
Appendix

Census of Agriculture Methodology

The purpose of a census is to enumerate all objects with a defined characteristic. For the census of agriculture, that goal is to account for “any place from which \$1,000 or more of agricultural products were produced and sold, or normally would have been sold, during the census year.” To do this, NASS creates a Census Mail List (CML) of agricultural operations that potentially meet the farm definition, collects agricultural information from those operations, reviews the data, corrects or completes the requested information, and combines the data to provide information on the characteristics of farm operations and farm operators at the national, State, and county levels. In this appendix, these census processes are described.

THE CENSUS POPULATION

The Census Mail List

The National Agricultural Statistics Service (NASS) maintains a list of farmers and ranchers from which the Census Mail List (CML) is compiled. The goal is to build as complete a list as possible of agricultural places that meet the farm definition. The CML compilation begins with the list used to define sampling populations for NASS surveys conducted for the agricultural estimates program. Each record on the list includes name, address, and telephone number plus additional information that is used to efficiently administer the census of agriculture and agricultural estimates programs.

NASS builds and improves the list on an ongoing basis by obtaining outside source lists. Sources include State and federal government lists, producer association lists, seed grower lists, pesticide applicator lists, veterinarian lists, marketing association lists, and a variety of other agriculture-related lists. NASS also obtains special commodity lists to address specific list deficiencies. These outside source lists are matched to the NASS list using record linkage programs. Most names on

newly acquired lists are already on the NASS list. Records not on the NASS list are treated as potential farms until NASS can confirm their existence as a qualifying farm. Staff in NASS field offices routinely contact these potential farms to determine whether they meet the farm definition. For the 2012 Census of Agriculture, NASS made a concerted effort to work with Community-Based Organizations not only to improve list coverage for minorities but also to increase census awareness and participation.

List building activities for developing the 2012 CML started in 2009 by updating list information from respondents to the 2007 Census of Agriculture. Between 2010 and 2012, NASS conducted a series of National Agricultural Classification Surveys (NACS) on approximately 1.7 million records, which included nonrespondents from the 2007 census and newly added records from outside list sources. The NACS report forms collected information that was used to determine whether an operation met the farm definition. If the definition was met, the operation was added to the NASS list and subsequently to the CML. Addressees that were nonrespondents to a NACS were also added to the CML and identified with a special status code.

Measures were taken to improve name and address quality. Additional record linkage programs were run to detect and remove duplicate records both within each State and across States. List addresses were processed through the United States Postal Service’s National Change of Address Registry and the Locatable Address Conversion System to ensure they were correct and complete. Records on the list with missing or invalid phone numbers were matched against a nationally available telephone database to obtain as many phone numbers as possible. To reduce costs, operations with characteristics that indicated they were unlikely to be farms, according to the farm definition, were removed from the list.

The official CML for the 2012 Census of Agriculture was established on September 1, 2012. The list contained 3,009,641 records. There were 2,387,326 records that were thought to meet the NASS farm definition and 622,315 potential farm records, which included NACS nonrespondents, other records added to the CML by the NASS field offices after the record linkage process, and late adds to the CML that were not included in any previous NACS or State screening survey.

Not on the Mail List (NML)

Extensive efforts are directed toward developing a CML that includes all farms in the U.S. However, some farms are not on the list, and some agricultural operations on the list are not farms. NASS uses its June Agricultural Survey (JAS) frame to quantify the number and types of farms not on the CML. The tracts in the JAS that are not on the CML are said to be in the Not on the Mail List (NML) domain. If a tract in the NML domain is determined to be a farm during the census, it is an NML farm. The NML farms are used to estimate the undercoverage associated with the census.

The NASS area frame, which is used for the JAS, covers all land in the U.S. and includes all farms. The land in the U.S. is stratified by characteristics of the land. A probability sample of segments is drawn within each stratum for the NASS annual area frame survey, (JAS). Segments of approximately equal size are delineated within each stratum and designated on aerial photographs. The JAS sample of segments is allocated to strata to provide accurate measures of acres planted to widely grown crops, farm numbers, and inventories of cattle. Sampled segments in the JAS are personally enumerated. Each operation identified within a segment boundary is known as a tract.

The 2012 JAS sample was increased to improve the farm counts for operations that produced specialty commodities or had socially disadvantaged or minority operators. The total sample consisted of 14,376 segments of which 3,291 were additional segments added to facilitate the use of the JAS as an Agricultural Coverage Evaluation Survey (ACES). The additional segments were added based upon multivariate sample allocations to target specific

items at the U.S. level that most improved the accuracy. The initial NML mailout consisted of 36,021 records. Census records that were Undeliverable as Addressed (UAA) were added to the NML domain and replicated from June 2012. A total of 36,424 NML records were summarized. Only 5,565 records were truly NML and classified as a farm (in-scope). The 2012 JAS consisted of sample segments from all States, with the exception of Alaska where NASS does not maintain an area frame.

During the JAS prescreening operation, each tract is identified as either agricultural or non-agricultural. Each JAS agricultural tract is identified as a farm or non-farm in June based on the farm definition. Non-agricultural tracts are further classified into categories – with farm potential, with unknown farm potential, or with no farm potential. The names and addresses collected in the 2012 JAS were matched to the CML. Those from the JAS 2012 survey that did not match a census record were determined to be in the NML domain and sent a yellow census report form so that they could be differentiated from the green report form sent to those on the CML. Instructions on the census report form directed any respondent who received duplicate forms to complete the CML form and to mail all duplicate forms back together. Those who returned a CML census form and an NML census form had been misclassified as NML and were removed from the NML classification.

The farm/nonfarm status of each NML domain operation was determined based on the reported data in the census form. An operation in the NML domain that was determined to be a farm is referred to as an NML farm. Characteristics of NML farms and their operators provided a measure of the undercoverage of farms on the CML. The percentage of farms not represented on the CML varied considerably by State. In general, NML farms tended to be small in acreage, production, and sales of agricultural products. Farm operations were missing from the CML for various reasons, including the possibility that the operation started after development of the CML, the operation was so small that it did not appear in any agriculture-related source lists, or the operation was misclassified as a nonfarm prior to

census mailout. The CML was used with the NML in a capture-recapture framework to represent all farming operations across all States in the JAS sample.

DATA COLLECTION OUTREACH AND PROMOTIONAL EFFORTS

NASS planned and executed a multi-phase strategic communications campaign for the 2012 Census of Agriculture, to increase the level of awareness and response among all U.S. agricultural producers.

- Phase 1 ran from October 2011 – July 2012. It raised awareness about the census and list building, encouraged producers to sign up in response to NASS mailings and at community, association, and other stakeholder meetings where NASS partners reached out.
- Phase 2 ran from July 2012 – December 2012. It notified farm operators and agricultural organizations that the census would be mailed in December, and encouraged communications regarding the census.
- Phase 3 ran from December 2012 – July 2013. It focused on census data collection with messaging urging response, reminding operators that it's-not-too-late-to-respond, and thank-you messaging.
- Phase 4 began in February 2014. It communicates information about the data release plan, which has four phases:
 - Phase A (November 2012 – December 2013) focused on thanking farmers for their participation in the census and partners for their leadership.
 - Phase B (January 2014 – February 2014) drew attention to the upcoming release.
 - Phase C (February 2014 through spring 2014) focuses on the census findings as they are released.
 - Phase D (ongoing) continues to focus on the census findings as they are released.

