

The 2015 Tanzania Enterprise Skills Survey (TESS)

I. Introduction

1. This document provides additional information on the data collected for the 2015 Tanzania Firm-Level Skills Survey between April 2015 and August 2015, conducted by the Enterprise Analysis (DECEA) and the Education Global Practice (GEDDR) of the World Bank Group.

The objective of the survey is to develop and test a methodological approach for a diagnostic of the composition and demand for skills and the relationship between skills (and/or skills constraints) and firm performance of selected economic sectors in Tanzania. A detailed skills module was developed as part of a larger firm-level survey collecting information, among others, on the characteristics of firms and their owners, innovation and export activities, and firm performance.

The report outlines and describes the sampling design of the data, the structure of dataset as well as additional information that may be useful when using the data, such as information on non-response cases and use of sampling weights.

II. Sampling Structure

2. The sample for the survey was selected using stratified random sampling, following a broadly similar methodology used in the World Bank's Enterprise Surveys (ES) – stratified random sampling¹. However, it is important to note that the universe of inference for the Tanzania Firm-Level Skills survey is not strictly comparable to that of the ES. For the ES, the universe of inference is private non-agricultural sectors in the country, excluding the following sectors: financial intermediation (group J²), real estate and renting activities (group K, except sub-sector 72, IT) and all public or utilities-sectors. For Tanzania Firm-Level Skills Survey, however, the universe is firms in eight selected economic activities, viz., food processing (ISIC15), textile and garments (ISIC 17 & 18), fabricated metal products (ISIC 28), furniture (ISIC 36), construction (ISIC 45), hotel and restaurant (ISIC 55), transport (ISIC 60 & 61) and Information technology (ISIC 72)).

¹ A stratified random sample is one obtained by separating the population elements into non-overlapping groups, called strata, and then selecting a simple random sample from each stratum. (Richard L. Scheaffer; Mendenhall, W.; Lyman, R., "Elementary Survey Sampling", Fifth Edition). The complete text of the World Bank Enterprise Survey's sampling methodology can be found at <http://goo.gl/EgMYXX>.

² ISIC refers to the ISIC code revision 3.1.

3. Three levels of stratification were used for this survey: industry, establishment size, and region. The original sample design with specific information of the economic activities and regions chosen is described in Appendix C.

4. The universe was stratified into *eight economic activities* (as noted above); *three size* stratification - small (5 to 19 employees), medium (20 to 99 employees), and large (more than 99 employees); and *five regions* (city and the surrounding business area): Arusha, Dar es Salaam, Mbeya, Mwanza, and Zanzibar.

III. Sampling Implementation

5. Given the stratified design, sample frames containing a complete and updated list of establishments as well as information on all stratification variables (number of employees, industry, and region) are required to draw the sample. Great efforts were made to obtain the best source for these listings. However, the quality of the sample frames was not optimal and, therefore, some adjustments were needed to correct for the presence of ineligible units. These adjustments are reflected in the weights computation (*see below*).

6. DataVision International Ltd was hired to implement the fieldwork.

7. Two sample frames were used:

- The first frame was the 2011/2012 Central Registry of Establishment (CRE) of the National Bureau of Statistics (NBS).
- The second frame was 2012 Central Registry of Establishment (CRE) of the Office of Chief Government Statistician (OCGS). The sample frame was used for the establishments in Zanzibar.

All databases contained the following information:

- a) Detailed stratification variables;
- b) Location identifiers- address, phone number, email; and
- c) Contact name(s).

8. The enumerated establishments with 5 employees or more were then used as the sample frame for the 2015 Tanzania Firm-Level Skills Survey with the aim of obtaining interviews of 390 establishments.

9. The quality of the frame was assessed at the onset of the project through visits to a random subset of firms and local contractor knowledge. The sample frame was not immune from the typical problems found in establishment surveys: positive rates of non-eligibility, repetition, non-existent units, etc.

10. Given the impact that non-eligible units included in the sample universe may have on the results, adjustments may be needed when computing the appropriate weights for individual observations.

