasticity-robust standard errors are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels.

Table A21: Endline Interviews: Mean Differences, with Inverse Probability Weights

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Operating	Stock & inventory	Fixed assets	Profit	Sales last week	Sales last 4 weeks	Total costs	Profit check
Monthly in-person	0.000 (0.027)	0.015 (0.031)	0.077 (0.040)*	0.035 (0.061)	0.016 (0.054)	0.071 (0.062)	-0.049 (0.067)	-0.021 (0.053)
Weekly by phone	-0.032 (0.028)	-0.011 (0.030)	0.013 (0.035)	-0.084 (0.059)	-0.012 (0.057)	0.020 (0.062)	0.032 (0.078)	-0.038 (0.054)
Observations	594	546	546	546	546	546	546	546
All treatments equal (<i>p</i>)	0.435	0.745	0.144	0.165	0.900	0.519	0.572	0.775

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Employees	Full-time	Paid	Money kept	Household takings	Hours yesterday	Honest	Careful	Written records
Monthly in-person	-0.138 (0.060)**	-0.081 (0.071)	-0.126 (0.060)**	-0.001 (0.035)	-0.003 (0.008)	-0.028 (0.047)	-0.054 (0.045)	-0.060 (0.048)	-0.009 (0.030)
Weekly by phone	-0.133 (0.065)**	-0.139 (0.074)*	-0.080 (0.069)	-0.018 (0.035)	-0.018 (0.007)***	-0.028 (0.050)	-0.076 (0.047)	-0.052 (0.049)	-0.062 (0.029)**
Observations	546	546	546	546	546	546	593	593	546
All treatments equal (<i>p</i>)	0.040**	0.162	0.108	0.870	0.017**	0.794	0.238	0.403	0.076*

Coefficients are from regressions of each outcome on a vector of data collection group indicators and randomisation stratum fixed effects. Regressions are weighted by the inverse probability of non-response. Continuous outcomes are standardised to have mean zero and standard deviation one in the monthly in-person group and winsorised at the 95th percentile. Owners who close their enterprises are included in regressions only for panel A column 1 and panel B columns 7 and 8. Heteroskedasticity-robust standard errors are shown in parentheses. ***, **, and * denote significance at the 1, 5, and 10% levels.

Table A22: Frequency and Medium Effects on Mean Outcomes in Endline Interviews with Multiple Testing Adjustments

Outcome	Monthly in-person			Weekly phone			All group equal	
	β	p	q	β	p	q	p	q
Operating	0.004	0.886	1.000	-0.030	0.281	0.682	0.437	1.000
Stock & inventory	0.010	0.755	1.000	-0.014	0.633	1.000	0.747	1.000
Fixed assets	0.059	0.123	1.000	0.008	0.820	1.000	0.270	0.781
Profit	0.016	0.782	1.000	-0.078	0.175	0.540	0.244	0.781
Sales last week	-0.002	0.975	1.000	-0.004	0.945	1.000	0.998	1.000
Sales last 4 weeks	0.055	0.362	1.000	0.025	0.678	1.000	0.659	1.000
Total costs	-0.058	0.384	1.000	0.042	0.588	1.000	0.424	1.000
Profit check	-0.029	0.567	1.000	-0.026	0.622	1.000	0.817	1.000
Employees	-0.143	0.016	0.327	-0.127	0.047	0.291	0.036	0.385
Full-time	-0.091	0.201	1.000	-0.139	0.060	0.291	0.157	0.687
Paid	-0.130	0.029	0.327	-0.070	0.304	0.682	0.091	0.468
Hours yesterday	-0.030	0.522	1.000	-0.027	0.591	1.000	0.785	1.000
Money kept	-0.011	0.757	1.000	-0.014	0.699	1.000	0.914	1.000
Household takings	-0.004	0.639	1.000	-0.018	0.011	0.185	0.031	0.385
Honest	-0.053	0.240	1.000	-0.082	0.077	0.300	0.191	0.706
Careful	-0.058	0.222	1.000	-0.053	0.273	0.682	0.397	1.000
Written records	-0.016	0.594	1.000	-0.068	0.018	0.185	0.049	0.385

This table shows, for each outcome, means in the endline interviews in the listed groups relative to the weekly in-person group (β 's). The table also shows a p -value for the null hypothesis that each difference equals zero and a sharpened q -value that controls the probability of incorrectly rejecting all null hypotheses (Benjamini, Krieger, and Yekutieli, 2006). The q -values adjust across outcomes for a specific null hypothesis – zero frequency effect, zero medium effect, zero frequency and medium effects – but do not adjust across null hypotheses.