As part of the plan, NASS targeted selective communications and outreach efforts on beginning and minority farm operators. All of these efforts were accomplished through an integrated communications program that focused on four primary areas: partnership building, local-level

outreach, public relations, and paid media. External support was provided by a private agricultural communications agency.

The unifying force behind the 2012 communications campaign was the theme “There’s Strength in Numbers.” This was accompanied by supporting messages and artwork that created a consistent look and feel for all census communications. All messages and materials served the purpose of inspiring action: *Grow Your Farm Future - Shape Your Farm Programs - Boost Your Rural Services - Fill out your Census of Agriculture - Do your part to be counted - There’s strength in numbers.*

Partnership and Local-Level Outreach

At the national level, NASS officials met with leaders from dozens of key agricultural organizations, State departments of agriculture, and other USDA agencies, to successfully secure their support in promoting the census among their constituencies. Stakeholders partnered with NASS to promote the 2012 Census of Agriculture through publications, special mailings, speeches, social media, websites, and other communications. In addition, through grassroots-level outreach and efforts, NASS partnered with a number of community-based organizations to reach minority and limited-resource farmers and ranchers. All national-level outreach was encouraged and mirrored at the regional, State, and local levels. Among the highlights of these partnership efforts was the production of more than 40 television and radio public service announcements (PSAs) featuring the U.S. Secretary of Agriculture, State secretaries, directors, and commissioners of agriculture and leaders from community-based organizations. The PSAs, available in both English and Spanish, encouraged farmers and ranchers to respond to the 2012 Census of Agriculture.

Public Relations

In the public relations arena, NASS and the contractor worked with internal and external stakeholders to equip them with communications tools and resources to deliver the census communications message to their audiences. NASS utilized its Intranet to deliver materials to the 12

regional and 46 field offices and created a “Partner to Promote the Census” portal on the census website to deliver public relations materials and tools to external stakeholders. The materials included, but were not limited to: customizable news releases, feature stories, newsletter articles, blogs; drop-in advertisements; website buttons and banners; PowerPoint templates; brochures; and more. In addition, at the national level NASS issued a dozen news releases citing department and agency spokespeople and published timely and relevant pieces to the USDA blog highlighting the census. These public relations efforts at the national, State, and local levels helped ensure that NASS’s message about the census was continually in the media, including print and online publications, a variety of social media, radio, and some television programs. Media outlets included both those specializing in agriculture and more general outlets.

Paid Media

For the 2012 Census of Agriculture, NASS placed special emphasis on reaching new and beginning farmers, while continuing efforts to improve its reach within previously under-represented populations. Even with increasingly limited budgets and resources, NASS was able to apply a portion of funds towards paid media. Strategically, NASS purchased limited print and online advertising in areas where there was the potential for high concentrations of under-represented populations and new and beginning farmers and ranchers.

DATA COLLECTION

Method of Enumeration

Data collection was accomplished primarily by mailout/mailback, but supplemented with Electronic Data Reporting (EDR) on the Internet, and personal enumeration for special classes of records in the census operations. Personal enumeration (interviewing) involved the use of both Computer-Assisted Telephone Interviewing (CATI) and Computer-Assisted Personal Interviewing (CAPI). Enumerators at the NASS National Operations Center in St. Louis, MO conducted CATI data collection. In addition, enumerators under contract with NASS through the National Association of State Departments of Agriculture (NASDA) conducted phone and personal interviews with

respondents. For the 2012 Census of Agriculture, NASS implemented a pre-notification strategy in an effort to increase awareness, improve overall responses, and encourage respondents to report early to avoid continued correspondence. All records in the initial mailout received either a postcard or pre-recorded voice message announcing the census mail packets were coming.

Report Forms

There were seven regionalized versions of the report forms used for the 2012 Census of Agriculture. The report form versions were designed to facilitate reporting crops most commonly grown within each report form region. Additionally, an American Indian report form was developed to facilitate reporting for operations on reservations in Arizona, New Mexico, and Utah. The regional report form numbers are: 12-A101, 12-A102, 12-A103, 12-A104, 12-A105, 12-A106 and 12-A107 (HI). The American Indian report form is 12-A200. All of the forms allowed respondents to write in specific commodities that were not listed on their form.

Report Form Mailings

Pre-notification by postcard or pre-recorded message began December 10, 2012. Approximately 3.0 million mail packets were mailed in December 2012. Each packet contained a cover letter, instruction sheet, a labeled report form, and a return envelope. The Census Bureau’s National Processing Center (NPC) in Jeffersonville, IN was contracted to perform mail packet preparation, initial mailout, and two follow-up mailings to nonrespondents.

The initial mailout was followed by a thank-you reminder postcard that was delivered in January 2013 to all operations that received mail packets. First follow-up mail packets were mailed in mid-February 2013 to approximately 1.0 million nonrespondents. Second follow-up mail packets were mailed in mid-March 2013 to approximately 750,000 nonrespondents.

Personal Follow-up

Operating concurrently with NPC’s mail data collection efforts, NASS telephone call centers targeted selected groups of census nonrespondents

for telephone enumeration. NASS field offices targeted selected groups of census records for in-person enumeration. These efforts were referred to as:

- Suspicious Out of Scope Follow-up
- Criteria Record Follow-up
- Must Case Follow-up
- American Indian and Alaska Native Farm Operator Follow-up
- Low Response County Follow-up
- Last Call Nonresponse Follow-up
- Not on Mail List (NML) Follow-up

Suspicious Out-of-Scope Follow-up. The Suspicious Out-of-Scope Follow-up was a phone follow-up that began in February 2013 and was conducted through May 2013. It included records that mailed their form back with a response that they were no longer farming. These operations had reported agricultural information in another survey during 2012. The operations were re-contacted with a CATI instrument to either verify the respondent was not farming or complete a census report form.

Criteria Record Follow-up. Nonrespondents and refusals to the National Agricultural Classification Surveys received unique coding on the CML and are referred to collectively as Criteria Records for follow-up data collection. These Criteria Records typically had a lower probability of meeting the farm definition and were less likely to respond. It was critical to identify those records in this group that represented farms to provide coverage of the small farm population. Small farms make up a significant portion of the overall U.S. farm population.

For the 2012 Census of Agriculture, 276,043 Criteria Records were included in the Census Mail List (CML). A sample of 23,739 Criteria Records was selected for targeted data collection efforts. The sampled records were first contacted by telephone using the census CATI instrument beginning in February 2013 after the initial mail returns were processed. Certified mail to 18,831 respondents was used for those who could not be contacted by telephone. Data collection resulted in 10,887 returns from both telephone and certified mail. The in-scope rate from the returns was applied to the remaining

criteria records during replication, which is described in the next sub-section.