Counts from sample frames are shown below:

Region Name	Sampling Size	Food	Textile & Garments	Fabricated Metals	Furniture	IT	Hotel & Restaurants	Construction	Transport	Total
Arusha	Small (5-19)	17	33	5	50	7	376	22	50	560
	Medium (20-99)	15	0	0	5	0	45	10	7	82
	Large (100+)	1	4	0	1	0	7	0	2	15
		33	37	5	56	7	428	32	59	657
Dar es Salaam	Small (5-19)	186	405	162	658	73	3047	367	680	5578
	Medium (20-99)	41	33	38	68	30	255	121	173	759
	Large (100+)	18	8	10	16	6	10	15	31	114
		245	446	210	742	109	3312	503	884	6451
Mbeya	Small (5-19)	43	197	12	81	9	747	47	49	1185
	Medium (20-99)	8	4	0	3	0	39	4	6	64
	Large (100+)	0	1	0	0	0	0	0	2	3
		51	202	12	84	9	786	51	57	1252
Mwanza	Small (5-19)	29	66	11	71	3	462	27	56	725
	Medium (20-99)	11	0	1	6	0	39	5	17	79
	Large (100+)	6	1	0	1	0	0	1	0	9
		46	67	12	78	3	501	33	73	813
Zanzibar	Small (5-19)	102	79	9	278	6	276	7	60	817
	Medium (20-99)	10	20	1	13	1	81	7	9	142
	Large (100+)	0	0	0	0	0	29	1	1	31
		112	99	10	291	7	386	15	70	990
Grand Total		487	851	249	1251	135	5413	634	1143	10163

Source: 2011/2012 National Bureau of Statistics (NBS) and 2012 Office of Chief Government Statistician (OCGS)

IV. Data Base Structure:

11. Data is collected using single and standardized questionnaire administered to all firms. The questionnaire has eight sections; six main sections and two sections on control information.

12. All variables are named using, first, the letter of each section and, second, the number of the variable within the section, i.e. *a1* denotes section A, question 1 (some exceptions apply). All variables are numeric with the exception of those variables with an “x” at the end of their names. The suffix “x” denotes that the variable is alpha-numeric.

13. There is a unique establishment identifiers, variable name *id*. The variables *a2* (sampling region), *a6a* (sampling establishment’s size), and *a4a* (sampling sector) contain the establishment’s classification into the strata chosen for each country using information from the sample frame. The strata were defined according to the guidelines described above.

14. All of the following variables contain information from the sampling frame. They may not coincide with the reality of individual establishments as sample frames may contain inaccurate information. The variables containing the sample frame information are included in the data set for researchers who may want to further investigate statistical features of the survey and the effect of the survey design on their results.

-*a2* is the variable describing sampling regions

-*a6a*: coded using the same standard for small, medium, and large establishments as defined above. The code -9 was used to indicate units for which size was undetermined in the sample frame.

-*a4a*: coded using ISIC codes for the chosen industries for stratification. These codes include food processing (ISIC³ 15), textile and garments (ISIC 17 & 18), fabricated metal products (ISIC 28), furniture (ISIC 36), construction (ISIC 45), hotel and restaurant (ISIC 55), transport (ISIC 60 & 61) and Information technology (ISIC 72)).

15. The surveys were implemented following a 2 stage procedure. Typically first a screener questionnaire is applied over the phone to determine eligibility and to make appointments. Then a face-to-face interview takes place with the Manager/Owner/Director of each establishment. In some cases, when the phone numbers were unavailable in the sample frame, the enumerators applied the screeners in person. The variables *a4b* and *a6b* contain the industry and size of the establishment from the screener questionnaire. Variables *a7* to *a11* contain additional information and were also collected in the screening phase.

16. Note that the fiscal years vary by firm as there is no standard for all firms in Tanzania. The start and end dates for the fiscal year for each firm can be found in the *fymonb*, *fyyearb*, *fymone* and *fyyeare* variables in the dataset

V. Universe Estimates

17. Universe estimates for the number of establishments in each cell (i.e., region-industry-size) were produced for the strict, weak and median eligibility definitions. The estimates were the multiple of the relative eligible proportions. Appendix B provides the estimates the universe based on the median eligibility assumption (*see below for median eligibility assumption*). Appendix E provides definition of eligibility codes.

