



# A GUIDE TO SAMPLING

DRAFT FOR INTERNAL DISCUSSION

This version: Oct 22, 2020

DRAFT

# Contents

Contents .....	2
I. Introduction.....	3
II. Overview of the measurement objectives and survey programs .....	3
II.A. Measurement objectives.....	3
II.B. 50x2030 Integrated Survey Programs .....	4
III. Populations of Interest .....	5
III.A. Households .....	6
III.B. Agricultural Holdings .....	6
IV. Sampling Frames .....	8
IV.A. Overview of multiple frame survey .....	8
IV.B. Frame for the Integrated Agricultural and Rural Survey Program .....	9
IV.C. Frame for the Agricultural Survey Program.....	11
V. Sampling designs .....	12
V.A. Sampling design for the Integrated Agricultural and Rural Survey Program .....	12
V.B. Sampling design for the Agricultural Survey Program.....	14
VI. Sample size .....	15
VI.A. Sample size for the household sector .....	15
VI.B. Sample size for non-household sector and special farms .....	18
VII. Stratification .....	18
VII.A. Stratification and allocation of PSUs .....	18
VII.B. Intra-PSU stratification and allocation .....	20
VIII. Estimations .....	20
VIII.A. Estimators and variances.....	21
VIII.B. Issue of difference between sampling and observation units in the household sector .....	24
VIII.C. Weighting using different observation units.....	25
IX. Subsampling and estimations.....	25
IX.A. Subsample size .....	26
IX.B. Subsampling for specific modules/questionnaires.....	26
IX.C. Case of crop cutting.....	27
IX.D. Case of farm level post-harvest losses .....	28
IX.E. Estimations with subsampling.....	28
X. Coverage of special farms (commercial, large, modern...) and dual frame estimation .....	29
X.A. Screening approach .....	30
X.B. Dual frame estimators.....	31
XI. Longitudinal data collection scheme and tracking procedures.....	34
XI.A. Length of panel/rotation design.....	35
XI.B. Handling sample attrition.....	35
XII. Integration in existing statistical systems.....	36
XII.A. Both agricultural and household surveys exist in the country .....	37
XII.B. Only an agricultural survey system exists.....	37
XII.C. Only a household survey exists .....	37
XIII. References.....	38

# I. Introduction

The 50x2030 Initiative to Close the Agricultural Data Gap aims to strengthen national data systems so that they are better equipped to meet the data demands coming from global, regional, and national data reporting systems and obligations. In particular, countries adopting the 50x2030 survey approach will be well-positioned to produce official statistics with sound methodology and report on critical agricultural-related SDG indicators, as well as to understand the drivers of agricultural productivity and income and their linkages with welfare and rural development. The approach integrates the collection of data on the basic features of the agricultural sector (including annual production figures) with a broader set of data on economic, environmental, and social factors of relevance to rural areas.

The 50x2030 Initiative proposes modular and integrated survey programs, whereby key agricultural data, namely production, is collected on an annual basis, while more in-depth agricultural data is collected every three years with an eye for understanding, not only monitoring, agricultural systems. The system builds on the experience of the FAO's Agricultural Integrated Surveys Programme (AGRIS) and the World Bank's Living Standards Measurement Study - Integrated Surveys on Agriculture (LSMS-ISA) programs and, just like those programs, will be an integral part of national statistical systems.

As any multipurpose survey system, the sampling design should be carefully elaborated to produce reliable estimates in a cost-effective way to fulfil the measurement objectives of the 50x2030 Initiative. The objective of this document is to discuss key technical features of a suitable sample design for the survey programs proposed by the Initiative. Starting with the definitions of the populations of interest, the document then discusses the development of sampling frames, stratification criteria, sampling size calculations, estimation procedures and sampling approaches over time. The final section covers sampling issues considering existing survey programs of countries adopting the 50x2030 Initiative's survey system.

## II. Overview of the measurement objectives and survey programs

### II.A. Measurement objectives

The measurement objectives of the 50x2030 Initiative consider in priority the estimations of official statistics of interest of the country and critical agricultural-related SDG indicators. In particular, the 50x2030 initiative focuses on SDG 2 (Zero Hunger) and 5 (Gender Equality), and data collection for the computation of four high priority SDG indicators:

- 2.3.1 – Volume of production per labour unit by classes of farming/pastoral/forestry enterprise size;
- 2.3.2 – Average income of small-scale food producers, by sex and indigenous status;
- 2.4.1 – Proportion of agricultural area under productive and sustainable agriculture;
- 5.a.1 – (a) Proportion of total agricultural population with ownership or secure rights over agricultural land, by sex; (b) Share of women among owners or rights-bearers of agricultural land, by type of tenure.

In addition to that the survey program may need to address additional indicators required at regional or international level. For instance, African countries are required to report the indicators of the Comprehensive Africa Agriculture Development Programme (CAADP). The survey instruments promoted

through the 50x2030 Initiative aimed to address several CAADP indicators and can be adapted and expanded to include national priority indicators as well as additional SDG indicators (50x30 Initiative, 2020). Covering so broad measurement requirements in a survey operation is a great challenge that the sampling design should address.

## II.B. 50x2030 Integrated Survey Programs

Two survey programs are supported by the Initiative: the Agricultural Survey Program (Agricultural Program) with only agricultural components and the Integrated Agricultural and Rural Survey Program (Integrated Program) that integrates the agriculture components with a household survey component.

**The 50x2030 Agricultural Survey Program** is a modular survey system with an annual core survey tool focused on crop, livestock, aquaculture, fishery, and forestry production ('CORE-AG'), and a set of specialized tools covering such topics as costs and farm income; labor and productivity; gender decision-making in agriculture; production practices and environmental aspects of farming ('ILP-AG', 'ILS-HH', 'PME-AG', 'MEA-AG'; illustrated in Figure 1). These specialized tools are administered at lower frequencies. Additional specialized instruments may be added according to country needs and demand. An example of the sequence in which the survey tools may be administered is presented in Figure 1, though this may be altered according to country needs.

FIGURE 1. SCHEMA OF THE 50x2030 AGRICULTURAL SURVEY PROGRAM

Years	1	2	3	4	5	6	7	8	9	10
Core Agricultural Module (CORE-AG)	█	█	█	█	█	█	█	█	█	█
Farm Income, Labor, and Productivity (ILP-AG)	█			█			█			
Production Methods and Environment (PME-AG)		█			█			█		
Machinery, Equipment, and Assets (MEA-AG)			█			█			█	

**The 50x2030 Integrated Agricultural and Rural Survey Program** follows the same logic as the Agricultural Program but integrates the agriculture tools with a household survey tool and broadens the target population to incorporate a sample of rural non-agricultural households into the system every three years (as illustrated

	Agricultural sector		Non-Agricultural sector
	Non-household sector holdings	Agricultural households	Non-agricultural households
Rural areas	█	█	█
Urban areas	█	█	

Agricultural Survey Program  
 Integrated Agricultural and Rural Survey Program

in Figure 1

). The Integrated model allows

partner countries to better understand, on the one hand, the drivers and dynamics of rural development,

structural transformation, and its linkages with agriculture; and on the other hand, the linkages between agricultural productivity and income with aspects of welfare and livelihoods, such as educational outcomes, non-agricultural income, or shocks and coping. The Integrated model achieves this through the combination of the Farm Income, Labor, and Productivity (ILP-AG) questionnaire and the Non-Farm Income and Living Standards Household (ILS-HH) questionnaire, which are administered together every three years (**Error! Reference source not found.2**).

FIGURE 2. SCHEMA OF THE 50x2030 INTEGRATED AGRICULTURAL AND RURAL SURVEY PROGRAM

Years	1	2	3	4	5	6	7	8	9	10
Core Agricultural Module (CORE-AG)	█	█	█	█	█	█	█	█	█	█
Farm Income, Labor, and Productivity (ILP-AG)	█			█			█			
Production Methods and Environment (PME-AG)		█			█			█		
Machinery, Equipment, and Assets (MEA-AG)			█			█			█	
Non-Farm Income and Living Standards (ILS-HH)	█			█			█			

### III. Populations of Interest

In line with the measurement objectives of the Initiative, the target populations of the survey program are (i) all households in rural areas and (ii) all agricultural holdings in the country. The Agricultural Program covers all agricultural holdings (in both household and non-household sectors) as population of interest while the Integrated Program integrates a household-based survey with the farm-based agricultural survey covering all households in rural areas and all agricultural holdings as population of interest.

FIGURE 3. COVERAGE AND POPULATIONS OF INTEREST OF THE AGRICULTURAL PROGRAM AND THE INTEGRATED PROGRAM IN 50x2030

	Agricultural sector		Non-Agricultural sector
	Non-household sector holdings	Agricultural households	Non-agricultural households
Rural areas	█	█	█
Urban areas	█	█	

Agricultural Survey Program  
 Integrated Agricultural and Rural Survey Program

### **III.A. Households**

The Initiative will consider, as a reference, the definition of households established by the UN World Population and Housing Census Programme 2020 (UN, 2017):

*“the concept of household is based on the arrangements made by persons, individually or in groups, for providing themselves with food or other essentials for living. A household may be either (a) a one person household, that is to say, a person who makes provision for his or her own food or other essentials for living without combining with any other person to form part of a multi-person household, or (b) a multi-person household, that is to say, a group of two or more persons living together who make common provision for food or other essentials for living. The persons in the group may pool their resources and may have a common budget; they may be related or unrelated persons, or constitute a combination of persons both related and unrelated”.*

If the national definitions differ from the international standards, the Initiative will help countries identify the implications of using a different definition and assist them in moving towards the recommended definition.

### **III.B. Agricultural Holdings**

Definition adopted for agricultural holdings is from the FAO World Programme for the Census of Agriculture (WCA) 2020 (FAO, 2015a):

*“Agricultural holdings are economic units of agricultural production under single management, comprising all livestock kept and all land used wholly or partly for agricultural production purposes, without regard to title, legal form or size. Single management may be exercised by an individual or household, jointly by two or more individuals or households, by a clan or tribe, or by a juridical person such as a corporation, cooperative or government agency. The holding’s land may consist of one or more parcels, located in one or more separate areas or in one or more territorial or administrative divisions, providing the parcels share the same production means, such as labor, farm buildings, machinery or draught animals.”*

The WCA distinguishes two types of agricultural holdings: (i) holdings in the household sector and (ii) holdings in the non-household sector.

#### **III.B.1. Holdings in the household sector**

Broadly speaking, agricultural holdings in the household sector are holdings operated by household members for their own account (either for sale or for own use). In some countries, a threshold (minimum size limit) is adopted for agricultural holdings (such as the area of agricultural land operated and number of livestock raised).

Households operating agricultural holdings for their own account are often called “agricultural production households” or simply “agricultural households”. As mentioned in the FAO World Programme for the Census of Agriculture 2020 (FAO, 2015a), there is usually a one-to-one correspondence between an agricultural

household and an agricultural holding. In other words, all the own-account agricultural production activities by members of a given agricultural household are usually undertaken under single management. It is unusual that different household members operate agricultural land or livestock completely independently, rather they pool together the income derived from these activities. Even if there is a degree of independence in the agricultural activities of individual household members, the income or produce generated by different household members is usually pooled.

However, although unusual, FAO (2015a) mentions two special cases where the agricultural holding and household concepts may diverge:

- If there are two or more units making up a household, such as where a married couple lives in the same dwelling as their parents, the two units may operate land independently but, as members of the same household, they make common arrangements for food and pool incomes. In this case there is a single household but two separate holdings.
- In addition to an individual household's agricultural production activities, a household may operate land or keep livestock jointly with another household or group of households. In this case, there are two agricultural holding units associated with the household and two sets of activities: (i) the agricultural production activities of the individual household itself; and (ii) the joint agricultural operations with the other household(s).

Section VII.B discusses how to treat such cases during data collection and estimation.