Must Case Follow-up. Must cases were known large operations, the absence of which could have significantly affected the accuracy of census results. For the 2012 Census of Agriculture, 118,533 records were categorized as Must cases. Each active Must operation was accounted for by mail receipt, phone interview, or personal enumeration; if an operation was no longer in operation, its nonfarm status was documented. CATI calling of nonrespondent Must cases was undertaken by call centers from March 2013 through May 2013, after the initial and first follow-up mailing. Following the CATI calling, the remaining nonresponse Must cases were assigned to field offices for personal enumeration. Because of the potential importance of Must cases, they were all accounted for and therefore not eligible for nonresponse weighting adjustment.

American Indian and Alaska Native Farm Operator Follow-up. The American Indian report form (12-A200) was mailed to all operations in Arizona, New Mexico and Utah thought to have an American Indian or Alaska Native operator. It was included in the initial mailout, but due to poor mail response a personal enumeration data collection strategy was utilized with no additional mail follow-up. A concerted effort was made to get individual reports from every American Indian and Alaska Native farm operator in the country. If this was not possible within a reservation, a single reservation-level census report was obtained from knowledgeable reservation officials. These reports covered agricultural activity on the entire reservation. The NASS reviewed these data and removed any duplicate data reported by American Indian or Alaska Native farm operators from that reservation who responded on an individual census report form. Additionally NASS obtained, from knowledgeable reservation officials, the count of American Indian and Alaska Native farm operators (on the reservations) who were not counted through individual census report forms, but whose agricultural activity was included in the reservation-level report form.

Low Response County Follow-up. The Low Response County (LRC) follow-up activity was used

to increase the response rate in all counties to at least 75 percent. CATI was used for this follow-up activity. NASS utilized an adaptive design technique to identify particular records for telephone contact, in an effort to increase coverage on minority operations and operations known to produce specialty commodities. In early April 2013, NASS identified nonresponse cases in counties with a response rate of less than 75 percent. Nonresponse records in these counties were then prioritized so that minority operations and specialty commodity producers were the primary records delivered to phone enumerators. Nonrespondent telephone contact information was transmitted electronically to NASS call centers and incorporated into their CATI instrument. CATI follow-up activities began in mid-April 2013 and continued through mid-June 2012. Automated procedures were employed biweekly to ensure that the record selection procedures were targeting counties that would meet the goals of increasing minority operation coverage and to monitor the number of respondents needed to reach the 75 percent county response rate. When the required number of completions was achieved for a given county, LRC activity was suspended in that county.

Last Call Nonresponse Follow-up. The Last Call Nonresponse Follow-up activity was utilized to increase the national response rate to 80 percent. All remaining nonresponse records with an expected value of sales greater than \$50,000 in counties that had not achieved a 75-percent response rate were eligible for this phone follow-up activity. CATI was used for this activity and began in mid-July 2013 and lasted until August 1, 2013. Automated procedures were employed to monitor the number of respondents needed and completed. When a 75 percent response rate was achieved for a given county, follow-up in that county was suspended. NASS achieved its goal of an 80-percent national response rate utilizing Last Call Nonresponse Follow-up.

Not on the Mail List (NML) Follow-up. To account for farming operations not on the CML, NASS used its 2012 JAS supplemented sample from the NASS area frame. The NASS area frame covers all land in the U.S. with the exception of Alaska and Rhode Island and includes all farms. As previously

described the NASS conducted a record linkage operation between the CML records and the records from the 2012 JAS. Those 2012 JAS records that did not match records on the CML were designated as “Not on the Mail List (NML)” records. These records were mailed a colored census form. The NML records were mailed at the same time as the census mailing and received the same follow-up procedures as the census mailing through the first follow-up in mid-February 2013. Beginning in March 2013, CATI was used for nonresponse follow-up for NML nonrespondents.

Replication

Replication is utilized to improve efficiency and reduce respondent burden. To adjust for nonresponse associated with criteria records in the 2007 Census of Agriculture, NASS replicated a set of respondents determined to be in-scope from the last mailing of the Agricultural Identification Survey (AIS), conducted in December 2006. The replicated records represented operations that were relatively small in size and homogeneous in nature. Replicated records were assumed to be in-scope, based on their AIS reported data.

For the 2012 Census of Agriculture, a first mailing was sent to the criteria records, a subpopulation consisting of all of the approximately 74,000 respondents to the 2011 NACS mailing. This included pre-notification using a pre-recorded message, the first mailing, and the thank-you reminder post card. No further follow-up efforts were conducted on this subpopulation. As in 2007, the agricultural operations in this subpopulation were relatively small in size and homogeneous in nature. The responses from the criteria records were used to estimate the in-scope rate for the 20,168 nonrespondents from this subpopulation.

Records were selected randomly for replication or coding as out-of-scope based on the estimated in-scope rate. The use of the in-scope rate after one mailing is supported by analysis of 2007 census data, which indicated the early in-scope rate was a reasonable proxy for the in-scope rate for the subpopulation of criteria records that did not respond to the NACS immediately preceding the census

mailing. Of the 20,168 NACS records with no response, 16,762 records were selected to be in-scope.

Data relationships between the 2012 responses and their respective NACS data were applied to the NACS data for the nonrespondents selected to be in-scope to derive values to seed replication. Then replication was conducted through imputation.

Criteria records with no response to the December 2011 NACS were excluded in the capture-recapture adjustments for coverage, response, or correct classification. The in-scope records were each given an initial weight of one. However, for calibration, the replicated in-scope records were eligible for a coverage adjustment.

REPORT FORM PROCESSING

Data Capture

The Census Bureau's National Processing Center (NPC) in Jeffersonville, IN was contracted to process returned mail packets. NASS staff on site at the NPC provided technical guidance and monitored NPC processing activities. All report forms returned to the NPC were immediately checked in, using bar codes printed on the mailing label, and removed from follow-up report form mailings. All forms with any data were scanned and an image was made of each page of a report form. Optical Mark Recognition (OMR) was used to capture categorical responses and to identify the other answer zones in which some type of mark was present.

Data entry operators keyed data from the scanned images using OMR results that highlighted the areas of the report forms with respondent entries. The keyer evaluated the contents and captured pertinent responses. Ten percent of the captured data were keyed a second time for quality control. If differences existed between the first keyed value and the second, an adjudicator handled resolution. The decision of the adjudicator was used to grade the performance of the keyers, who were required to maintain a certain accuracy level.

The images and the captured data were transferred to NASS's centralized network and became available to field offices and headquarters on a flow basis. The

images were available for use in all stages of review. Images were computer generated for reports obtained from the telephone interviews and the Internet.

Editing Data

Captured data were processed through a computer formatting program, which verified that records were valid – that the record identification number was on the list of census records, that the reported counties of operation and production were valid, and other related criteria. Rejected records were referred to analysts for correction. Accepted records were sent to a complex computer batch edit process. Each execution of the computer edit in batch mode consisted of records from only one State and flowed as the data were received from the NPC, the NASS Electronic Data Reporting (EDR) web utility, or the Computer-Assisted Telephone Interview (CATI) applications.