18. For some establishments where contact was not successfully completed during the screening process (because the firm has moved and it is not possible to locate the new location, for example), it is not possible to directly determine eligibility. Thus, different assumptions about the eligibility of establishments result in different adjustments to the universe cells and thus different sampling weights.

19. Three sets of assumptions on establishment eligibility are used to construct sample adjustments using the status code information. Appendix...provides the definition of eligibility.

20. Strict assumption: eligible establishments are only those for which it was possible to directly determine eligibility. The resulting weights are included in the variable

³ ISIC refers to the ISIC code revision 3.1.

wstrict.

Strict eligibility = (Sum of the firms with codes 1,2,3 & 4) / Total

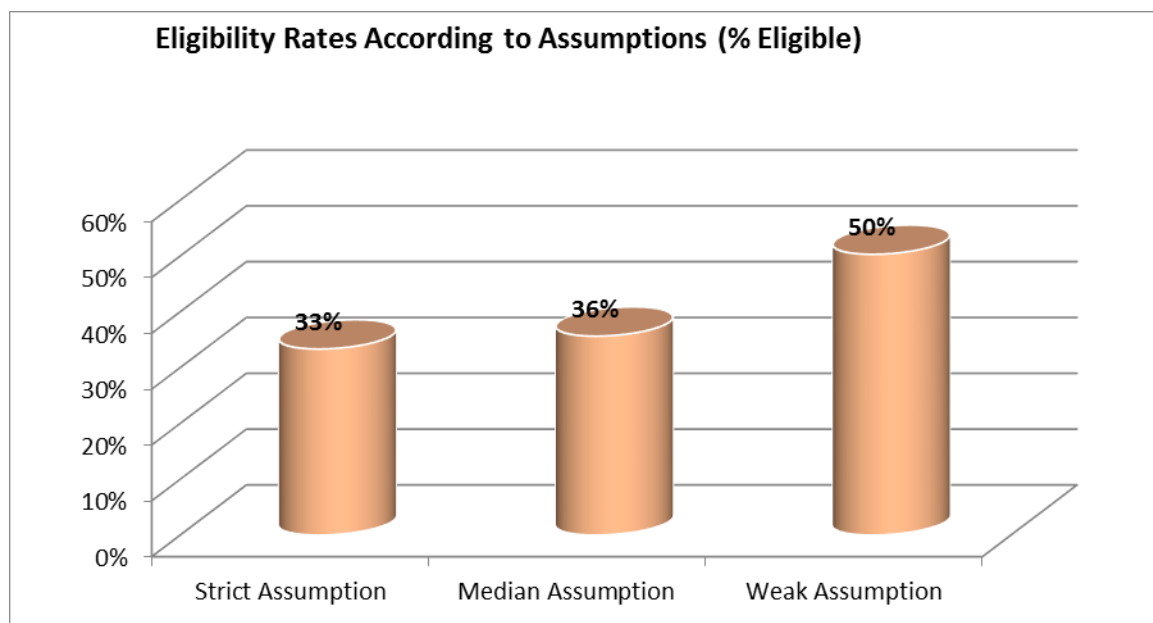
21. Median assumption: eligible establishments are those for which it was possible to directly determine eligibility and those that rejected the screener questionnaire or an answering machine or fax was the only response. The resulting weights are included in the variable *wmedian*.

Median eligibility = (Sum of the firms with codes 1,2,3,4,10,11, & 13) / Total

22. Weak assumption: in addition to the establishments included in points a and b, all establishments for which it was not possible to contact or that refused the screening questionnaire are assumed eligible. This definition includes as eligible establishments with dead or out of service phone lines, establishments that never answered the phone, and establishments with incorrect addresses for which it was impossible to find a new address. Under the weak assumption only observed non-eligible units are excluded from universe projections. The resulting weights are included in the variable *wweak*.

Weak eligibility = (Sum of the firms with codes 1,2,3,4,91,92,93,10,11,12,&13) / Total

23. The following graph shows the different eligibility rates calculated for firms in each sample frame under each set of assumptions.



24. Once an accurate estimate of the universe cell projection was made, weights for the probability of selection were computed using the number of completed interviews for each cell.