### **III.B.2. Holdings in the non-household sector**

The WCA 2020 considers the holdings of the non-household sector as agricultural holdings operated by entities such as corporations, government institutions, cooperatives, etc. A clear definition of non-household holdings and their demarcation from holdings in the household sector is essential because it has direct implications on the sampling design. In some contexts, confusion may arise between the household and the non-household sectors especially with respect to households that operate large/modern, commercial farms that could be classified as quasi-corporations in the System of National Accounts. Provided that business registers are among the primary sources for establishing the sampling frame for non-households holdings, FAO advises the use of such registration as a criterion for classifying the quasi-corporations in the non-household sector. Therefore, holdings in the non-household sector could be considered as agricultural holdings operated by:

- Corporations as defined by the System of National Account
- Government institutions: agricultural production entities operated by a central or local government directly or through a special body.
- Cooperatives
- Institutional households such as hospitals, schools, prisons, religious institutions, etc.
- Non-profit institutions
- Registered quasi-corporations: households with large/modern or specific agricultural operations in which income and expenditure flows from agricultural activities can be separated from the other household activities, and the operations are registered in the nation's business registry.

Corporations and government institutions may have complex structures, in which different activities are undertaken by different parts of the organization. In some cases, corporations may consist of different establishments that would constitute different agricultural holdings. The WCA advises using the national accounting concept of establishment when dealing with corporations in the listing of non-household farms, whereby an establishment is an economic unit engaged in one main productive activity, operating in a single location (FAO, 2015a). For instance, a corporation that owns two distinct establishments in different locations would be considered as two different agricultural holdings.

### III.B.3. Special farms

Some countries may have particular interest in special farms (commercial farms, large farms, modern farms, specialized farms, etc.) that could be in both household and non-household sectors. In that framework, it would be important to develop a clear definition of these farms to be considered in sampling and estimations procedures. Estimations procedures with two overlapping frames (dual frame) are discussed in section VII.E.

## IV. Sampling Frames

A quality sampling frame is necessary for producing reliable survey estimates. In the framework of agricultural surveys, FAO (2017) mentions the following issues with sampling frames that are often faced and must be avoided:

- Under-coverage or incomplete frame: failure to include some holdings in the sampling frame.
- Over-coverage: some units that are not agricultural holdings are included in the frame. A particular case is constituted by holdings that are listed in the frame but no longer exist.
- Multiplicity: some holdings are duplicated, resulting in greater probabilities of being included in the sample.
- Clusters of elements: some holdings in the frame are in fact clusters of holdings rather than individual holdings. This is usually due to a misunderstanding of the holding's definition on part of some enumerators, who may inventory two or more holdings as a single one.

Research projects on sampling frames for agricultural statistics implemented by FAO in the framework of the Global Strategy to Improve Agricultural and Rural Statistics (Global Strategy) highlighted the importance of using master sampling frames for cost-effectiveness, consistency and integration of agricultural statistics in countries. A master sampling frame is a frame that enables selection of different samples (including from different sampling designs) for specific purposes: agricultural surveys, household surveys, and farm management surveys. Such a frame enables samples to be drawn for several different surveys or different rounds of the same survey, which makes it possible to avoid building an ad hoc frame for each survey (FAO, 2015b).

### IV.A. Overview of multiple frame survey

Broadly speaking, multiple-frame survey refers to surveys where two or more frames (dual or multiple-frame) are used and independent samples are selected from each frame. In the literature, the use of many frames in a survey is motivated by various reasons including:

- **Absence of a frame with a complete coverage of the population of interest:** there are situations where it is difficult (or sometimes impossible for some populations) to get a frame with complete



coverage. A solution is to use many frames on the same populations for improved coverage, exploiting the strengths and offsetting the weaknesses of each type of frame.

- **Production of reliable statistics on subpopulation, rare and hard to reach populations** depending on the country's interest.

Advantages associated with the use of multiple frames include (i) saving cost in under sampling expensive frame and over sampling cheaper frames (e.g. area and list frames) and (ii) having more flexibility in the survey design to better control survey costs, coverage, response rates, accuracy (use of different sampling designs for different frames, use of different modes of data collection: face-to-face, phone, web interviews). In the framework of agricultural survey, this approach facilitates (i) the coverage of all agricultural holdings in both household and non-household sector, (ii) the use of the suitable sampling design recommended for each type of holdings, (iii) the production of reliable statistics on special farms whom countries are interested in. However, estimations procedures can be quite complex in particular when the number of frames become high (typically more than three).

There are two important requirements for the use of multiple frame:

- **Completeness:** the union of all frames should provide full coverage for the target population. In this way, every element should be listed in at least one of the frames.
- **Identifiability:** for any sampled unit, it should be possible to understand whether or not it belongs to one of the frames.

When there is no overlap, estimations procedures are straightforward as independent samples are selected from each of them for survey implementation. In presence of overlap, methods proposed in the literature for dual-frame estimations are discussed below.

Methodological work on the sampling strategy of the Agricultural Integrated Survey (AGRIS) recommended two types of master sampling frames for integrated agricultural surveys (FAO, 2017):

1. A multiple frame consisting in two list frames: list agricultural holdings (i) in the household sector and (ii) in the non-household sector.
2. A multiple frame consisting of an area frame and the following two list frames: (i) landless holdings raising livestock and (ii) large commercial agricultural holdings).

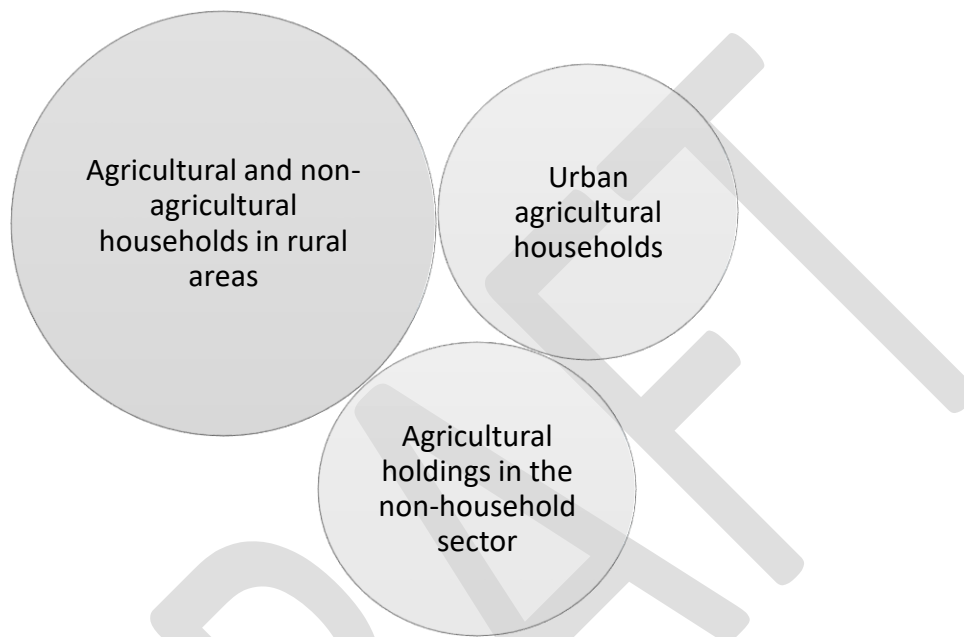
Both recommendations of multiple frames are valid for the Agricultural Survey Program (without the rural non-farm household sample) of the 50x2030 Initiative. However, for the Integrated Agricultural and Rural Survey Program, the first recommendation (to be combined with a frame of rural households) is the most appropriate. In fact, it complies with the necessity of the selection of a representative sample of households (irrespective of agricultural engagement) in rural areas. Using an area frame will lead to an indirect sample of agricultural households with no guarantee of representativeness regarding some households' typologies (more details in the section on sampling methods). However, to further improve the estimations of agricultural area and production, an area frame can be associated to the list frame of households in the multiple frame scheme.

#### **IV.B. Frame for the Integrated Agricultural and Rural Survey Program**

The frame for the Integrated Agricultural and Rural Survey Program should cover rural households (irrespective of agricultural engagement) and all agricultural holdings in the country. Accordingly, an ideal

frame would be complete lists of rural households and all agricultural holdings in both the household and non-household sectors. In the household sector, agricultural holdings are generally indirectly sampled through agricultural households because a direct listing of agricultural holdings is practically problematic in most contexts as holdings may include many parcels in different locations. Therefore, a suitable master sampling frame is multiple frame composed of (i) the list of all agricultural and non-agricultural households in rural areas, (ii) the list of urban agricultural households and (iii) the list of agricultural holdings in the non-household sector.

Fig.2 Sampling frame for the Integrated Agricultural and Rural Survey Program



#### **IV.B.1. Developing the frame of households (urban and rural)**

A complete list of agricultural and non-agricultural households in rural areas and urban agricultural households could be established from data of the Population and Housing Census (PHC). However, countries should be aware of the potential challenges with the use of PHC data, especially the identification of agricultural households. In particular in urban areas, there is a risk of over coverage of agricultural households if households whose members practice agriculture as paid employees, for examples, are considered as agricultural households in the frame. When the PHC is implemented taking into account the United Nations Principles and Recommendations for the Population and Housing Censuses, FAO (2015b) advises that households are considered to be agricultural households if at least one household member's *Economic Activity Status* is *Own-account worker* or *self-employed* or *employer*, and if the main occupation is an *occupation related to agriculture* (growing of crops, animal production, plant propagation, mixed farming).

However, it is very unusual to have on hand current PHC data to directly select a sample of households for survey implementation. A standard practice, where multistage sampling design is suitable (see section below on sampling design), is to select a sample of primary sampling units (e.g. enumeration areas) from the PHC data and perform a fresh household listing (micro-census) in each of them to select the final sample of households. For the specific case of urban areas, it may be important to assess the suitability of a multistage sampling design before using such approach, in particular for countries that do not have previous experience

in implementing agricultural surveys. This may be performed through using information from previous PHC or agricultural census to estimate the design effect. In presence of high design effect, a simple random sampling would be more cost effective. That would require the development of an exhaustive list of agricultural urban households through a field listing operation or based on available administrative sources. Big data could be also explored (see Young et al. 2018).

#### **IV.B.2. Coverage of urban agricultural households**

The inclusion of the population of urban agricultural households is optional, depending on the importance of urban agriculture in the household sector in the country. If the incidence of urban agricultural households is particularly low, it would be cost effective to ignore these populations in data collection and consider their contribution in the estimation stage through. For example, final aggregates may be adjusted using the proportions of urban agricultural households in the inference domains. In the sampling design for the 50x2030 Initiative, these populations are treated as an optional domain that countries may decide to take into account considering their importance and survey operational costs. The proportion of agricultural households in urban areas out of all agricultural households is rather low in several countries. However, these households may be responsible for a large share of the overall agricultural production, or the production of specific products (e.g. horticulture, livestock...). In such case, it may be important to consider them in the survey sample.

#### **IV.B.3. Developing the frame of holdings in the non-household sector**

The ideal situation for developing a complete list of agricultural holdings in the non-household sector would be using data of a recent census of agriculture that covered all farms in the country. Otherwise, to build the frame for these holdings, business registers of farms, including the national business register and informal business registers of farmers' organisations, are often used as a starting point. Efforts must be made to handle the probable overlap between the formal and informal registers. In addition, all other relevant registers should be considered including the list of government institutions (agricultural research centres, schools, hospitals, prisons etc.) and non-government organisations operating farms. Local knowledge and information from extension agents and local authorities about large specialty-type farms are useful in the process of developing the non-household sector frame.