The computer edit determined whether a reporting operation met the qualifying criteria to be counted as a farm (in-scope). The edit examined each in-scope record for reasonableness and completeness and determined whether to accept the recorded value for each data item or to take corrective action. Such corrective actions included removing erroneously reported values, replacing an unreasonable value with one consistent with other reported data, or providing a value for an overlooked item. To the extent possible, the computer edit determined a replacement value. Strategies for determining replacement values are discussed in the next section. Operations failing to meet the qualifying criteria were categorized as out-of-scope for the census; that is, they were classified as being a nonfarm. Out-of-scope records that NASS had reason to believe might be in-scope (indications of recent and/or significant agricultural activity reported on NASS surveys, for example) were referred to analysts for verification.

The edit systematically checked reported data section-by-section with the overall objective of achieving an internally consistent and complete report. NASS subject-matter experts had previously defined the criteria for acceptable data. Problems that could not be resolved within the edit were referred to an analyst for intervention. Prior to the

census mailout, NASS established a group of 90 analysts in a Census Editing Unit in the National Operations Center in St. Louis, MO who examined the scanned images, consulted additional sources of information, and determined an appropriate action. Field office analysts also participated using an interactive version of the edit program to submit corrected data and immediately re-edit the record to ensure a satisfactory solution.

Imputing Data

The edit determined the best value to impute for reported responses that were deemed unreasonable and for required responses that were absent. If an item could not be calculated directly from other current responses, the edit determined whether acreage, production or inventory items had been reported for that farm on a recent NASS crop or livestock survey. For operators who had not changed in five years, demographic variables such as race and gender were taken from the previous census. Administrative data from the Farm Service Agency were used for a few items, such as Conservation Reserve Program acreage. When deterministic edit logic and previously-reported data sources proved inadequate, data from a reporting farm of similar type, size, and location (a donor farm) were considered. In cases where automated imputation was unable to provide a consistent report, the record was referred to an analyst for resolution.

Separate system processes were established to efficiently provide data from a similar farm to the edit when donor imputation was required. The farm characteristics used to define similarity between a recipient record and its donor record were determined dynamically by the edit logic. Euclidean distance was used for similarity computations, with each contributing similarity characteristic scaled appropriately. The most similar farm based on this criterion (the “nearest neighbor”) was identified and returned to the edit for use as a donor. The calculated distance between the centroids of the principal counties of production of the donor and recipient was always included as one of the measures of similarity.

To provide donors to the automated edit, a pool of successfully edited records was maintained for each

section of the report form. These donor pools began with 2007 census data, reconfigured to emulate 2012 data and then edited using 2012 logic. Data from the 2010 Census Content Test were similarly remapped and edited before being added to the original donor pools. As 2012 records were successfully processed, they were added to the donor pools, which maintained the most recent data for each farm. Donor pools were updated approximately every other week, as determined by edit processing schedules. After several updates, all initial data records were dropped, leaving only 2012 records in the donor pools. After each update, donor pool records were grouped into strata containing farms in the same state of similar type and size, using a data-driven algorithm to define strata. Certain American Indian farms were treated as a separate group, effectively having their own donor pool.

In response to each donor request issued by the edit, a dedicated system process would search the appropriate stratum and respond with the most similar donor, while giving preference to more recent donors. In relatively rare instances where it was unable to provide a donor, the donor selection process issued an appropriate failure message to the edit. Imputation failures occurred for several different reasons. The requirement that an imputed value be positive could have ruled out all available donors, as could have the necessity for the donor record to satisfy a particular constraint – say, that the donor record has cattle, but no milk cows. In general, an imputation failure occurred if there was no satisfactory donor in the same profile as the report being edited. Records with imputation failures were either held until more records were available in the donor pool or referred to an analyst. In addition, when such a failure occurred in finding a donor for expenditure data, a program provided values from a table of donor pool averages in lieu of values from an individual donor, wherever possible. This ‘failover’ utility was new for the 2012 census imputation process, and significantly reduced the number of imputation failures among the expenditure and labor variables. During the early stages of editing, records requiring imputation for production (and hence yields) of field crops or hay, land values, or certain expenditure variables were set aside or “parked.” These records were edited when the donor pools contained only 2012 records,

ensuring that 2012 data were used in imputations for these variables.

After receiving a donor's data, the edit substituted the values into the edited record. In many cases, the donor record's data value was scaled using another data field specified in the edit logic. In such cases, the size of the auxiliary field's value in the edited record, relative to its value in the donor record, was used to inflate or reduce the donor record's value for the imputed field. The imputed data were then validated by the same edit logic to which reported data were subject. Since imputation was conducted independently for each occurrence, reports requiring multiple imputations may have drawn from multiple donors.

Data Analysis

The complex edit ensured the full internal consistency of the record. Successfully completing the edit did not provide insight as to whether the report was reasonable compared to other reports in the county. Analysts were provided an additional set of tools, in the form of listings and graphs, to review record-level data across farms. These examinations revealed extreme outliers, large and small, or unique data distribution patterns that were possibly a result of reporting, recording, or handling errors. Potential problems were researched and, when necessary, corrections were made and the record interactively edited again.

When NASS summarizes the census of agriculture, it assigns the data from an individual report to the "principal" county. The principal county is based on the operator's response to a census question and is the one county in which the majority of agricultural products are produced. Because some large operations have significant production in multiple counties, some reports were broken up into multiple source counties, to more accurately allocate the data. Similarly, large farms operating in more than one State were treated as distinct, state-specific operations. A separate report form was completed for each county or State and a separate record was added.

ACCOUNTING FOR UNDERCOVERAGE, NONRESPONSE, AND MISCLASSIFICATION

Although much effort was expended making the CML as complete as possible, the CML did not include all U.S. farms, resulting in list undercoverage. Some farm operators who were on the CML did not respond to the census, despite numerous attempts to contact them. In addition, although each operation was classified as a farm or a nonfarm based on the responses to the census report form, some were misclassified; that is, some nonfarms were classified as farms and some farms were classified as nonfarms. NASS's goal was to produce agricultural census totals for publication that were fully adjusted for list undercoverage, nonresponse and misclassification at the county level.

In the 2007 Census of Agriculture, adjustments for undercoverage and nonresponse were estimated independently. In 2007, as in earlier censuses, the NASS area frame was used to adjust for undercoverage. This process assumed that the area frame provided complete coverage and that all operations are correctly classified as farm/nonfarm. To determine the extent of undercoverage in 2007, the CML records were matched to the area-frame tracts designated as agricultural, non-agricultural with potential, or non-agricultural with potential unknown in June. The area-frame tracts that did not match a CML record were designated as being in the Not on the Mail List (NML) domain. In 2007, tracts that were determined to be non-agricultural without potential during the pre-screening phase of the June Agricultural Survey (JAS) were not considered in the NML domain construction. The NML domain tracts were sent a census form and, if a tract was associated with a farm, then that farm contributed to the correction for undercoverage.

