



January 19, 2018

MEMORANDUM FOR: Wilbur L. Ross, Jr.
Secretary of Commerce

Through: Karen Dunn Kelley
Performing the Non-Exclusive Functions and Duties of the Deputy
Secretary

Ron S. Jarmin
Performing the Non-Exclusive Functions and Duties of the Director

Enrique Lamas
Performing the Non-Exclusive Functions and Duties of the Deputy
Director

From: John M. Abowd
Chief Scientist and Associate Director for Research and Methodology

Subject: Technical Review of the Department of Justice Request to Add
Citizenship Question to the 2020 Census

The Department of Justice has requested block-level citizen voting-age population estimates by OMB-approved race and ethnicity categories from the 2020 Census of Population and Housing. These estimates are currently provided in two related data products: the PL94-171 redistricting data, produced by April 1st of the year following a decennial census under the authority of 13 U.S.C. Section 141, and the Citizen Voting Age Population by Race and Ethnicity (CVAP) tables produced every February from the most recent five-year American Community Survey data. The PL94-171 data are released at the census block level. The CVAP data are released at the census block group level.

We consider three alternatives in response to the request: (A) no change in data collection, (B) adding a citizenship question to the 2020 Census, and (C) obtaining citizenship status from administrative records for the whole 2020 Census population.

We recommend either Alternative A or C. Alternative C best meets DoJ's stated uses, is comparatively far less costly than Alternative B, does not increase response burden, and does not harm the quality of the census count. Alternative A is not very costly and also does not harm the quality of the census count. Alternative B better addresses DoJ's stated uses than Alternative A. However, Alternative B is very costly, harms the quality of the census count, and would use substantially less accurate citizenship status data than are available from administrative sources.

<i>Summary of Alternatives</i>			
	<i>Alternative A</i>	<i>Alternative B</i>	<i>Alternative C</i>
Description	No change in data collection	Add citizenship question to the 2020 Census (i.e., the DoJ request), all 2020 Census microdata remain within the Census Bureau	Leave 2020 Census questionnaire as designed and add citizenship from administrative records, all 2020 Census microdata and any linked citizenship data remain within the Census Bureau
Impact on 2020 Census	None	Major potential quality and cost disruptions	None
Quality of Citizen Voting-Age Population Data	Status quo	Block-level data improved, but with serious quality issues remaining	Best option for block-level citizenship data, quality much improved
Other Advantages	Lowest cost alternative	Direct measure of self-reported citizenship for the whole population	Administrative citizenship records more accurate than self-reports, incremental cost is very likely to be less than \$2M, USCIS data would permit record linkage for many more legal resident noncitizens
Shortcomings	Citizen voting-age population data remain the same or are improved by using small-area modeling methods	Citizenship status is misreported at a very high rate for noncitizens, citizenship status is missing at a high rate for citizens and noncitizens due to reduced self-response and increased item nonresponse, nonresponse followup costs increase by at least \$27.5M, erroneous enumerations increase, whole-person census imputations increase	Citizenship variable integrated into 2020 Census microdata outside the production system, Memorandum of Understanding with United States Citizen and Immigration Services required to acquire most up-to-date naturalization data

Approved: _____ Date: _____

John M. Abowd, Chief Scientist
and Associate Director for Research and Methodology

Detailed Analysis of Alternatives

The statistics in this memorandum have been released by the Census Bureau Disclosure Review Board with approval number CBDRB-2018-CDAR-014.

Alternative A: Make no changes

Under this alternative, we would not change the current 2020 Census questionnaire nor the planned publications from the 2020 Census and the American Community Survey (ACS). Under this alternative, the PL94-171 redistricting data and the citizen voting-age population (CVAP) data would be released on the current schedule and with the current specifications. The redistricting and CVAP data are used by the Department of Justice to enforce the Voting Rights Act. They are also used by state redistricting offices to draw congressional and legislative districts that conform to constitutional equal-population and Voting Rights Act nondiscrimination requirements. Because the block-group-level CVAP tables have associated margins of error, their use in combination with the much more precise block-level census counts in the redistricting data requires sophisticated modeling. For these purposes, most analysts and the DoJ use statistical modeling methods to produce the block-level eligible voter data that become one of the inputs to their processes.

If the DoJ requests the assistance of Census Bureau statistical experts in developing model-based statistical methods to better facilitate the DoJ's uses of these data in performing its Voting Rights Act duties, a small team of Census Bureau experts similar in size and capabilities to the teams used to provide the Voting Rights Act Section 203 language determinations would be deployed.

We estimate that this alternative would have no impact on the quality of the 2020 Census because there would be no change to any of the parameters underling the Secretary's revised life-cycle cost estimates. The estimated cost is about \$350,000 because that is approximately the cost of resources that would be used to do the modeling for the DoJ.

Alternative B: Add the question on citizenship to the 2020 Census questionnaire

Under this alternative, we would add the ACS question on citizenship to the 2020 Census questionnaire and ISR instrument. We would then produce the block-level citizen voting-age population by race and ethnicity tables during the 2020 Census publication phase.

Since the question is already asked on the American Community Survey, we would accept the cognitive research and questionnaire testing from the ACS instead of independently retesting the citizenship question. This means that the cost of preparing the new question would be minimal. We did not prepare an estimate of the impact of adding the citizenship question on the cost of reprogramming the Internet Self-Response (ISR) instrument, revising the Census Questionnaire Assistance (CQA), or redesigning the printed questionnaire because those components will not be finalized until after the March 2018 submission of the final questions. Adding the citizenship question is similar in scope and cost to recasting the race and ethnicity questions again, should that become necessary, and would be done at the same time. After the 2020 Census ISR, CQA and printed questionnaire are in final form, adding the citizenship question would be much more expensive and would depend on exactly when the implementation decision was made during the production cycle.