### **IV.C. Frame for the Agricultural Survey Program**

The main requirement for this frame is suitable coverage of all agricultural holdings in the country, i.e. all agricultural households and all agricultural holdings in the non-household sector. The directives discussed in section II.A.1 for developing the frame of households can be used as well for the specific case of agricultural households. However here, apart from the PHC data, data from recent census of agriculture may be naturally used to develop the sampling frame for not only agricultural households but also non-household farms if they were covered by the census. As explained in section II.A.1, in most cases, census data will be outdated and the selection of a sample of primary sampling units from the data of the most recent PHC or census of agriculture followed by fresh listing operations in each of them would be necessary.

Recommendations for developing master sampling frame for integrated agricultural surveys are provided in the Handbook on the Agricultural Integrated Survey.<sup>1</sup> The Agricultural Survey Program of the 50x2030

---

<sup>1</sup>chapter 5 of FAO (2017)

initiative has quite a similar structure, and so countries are referred to these recommendations for the development of sampling frames for the program. The two main recommendations are summarised below:

- 1) Building a complete list of agricultural holdings (i) in the household sector and (ii) in the non-household sector. The list for the non-household sector can be developed as explained above in section II.A.3. For the household sector, a cost-effective approach is to link the population and agricultural censuses as suggested by the WCA 2020 (FAO, 2015a).
- 2) Using a multiple frame consisting of an area frame and two list frames (landless holdings raising livestock and large commercial agricultural holdings). The operational process for building an area frame for agricultural surveys is explained in FAO (2015b). Considering that an area frame does not cover landless holdings that raise livestock, a complementary listing of these holdings is recommended. In addition, if large agricultural holdings happen to be sampled from an area frame, they may behave like outliers. A second list of large commercial agricultural holdings is therefore also recommended (FAO, 2017). Procedures for developing these lists (landless holdings and large farms) are quite similar to the ones proposed for holdings in the non-household sector (see section III.B.3 above). The starting point would be the consultation of agricultural registers (formal and informal) in the country and other administrative sources to establish a preliminary list that should be updated with information from extension agents, local authorities and ad hoc field listing operations. Estimation approaches with overlapping sample frames are discussed in section VII.

## V. Sampling designs

The sampling strategy of the 50x2030 Initiative combines the experiences of the FAO AGRISurvey Programme and WB LSMS-ISA. Accordingly, the sampling design of the 50x2030 survey program should integrate the main features of the sampling strategies of the two survey systems.

### V.A. Sampling design for the Integrated Agricultural and Rural Survey Program

The integrated sampling design should ensure reliable estimates of the main variables of interest on the populations of interest at the level of the estimation domains (also called domains of inference or analytical strata). These domains are usually administrative zones for which countries expect reliable statistics from the survey data: regions, provinces, districts, etc. The domains are country specific and the sampling design should be adapted to each of them. Countries should be advised that the choice of the domains has a direct impact on the overall sample size and accordingly on the budget: the lower the geographical level estimation domain, the higher the overall budget. The sampling design for the Integrated Agricultural and Rural Survey Program is summarised in Table 1 below for the different sectors considered. These recommendations could need to be adjusted to meet the needs of countries' specific situations

Table 1. Summary of the major elements of the sampling design for the Integrated Program

Items	Populations of interest		
	Household (rural)	Household (urban)	Non-household sector
Observation units	<ul style="list-style-type: none"> <li>▪ Households</li> <li>▪ Agricultural holdings</li> </ul>	Agricultural holdings	Agricultural holdings
Final sampling units	Households	Households	Agricultural holdings

Frames	List of households from population census or list of EAs from population census and micro censuses in sampled EAs	List of households from population census	List of non-household farms developed from registers and/or field operation
Sampling method	Stratified two-stage	Country specific: Stratified one-stage or two-stage	Stratified one-stage
Stratification	Country specific: PSU level strata: Administrative zones, Agro-Ecological Zones Intensity of agricultural activity using land use data, proportion of agricultural households SSU level strata (intra PSU): practice of agriculture	Country specific: Administrative zones, Agro-Ecological Zones	Country specific: Production systems (crop/livestock/mixed), ad-hoc categorization
Sampling scheme	1 <sup>st</sup> stage: PPS of PSUs (EAs) 2 <sup>nd</sup> stage: Systematic or Simple random sampling without replacement of Households	Country specific: depending on the sampling method adopted	Systematic or Simple random sampling without replacement within each stratum

### V.A.1. Household sector (rural areas)

A standard stratified two-stage sampling design is recommended for this sector. This sampling scheme is widely used in households and agricultural surveys. The ideal frame, as discussed above, should be the list of rural households from a recent population and housing census (PHC).

The primary sampling units (PSU) are preferably enumeration areas (EAs) defined for the PHC because they are generally quite homogeneous in terms of population size. Homogeneity in population will serve to improve final estimates. In some cases, it may be useful to merge some small EAs and split up big ones to improve the overall size homogeneity. FAO (2017) discusses issues related to the choice of the PSU. The PSUs should be stratified as discussed in section VI.A. In each stratum, a sample of PSUs is drawn with a probability-proportional-to-size (PPS) (without replacement).<sup>2</sup> The measure of the size of PSUs (EAs) is usually equal to the number of households within that enumeration area (as results from the sampling frame).

The secondary sampling units (SSU) are households. Within each sampled PSU, the sample of SSUs may be selected by means of stratified simple random sampling (or systematic sampling) without replacement. The intra-stratification in the PSU is discussed in section VI.B.

### V.A.2. Household sector (urban areas)

Given the rather low proportion of urban agricultural households in most contexts (as discussed above), a stratified simple random sampling would be suitable for these populations in most countries. In most cases, that operation would be cost effective if the recommended frame (list of households from population census) is recent, allowing a direct selection of households without making a fresh listing in the field. In any case, a

<sup>2</sup> Probability proportional to size is a sampling procedure whereby the probability of selection of each unit in the universe is proportional to the size of some known relevant variable, in this case the population of number of households.

relatively high geographical dispersion of the sample should be expected and that would increase the data collection cost. If the total number of urban agricultural households is relatively important and depending on the repartition of these households in the urban area, a stratified two-stage sampling may be explored.

### V.A.3. Non-household sector

A stratified one-stage design is appropriate for holdings in the non-household sector (FAO, 2017). The stratification criteria may be the agricultural production systems (crop/livestock/mixed) or other ad-hoc typology.

### V.B. Sampling design for the Agricultural Survey Program

The main difference between the Agricultural Survey Program and the Integrated Agricultural and Rural Survey Program is the exclusion of non-agricultural households in rural areas. The Agricultural Survey Program covers the same populations of interest as the FAO's Agricultural Integrated Survey (AGRIS) Program. Two types of master sampling frames proposed in FAO (2017) could be used in this framework: either multiple frame of list frames or multiple frame of area frame and list frame. The basic features of the sampling designs for each type of frame, discussed in detail in FAO (2017) are presented in Table 2 below. Further customization could be needed depending on countries specific situations.

Table 2. Summary of the major elements of the sampling design for the Agricultural Program

Type of frame	Multiple frame of list frames		Multiple frame of area frame and list frame	
Sub-frame	Holdings of the household sector	Holdings of the non-household sector	Area frame	Lists
Observation units	Agricultural holdings (AH)	Agricultural holdings	Agricultural holdings	Agricultural holdings
Sampling units	AH of the household sector	AH of the non-household sector	Segments or points	-Landless AH raising livestock; -Large commercial holdings
Sampling method	Stratified two-stage	Stratified one-stage	Country specific: (i)Stratified cluster sampling or (ii) Stratified two-stage	Stratified one-stage
Stratification	Country specific: Administrative zones, Agro-Ecological Zones; Intensity of agricultural activity using land use data	Country specific: Production systems (crop/livestock/mixed), ad-hoc categorization	Land cover/use	ad-hoc categorization e.g based on size, activities
Sampling scheme	1 <sup>st</sup> stage: PPS of PSUs (EAs); 2 <sup>nd</sup> stage: SRSWOR of agricultural households	SRSWOR of agricultural holdings	Country specific: (i) PPS selection of segments and full coverage of all farms in the selected segments or (ii) two stage sampling: 1 <sup>st</sup> stage: PPS of PSUs (segments or grids/clusters of points); 2 <sup>nd</sup> stage: SRSWOR of segments or points	SRSWOR of agricultural holdings

## VI. Sample size

The calculation of the sample size is a critical step of the survey design to ensure the production of reliable statistics by keeping the sampling error to the minimum possible. The recommended approach to calculate the sample size is the one that considers the analytical requirements of the survey, i.e. ensuring reliable estimations of key variables of interest. The choice of the variable of interest can be performed considering key variables necessary for the calculation of the most important indicators expected from the survey operation.

A measure of statistical dispersion (coefficient of variation, variance, standard error, etc.) of the variable of interest in the population is required in the sample size formula. It should be ideally calculated from census data but that is not always possible given the limited scope of census questionnaires. A common alternative is to estimate it from data collected in previous surveys. The measure of statistical dispersion of any variable strongly correlated to the variable of interest may also be used. Suggestions of key variables of interest are proposed in the table below considering the measurement objectives of Initiative (see section II) and the recommended sampling frames.

Table 2. Variables of interest suggested for determining sample size

<b>Sampling units</b>	<b>Variables of interest suggested for determining sample size</b>
Households (agricultural and non-agricultural)	Income Agricultural area Agricultural production value
Agricultural holding in the non-household sector	Agricultural area Agricultural production value
Landless livestock holdings	Agricultural production value Livestock size unit
Large agricultural holdings	Agricultural area
Segments/points (area frame)	Agricultural area

In case there is a need to consider many variables of interest, the maximum of the minimum sample sizes required for each of them can be considered. This is what is proposed below for sample size of households for the Integrated Agricultural and Rural Survey Program.

If the key variable of interest is very heterogeneous in the population, the required sample size will be very high and may be unaffordable considering budget constraints. The solution in such situation is to consider the maximum sample size that the budget can support. Of course, that would not guarantee reliable estimations but the use of advanced statistical methods for optimal stratification and allocation may be helpful. One of such method is the optimal multivariate stratification and allocation procedures proposed by Barcaroli et al (2020) through the R Package 'SamplingStrata'.

### VI.A. Sample size for the household sector

In this sector, a stratified two-stage sampling design is recommended. The primary sampling units (PSU) recommended are usually the enumeration areas and the secondary sampling units (SSU) are either households (Integrated Agricultural and Rural Survey Program) or agricultural households (Agricultural

Survey Program). Approaches for calculating the sizes of both samples of PSU and SSU are discussed in this section.

### VI.A.1. Number of secondary sampling units

#### **Integrated Agricultural and Rural Survey Program**

In rural areas, the Integrated Agricultural and Rural Survey Program has two main estimations goals: producing estimates on the whole population of rural households and estimates for the subset of agricultural households at both national and sub-national level. The optimal sampling strategy to meet these objectives would require a complete list of rural households from a recent PHC, with their type, agricultural (denoted as A from now on) and non-agricultural households (denoted as B).

In the integrated survey, the household-sector sample size should ensure a reliable estimation of a key household related variable (e.g. income) in the population of rural households (A and B) and reliable estimation of a key agricultural variable (e.g. agricultural area) from the sub population of agricultural households (A) as households in the sub population B do not operate agricultural land. To calculate the minimum sample size of households to fulfil this goal, the usual approximate formula based on the coefficient of variation can be used.

Let's consider for each estimation domain  $U_d$ :

- $M_{Ad}$  and  $M_{Bd}$  total number of households respectively of type A and B.
- $cv_{Aincd}^2$  and  $cv_{Bincd}^2$  coefficient of variation of income of households of type A and B, respectively
- $cv_{Aland}^2$  coefficient of variation of agricultural area of the agricultural household.
- $cv_d^{*2}$  maximum relative error accepted for estimating the total (average) of income and agricultural area.
- $\widehat{deff}_{Aincd}$ ,  $\widehat{deff}_{Bincd}$  and  $\widehat{deff}_{land}$  estimates of the design effect for income of households of type A and B and agricultural area, respectively.
- $g$  is the expected response rate.