VI. Weights

25. Since the sampling design was stratified and employed differential sampling, individual observations should be properly weighted when making inferences about the population. Under stratified random sampling, unweighted estimates are biased unless sample sizes are proportional to the size of each stratum. With stratification the probability of selection of each unit is, in general, not the same. Consequently, individual observations must be weighted by the inverse of their probability of selection (probability weights or *pw* in Stata.)⁴

26. Three versions of sampling weights are provided based on the three eligibility assumptions noted above, i.e., strict, median and weak weights. Special care was given to the correct computation of the weights. It was imperative to accurately adjust the totals within each region/industry/size stratum to account for the presence of ineligible units (the firm discontinued businesses or was unattainable, education or government establishments, establishments with less than 5 employees, no reply after having called in different days of the week and in different business hours, no tone in the phone line, answering machine, fax line⁵, wrong address or moved away and could not get the new references) The information required for the adjustment was collected in the first stage of the implementation: the screening process. Using this information, each stratum cell of the universe was scaled down by the observed proportion of ineligible units within the cell. Once an accurate estimate of the universe cell (projections) was available, weights were computed using the number of completed interviews.

27. Appendix B shows the *median cell weights* for registered establishments for the 2015 Tanzania Firm-Level Skills Survey. (The three weights - strict, median and weak weights - are all in the dataset.)

VII. Appropriate use of the weights

28. Under stratified random sampling weights should be used when making inferences about the population. Any estimate or indicator that aims at describing some feature of the population should take into account that individual observations may not represent equal shares of the population. For estimations with weighting, we recommend the use of the median weights.

29. However, there is some discussion as to the use of weights in regressions (see Deaton, 1997, pp.67; Lohr, 1999, chapter 11, Cochran, 1953, pp.150). There is not a strong large sample econometric argument in favor of using weighted estimation for a common population coefficient if the underlying model varies per stratum (stratum-specific coefficient): both simple OLS and weighted OLS are inconsistent under regular conditions.

⁴ This is equivalent to the weighted average of the estimates for each stratum, with weights equal to the population shares of each stratum.

⁵ For the surveys that implemented a screener over the phone.

However, weighted OLS has the advantage of providing an estimate that is independent of the sample design.⁶

VIII. Non-response

30. Survey non-response must be differentiated from item non-response. The former refers to refusals to participate in the survey altogether whereas the latter refers to the refusals to answer some specific questions. The 2015 Tanzania Firm-Level Skills Survey suffers from both problems and different strategies were used to address these issues.

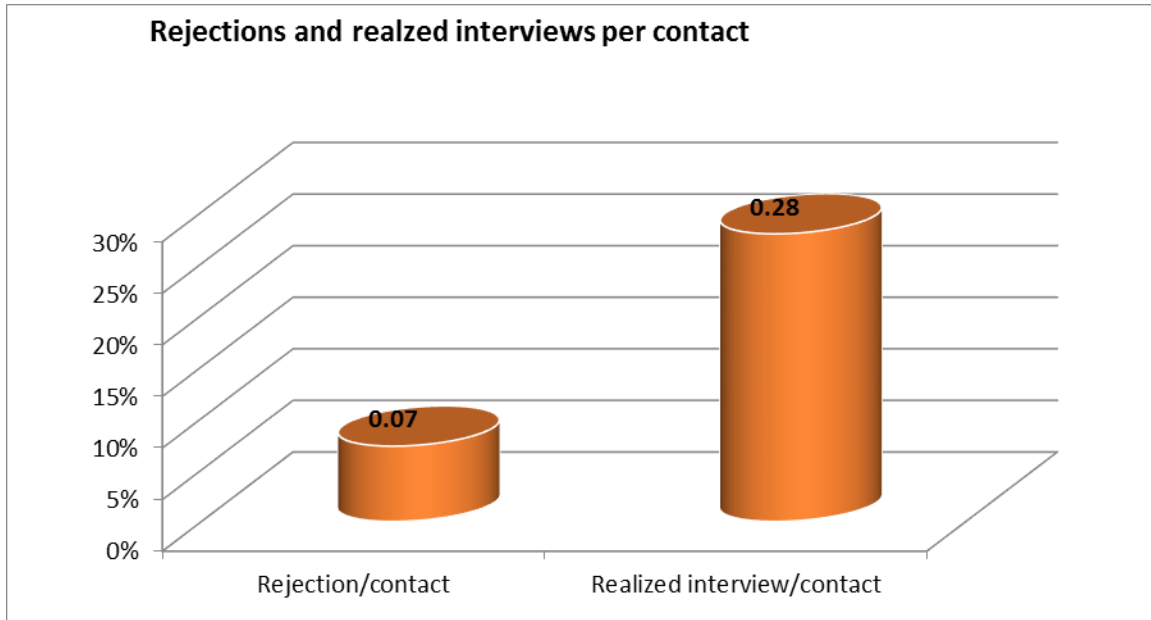
31. Item non-response was addressed by re-contacting firms. That is, establishments with incomplete information were re-contacted in order to complete this information, whenever necessary. However, there were clear cases of low response. The response rates are particularly low for questions about the names and locations of the main universities and schools attended by the establishment's recent hires (questions 113, 114, 115 and 125 in the questionnaire). Despite repeated callbacks, respondents note that they just do not know the names and locations of schools attended by their employees.