Figure A4: Frequency and Medium Effects on Outcome Distributions in Endline Interviews

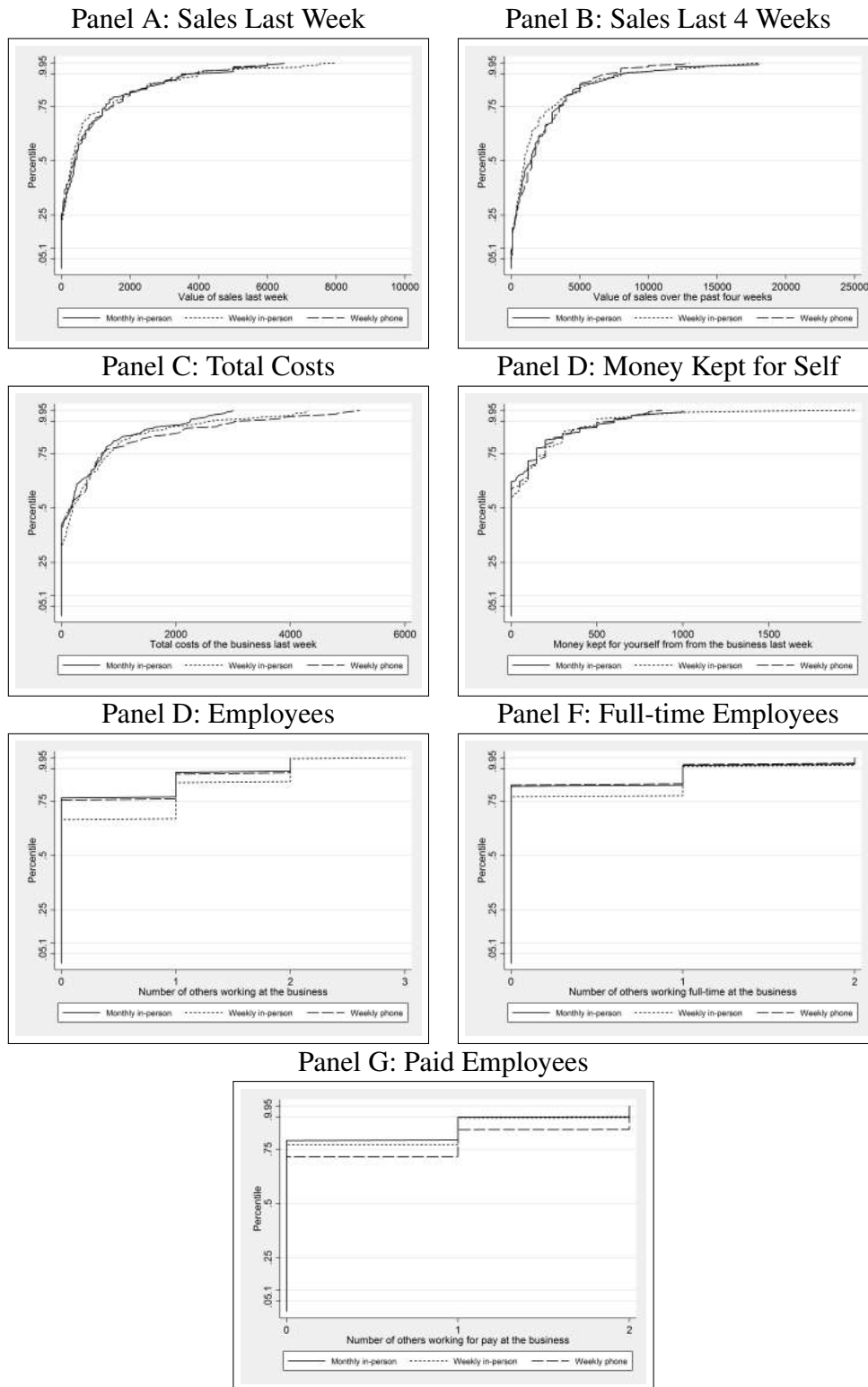


Figure shows empirical CDFs of outcomes in endline interviews. We use quantile regression to test for differences at each of the quantile shown on the y -axis. We cluster by enterprise (Parente and Silva, 2016) and use the false discovery rate (Benjamini, Krieger, and Yekutieli, 2006) to control for multiple testing across quantiles. $+$ indicates a medium effect: rejection of the null hypothesis that the coefficients for weekly in-person and phone interviews are equal. $*$ indicates a frequency effect: rejection of the null hypothesis that the coefficients for weekly and monthly in-person interviews are equal. $+++/**$, $++/*$, and $+/*$ denote significance at the 1, 5, and 10% levels.

Figure A5: Frequency and Medium Effects on Outcome Distributions in Endline Interviews