The minimum sample size of households ( $m_d$ ) in the domain  $U_d$  is:

$$m_d = \frac{1}{g} \left[ \text{Max} \left( \widehat{deff}_{land} \frac{cv_{Aland}^2}{cv_d^{*2} + \frac{cv_{Aland}^2}{M_{Ad}}}, \widehat{deff}_{Aincd} \frac{cv_{Aincd}^2}{cv_d^{*2} + \frac{cv_{Aincd}^2}{M_{Ad}}} \right) + \widehat{deff}_{Bincd} \frac{cv_{Bincd}^2}{cv_d^{*2} + \frac{cv_{Bincd}^2}{M_{Bd}}} \right]$$

Or

$$m_d = \max(m_{dA,inc}, m_{dA,land}) + m_{dB,inc} = m_{dA} + m_{dB,inc}$$

This procedure requires having all the variables in the formula for A and B type households (agricultural and non-agricultural rural households) in each domain  $d$ . However, it may happen for instance that the coefficient of variation of the income cannot be estimated for each sub population if the exercise is done with data from a household survey that did not cover agricultural activities. In such case if  $m_{d,inc}$  is the overall minimum size of rural households for a reliable estimate of the income, we have:

$$m_{d,inc} = \frac{1}{g} \widehat{deff}_{incd} \frac{cv_{incd}^2}{cv_d^{*2} + \frac{cv_{incd}^2}{M_{Ad} + M_{Bd}}}$$



And

$$m_d = \max(\tilde{W}_{Ad}m_{d,inc}, m_{d,land}) + (1-\tilde{W}_{Ad})m_{d,inc}$$

Where:

- $cv_{incd}^2$  is the coefficient of variation of the income of rural households in the domain  $d$
- $\widetilde{def}_{incd}$  is an estimate of the design effect for the income of rural households
- $\tilde{W}_{Ad}$  is an estimate of the proportion of agricultural households in the domain  $d$ .

### **Agricultural Survey Program**

This program being interested only in the agricultural households (A), the sample size will be simply:

$$m_d = \frac{1}{g} \widetilde{def}_{land} \frac{cv_{Aland}^2}{cv_d^{*2} + \frac{cv_{Aland}^2}{M_{Ad}}}$$

### **VI.A.2. Number of PSU**

When PSUs are selected with probability proportional to their size, selecting a fixed number of  $m_0$  households per PSU will allow having constant weights. This means that the number of PSUs to be selected in  $d$  would be given by dividing the sample size of households by  $m_0$ . With this approach, the number of PSUs to be selected in the domain  $d$  is given by

$$n_d = \left[ \frac{m_d}{m_0} \right] + 1$$

Where  $\left[ \frac{m_d}{m_0} \right]$  is the integer part.

The value of  $m_0$  can be determined considering both costs and homogeneity of households in the PSUs (intra-class correlation  $\bar{\rho}$ ),

$$m_0 = \sqrt{\frac{c_p \times (1 - \bar{\rho})}{c \times \bar{\rho}}}$$

where  $c_p$  and  $c$  are respectively the cost of adding an additional PSU into the sample and the unit cost of an interview. The intraclass correlation  $\bar{\rho}$  can be estimated from previous surveys; since two variables are considered, income and agricultural land area, the minimum value,  $\bar{\rho} = \min(\bar{\rho}_{cons}, \bar{\rho}_{land})$ , should be considered (it is a conservative choice). It is worth noting that this formula is an approximation based on two stage simple random sampling of both PSUs and SSUs, when PSUs size does not vary greatly.

Alternatively, a common practice to facilitate the organization of the field work consists of fixing the number of households  $m_0$  to select in each PSU, considering the maximum enumerator workload during survey implementation. An arbitrary value generally varying between 10 and 15 is often considered.

## VI.B. Sample size for non-household sector and special farms

The sampling design recommended for the agricultural holdings in the non-household sector is a stratified simple random sampling. The same design is advised for complementary lists of special farms recommended in some cases (see section 4) to improve the coverage of the sampling frames like the landless agricultural holdings raising livestock or the large commercial holdings. For these populations, the agricultural production value would be a suitable variable to be considered for sample size calculations. However, for agricultural area and livestock size units may also be suitable to calculate the size of the sample for large holdings and landless livestock holdings respectively.

Let's consider for each estimation domain  $U_d$ :

- $M_d$  total number of holding.
- $cv_{yd}$  coefficient of variation of the key variable of interest
- $cv_d^*$  maximum error accepted
- $g$  is the expected response rate.

The sample size  $m_d$  can be calculated using the following formula:

$$m_d = \frac{1}{g} \frac{cv_{yd}^2}{cv_d^{*2} + \frac{cv_{yd}^2}{M_d}}$$

## VII. Stratification

Stratification consists of dividing the population into subpopulations (strata) and performing independent sample selection in each of them. Stratification can contribute to a great extent to decreasing sampling errors in particular when the subpopulations are homogeneous and if the sample size is well allocated in the strata. Other important advantages of stratification include the control of sample sizes for different sub-populations for reporting and analysis purposes and the possibility to adopt different sampling strategies for each stratum. In many countries, households and agricultural surveys are expected to produce reliable estimates at sub-national levels (regions, provinces, districts, etc.) that are the estimation domains or analytical strata. To improve estimations in the estimation domains, additional stratification, discussed here, is usually performed within them. Discussion is focused on stratification in the framework of a two-stage design with a list frame as stratification procedures are quite straightforward with area frames and in contexts where a stratified single-stage sampling design is recommended.

### VII.A. Stratification and allocation of PSUs

When implementing two-stage sampling, FAO (2017) recommends a stratification of the EAs by administrative zones (e.g. regions, provinces, etc.) and agro-ecological zones before the first stage selection in order to improve the estimates of agricultural statistics. Stratification of PSUs should be carefully controlled since having too many strata is not desirable (an independent sample has to be selected in each stratum). To avoid too many strata, explicit stratification can be coupled with an implicit stratification. That consists of sorting the sampling frame by relevant criteria (usually geographical) in each stratum and selecting an independent sample in each stratum with systematic sampling.

### Stratification in the Integrated Agricultural and Rural Survey Program

When the list of households from the PHC is outdated, the actual structure of the households within the sampled PSUs can be known only after a fresh listing of households in these PSUs. A major drawback is the lack of control of the final sample especially the number of agricultural households required in the domain. Since the selection is done at level of PSUs, that may show a varying situation in terms of proportion of agricultural households.

In the context of the Integrated Agricultural and Rural Survey, to maintain control of the final sample size by household type (A and B) it is preferable to make a first level stratification of the EAs in terms of proportion of agricultural households in each of them estimated from the latest PHC or other suitable sources. Even if the PHC data are considered outdated, this structural information (proportion of agricultural households) may not likely vary too much in all PSUs and could be helpful for stratification purposes. The first level stratification below may be considered using a proportion threshold  $\rho$  ( $\frac{1}{2} < \rho < 1$ ).

First level PSU strata	Definition
Agricultural	<i>Proportion of agriculture households in the PSU <math>\geq \rho</math></i>
Mixed	<i><math>1 - \rho &lt; \text{Proportion of agriculture households in the PSU} &lt; \rho</math></i>
Non-agricultural	<i>Proportion of agricultural households in the PSU <math>\leq 1 - \rho</math></i>

The sample of PSUs in the domain  $d$  ( $n_d$ ) can be allocated using parameters  $\theta_a$ ,  $\theta_m$  and  $\theta_{na}$  with  $\theta_a + \theta_m + \theta_{na} = 1$

First level allocation	
First level PSU strata	Allocation of the sample of PSU
Agricultural	$\theta_a n_d$
Mixed	$\theta_m n_d$
Non-agricultural	$\theta_{na} n_d$

If  $m_0$  households will be selected in each sampled PSU using a systematic or simple random sampling without replacement, the expected number of agricultural households in the final sample ( $m_{dAexp}$ ) is:

$$m_{dAexp} = \rho m_0 \theta_a n_d + (1 - \rho) m_0 \theta_m n_d + \delta m_0 \theta_{na} n_d = (\rho \theta_a + (1 - \rho) \theta_m) m_0 n_d + \delta \theta_{na} m_0 n_d$$

$\delta < 1$  is unknown before the selection of the sample of households contrary to the other parameters that are fixed by the sample designer.

Let's consider  $\tau$  the proportion of agricultural households in the planned sample.

$$\tau = \frac{m_{dA}}{m_d} = \frac{m_{dA}}{m_0 n_d} \Rightarrow m_{dA} = \tau m_0 n_d$$

To ensure the achievement of the planned sample of agricultural households in the final sample of households, parameters  $\theta_a$ ,  $\theta_m$  and  $\theta_{na}$  could be fixed to have  $m_{dAexp} \geq m_{dA}$ . That corresponds to:

$$(\rho \theta_a + (1 - \rho) \theta_m) m_0 n_d + \delta \theta_{na} m_0 n_d \geq \tau m_0 n_d$$

$\delta$  being unknown, parameters  $\theta_a$ ,  $\theta_m$  and  $\theta_{na}$  can be therefore fixed under the following conditions:

$$\rho \theta_a + (1 - \rho) \theta_m \geq \tau$$

$$\theta_{na} = 1 - (\theta_a + \theta_m)$$

A second level stratification of PSUs may be performed inside the first level strata. Common stratification criteria for improving estimates in agricultural and households surveys are: agro-ecological zones, land use classes, size categories based on population, agricultural area, intensity of agricultural activity.... For countries having a strong interest in gender-disaggregated statistics, a classification of PSUs based on the number of female headed households can be considered for stratification. The allocation in these second level strata can follow different criteria. Typically, in household surveys an allocation proportional to the population in the strata is considered. FAO (2017) recommends the optimal allocation of Neyman for agricultural surveys. Kish (1987, page 228) suggests a compromise solution between equal and proportional allocation:

$$n_{dh} = n_d \times \frac{\theta_{dh}}{\sum_{h=1}^{H_d} \theta_{dh}}$$

Where

$$\theta_{dh} = \sqrt{\left(W_{dh}^2 + \frac{1}{H_d^2}\right)}$$

$H_d$  is the number of strata in the domain  $d$ , while  $W_{dh}$  is the relative size of stratum  $h$  in domain  $d$ , it can be the proportion of PSU in stratum  $h$  compared to the domain total,  $W_{dh} = N_{dh}/N_d$ , (relative size in terms of population). A multivariate stratification and allocation (Barcaroli, 2020) or compromise power allocation (Bankier, 1988) could also be explored if the frame contains relevant variables correlated with households' income or agricultural area (household size, livestock, agricultural production, etc.) at PSU level

### **Stratification in the Agricultural Survey Program**

In the context of the Agricultural Survey Program, the same PSU stratification strategy proposed above could be used but the first level allocation is not necessary as this survey program does not cover non-agricultural households.

### **VII.B. Intra-PSU stratification and allocation**

Although intra-PSU stratification has a limited effect on the total variance, it can be helpful to ensure good coverage of specific populations of interest. For the Integrated Agricultural and Rural Survey Program, a relevant stratification criterion could be the practice of agriculture or not. Countries interested particularly in gender-disaggregated statistics may consider also the gender of the head of the household. For the Agricultural Survey Program, production systems (crop/livestock/mixed) should be explored if such information are available in the frame. A proportional allocation of the sample should be used in the PSU to ensure the final units have equal probability of selection in the PSU. Alternatively, an implicit stratification could be considered, i.e. ordering of households by type and subsequent selection of the random sample of  $m_0$  by means of systematic criterion.

## **VIII. Estimations**

Estimators, variances and specific issues related to estimations are discussed in this section.

## VIII.A. Estimators and variances

The main estimators and variances from the sampling schemes proposed are presented here.