32. Survey non-response was addressed by maximizing efforts to contact establishments that were initially selected for interview. Attempts were made to contact the establishment for interview at different times/days of the week before a replacement establishment (with similar strata characteristics) was suggested for interview. Survey non-response did occur but substitutions were made in order to potentially achieve strata-specific goals.

33. As the following graphs show, the number of realized interviews per contact contacted establishments was 0.28.⁷ This number is the result of two factors: explicit refusals to participate in the survey, as reflected by the rate of rejection (which includes rejections of the screener and the main survey) and the quality of the sample frame, as represented by the presence of ineligible units. The number of rejections per contact was 0.07.

⁶ Note that weighted OLS in Stata using the command `regress` with the option of weights will estimate wrong standard errors. Using the Stata survey specific commands `svy` will provide appropriate standard errors.

⁷ The estimate is based on the total no. of firms contacted including ineligible establishments.



34. Details on the rejection rate, eligibility rate, and item non-response are available at the level of strata. This report summarizes these numbers to alert researchers of these issues when using the data and when making inferences. Item non-response, selection bias, and faulty sampling frames are not unique to the 2015 Tanzania Firm-Level Skills Survey. All firm-level surveys suffer from these shortcomings, but in very few cases they have been made explicit.

References:

Cochran, William G., Sampling Techniques, 1977.

Deaton, Angus, The Analysis of Household Surveys, 1998.

Levy, Paul S. and Stanley Lemeshow, Sampling of Populations: Methods and Applications, 1999.

Lohr, Sharon L. Sampling: Design and Techniques, 1999.

Scheaffer, Richard L.; Mendenhall, W.; Lyman, R., Elementary Survey Sampling, Fifth Edition, 1996.

Appendix A

Median Cell Weights: Tanzania Firm-Level Skills Survey

Region	Size	Food	Textile & Garments	Fabricated Metals	Furniture	IT	Hotel & Restaurants	Construction	Transport
Arusha	Small (5 to 19)	1.2	1.9	1	8.1	1	31.7	1.4	8.7
	Medium (20 to 99)	1.1			1.1		20	1.1	1.6
	Large (100+)	1	2		1		1.4		1.2
Dar es Salaam	Small (5 to 19)	15.1	9.8	2.1	9.2	1.3	33.5	12.7	9.3
	Medium (20 to 99)	5.9	2.3	1	5.5	1.7	17.9	7.7	25
	Large (100+)	2	3.7	1.8	2.8	1.3	2.8	2.1	2.2
Mbeya	Small (5 to 19)	2.9	5.5	1.1	2.6	2.8	16.6	3.2	13.4
	Medium (20 to 99)	2.9	1.2		1		3.4	1	2.2
	Large (100+)								1
Mwanza	Small (5 to 19)	1.9	4.4	1.2	3.1	1.1	24.7	1.6	18.5
	Medium (20 to 99)	1		1	1		16.5	1.6	2.5
	Large (100+)	1.7	1		1			1	
Zanzibar	Small (5 to 19)	8.2	10.4	1.4	21.1	1.7	130.9	1	29.3
	Medium (20 to 99)	2.1	1.7	1	3.9	1	12.7	1.1	5.8
	Large (100+)						8	1	1