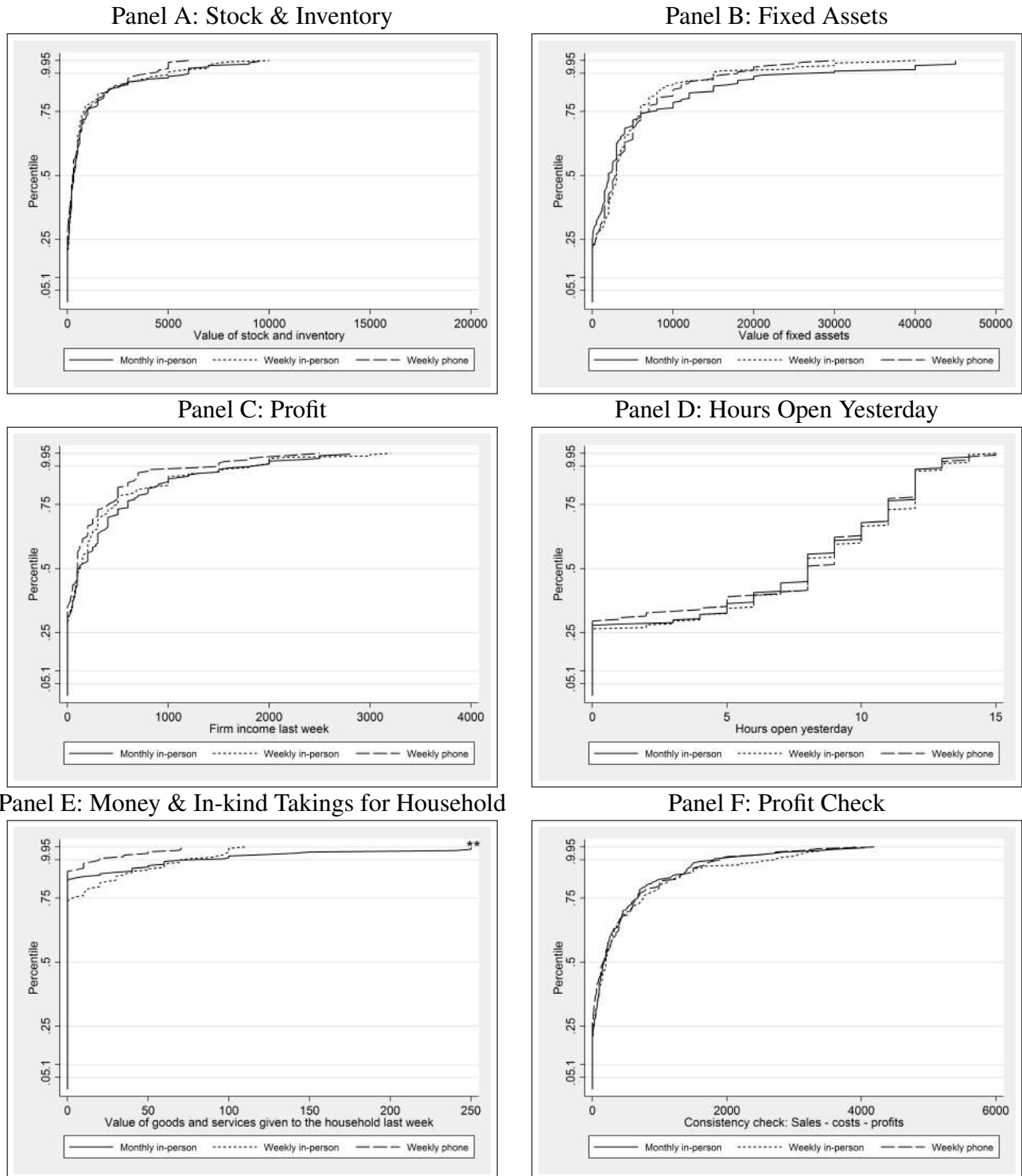


Figure shows empirical CDFs of outcomes in endline interviews. We use quantile regression to test for differences at each of the quantile shown on the y-axis. We cluster by enterprise (Parente and Silva, 2016) and use the false discovery rate (Benjamini, Krieger, and Yekutieli, 2006) to control for multiple testing across quantiles. + indicates a medium effect: rejection of the null hypothesis that the coefficients for weekly in-person and phone interviews are equal. * indicates a frequency effect: rejection of the null hypothesis that the coefficients for weekly and monthly in-person interviews are equal. +++/***, ++/**, and +/* denote significance at the 1, 5, and 10% levels.

References

- ARTHI, V., K. BEEGLE, J. DE WEERDT, AND A. PALACIOS-LOPEZ (2018): “Not Your Average Job: Measuring Farm Labor in Tanzania,” *Journal of Development Economics*, 130, 160–172.
- BEAMAN, L., J. MAGRUDER, AND J. ROBINSON (2014): “Minding Small Change: Limited Attention among Small Firms in Kenya,” *Journal of Development Economics*, 108, 69–86.
- BENJAMINI, Y., A. M. KRIEGER, AND D. YEKUTIELI (2006): “Adaptive Linear Step-Up Procedures that Control the False Discovery Rate,” *Biometrika*, 93(3), 491–507.
- BLUNDELL, R., AND S. BOND (1998): “Initial Conditions and Moment Restrictions in Dynamic Panel Data Models,” *Journal of Econometrics*, 87(1), 115–143.