### VIII.A.1. Stratified two-stage sampling design

#### Notation

$h$  = stratum

$H$  = total number of strata

$i$  = PSU

$N$  = total number of PSUs

$I_h$  = total number of PSUs in the  $h^{th}$  stratum

$j$  = SSU

$M_{hi}$  = total number of SSUs found in the  $i$ -th PSU in stratum  $h$  ( $j = 1, 2, \dots, M_{hi}$ )

$M = \sum_h \sum_i M_{hi}$  = total number of SSUs in the country

$F_{hi}$  = total number of SSUs listed in the sampling frame as belonging to the  $i$ -th PSU in stratum  $h$

$F_h = \sum_i F_{hi}$ , is the total number of SSUs listed in the sampling frame in stratum  $h$

$n_h$  = number of sample PSUs selected in stratum  $h$

$m_{hi}$  = number of sample SSUs selected in  $i$ -th PSU in stratum  $h$

$y_{hij}$  = value of the target variable  $Y$  observed on the  $j$ -th SSU, in  $i$ -th PSU in stratum  $h$

#### Estimators

The probability of selecting the SSU  $j$  in the sample is the product of the probability of selection of the PSU  $i$  in which it is located ( $n_h \frac{F_{hi}}{F_h}$ ) and its probability of selection in the PSU  $i$  ( $\frac{m_{hi}}{M_{hi}}$ ).

Thus, the *weight* assigned to the SSU  $f$  selected in the  $i$ -th PSU in stratum  $h$  is:

$$w_{hij} = \left( n_h \frac{F_{hi}}{F_h} \right) * \left( \frac{M_{hi}}{m_{hi}} \right)$$

The weights are roughly constant within each stratum when  $m_{hi}$  is constant ( $m_{hi} = m_{h0}$ ) and when the number of SSUs found in a sampled PSU is approximately equal to the number of SSUs resulting from the sampling frame ( $M_{hi} = F_{hi}$ ).

An estimate of the total amount of  $Y$  for the entire population may be computed with the following formula:

$$\hat{Y} = \sum_h \sum_i \sum_j w_{hij} y_{hij}$$

The mean of Y is can be estimated with two different estimators:

- *Simple mean*

$$\hat{\bar{Y}} = \hat{Y}/M$$

- *Weighted sample mean*

$$\tilde{\bar{Y}} = \frac{\hat{Y}}{\sum_h \sum_i \sum_j w_{hij}}$$

This latter estimator tends to be preferable to the simple mean when the total size of the population is unknown or uncertain (as occurs, for example, if the frame is obsolete) and when estimating the mean of Y for an unplanned domain of interest (sub population).

### **Variance**

The variance of the estimator of the total amount Y in the population is a rather complex formula and can be found, for instance, in Cochran (1977, equation 11.42). A simple approximate estimation of such a variance involving the PSUs alone can be obtained with the following estimator, provided by Särndal, Swensson, and Wretman (1992, p. 154), which overestimates this variance.

$$\tilde{v}(\hat{Y}) = \sum_{h=1}^H M_h^2 \frac{1}{m_h(m_h - 1)} \sum_{i=1}^{I_h} (\hat{Y}_{hi} - \frac{1}{m_h} \sum_{i=1}^{I_h} \hat{Y}_{hi})^2$$

where  $\hat{Y}_{hi}$  and  $\hat{Y}_h$  are the estimates of the total amount of Y at PSU and stratum levels, respectively.

An approximate estimator of the variance of the mean is:

$$\tilde{v}(\hat{\bar{Y}}) = \frac{1}{M^2} \tilde{v}(\hat{Y})$$

### **Coefficient of variation of the total**

$$\tilde{cV}(\hat{Y}) = \frac{\sqrt{\tilde{v}(\hat{Y})}}{\hat{Y}}$$

### **Coefficient of variation of the mean**

$$\tilde{cV}(\hat{\bar{Y}}) = \frac{\sqrt{\tilde{v}(\hat{\bar{Y}})}}{\hat{\bar{Y}}}$$

## VIII.A.2. Stratified one-stage sampling design

Let us consider the following notation:

$h$  = stratum

$H$  = total number of strata

$i$  = SSU

$N_h$  = total number of SSU in stratum  $h$

$N$  = total number of SSU

$n_h$  = number of sample SSU selected in stratum  $h$

$y_{hi}$  = value of the target variable  $Y$  observed on the  $i$ -th SSU, in stratum  $h$

The weight of the  $i$ -th SSU in stratum  $h$  is simply the inverse of its probability of selection:

$$w_{hi} = \frac{N_h}{n_h}$$

An estimate of the total amount of  $Y$  for the population of SSUs is:

$$\hat{Y} = \sum_h \sum_i w_{hi} y_{hi}$$

An estimate of the sampling variance is provided by:

$$\tilde{V}(\hat{Y}) = \sum_{h=1}^H N_h(N_h - n_h) \frac{1}{n_h(n_h - 1)} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2$$

where

$$\bar{y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi}$$

An estimate of the population mean of  $Y$  is provided by

$$\hat{\bar{Y}} = \hat{Y}/N$$

And the corresponding sampling variance is estimated through

$$\tilde{V}(\hat{\bar{Y}}) = \frac{1}{N^2} \tilde{V}(\hat{Y})$$

The coefficient of variation of the total and the mean can be estimated using the same formulas provided for the stratified two-stage sampling design.

### **VIII.B. Issue of difference between sampling and observation units in the household sector**

The sampling unit proposed for the survey programs in the household sector is the household. However, the observation unit for the agricultural component is the agricultural holding. Therefore, it is important to discuss the relationship between households and holdings to be considered in computing the sampling weights of the holdings. Four types of links are expected between observation and sampling units (Lessler and Kalsbeek, 1992). In the framework of integrated surveys, they can be described as follow:

- One-to-one: each household is associated with a unique agricultural holding and each agricultural holding is associated with a unique household.
- One-to-many: one household can be associated with many agricultural holdings, but each agricultural holding is associated with only one household.
- Many-to-one: many households can be associated with one agricultural holding, but each household is associated with only one agricultural holding
- Many-to-many: each agricultural holding may be associated with many households and each household may be associated with many agricultural holdings.

The number of sampling units that would lead to the collection of data from the same reporting unit is referred to as “multiplicity” (FAO, 2015b). Here, multiplicity arises when an agricultural holding is operated by two distinct households or more or when a single household operates two or more distinct holdings. In the household-sector, the latter case is very rare and it is operationally complex to capture in the survey considering the definition of the holding. Acknowledging the rarity of this situation and the complexity of treating the holdings present in the household as two separate entities, the 50x2030 survey instruments collect the information as if the holdings were a unique entity.

Existence of multiplicities leads to biased estimates (Lavalée, 2007). Falorsi et al. (2015) recommends using the Generalized Weight Share Method (GWSM), proposed by (Lavalée, 2007), when dealing with multiplicities between holdings and households.

The following operational recommendations can be made for the use of the GWSM in the framework of the 50x2030 integrated sampling design:

- (i). Identifying multiplicities during households’ listing: when listing households in PSU, include questions to identify multiple-households holdings.
- (ii). After sampling and during the actual survey, identify the sampled households linked to each multiple households-holding
- (iii). During data processing compute sampling weights of the multiple-households holdings using the Generalized Weight Share Method

Let’s consider

$w_i$  the sampling weight of the household  $i$  ( $w_i = 0$  if the household is not sampled)

$l_{ij} = 1$  if the household  $i$  operates the agricultural holding  $j$  and  $l_{ij} = 0$  otherwise



The sampling weight  $w_j'$  of the agricultural holding  $j$  can be expressed as:

$$w_j' = \frac{\sum_i l_{ij} w_i}{\sum_i l_{ij}}$$

### **VIII.C. Weighting using different observation units**

To fulfil the main objectives of the Initiative, two observation units are considered: households and agricultural holdings. However, as indicated above, there is evidence of the existence of multiple-household holdings and they should be identified in the sampling frame. In countries where there are many cases of households operating agricultural holding in partnership, given the requirement of collecting data at both household and holding levels, the following actions should be taken:

- Identify agricultural information that should be captured at both household and holding levels including revenues and expenses, assets, investments etc.
- When interviewing households operating a multiple-household holding, collect these information separately at household level and at holding level. Obviously, that will increase the burden of the interview for these respondents but hopefully as mentioned below, such cases are unusual in the field.
- Finally:
  - weighted estimations of household level data at national level will be done using the households' direct sampling weights as design weights. Design weights are then adjusted through calibration, post stratification, non-responses adjustments etc. to calculate the final weights.
  - weighted estimations of holding level data at national level will be done using the holdings' sampling weights ( $w_j'$ ) calculated as explained in the previous section and using the final weights of households.

Therefore, final household and holding weights will be different only in the case of holdings operated by more than one household.

### **VIII.D. Estimation for rare crops**

Covering rare items is a challenge in all surveys as it could require a very large sample. With the 50x2030 survey methodology, reliable estimations for most crops should be possible at national level but not for sub national domains. For very rare crops, there is no such guarantee even at national level but upon country request, advanced estimation approaches (e.g. small area estimation models) could be explored using available auxiliary information from administrative data or previous agricultural census. Such auxiliary information could be also used at the design stage to explore ad-hoc stratification to improve the coverage of rare crops for reliable estimates.

## **IX. Subsampling and estimations**

In agricultural and household surveys, subsampling can be used as a cost-effective tool for various purposes including the following.

### ***Use of different estimation domains for specific information***

If the main domains of the survey are, as usual, sub-national administrative areas (regions, provinces, districts etc.), the country could consider that some information are just necessary for estimation at the national level. Therefore, it will not be necessary to collect them in the full sample in each estimation domain. The sample size for country level estimation can be calculated for this information and the corresponding questions will be administered only to a subsample in each estimation domain. For instance, if most information collected by a rotating module are not requested for sub-national policy making, the rotating questionnaires could be administered to a subsample.

### ***Collecting information with higher operation cost and/or time requirements***

Some data collection methods provide high quality information but their implementation may be too expensive to implement at full scale. Objective measurements of land (e.g. using GPS) or of yield (crop-cutting) may be considered as an example. Such operations may be performed on a subsample when full scale implementation is not feasible, and the results may subsequently be used to correct measurement error for the whole sample. Successful construction of correction factors depends on the nature of the measurement error, and may not be possible in all cases.

#### **IX.A. Subsample size**

The size of the subsample can be calculated using the variability of key variables aimed to be collected through the subsample. One possibility with the subsampling is to stratify the main sample before selecting the subsample. In such case, Fuller (2009) proposes a formula to calculate the optimal subsample size. For instance let us suppose that the subsampling is performed for objective measurement of the yield and let  $y$  be the yield collected through declaration. If the sample is stratified into  $H$  strata, the optimal sample size of the subsample following Fuller (2009) will be

$$m_{\text{subsample}} = m \sqrt{\frac{\sigma_w^2 C_1}{\sigma_b^2 C_2}}$$

Where:

$$\sigma_w^2 = \frac{\sum_{h=1}^H \sigma_{yh}^2}{H}$$

$$\sigma_b^2 = \sigma_y^2 - \sigma_w^2$$

$\sigma_y^2$  and  $\sigma_{yh}^2$  are the variances of  $y$  in the sample and in the stratum  $h$ , respectively.

$m$  is the sample size of the survey

$C_1$  and  $C_2$  are the costs of the interview by declaration and objective measurement, respectively.

#### **IX.B. Subsampling for specific modules/questionnaires**

The 50x2030 Initiative broadly proposes a modular survey system with an annual core survey tool focused on crop, livestock, aquaculture, fishery, and forestry production (CORE-AG), and a set of specialized modules (see section II):

- Farm Income, Labor, and Productivity (ILP-AG)
- Production Methods and Environment (PME-AG)
- Machinery, Equipment, and Assets (MEA-AG)
- Non-Farm Income and Living Standards Household (ILS-HH)

These specialized modules are administered at lower frequencies and all of them are part of the Integrated Agricultural and Rural Survey Program. However, the ILS-HH is not administered in the Agricultural Survey Program.