Appendix B

Universe Estimates based on the median Eligibility Assumption

Region	Size	Food	Textile & Garments	Fabricated Metals	Furniture	IT	Hotel & Restaurants	Construction	Transport	Grand Total
Arusha	Small (5 to 19)	6	9	2	16	3	127	6	17	186
	Medium (20 to 99)	7	0	0	2	0	20	3	3	36
	Large (100+)	1	2	0	1	0	4	0	1	9
Dar Es Salaam	Small (5 to 19)	61	108	51	202	28	972	89	224	1734
	Medium (20 to 99)	18	12	16	27	15	107	39	75	309
	Large (100+)	10	4	6	9	4	6	6	18	62
Mbeya	Small (5 to 19)	12	44	3	21	3	199	9	13	304
	Medium (20 to 99)	3	1	0	1	0	14	2	2	23
	Large (100+)	0	1	0	0	0	0	0	1	2
Mwanza	Small (5 to 19)	10	18	4	22	1	148	7	19	227
	Medium (20 to 99)	5	0	1	3	0	17	2	7	35
	Large (100+)	3	1	0	1	0	0	1	0	6
Zanzibar	Small (5 to 19)	49	31	4	127	3	131	3	29	378
	Medium (20 to 99)	6	10	1	8	1	51	3	6	86
	Large (100+)	0	0	0	0	0	24	1	1	26
Total		190	241	88	439	58	1819	171	417	3422

Appendix C

Original Sample Design, Tanzania Firm-Level Skills Survey

Region Name	Sampling Size	Food	Textile & Fabricate Garments d Metals	Furniture	IT	Hotel & Restaurants	Construction	Transport	Total	
Arusha	Small (5-19)	5	5	2	2	3	4	6	2	29
	Medium (20-99)	5	0	0	2	0	2	3	2	14
	Large (100+)	1	2	0	1	0	3	0	1	8
		11	7	2	5	3	9	9	5	51
Dar es Salaam	Small (5-19)	4	11	18	16	19	20	12	18	118
	Medium (20-99)	3	3	12	2	9	3	5	3	40
	Large (100+)	5	3	3	4	2	2	5	7	31
		12	17	33	22	30	25	22	28	189
Mbeya	Small (5-19)	2	7	4	3	3	12	3	2	36
	Medium (20-99)	3	2	0	1	0	2	2	2	12
	Large (100+)	0	1	0	0	0	0	0	1	2
		5	10	4	4	3	14	5	5	50
Mwanza	Small (5-19)	5	4	4	7	1	5	3	2	31
	Medium (20-99)	4	0	1	2	0	2	2	3	14
	Large (100+)	2	1	0	1	0	0	1	0	5
		11	5	5	10	1	7	6	5	50
Zanzibar	Small (5-19)	6	3	3	5	2	3	2	2	26
	Medium (20-99)	3	6	1	2	1	2	3	2	20
	Large (100+)	0	0	0	0	0	2	1	1	4
		9	9	4	7	3	7	6	5	50
Grand Total		48	48	48	48	40	62	48	48	390

Appendix D

Completed Interviews, Tanzania Firm-Level Skills Survey

Region	Size	Food	Textile & Garments	Fabricated Metals	Furniture	IT	Hotel & Restaurants	Construction	Transport	Total
Arusha	Small (5 to 19)	5	5	2	2	3	4	4	2	27
	Medium (20 to 99)	6			2		1	3	2	14
	Large (100+)	1	1		1		3		1	7
Dar Es Salaam	Small (5 to 19)	4	11	25	22	22	29	7	24	144
	Medium (20 to 99)	3	5	16	5	9	6	5	3	52
	Large (100+)	5	1	3	3	3	2	3	8	28
Mbeya	Small (5 to 19)	4	8	3	8	1	12	3	1	40
	Medium (20 to 99)	1	1		1		4	2	1	10
	Large (100+)								1	1
Mwanza	Small (5 to 19)	5	4	3	7	1	6	4	1	31
	Medium (20 to 99)	5		1	3		1	1	3	14
	Large (100+)	2	1		1			1		5
Zanzibar	Small (5 to 19)	6	3	3	6	2	1	3	1	25
	Medium (20 to 99)	3	6	1	2	1	4	3	1	21
	Large (100+)						3	1	1	5
Total		50	46	57	63	42	76	40	50	424