- BRUHN, M., AND D. MCKENZIE (2009): “In Pursuit of Balance: Randomization in Practice in Development Field Experiments,” *American Economic Journal: Applied Economics*, pp. 200–232.
- CHO, W., AND B. GAINES (2007): “Breaking the (Benford) Law: Statistical Fraud Detection in Campaign Finance,” *The American Statistician*, 61(3), 1–6.
- CROKE, K., A. DABALEN, G. DEMOMBYNES, M. GIUGALE, AND J. HOOGEVEEN (2014): “Collecting High Frequency Panel Data in Africa using Mobile Phone Interviews,” *Canadian Journal of Development Studies*, 35(1), 186–207.
- DILLON, B. (2012): “Using Mobile Phones to Collect Panel Data in Developing Countries,” *Journal of International Development*, 24, 518–27.
- GALLUP (2012): “The World Bank Listening to LAC (L2L) Pilot: Final Report,” *Gallup Report*.
- HEATH, R., G. MANSURI, D. SHARMA, B. RIJKERS, AND W. SEITZ (2017): “Measuring Employment: Experimental Evidence from Ghana,” Working paper.
- IMBENS, G. (2015): “Matching Methods in Practice: Three Examples,” *Journal of Human Resources*, 50(2), 373–419.
- PAPKE, L., AND J. WOOLDRIDGE (1996): “Econometric Methods for Fractional Response Variables with an Application to 401(k) Plan Participation Rates,” *Journal of Applied Econometrics*, 11, 619–632.
- PARENTE, P. M., AND J. M. S. SILVA (2016): “Quantile Regression with Clustered Data,” *Journal of Econometric Methods*, 5(1), 1–15.
- SCHEINER, S., AND J. GUREVICH (2001): *Design and Analysis of Ecological Experiments*. Oxford University Press, 3rd edn.