As the specialized modules are generally considered for indicators and analyses that are required at the national level, it is suggested to consider for them for larger domains of estimation, potentially even at the national level only, instead of sub-national administrative areas. For example, while the CORE-AG questionnaire may seek data representative at the district level, and potentially for specific crops, questionnaires such as the ILS-HH may be implemented with a sample representative at the regional or national level, depending on country needs. This recommendation has some advantages as it will reduce the respondents' burden of an important part of the sample. The reduction of interview length has a direct effect on the data collection cost and the reduction of measurement errors due to respondent burden. However, the use of a subsample in this framework may affect the precision of estimators as estimations will be performed in a two-phase sampling scheme.

In order to improve the estimation of the total income (farm and non-farm), it is recommended to consider a unique subsample of agricultural households for both ILS-HH and ILP-AG (the same agricultural households will receive the ILS and ILP). As the non-agricultural households will be administered only ILS-HH, there is obviously no need to select a subsample of these populations for the survey.

The subsample can be selected simply through a simple random sampling without replacement using the large sample. However, if there are relevant information on the large sample for stratification, it can be stratified before sample selection. For instance, in the framework of a two-stage sampling, information collected with the listing questionnaire in the primary sampling units may be helpful for such stratification.

### **IX.C. Case of crop cutting**

Crop-cutting is an operationally demanding activity, but one that adds significant value to survey operations. It requires the acquisition of specific equipment and at least two visits of the enumerator of the holding (during planting and harvesting periods). Some countries, such as Niger and Burkina Faso, implement crop-cutting in their agricultural survey operations on the full sample of holdings and all plots covered by the operation. However, a subsample of holdings and/or plots may be used especially when only national level estimates of crop yields are expected from the survey or if the results of crop-cutting are aimed to be used for correcting yield estimates collected by farmers' declarations.

#### ***Subsampling approaches***

Options for subsampling for crop-cutting include:

- Directly selecting a subsample of plots (bypassing the holding level) for implementing crop-cutting. This is an efficient option (in terms of quality of estimations) but has some operational constraints. In fact, it can be performed only after the listing of all plots of the holdings. Therefore, this listing should be completed in a timely manner to allow the processing of the data and selection of the

subsample of plots. In addition, large plots or plots very far from the holding dwelling may appear in the sample increasing operation costs.

- Selecting a subsample of holdings and covering all or some plots of each subsampled holding by the crop-cutting operation. This option may allow coverage of more parcels than the previous at a lower cost as the holding is a cluster of plots. However, it is less efficient than the previous as cluster sampling leads to higher variance.

#### **IX.D. Case of farm level post-harvest losses**

The FAO's Guidelines on the measurement of harvest and post-harvest losses discusses the main postproduction operations during which harvest losses occur at the different stages of the value chain. At the farm level, in the case of grains (cereals and pulses) losses occur mainly during the following operations: harvesting, threshing or shelling, cleaning or winnowing, drying and storage in the holding (FAO, 2018a). The Guidelines recommend using probability sample surveys as the backbone of any loss assessment, complemented by other methods that may be used mainly as preliminary assessments or to further analyze certain aspects related to PHL. With such surveys, loss measurements can be (i) objective – drawn from crop-cutting on the field or laboratory analysis of grain sampled from storage facilities – or (ii) subjective, by asking the respondent (farmer, storage facility manager, etc.) to provide his or her own estimate of loss.

In the framework of the initiative, it is recommended to make post-harvest losses assessments using a subsample for a number of reasons including reducing operation costs and respondent burden. In fact, both options of loss measurements (objective and subjective) are relatively expensive and time-consuming, and require well-trained personnel (FAO, 2018a). Objective measurements are particularly expensive and require many visits of the enumerators to the farm. Subjective assessments are cheaper but require additional visits to the farm in particular for collecting information on storage losses; thus, additional cost may be important. In addition, field tests performed by FAO in the framework of the Global Strategy showed important measurement errors when comparing objective and subjective measurements in Ghana, Malawi, Namibia and Zimbabwe (FAO, 2018b).

Another reason for subsampling is that this information is generally expected at national level. Subnational domain estimations are therefore not particularly required. In addition, regarding international demand, SDG 12.3.1 on Global Food Loss Index is expected at country level.

Important indicators related to post-harvest losses are the proportions of losses at crop and operation level. The crop of interest should be specified, they are generally cereals and pulses. Therefore, the subsampling should be performed among holdings producing the target crops.

#### **IX.E. Estimations with subsampling**

If subsampling is used to collect specific information, estimations can be done using (i) a three-stage sampling scheme or (ii) a two-phase sampling scheme.

If all the sample of holdings are covered and a subsample of plots selected in each of them, that corresponds to a three-stage sample selection. New sampling weights should be calculated for the subsampled plots by multiplying the inverse of their probability of selection and the sampling weights of holdings. Estimators of variances provided for the two-stage design can still be used here with few adaptations.

However, selecting a subsample of holdings would correspond to a two-phase sampling. Estimations can be done using regression or ratio estimators. Regression estimators are considered more efficient in this context (Cochran, 1977). Let's consider the use of subsampling for objective measurements (GPS measurement or crop cutting) in a two-stage scheme (subsample of holdings). For simplicity we will consider the case of a sample holdings selected through a simple random without replacement from which a simple random subsample of holdings is selected for the objective measurements. Let's consider  $y$  the yield measured using the crop-cutting in the subsample and  $x$  the yield collected by declaration on the whole sample. From Sitter (1997), the regression estimator used to estimate a more accurate average yield is:

$$\bar{y}_{reg} = \bar{y} + \beta(\bar{x} - \bar{x}_{subsample})$$

Where:  $\bar{x}_{subsample}$  is the average of  $x$  in the subsample and  $\beta$  the least squares regression coefficient of  $y$  on  $x$  in the subsample.

An estimator of the variance of  $\bar{y}_{reg}$  is

$$\tilde{v}(\bar{y}_{reg}) = \left( \frac{1}{m_{subsample}} - \frac{1}{m} \right) s_d^2 + \left( \frac{1}{m} - \frac{1}{M} \right) s_y^2$$

Where:

- $m_{subsample}$  : size of the subsample
- $m$  : size of the whole sample
- $M$  : size of the population of holdings
- $s_d^2$  is the sample variance of the quantities  $d_i = y_i - \bar{y}_i - \beta(x_i - \bar{x}_i)$

## X. Coverage of special farms (commercial, large, modern...) and dual frame estimation

The definition adopted by the initiative for agricultural holdings of the non-household sector (non-households' holdings) allows to avoid any overlap of this population with the population of agricultural holdings of the household sector (households' holdings). However, it usually happens that countries are interested in some special farms (e.g. commercial farms, modern farms, large farms, organic farms, farms producing specific crops, etc.) for various reasons including the need of information for making decisions related to specific types of holdings or the monitoring of agricultural programs supporting particularly those holdings. In addition, complementary lists of holdings with special farms are recommended in some cases (see section IV) to improve the coverage of the sampling frames like the landless agricultural holdings raising livestock or the large commercial holdings. The agricultural holdings operating such special farms can belong to both household sector and non-household sector.

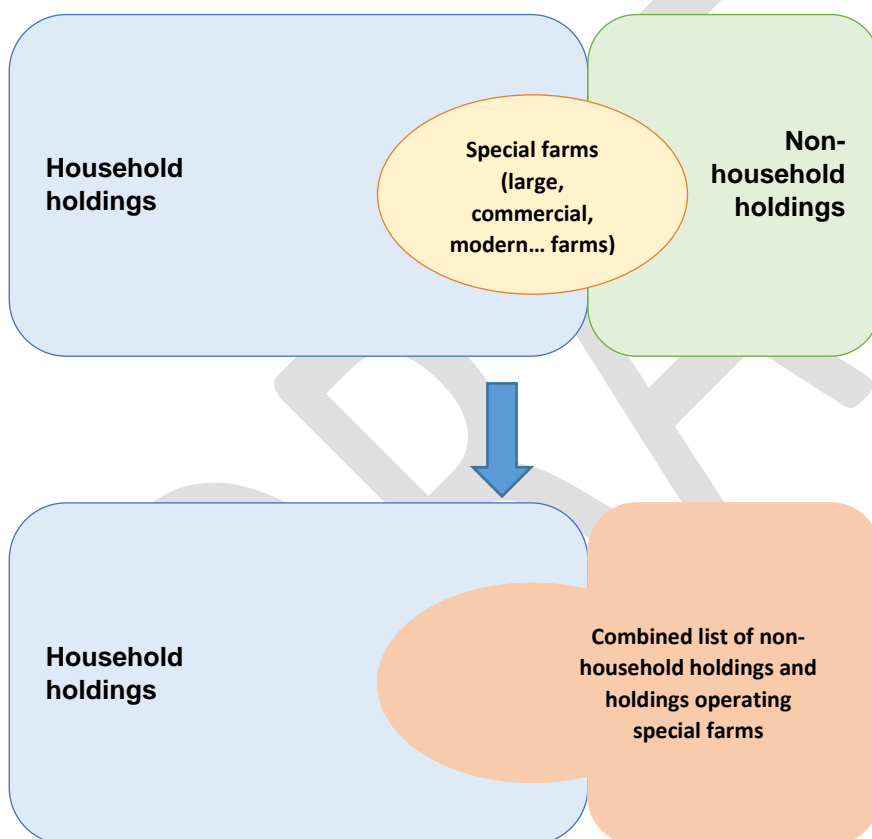
When a country has an interest in special farms, a clear definition of these farms should be elaborated and a list of special farms should be established through administrative data and/or listing operations as a sampling frame for them. This list will overlap with either the frame of households' holdings or the frame of non-households' holdings or both leading inevitably to a dual frame survey.

Regarding the sampling approach of special farms, it can be observed that some countries opt for the creation of an explicit stratum of agricultural households operating special farms after listing operations in the sample of primary sampling units used for the selection of households' holdings. However, this approach presents

some weaknesses affecting the efficiency of the sample especially when the proportion of households operating special farms is low and/or the geographical distribution of these farms is asymmetric. This does not guarantee reliable estimations on the population of special farms as the explicit stratification is performed in a sample and may lead to an under-sampling of households' holdings affecting the precision of estimates on them.

In practice, when special farms are of interest and list frames used as sampling frame, a single list can be developed for both populations of holdings operating special farms and non-households' holdings as the same sampling design is recommended for both (stratified one-stage), and listing strategies are quite similar for both. Thus, in the end, the latter list and the list of households' holdings will constitute an overlapping dual frame that will represent the sampling frame covering all holdings and allowing the production of reliable statistics on the special farms. The overlapping units are simply households' holdings operating special farms.

Fig.3 procedure of frame development for covering special farms



### X.A. Screening approach

The screening estimator approach consist in removing the overlapping units either from one frame before sample selection or from one sample before data collection. It can be considered as a special case of stratified sampling and accordingly estimations are straightforward. Common drawbacks of the screening approaches

are that they can be resource-consuming, error-prone and amount to missed opportunity to collect data from a willing participant.

### **X.A.1. Frame level screening**

This screening approach consists in removing the overlap between frames (de-duplicating the frames) before the sample selection. For instance, in the case of dual frame of household's holdings and special farms, the screening will consist in removing households operating special farms from the complete lists of agricultural households in EAs before selecting the sample of agricultural households. Before such removals, it is important to ensure that the special farms screened in the EAs are included in the frame of special farms. This may increase the listing time in EAs when the identification of special farms required many additional questions in the listing form.