Appendix E: Eligibility Code, Tanzania Enterprise Skills Survey

0	Screening in process	I4. In process (the establishment is being called/ is being contacted - previous to ask the screener)	0
507	Eligible	1. Eligible establishment (Correct name and address) 2. Eligible establishment (Different name but same address - the new firm/establishment bought the original firm/establishment) 3. Eligible establishment (Different name but same address - the firm/establishment changed its name) 4. Eligible establishment (Moved and traced) 16. Eligible establishment (Panel Firm - now less than five employees; this code applies only to panel firms.)	430 23 53 1 0
29	Screener refusal	I3. Refuses to answer the screener	29
409	Ineligible	5. The establishment has less than 5 permanent full time employees 616. The firm discontinued businesses - (Establishment went bankrupt) 617. 618. The firm discontinued businesses - (Original establishment disappeared and is now a different firm) 619. The firm discontinued businesses - (Establishment was bought out by another firm) 620. The firm discontinued businesses - (It was impossible to determine for what reason) 621. The firm discontinued businesses - (Other) 7. Not a business: Private household 8. Ineligible activity: Education, Agriculture, Finances, Government, etc.	225 26 0 15 4 72 3 20 44
27	Out of target	I51. Out of target - outside the covered regions I52. Out of target - moved abroad I53. Out of target - Not registered with Statistical Authority I54. Out of target - establishment is HQ without production or sales of goods or services I55. Out of target - establishment was not in operation for the entirety of last fiscal year I56. Duplicated firm within the sample	7 0 1 4 2 13
549	Unobtainable	91. No reply after having called in different days of the week and in different business hours 92. Line out of order 93. No tone 94. Phone number does not exist I0. Answering machine I1. Fax line- data line I2. Wrong address/ moved away and could not get the new references	195 23 5 67 4 2 253
1521	Total contacted		

Appendix F

Local Agency team involved in the study:

Local Agencies	Name: DataVision International Ltd. Country: Tanzania Activities since: 1998
Enumerators involved:	Enumerators: 23 Recruiters: 23
Other staff involved:	Fieldwork Coordinators: 2 Data Processing: 1

Sample Frame:

Characteristic of sample frame used:	Census of registered businesses operating in Tanzania Census of registered businesses operating in Zanzibar
Source:	2011/12, Tanzania National Bureau of Statistics (NBS) 2012, Office of Chief Government Statistician (OCGS)
Year:	2011/12 and 2012

Sectors included in the Sample:

Original Sectors	Eight selected economic activities, viz., food processing (ISIC ⁸ 15), textile and garments (ISIC 17 & 18), fabricated metal products (ISIC 28), furniture (ISIC 36), construction (ISIC 45), hotel and restaurant (ISIC 55), transport (ISIC 60 & 61) and Information technology (ISIC 72)).
Added (top up) Sectors	None

⁸ ISIC refers to the ISIC code revision 3.1.

Fieldwork and country situation:

Date of Fieldwork	22 nd April 2015 to 30 th June 2015.
Country	Tanzania
Use of CAPI	<ul style="list-style-type: none">• Computer-assisted personal interviewing (CAPI) was used to collect data. The TIKITI mobile research app was used.
Problems found during fieldwork:	<ul style="list-style-type: none">• Some of the respondents refused to participate in the survey, mostly noting that they have been participating in a lot of similar surveys with no benefits to them and for which they never obtained any feedback.• Some respondents were less transparent, particularly on financial information. This necessitated repeated callbacks to convince respondents to answer these questions.• The questionnaire was a bit long for managers to complete. Consequently, in many cases managers either asked enumerators to return on another date, or referred them to another person who did not have enough information about the company as the managers do. Further, respondents generally found some of the questions difficult to recall/answer. In addition to question on financial related information, respondents have hard time on questions relating to names of school and university attended by employees since that requires digging the human resource record for all employees.• Although the sampling frame (the Tanzania CRE) is the most comprehensive frame for the country, it is less up-to-date as determined during the fieldwork. The contact addresses were found out to be incorrect and in some cases the firm does not exist or is now a private household. This made the screening process to locate a firm a bit time consuming.