To adjust for nonresponse in 2007, each responding CML record was given a probability of being a farm using a classification tree. The inverse of this probability became the nonresponse weight for that record. For undercoverage, the adjustment provided State-level values. A State-level estimate was based on the weighted sum of the responders with an adjustment for the non-responders within that State plus the State-level undercoverage adjustment.

Because State-level farm count estimates based on this two-step process sometimes had high standard errors and apparent biases, the national-level adjusted estimates were smoothed across States, producing initial State-level farm operation coverage targets.

Research following the 2007 Census of Agriculture led to the realization that some area-frame operations were misclassified as farm/nonfarm, which was in conflict with the previous assumption that the JAS farm classification was the accurate classification. Further, because nonresponse could only occur if the operation was on the CML, undercoverage and nonresponse were dependent. Thus in 2012, NASS used capture-recapture methodology to adjust for undercoverage, nonresponse, and misclassification. To implement capture-recapture methods, two independent surveys were required. The 2012 Census of Agriculture (based on the CML) and the 2012 JAS (based on the area frame) were those two surveys. Historically, NASS has been careful to maintain the independence of these two surveys.

A second assumption was that the proportion of JAS farms with a given set of characteristics captured by the census was equal to the proportion of U.S. farms with those same characteristics captured by the census.

For a farm to be identified as a farm, and thus captured by the census, it must be on the CML, respond to the census report form and, based on the census response, be classified as a farm; that is, the capture probability π_C is of interest:

$$\pi_C = \pi(\text{CML, Responded, Farm on Census} | \text{Farm})$$

Two types of classification error can occur. First, a farm can be misclassified as a nonfarm. This type of misclassification is accounted for in determining the probability of capture π_C . The second type of classification error results when a response to the census is classified as a farm operation when it does not meet the definition of a farm. That is, some farms on the CML may be misclassified from their census report response and may be nonfarms. To account for the misclassification of nonfarms as farms, the probability of a farm on the census being classified correctly must be estimated; that is,

$$\pi_{CCFC} = \pi(\text{Farm} | \text{Farm on Census})$$

where *CCFC* represents Correct Census Farm Classification. To adjust for undercoverage, nonresponse, and misclassification, each CML record classified as a farm based on its response to the census report form was given a weight of the ratio of the estimated probability of correct classification of a farm on the census and the estimated probability of capture ($\hat{\pi}_{CCFC} / \hat{\pi}_C$ where the hat symbol ($\hat{\cdot}$) denotes an estimate). To estimate the number of farms with a given set of characteristics, the weights of CML records responding as farms on the census and having that set of characteristics were summed. This estimator is referred to as the capture-recapture estimator (*CR*):

$$CR = \sum_{i \in F} \frac{\hat{\pi}_{CCFC,i}}{\hat{\pi}_{C,i}}$$

where *F* is the set of all CML records classified as farms based on their responses to the census questionnaire.

To estimate the capture and correct census farm classification probabilities, a matched dataset consisting of JAS records and census records was created. Records in the 2012 JAS sample were matched to the 2012 census using probabilistic record linkage. The CML records that matched with JAS tracts represent the Census sample. Note: The Census Sample is a subset of the CML records and includes only those records matching a JAS tract. Both agricultural and non-agricultural tracts were included in the matched dataset. (This differs from the 2007 processes, which considered only the agricultural tracts and non-agricultural tracts with potential or with potential unknown. It also included CML records that responded to the census as a farm or nonfarm and CML records that did not respond to the census.)

Resolving Farm Status

The farm status based on census responses to either the CML or NML census data collection and the JAS agreed in most cases; these records are referred to as having resolved farm status. However, in other cases, a record was identified as a farm (nonfarm) on the JAS and as a nonfarm (farm) by the census

through either the CML or the NML. Such records are said to have conflicting or unresolved farm status. An operation identified as a farm is referred to as in-scope; one identified as a nonfarm is referred to as out-of-scope. From the set of matched records, three groups with conflicting farm status were identified: 1) in-scope JAS records that were out-of-scope on the census and 2) census in-scope and JAS out-of-scope records, and 3) in-scope JAS records that did not have a census response. The records with conflicting farm status were sent to regional field offices for review. In each case, efforts were made to determine whether (1) the status had changed between June and December when the census was conducted, (2) the JAS farm status was correct, (3) the census farm status was correct, (4) the records were incorrectly matched, or (5) the farm status could not be resolved. Not all of the records with conflicting farm status could be resolved. In 2012, 11.6 percent of the records in the Census Sample had unresolved farm status. Of these, 18.9 percent were from nonresponse to the census report form.

The probability an operation is a farm was estimated for the records with unresolved farm status. Using the 2012 matched dataset, a logistic model of the probability an operation is a farm based on the records with resolved farm status was developed; that is, the operations where the farm (or nonfarm) status agreed between the JAS and the census were used to develop a missing data model, which was then used to resolve farm status. The final missing data model was used to impute the probability that each of the agricultural operations with unresolved farm status is a farm. For the resolved farms and nonfarms, the probability of the operation being a farm was 1 and 0, respectively. Five-fold cross-validation was used to develop and to compare competing models. The accuracy of the model was thereby not overstated due to fitting and evaluating the model on the same set of data. To ensure that each of the cross-validation samples covered the U.S., the five cross-validation samples of JAS segments were drawn within State-stratum combinations. Characteristics of the JAS tracts were considered as potential covariates in the model. Because limited information is available for JAS nonfarm tracts, county-level socio-demographic

variables from the most recent U.S. population census were also considered. The sample weight associated with each JAS tract was multiplied by the probability of being a farm. This adjusted weight was used in all subsequent modeling.

Capture Probabilities

Recall that, for a farm to be identified as a farm, and thus captured, by the census, it must be on the CML, respond to the census report form and, based on the census response, be classified as a farm. These adjustments are dependent so that the probability of capture π_C may be written as

$$\begin{aligned}\pi_C &= \pi(\text{CML, Responded, Farm on Census}|\text{Farm}) \\ &= \pi(\text{CML}|\text{Farm})\pi(\text{Responded}|\text{CML, Farm})\pi(\text{Farm on Census}|\text{CML, Responded, Farm})\end{aligned}$$

The probability of capturing a farm depends on the characteristics of the farm. Using five-fold cross-validation, three logistic models were developed based on the matched dataset. The first model estimated the probability of a farm being on the CML. The second model estimated the probability that a farm on the CML responded to the census report form. The final model estimated the probability that a farm that was on the CML and responded to the census was identified as a farm based on its response. The probability that a farm is captured by the census of agriculture is then the product of the three conditional probabilities that a farm is on the CML, responds, and is identified as a farm.

Note 1: Responses were required for Must cases. These operations were only included in modeling the probability of a farm being on the CML. Consequently, the weight associated with a Must record was the reciprocal of the probability of a farm being on the CML.

Note 2: Two sets of models were created. One set estimated the probability of capture for Texas farms. The other set provided estimated capture probabilities for farms in the remaining States, except for Alaska.