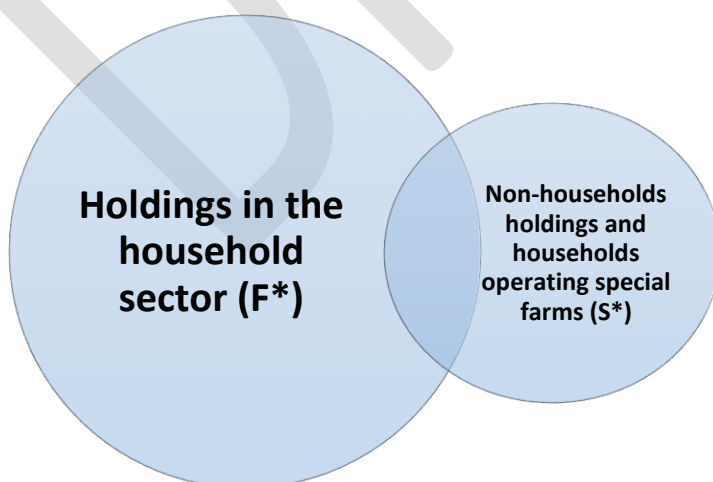
### **X.A.2. Sample level screening**

In this screening approach, the overlap between one frame and the sample selected from the other frame (usually the most expensive frame) is removed. Example: when a dual frame of area and list frame is used, agricultural holdings identified in the area sample (e.g. sample of segments) are removed from the list frame before selecting the sample from the list. The sample level screening is considered acceptable when the overlap is relatively small. A specific drawback of this approach is that it may increase non-response errors.

### **X.B. Dual frame estimators**

Consider a typical situation when a country is interested in special farms (commercial, large farms...) and is using a dual-frame consisting in a list of holdings in the household sector and a combined list of non-households holdings and households operating special farms.

Fig.4 illustration of a dual frame covering agricultural holdings and special farms



Let  $N$  be the total number of agricultural holdings in the estimation domain (or domain of inference),  $F^*$  the population of holdings in the household sector and  $S^*$  the population of non-household holdings and

households operating special farms. The population of agricultural holdings from these two frames ( $F^* \cup S^*$ ) can be divided into three mutually exclusive subpopulations:

$F$ : the population of household holdings not operating special farms, with a size  $N_f$

$FS$ : the overlapping population of household holdings operating special farms, with a size  $N_{fs}$

$S$ : the population of non-household holdings with a size  $N_s$

Therefore  $N = N_f + N_{fs} + N_s$

Let  $Y$  be a variable of interest (e.g. agricultural planted area) in the population of agricultural holdings and let  $y_k$  be its value on unit  $k$ , for  $k = 1, \dots, N$ .

The objective is to estimate from the data of two independent samples surveys (households and registered farms) the population total  $Y = \sum_{k=1}^N y_k$  that can be written as:

$$Y = Y_f + Y_{fs} + Y_s = \sum_{k=1}^{N_f} y_k + \sum_{k=1}^{N_{fs}} y_k + \sum_{k=1}^{N_s} y_k$$

The simple summation of the two totals of  $Y$  estimated from the two samples is biased because of the overlap between the two sampling frames used to select the samples. Different methods are proposed in the literature for estimations from dual-frame surveys. The most famous were proposed by Hartley (1962; 1974), Fuller and Burmeister (1972), Bankier (1986), Kalton and Anderson (1986), Skinner (1991), Skinner and Rao (1996), Singh and Wu (1996; 2003), Lohr and Rao (2006), Mecatti (2007) and Singh and Mecatti (2011).

Estimations in general and variance estimations in particular are not always straightforward in dual-frame survey. Arcos et al., (2015) note that standard software packages for complex surveys cannot be used directly when the sample is obtained from a dual frame survey because the classical design-based estimators are severely biased and there is an underestimation of standard errors. The authors developed the R package *Frames2* that can be used for dual-frame estimations using most of the methods mentioned above.

The dual frame estimation method proposed by Hartley (1962) consists basically in weighting the two estimates of  $Y$  for the overlapping domain (here  $FS$  the overlapping population) to avoid the overestimation bias. The Hartley's general class of dual frame estimators is given by:

$$\hat{Y}_H = \hat{Y}_f + \theta \hat{Y}_{fs}^f + (1 - \theta) \hat{Y}_{fs}^s + \hat{Y}_s$$

Where

- $\hat{Y}_{fs}^f$  is the estimator on the domain  $FS$  based only on the data of the sample of household farms
- $\hat{Y}_{fs}^s$  is the estimator on the domain  $FS$  based only on the data of the sample of non-household or special farms
- $\theta$  is an arbitrary constant, such that  $0 \leq \theta \leq 1$

It can be easily noted that using  $\theta = 0$  or  $\theta = 1$  is equivalent to a sample level screening approach (see section 5.3.2.). Hartley (1974) proposed an optimal value for  $\theta$  minimising the variance of  $\hat{Y}_H$ :

$$\theta_{opt} = \frac{V(\hat{Y}_{fs}^s) + Cov(\hat{Y}_s, \hat{Y}_{fs}^s) - Cov(\hat{Y}_f, \hat{Y}_{fs}^f)}{V(\hat{Y}_{fs}^f) + V(\hat{Y}_{fs}^s)}$$



This optimal value can be estimated using the R package Frames2 (Arcos et al., 2015). The main drawback is that it is variable-specific i.e. shall be estimated for each indicator.

The Hartley's estimator with  $\theta = 1/2$  (also called average estimator) is equivalent to the multiplicity estimator proposed by Mecatti (2007) for multiple-frame surveys and presents interesting features that could interest countries. Lohr (2011) mentioned that and the value of  $\theta = 1/2$  is frequently recommended with the Hartley's estimator. The average estimator may not be the most efficient one in many contexts compare to other dual-frame estimators but it can be recommended to countries because it offers operational advantages. First, it is straightforward to compute and implement because the value of  $\theta$  does not depend on the quantity of interest (Baffour et al., 2016) and its variance is quite easy to estimate. In addition, it provided good results in a number of experiences in the literature including Brick et al. (2006) dual frame survey of landline and cell phone numbers and Ferraz et al. (2017) with a dual frame of area and list frames in agricultural survey. Ferraz et al. (2017) notes that the average estimator is more efficient than the screening estimator. Chauvet and Tandeau de Marsac (2014) performed simulations comparing the performance of several dual-frame estimators in two-stage sampling designs and conclude that a simple estimator is sometimes preferable, even if it uses only part of the information collected.

Below is described how estimations are straightforward with the average (or multiplicity) estimator that can be presented as follows:

$$\hat{Y} = \hat{Y}_f + \frac{1}{2}\hat{Y}_{fs}^f + \frac{1}{2}\hat{Y}_{fs}^s + \hat{Y}_s$$

Let's consider:

$S^f$  and  $S^s$  the samples selected from the populations  $F^*$  and  $S^*$  respectively

$w_i^f$ : sampling weight of unit  $i$  selected in  $S^f$

$w_i^s$ : sampling weight of unit  $i$  selected in  $S^s$

The use of the average estimator will consist simply in calculating adjusted weight as follow:

$$w_i^{*f} = \begin{cases} w_i^f & \text{if } i \in F \\ \frac{1}{2}w_i^f & \text{if } i \in FS \end{cases}$$

$$w_i^{*s} = \begin{cases} w_i^s & \text{if } i \in S \\ \frac{1}{2}w_i^s & \text{if } i \in FS \end{cases}$$

The adjusted weights should be used for estimations on the population of agricultural holdings (e.g. area, production...) using the standard Horvitz-Thompson estimator. However, the initial weights  $w_i^f$  and  $w_i^s$  are kept in the two samples to be used for estimations specific for population  $F$  and  $S$  respectively. For instance, the estimation of the production of special farms would be performed using the sample  $S^s$  and the initial weights  $w_i^s$ .

Given that samples were selected independently from the two frames, the variance estimations are straightforward using any standard statistical software:  $V(\hat{Y}) = V(\hat{Y}_{hf} + \frac{1}{2}\hat{Y}_{hfr}^f) + V(\frac{1}{2}\hat{Y}_{hfr}^s + \hat{Y}_r)$ .

## XI. Longitudinal data collection scheme and tracking procedures

The 50x2030 survey programs recommend annual data collection on the agricultural sector. From one year to another, there are three alternatives regarding the samples for such repeated survey: (i) selecting a new sample every year (often called “repeated cross-section”), (ii) using the same sample during a number of years (panel) or (iii) changing a proportion of the sample from one year to another (partial rotation).

The panel approach generally presents lower operation costs as the same sample is surveyed every year over a period of time, especially for surveys that do not require heavy tracking operations. The panel is also well suited for estimating changes but panel sample may not be representative after a number of years because of sample attrition and structural changes in the population.

The partial rotation scheme is therefore a good alternative especially for a survey plan with a relatively long period of implementation, although it could also suffer from sample attrition. It is less expensive than the repeated cross section approach, allows longitudinal analyses and facilitates more precise estimates of changes than the repeated cross section approach.

The first option would improve annual cross-sectional estimates if the sampling frame is well updated every year before sample selection. However, compared to the other options, it will require higher operational costs of the survey program including yearly additional costs for updating the sampling frame and locating the sampling units for survey implementation. In addition, because there is no or little overlap between successive samples, repeated cross-sections usually present more discrepancies in time series, estimates of changes are less precise and longitudinal analyses are very limited or sometimes impossible.

For the 50x2030 Initiative, the panel and the partial rotation approaches are advised as cost effective sampling approaches over time. In countries where the rate of unit non-responses is usually high (as observed in previous surveys), panel should be avoided because of the risk of high non-response rates over time due to respondent burden.

Table 3: Pros and cons of sampling approaches over time

Approach	Pros	cons
<b>Repeated cross-sections</b>	<ul style="list-style-type: none"> <li>▪ Better sample representativeness (updated frame and sample)</li> <li>▪ More precise cross-sectional estimates</li> </ul>	<ul style="list-style-type: none"> <li>▪ High annual operational costs: update of the frame, new sample to be interviewed</li> <li>▪ Less precise estimates of change/not allows longitudinal studies</li> <li>▪ May need data reconciliation from one year to another</li> </ul>
<b>Panel</b>	<ul style="list-style-type: none"> <li>▪ Reduce the variation of the estimates</li> <li>▪ Precise estimates of change</li> <li>▪ Smoother time series data</li> <li>▪ Low operation cost (in case without heavy tracking)</li> </ul>	<ul style="list-style-type: none"> <li>▪ Low sample representativeness if there are important structural change in the population</li> <li>▪ Sample attrition: respondent burden, change or movement of units</li> </ul>

<b>Rotation</b>	<ul style="list-style-type: none"> <li>▪ Compared to repeated cross-sections:</li> <li>▪ Improved precision of estimates of change</li> <li>▪ Lower operation cost</li> </ul>	<ul style="list-style-type: none"> <li>▪ Affects sample representativeness depending on sample fraction</li> </ul>
-----------------	---	--

### **XI.A. Length of panel/rotation design**

For the Integrated Agricultural and Rural Survey Program, a three year panel or rotation would be cost effective because non-agricultural households are considered every three years in that survey program. For the Agricultural Survey Program, there is more flexibility in the choice of the length of the panel/rotation. However, as the efficiency of the panel/rotation sample will decrease over time, it is important to avoid lengthy periods for the panel/rotation. Three to five years would be suitable depending on countries' specific contexts.

### **XI.B. Handling sample attrition**

Attrition in a repeated survey is the dropout of sample units from one survey round to another. Both panel and rotating samples could suffer from attrition bias. Attrition affects both cross sectional and longitudinal estimates because of the reduction of sample size and potential systematic nature of the attrited units. Attrition bias is usually well addressed with weights adjustments for cross sectional estimations when the attrition rate is not very high but estimates of longitudinal changes are definitively less precise in presence of attrition. Therefore, it is best to develop strategies to avoid this issue as much as possible.

There are a number of approaches in the literature to handle sample attrition including sensitizing, giving incentives to respondents, replacement, oversampling, tracking and follow up interviews. In the framework of the 50x2030 Initiative, replacements are not advised because of selection bias issues (Witoelar 2011). Oversampling and tracking approaches, explained below, are widely used solutions widely used.