Note 3: Because Alaska is not included in the JAS and thus has no area frame, the Alaskan agricultural

operations were not included in the capture-recapture process. No adjustments were made for undercoverage or misclassification. To account for nonresponse, the CML records were divided into three groups: (1) the Must records, (2) the Criteria Records, and (3) the remaining CML records. The must records received a weight of one, thereby receiving no adjustment for nonresponse. The probability of response for each of the other two groups was the proportion of responders within the group. Each record within the group was then given a weight equal to the reciprocal of the probability of response.

Misclassification

An operation is misclassified if (1) it meets the definition of a farm, but is classified as a nonfarm on the census or (2) it does not meet the definition of a farm, but is classified as a farm on the census. The first type of misclassification is accounted for when modeling the probability of capture. An adjustment is still needed for the misclassification of nonfarms as farms. As with farm status and capture, the probability of this misclassification depends on an operation's characteristics. Thus, a final logistic model was developed. Given that an operation was classified as a farm on the CML, the probability of its being a farm was modeled based on its characteristics. Five-fold cross-validation was used to ensure that the model was not over-fitted.

CALIBRATION

Each operation identified as being in-scope on the CML was given a weight equal to the probability of misclassification divided by the probability of capture. This weight accounted for undercoverage, nonresponse, and both types of misclassification.

The record weighting processes were initially applied at the State level to produce adjusted estimates of farm numbers and land in farms for 63 different categories of 8 characteristics of the farm operation or the farm operator -- value of agricultural sales (8); age (2); female; race (4); Hispanic origin of principal farm operator ; 4 sales categories for each of 10 major commodities (40); and farm type groups (7). The State-level number of farms and land in farms were two additional adjusted estimates,

resulting in 65 categories. To reduce the intercensal variation at the State level, the State targets were smoothed by averaging the 2012 estimates from capture-recapture and the published 2007 state estimates with the restrictions that the smoothed targets were within one standard error of the capture-recapture estimates. The smoothed State targets were rescaled so that they summed to the national capture-recapture estimates.

These State estimates were general purpose in that they did not provide any control over expected levels of commodity production of the individual farm operation. As a result of this limitation, the procedures could have over-adjusted or under-adjusted for commodity production. To address this, a second set of variables, known as commodity targets, was added to the calibration algorithm. These targets were commodity totals from administrative sources or from NASS surveys of nonfarm populations (e.g. USDA Farm Service Agency program data, Agricultural Marketing Service market orders, livestock slaughter data, cotton ginning data). The introduction of these commodity coverage targets strengthened the overall adjustment procedure by ensuring that major commodity totals remained within reasonable bounds of established benchmarks. Commodity coverage targets with acceptable ranges were established by subject-matter experts for each State, with New England treated as a State.

Each State was calibrated separately. The calibration algorithm addressed commodity coverage. The algorithm was controlled by the 65 State farm operation coverage targets and the State commodity coverage targets. To ensure that the calibration process converged with so many constraints, it was desirable to provide some tolerance ranges for each target. Although full calibration to a single point estimate would assure that the weighted total among census respondents equaled its target for each calibration variable in either set, it was not always possible to calibrate to such a large number of target values while ensuring that farm weights were within a reasonable range and not less than one. Because of this and because calibration targets are estimates themselves subject to uncertainty, NASS allowed some tolerance in the determination of the adjusted weights. Rather than forcing the total for each

calibration variable computed using the adjusted weights to equal a specific amount, NASS allowed the estimated total to fall within a tolerance range. This tolerance strategy made it possible for the calibration algorithm to produce a set of satisfactory, adjusted weights.

Ranges for the farm operation coverage targets were determined differently from the commodity targets. The State target for number of farms had no tolerance range. The tolerance range for the 64 other State farm operation coverage targets was the estimated smoothed State total for the variable plus or minus one-half of one estimated standard error of the capture-recapture estimate. This choice limited the cumulative deviation from the estimated total for a variable when State totals were summed to a U.S. level total. The commodity target tolerance ranges were determined by subject-matter experts, based on the amount of confidence in the source, and usually were less than plus or minus two percent of the target. Ranges were not necessarily symmetric around the target value.

Census data collection was assumed to be complete for very large and unique farms with their weight being controlled to 1 during the calibration adjustment process. For all other farms, adjustment weights were obtained using truncated linear calibration which forced the final census record weights to fall in the interval [1,6]. Adjustments began with the nonresponse and misclassification adjusted weights. Through calibration, a second stage weight that simultaneously satisfied all farm operation coverage and commodity coverage calibration targets was obtained. Calibration was seldom able to adjust weights so that all State targets were met. Within the calibration process, the highest priority for meeting a target was given to the number of farms, total land in farms, and top cash-receipt commodities accounting for 80 percent of the State's production. All remaining targets associated with commodities and characteristics of farms and farm operators had equal priority. If a value within the tolerance range of any variable could not be achieved in a given State, the variable was removed as a target in that State and the calibration algorithm was rerun.

Weight computations in the final algorithms were performed to several decimals. Thus, the fully-adjusted weights were non-integer numbers. To ensure that all subdomains for which NASS publishes summed to their grand total, fully-adjusted weights were integerized. This eliminated the need for rounding individual cell values and ensured that marginal totals always added correctly to the grand total. As an example of how the integerization process worked, assume there were five census records in a county with final noninteger coverage weights of 2.2, for a total of 11. The integerization process randomly selected four of these records and rounded their final weight down to 2.0 and rounded the fifth record up to 3.0, for a total of 11.

DISCLOSURE REVIEW

After tabulation and review of the aggregates, a comprehensive disclosure review was conducted. NASS is obligated to withhold, under Title 7, U.S. Code, any total that would reveal an individual's information or allow it to be closely estimated by the public. Cell suppression was used to protect the cells that were determined to be sensitive to a disclosure of information. Farm counts are not considered sensitive and are not subject to disclosure controls.

Based on agency standards, data cells were determined to be sensitive to a disclosure of information if they violated either of two criteria rules. The threshold rule was violated if the data cell contained less than three operations. For example, if only one farmer produced turkeys in a county, NASS could not publish the county total for turkey inventory without disclosing that individual's information. The dominance rule was violated if the distribution of the data within the cell allowed a data user to estimate any respondent's data too closely. For example, if there are many farmers producing turkeys in a county and some of them were large enough to dominate the cell total, NASS could not publish the county total for turkey inventory without risking disclosing an individual respondent's data. In both of these situations, the data were suppressed and a "(D)" was placed in the cell in the census publication table. These data cells were referred to as primary suppressions.

Since most items were summed to marginal totals, primary suppressions within these summation relationships were protected by ensuring that there were additional suppressions within the linear relationship that provided adequate protection for the primary. A detailed computer routine selected additional data cells for suppression to ensure all primary suppressions were properly protected in all linear relationships in all tables. These data cells were referred to as complementary suppressions. These cells were not themselves sensitive to a disclosure of information but were suppressed to protect other primary suppressions. A “(D)” was also placed in the cell of the census publication table to indicate a complementary suppression. A data user could not determine whether a cell with a (D) represented a primary or a complementary suppression.

Field office analysts reviewed all complementary suppressions to ensure no cells had been withheld that were vital to the data users. In instances where complimentary suppressions were deemed critically important to a State or county, analysts requested an override and a different complementary cell was chosen.

CENSUS QUALITY

The purpose of the census of agriculture is to account for “any place from which \$1,000 or more of agricultural products were produced and sold, or normally would have been sold, during the census year.” To accomplish this, NASS develops a CML that contains identifying information for operations that have an indication of meeting the census definition, develops procedures to collect agricultural information from those records, establishes criteria for analyst review of the data, creates computer routines to correct or complete the requested information, and provides census estimates of the characteristics of farms and farm operators with associated measures of uncertainty.

It is not likely that either the CML includes all operations that meet the definition of a farm or that all those that do meet the definition of a farm respond to the census inquiry. The goal is to publish data with a high level of quality. There are many ways to measure the quality of a census.

One of the first indicators used is a measure of the response to the census data collection as it has generally been thought that a high response rate indicates more complete coverage of the population of interest. This is a valid assumption if the enumeration list, the CML here, has complete coverage of the population of interest. In the case of the census of agriculture, the definition requiring advance knowledge of sales makes achieving a high level of coverage difficult. To ensure that the census of agriculture is as complete as possible, records are included that might not meet the census definition of a farm – in fact, almost 50 percent more records than the anticipated number of qualifying farm operations were included in the 2012 CML. A second indicator of quality then is the coverage of the farm population by the CML. Other indicators of quality relate to the accuracy and completeness of the data, and the validity of the procedures used in processing the data.

In some cases, NASS was able to produce measures of quality – such as the response rate to the data collection, the coverage of the census mail list, and the variability of the final adjusted estimates. In other cases, measures were not produced but descriptions of procedures that NASS used to reduce errors from the procedures were subsequently provided.

Census Response Rate

The response rate is one indicator of the quality of a data collection. It is generally assumed that if a response rate is close to a full participation level of 100 percent, the potential for nonresponse bias is small, although this has been questioned recently in the literature. Because the CML contains both farm and nonfarm records, the response rate is an indicator of replying to the census data collection effort, but does not reflect whether those responding met the farm definition. The response rate for the 2012 Census of Agriculture CML is 80.1 percent as compared with a response rate of 85.2 percent for the 2007 Census of Agriculture and 88.0 for the 2002 Census of Agriculture.

The 2012 Census of Agriculture response rate used the fourth response rate formula from the American

Association of Public Opinion Research Response Rate Standard Definitions manual:

$$RR4 = \frac{C_{adj}}{C_{adj} + R + NC + O + Replicated + e(U)} (100)$$

where

- C_{adj} = number of fully and partially completed records, excluding replicated records
- R = number of explicit refusals
- NC = number of non-contacted operations
- O = number of other types of nonrespondents
- $Replicated$ = number of replicated records
- U = number of operations of unknown eligibility
- $e(U)$ = estimated number of operations of unknown eligibility assumed to be eligible

Records were classified into the above variables based on the combination of their active status (AS) codes, in-scope status, and replication status. Active status refers to the eligibility status of records for selection on the CML. All replicated records were considered to be a form of nonresponse and were classified into other nonrespondents; in-scope status was considered immaterial.

Certain active status classifications indicated records of unknown agricultural status. These classifications included records to be removed from the CML but had data from outside sources indicating agricultural activity, new records from outside data sources, nonrespondents and refusals to the NACS, records for regional office handling only, and records with Farm Service Agency or Conservation Reserve Program data on operations that are not owned by the principal operator. These records were stratified (grouped) based on their probabilities of being in-scope had they responded. The estimated number of in-scope nonrespondents was calculated for the h th stratum (group) by the following formula:

$$e(U_h) = \left(\frac{C_{in-scope,h}}{C_h} \right) U_h$$

where

- $e(U_h)$ = estimated number of operations of unknown eligibility assumed to be eligible in the h th group

$C_{in-scope,h}$ = the number of completed and in-scope census records in the h th group

C_h = the number of completed census records in the h th group

U_h = number of operations of unknown eligibility in the h th group

Census Coverage

As a side-product of the statistical adjustment used to account for undercoverage, nonresponse of farms on the CML, and misclassification of responses to the census, the proportion of the adjustments due to each of those factors can be derived. The final census numbers reflected an adjustment of 12 percent for undercoverage, 16 percent for nonresponse to the census data collection, and 6 percent due to incorrect classification. Tables for State estimates of these adjustments will be provided in the final census publication.

MEASURED ERRORS IN THE CENSUS PROCESS

Although the census of agriculture does not inherently rely on a sample, it uses statistical procedures in compiling the CML, in its data collection procedures, in data editing and processing, and in compiling the final data. Additionally, it uses statistical procedures to both measure errors in the various processes and in making adjustments for those errors in the final data. One example is the statistical process used to account for undercoverage, nonresponse of farms on the CML, and misclassification of responses to the census. The basis of the undercoverage adjustment is the capture-recapture procedure that uses the area sample enumeration from the June Agricultural Survey. The largest contribution to error in the census estimates is due to the adjustments for nonresponse, undercoverage, misclassification, calibration and integerization.

Variability in Census Estimates due to Statistical Adjustment

In conducting the 2012 Census of Agriculture, efforts were initiated to measure error associated with the adjustments for farm operations that were not on the CML, for farm operations that were on the

CML but did not respond to the census report form , for farms and nonfarms that were misclassified as nonfarms and farms, respectively, for calibration, and for integerization. These error measurements were developed from the standard error of the estimates at the national, State, and county levels and were expressed as coefficients of variation (CVs) at the national and State levels and as generalized coefficients of variation (GCVs) at the county levels.

The standard error of an estimate is an estimate of the standard deviation of the sampling distribution of the estimator. Because Texas and Alaska were modeled separately from the other States, the variances of a national-level data item for these two States were computed separately and added to the variance of that data item for the rest of the U.S. The standard error was then the square root of the total variance. In each case, standard errors were computed using the group jackknife approach. To conduct the jackknifing, k mutually exclusive and exhaustive groups of JAS segments were formed. The groups were selected using a stratified random design so that each group reflected the survey design, including State and agricultural strata within a State. In turn, each group, $j = 1, 2, \dots, k$, was deleted and the capture-recapture estimate $CR_i^{(j)}$ was computed for each data item i at the specified geographical level, such as nation, State, or county, using the remaining $(k - 1)$ groups. Estimate of the variance and standard error associated with the capture-recapture estimate CR_i are then, respectively,

$$\sigma_i^2 = \frac{k-1}{k} \sum_{j=1}^k (CR_i^{(j)} - CR_i)^2; \quad SE(CR_i) = \sqrt{\sigma_i^2}$$

Increasing k improves the estimate of the variance but, as k increases, the observations become too sparse to reflect the survey design and to provide country-wide coverage. Based on 2007 data, $k = 10$ was determined to be the largest number of groups that could be formed and still have each group provide adequate coverage within all States and agricultural strata. Thus, 10 jackknife groups were used to provide standard errors for 2012 State and national estimates. To capture the additional variability from calibration and integerization, the standard errors were computed using the calibrated, integerized capture-recapture estimates from the jackknife groups. For the estimate of the number of

farms with a given set of characteristics, only the CML records with those characteristics were used to obtain the overall estimate as well as the estimates from each jackknife group.