However, in presence of high attrition, it would be suitable to update the sample in the successive survey round. In fact, the adoption of panel/rotation design supposes implicitly that significant change in the population of interest is not expected in the very short term (1-3 years). Such assumption would not hold if high attrition is recorded.

#### **XI.B.1. Oversampling**

Oversampling is a way of anticipating attrition in increasing the sample size accordingly. It can be performed in multiplying the sample size by the inverse of the expected response rate. Response rates from previous survey rounds or similar surveys in the country are generally considered. That attenuates the effect of attrition on sampling error due to lower sample size

#### **XI.B.2. Tracking**

In longitudinal surveys, tracking refers to following and interviewing respondents who moved from the location where they were first interviewed. Tracking procedures aim to address not only the issue of missing

units but also the issue related to important changes in the observation units. Tracking in the framework of 50x2030 Initiative concerns agricultural holdings (i.e. agricultural households and holdings in the non-household sector).

#### *a) When to track?*

When deciding on conditions for tracking, it is important to consider the objectives of survey (Witoelar, 2011). In the framework of agricultural surveys, tracking of agricultural holdings may be necessary in case of disappearance or important changes of the unit that could affect at an important extent the agricultural outputs of the holding.

#### **Disappearance of a holding**

- an agricultural household moving from its initial dwelling where it was surveyed to another place (for migration or other purpose)
- an agricultural enterprise changing the location of its headquarters
- an agricultural holding (household/non-household) completely merging with another holding. In the household sector, for instance, two individual households may merge into a single one through marriage.

#### **Specific change of the holding**

- An agricultural holding splitting into two or more different holdings. Household members may leave the household with agricultural parcel or other important assets and constitute independent households.
- An agricultural holding splitting into many entities: some entities being independent holdings and other entities merging with other holdings. A member of an agricultural household may join another household through marriage.
- Similarly, an agricultural holding may welcome parts of another holding joining with agricultural assets such as land, livestock, etc.

#### *b) How far to track?*

To avoid biases in cross sectional estimations, it is preferable to limit the geographical area of the tracking to the geographical domain in which the holding was sampled. For instance, in the framework of a two-stage design, the tracking areas would be the primary sampling units (PSU). If the unit was sampled in an enumeration area considered as PSU, tracking will not be performed if the unit moved out of the PSU. This approach could affect the quality of estimates of change or other thematically analyses but cross-sectional estimates are priority for the survey programs.

## **XII. Integration in existing statistical systems**

Countries seeking to join the 50x2030 Initiative may already have an agricultural survey and/or household survey that are conducted regularly as part of the national statistical system. Depending on countries' situation, some recommendations are provided below.

### **XII.A. Both agricultural and household surveys exist in the country**

If the country wants to keep separate surveys for agriculture and households, then they may use the data integration model with the two sources of microdata to develop integrated agricultural and household microdata for analyses. A separate document is being developed on microdata integration.

If the country is willing to integrate the two survey operations, then two options could be explored:

- (i) the two samples may be merged (at least in rural areas) every three years for a single survey operation using the integrated questionnaire (which is to be adapted to the country context). This supposes that the household survey covers a representative sample of households in the country and that the agricultural survey covers both household and non-household sectors. This option may not be cost effective as the new dispersion of the final sample will increase operation costs compared to the integrated sampling design proposed here. The probability of selection of the units in the final sample may be complex (especially if the sampling designs of the two samples are different) to calculate and alternative estimation options like using multi-frame estimators could be explored
- (ii) a new integrated agricultural and household sample may be selected following the sampling strategy proposed here. This is the most cost-effective, preferred option.

It is important to note that integrating the two surveys operations may encounter institutional challenges if the two surveys are implemented by two different departments.

### **XII.B. Only an agricultural survey system exists**

A complementary sample of non-agricultural households may be selected for the years in which the ILS-HH is administered. The final sample of households may not be fully representative. In fact, the size of the existing sample of agricultural households may not be enough for reliable estimations of some household related indicators especially if the country is using an area frame for agricultural surveys.

Here again the best option is to select a new integrated agricultural and household sample.

### **XII.C. Only a household survey exists**

Depending on the number of agricultural households in the sample of households compared to the required number of agricultural households for reliable estimations of agricultural statistics, the existing sample may be used for the data collection in the years in which the ILS-HH is administered, or else a new integrated sample should be designed. If the existing sample can be used, the subpopulation of agricultural households in the sample can also be used for annual agricultural surveys, supplemented by a sample of non-household sector holdings.

### XIII. References

- 50x30 Initiative. (2020). A Guide to the 50x2030 Data Collection Approach: Questionnaire Design. Technical Paper Series #2. Rome.
- Abreu, D.A., Lawson, L. A., and Hickman, S. (2018). *Assessment of a Review Process for the 2017 Census of Agriculture*. Proceedings of the Joint Statistical Meetings Survey Research Methods Section.
- Arcos, A., Molina, D., Rueda, M., & Ranalli, M. G. (2015). *Frames2: A Package for Estimation in Dual Frame Surveys*. The R Journal, 7-1, 52-72.
- Baffour, B., Haynes, M., Western, M., Pennay, D., Misson, S., Martinez, A. (2016). *Weighting strategies for combining data from dual-frame telephone surveys: Emerging evidence from Australia*. Journal of Official Statistics, 32, 549-578.
- Bankier, M. D. (1986). *Estimators based on several stratified samples with applications to multiple frame surveys*. Journal of the American Statistical Association, 81(396):1074–1079.
- Bankier, M.D. (1988). "Power allocations: determining sample sizes for subnational areas." *The American Statistician* 42, 174-177.
- Barcaroli, G. Ballin, M. Odendaal, H. Pagliuca, D. Willighagen, E. Zardetto D. (2020). *SamplingStrata: optimal stratification of sampling frames for multipurpose sampling surveys*, R package. Version 1.5-1
- Bethel, J. (1989). *Sample Allocation in Multivariate Surveys*, Survey Methodology, 15,47-57
- Brick, J. M., Dipko, S., Presser, S., Tucker, C. and Yuan, Y. (2006). Nonresponse bias in a dual frame sample of cell and landline numbers. *Public Opinion Quarterly*, 70(5), 780—793.
- Chauvet G., Tandeau de Marsac, G. (2014). Estimation methods on multiple sampling frames in two-stage sampling designs. *Survey Methodology*, Vol. 40, No. 2, p. 335-346.
- Cochran, W.G. (1977). *Sampling Techniques*. 3rd Edition. John Wiley & Sons: New York, USA.
- Falorsi, P.D. Bako, D. Righi, P. Piersante, A. (2015). *Integrated Survey Framework*. FAO Publication. Rome
- FAO (2015a). *World Census of Agriculture 2020. Volume 1: Programme, concepts and definitions*. FAO Publication. Rome.
- FAO (2015b). *Handbook on Master Sampling Frames for Agricultural Statistics*. FAO Publication. Rome.
- FAO (2017). *Handbook on the Agricultural Integrated Survey (AGRIS)*. FAO Publication. Rome
- FAO (2018a). *Guidelines on the measurement of harvest and post-harvest losses. Recommendations on the design of a harvest and post-harvest loss statistics system for food grains (cereals and pulses)*. FAO Publication. Rome
- FAO (2018b). *Accelerated Technical Assistance Plan for Africa. Global Office Final report*. Technical Report Series GO-44-2018. FAO Publication. Rome. Available at <http://gsars.org/en/accelerated-technical-assistance-plan-for-africa-global-office-final-report/>
- Ferraz, C. 2018. Master Sampling Frame: The Field Experiments conducted in Brazil. GSARS Technical Report: Rome

- Fuller, W. A. Burmeister, L. F. (1972). *Estimation for samples selected from two overlapping frames*. In ASA Proceedings of the Social Statistics Sections, pages 245–249.
- Fuller, W.A. (2009). *Sampling Statistics*. Wiley
- Hartley, H. O. (1962). *Multiple frame surveys*. In Proceedings of the American Statistical Association, Social Statistics Sections, pages 203–206.
- Hartley, H. O. (1974). *Multiple frame methodology and selected applications*. Sankhya C, 36(3):99–118.
- Kalton, G. Anderson, D. W. (1986). *Sampling rare populations*. Journal of the Royal Statistical Society A, 149 (1):65–82.
- Kish, L. (1965). *Survey Sampling*. John Wiley & Sons: New York, USA.
- Kish, L. (1987). *Statistical design for research*. New York, NY: John Wiley & Sons.
- Kokan, A. R. (1963). *Optimum Allocation in Multivariate Surveys*. Journal of the Royal Statistical Society. Series A (General), 126(4), 557.
- Kozak, M. (2006). *On Sample Allocation in Multivariate Surveys*. *Communications in Statistics - Simulation and Computation*, 35(4), 901–910.
- Lavallée, P. (2007). *Indirect Sampling*. Springer: Ottawa
- Lessler, J.T. Kalsbeek, W.D. (1992). *Nonsampling errors in surveys*. John Wiley & Sons Inc.: New York, US
- Lohr, S.L., and Rao, J.N.K. (2006). *Estimation in multiple-frame surveys*. Journal of the American Statistical Association, 101, 1019-1030.
- Lohr, S. L. (2009). *Sampling: design and analysis*. Nelson Education
- Lohr, S.L. (2011). *Alternative survey sample designs: Sampling with multiple overlapping frames*. *Survey Methodology*, Vol.37 no.2, p. 197-213.
- Mecatti, F. (2007). *A single frame multiplicity estimator for multiple frame surveys*. *Survey Methodology*, 33, 151-157.
- Singh, A. and Wu, S. (1996). *Estimation for multiframe complex surveys by modified regression*. Proceedings of the Statistical Society of Canada, Survey Methods Section, n, pages 69–77.
- Singh, A. and Wu, S. (2003). *An extension of generalized regression estimator to dual frame surveys*. Proceedings of the Joint Statistical Meeting - Section on Survey Research Methods, pages 3911–3918.
- Singh, A. and Mecatti, F. (2011) *Generalized Multiplicity-Adjusted Horvitz-Thompson Type Estimation as a Unified Approach to Multiple Frame Surveys*. *Journal of Official Statistics*, 27, 633-650.
- Sitter, R. R. (1997). *Variance Estimation for the Regression Estimator in Two-Phase Sampling*, *Journal of the American Statistical Association*, 92:438, 780-787
- Skinner, C. J. (1991). *On the efficiency of raking ratio estimation for multiple frame surveys*. *Journal of the American Statistical Association*, 86(415):779–784.
- Skinner, C. J. Rao J. N. K. (1996). *Estimation in dual frame surveys with complex designs*. *Journal of the American Statistical Association*, 91(443):349–356.

U.N. (2017). *Principles and Recommendations for Population and Housing Censuses- Revision 3*. United Nations Publication. New York, US

Valliant, R., Dever, J. A., & Kreuter, F. (2015). Effects of cluster sizes on variance components in two-stage sampling. *Journal of Official Statistics*, 31(4), 763-782.

Valliant, R., Dever, J. A., and Kreuter, F. (2018). *Practical tools for designing and weighting survey samples*. 2<sup>nd</sup> edition. New York: Springer.

Witoelar, F. (2011). *Tracking in Longitudinal Household Surveys*. The World Bank. Washington, DC:

Young, L.J., Lamas, A.C., & Abreu, D. A. (2017). The 2012 Census of Agriculture: a capture–recapture analysis. *Journal of Agricultural, Biological and Environmental Statistics*, 22(4), 523-539.

Young, L. J., Hyman, M., & Rater, B. R. (2018). Exploring a big data approach to building a list frame for urban agriculture: A pilot study in the city of Baltimore. *Journal of Official Statistics*, 34(2), 323-340.

DRAFT