When the constraints of the calibration process produced an artificially small standard error, the more conservative capture-recapture standard error was used. Note that the jackknife groups must only be constructed once, and different subsets of the records were used to compute estimates and standard errors for the data items.

The CV is a measure of the relative amount of error associated with the sample estimate:

$$CV = \frac{SE(CR_i)}{CR_i} 100\%$$

where $SE(CR_i)$ is the standard error of the capture-recapture estimate for data item i . This relative measure allows the reliability of a range of estimates to be compared. For example, the standard error is often larger for large population estimates than for small population estimates, but the large population estimates may have a smaller CV, indicating a more reliable estimate. For county-level estimates, a generalized coefficient of variation (GCVs) was determined for each estimate within a State. A generalized variance function relates a function of the variance of an estimator to a function of the estimator. Within a State, the standard error of an estimate for a data item was often found to be linearly related to the estimate of that item with an intercept of zero. Based on this modeled relationship, the GCV is the slope of the line relating the standard error to the estimate, multiplied times 100 to represent the GCV as a percentage.

The standard error is the product of the CV (or GCV for county estimates) and the estimate divided by 100. As an example, if the GCV for a State is 25 percent and a county's estimate is 4, then the standard error is $25(4)/100 = 1$. The standard error of an estimated data item from the census provides a measure of the error variation in the value of that estimated data item based on the possible outcomes of the census collection, including variants as to who was on the CML, who returned a census form, who was misclassified either as a farm or as a nonfarm,

and the uncertainty associated with calibration and integerization. With 95 percent confidence, an estimate is within two standard errors of the true value being estimated. For this example, with 95 percent confidence, the estimate of 4 is within $2(1) = 2$ of the true county value.

NONMEASURED ERRORS IN THE CENSUS PROCESS

As noted in the previous section, sampling errors can be introduced from the coverage, nonresponse and misclassification adjustment procedures. This error is measureable. However, nonsampling errors are imbedded in the census process that cannot be directly measured as part of the design of the census but must be contained to ensure an accurate count. Extensive efforts were made to compile a complete and accurate mail list for the census, to elicit response to the census, to design an understandable report form with clear instructions, to minimize processing errors through the use of quality control measures, to reduce matching error associated with the capture-recapture estimation process, and to minimize error associated with identification of a respondent as a farm operation (referred to as classification error). The weight adjustment and tabulation processes recognize the presence of nonsampling errors; however, it is assumed that these errors are small and that, in total, the net effect is zero. In other words, the positive errors cancel the negative errors.

Respondent and Enumerator Error

Incorrect or incomplete responses to the census report form or to the questions posed by an enumerator can introduce error into the census data. Steps were taken in the design and execution of the census of agriculture to reduce errors from respondent reporting. Poor instructions and ambiguous definitions lead to misreporting. Respondents may not remember accurately, may give rounded numbers, or may record an item in the wrong cell. To reduce reporting and recording errors, the report form was tested prior to the census using industry accepted cognitive testing procedures. Detailed instructions for completing the report form were provided to each respondent. Questions were phrased as clearly as possible based on previous tests

of the report form. Computer-assisted telephone interviewing software included immediate integrity checks of recorded responses so suspect data could be verified or corrected. In addition, each respondent's answers were checked for completeness and consistency by the complex edit and imputation system.

Processing Error

Processing of each census report form was another potential source of nonsampling error. All mail returns that included multiple reports, respondent remarks, or that were marked out of business and report forms with no reported data were sent to an analyst for verification and appropriate action. Integrity checks were performed by the imaging system and data transfer functions. Standard quality control procedures were in place that required that randomly selected batches of data keyed from image be re-entered by a different operator to verify the work and evaluate key entry operators. All systems and programs were thoroughly tested before going on-line and were monitored throughout the processing period.

Developing accurate processing methods is complicated by the complex structure of agriculture. Among the complexities are the many places to be included, the variety of arrangements under which farms are operated, the continuing changes in the relationship of operators to the farm operated, the expiration of leases and the initiation or renewal of leases, the problem of obtaining a complete list of agriculture operations, the difficulty of contacting and identifying some types of contractor/contractee relationships, the operator's absence from the farm during the data collection period, and the operator's opinion that part or all of the operation does not qualify and should not be included in the census. During data collection and processing of the census, all operations underwent a number of quality control checks to ensure results were as accurate as possible.

Item Nonresponse

All item nonresponse actions provide another opportunity to introduce measurement errors. Regardless of whether it was previously reported data, administrative data, the nearest neighbor

algorithm, or manually imputed by an analyst, some risk exists that the imputed value does not equal the actual value. Previously reported and administrative data were used only when they related to the census reference period. A new nearest neighbor was randomly selected for each incident to eliminate the chance of a consistent bias.

Record Matching Error

The process of building and expanding the CML involves finding new list sources and checking for names not on the list. An automated processing system compared each new name to the existing CML names and “linked” like records for the purpose of preventing duplication. New names with strong links to a CML name were discarded and those with no links were added as potential farms. Names with weak links, possible matches, were reviewed by staff to determine whether the new name should be added. Despite this thorough review, some new names may have been erroneously added or deleted. Additions could contribute to duplication (overcoverage) whereas deletions could contribute to undercoverage. As a result, some names received more than one report form, and some farm operators did not receive a report form. Respondents were instructed to complete one form and return all forms so the duplication could be removed.

Another chance for error came when comparing June Agricultural Survey tract operator names to the CML. Area operators whose names were not found on the CML were part of the measure of list incompleteness, or NML. Mistakes in determining overlap status resulted in overcounts (including a tract whose operator was on the CML) or undercounts (excluding a tract whose operator was not on the CML). All tracts determined to not be on

the list were triple checked to eliminate, or at least minimize, any error. NML tract operators were mailed a report form printed in a different color. In order to attempt to identify duplication, all respondents who received multiple report forms were instructed to complete the CML version and return all forms so duplication could be removed. Records in the 2012 JAS were matched to the 2012 census using probabilistic record linkage. The records of operations with unresolved farm status were reviewed by the field offices. If farm status could not be resolved, the probability of an operation being a farm was imputed using a missing data model. The uncertainty associated with this estimate, with the exception of model uncertainty, was accounted for, but errors not found through this process were not.

Model Uncertainty Error

Five logistic models were developed in the process of adjusting the farm numbers for undercoverage, nonresponse, and misclassification. One model estimated the probability of an agricultural operation with unresolved farm status being a farm. The remaining four models estimated the probability of coverage, response, and correct classification of farms and of nonfarms. Each model was fit independently by two people. For some models, both statisticians obtained the same model. Although the covariates in the two selected models differed from some for the other logistic models, the estimated probabilities were similar, but not identical. The reported standard errors account for the variability in the parameter estimates of the selected models, but not for the additional variation due to model uncertainty. They also do not account for any bias associated with a model.