For these reasons, we analyzed Alternative B in terms of its adverse impact on the rate of voluntary cooperation via self-response, the resulting increase in nonresponse followup (NRFU), and the consequent effects on the quality of the self-reported citizenship data. Three distinct analyses support the conclusion of an adverse impact on self-response and, as a result, on the accuracy and quality of the 2020 Census. We assess the costs of increased NRFU in light of the results of these analyses.

B.1. Quality of citizenship responses

We considered the quality of the citizenship responses on the ACS. In this analysis we estimated item nonresponse rates for the citizenship question on the ACS from 2013 through 2016. When item nonresponse occurs, the ACS edit and imputation modules are used to allocate an answer to replace the missing data item. This results in lower quality data because of the statistical errors in these allocation models. The analysis of the self-responses responses is done using ACS data from 2013-2016 because of operational changes in 2013, including the introduction of the ISR option and changes in the followup operations for mail-in questionnaires.

In the period from 2013 to 2016, item nonresponse rates for the citizenship question on the mail-in questionnaires for non-Hispanic whites (NHW) ranged from 6.0% to 6.3%, non-Hispanic blacks (NHB) ranged from 12.0% to 12.6%, and Hispanics ranged from 11.6 to 12.3%. In that same period, the ISR item nonresponse rates for citizenship were greater than those for mail-in questionnaires. In 2013, the item nonresponse rates for the citizenship variable on the ISR instrument were NHW: 6.2%, NHB: 12.3% and Hispanic: 13.0%. By 2016 the rates increased for NHB and especially Hispanics. They were NHW: 6.2%, NHB: 13.1%, and Hispanic: 15.5% (a 2.5 percentage point increase). Whether the response is by mail-in questionnaire or ISR instrument, item nonresponse rates for the citizenship question are much greater than the comparable rates for other demographic variables like sex, birthdate/age, and race/ethnicity (data not shown).

B.2. Self-response rate analyses

We directly compared the self-response rate in the 2000 Census for the short and long forms, separately for citizen and noncitizen households. In all cases, citizenship status of the individuals in the household was determined from administrative record sources, not from the response on the long form. A noncitizen household contains at least one noncitizen. Both citizen and noncitizen households have lower self-response rates on the long form compared to the short form; however, the decline in self-response for noncitizen households was 3.3 percentage points greater than the decline for citizen households. This analysis compared short and long form respondents, categories which were randomly assigned in the design of the 2000 Census.

We compared the self-response rates for the same household address on the 2010 Census and the 2010 American Community Survey, separately for citizen and noncitizen households. Again, all citizenship data were taken from administrative records, not the ACS, and noncitizen households contain at least one noncitizen resident. In this case, the randomization is over the selection of household addresses to receive the 2010 ACS. Because the ACS is an ongoing survey sampling fresh households each month, many of the residents of sampled households completed the 2010 ACS with the same reference address as they used for the 2010 Census. Once again, the self-response rates were lower in the ACS than in the 2010 Census for both citizen and noncitizen households. In this 2010 comparison, moreover, the decline in self-response was 5.1 percentage points greater for noncitizen households than for citizen households.

In both the 2000 and 2010 analyses, only the long-form or ACS questionnaire contained a citizenship question. Both the long form and the ACS questionnaires are more burdensome than the shortform. Survey methodologists consider burden to include both the direct time costs of responding and the indirect costs arising from nonresponse due to perceived sensitivity of the topic. There are, consequently, many explanations for the lower self-response rates among all household types on these longer questionnaires. However, the only difference between citizen and noncitizen households in our studies was the presence of at least one noncitizen in noncitizen households. It is therefore a reasonable inference that a question on citizenship would lead to some decline in overall self-response because it would make the 2020 Census modestly more burdensome in the direct sense, and potentially much more burdensome in the indirect sense that it would lead to a larger decline in self-response for noncitizen households.

B.3. Breakoff rate analysis

We examined the response breakoff paradata for the 2016 ACS. We looked at all breakoff screens on the ISR instrument, and specifically at the breakoffs that occurred on the screens with the citizenship and related questions like place of birth and year of entry to the U.S. Breakoff paradata isolate the point in answering the questionnaire where a respondent discontinues entering data—breaks off—rather than finishing. A breakoff is different from failure to self-respond. The respondent started the survey and was prepared to provide the data on the Internet Self-Response instrument, but changed his or her mind during the interview.

Hispanics and non-Hispanic non-whites (NHNW) have greater breakoff rates than non-Hispanic whites (NHW). In the 2016 ACS data, breakoffs were NHW: 9.5% of cases while NHNW: 14.1% and Hispanics: 17.6%. The paradata show the question on which the breakoff occurred. Only 0.04% of NHW broke off on the citizenship question, whereas NHNW broke off 0.27% and Hispanics broke off 0.36%. There are three related questions on immigrant status on the ACS: citizenship, place of birth, and year of entry to the United States. Considering all three questions Hispanics broke off on 1.6% of all ISR cases, NHNW: 1.2% and NHW: 0.5%. A breakoff on the ISR instrument can result in follow-up costs, imputation of missing data, or both. Because Hispanics and non-Hispanic non-whites breakoff much more often than non-Hispanic whites, especially on the citizenship-related questions, their survey response quality is differentially affected.

B.4. Cost analysis

Lower self-response rates would raise the cost of conducting the 2020 Census. We discuss those increased costs below. They also reduce the quality of the resulting data. Lower self-response rates degrade data quality because data obtained from NRFU have greater erroneous enumeration and whole-person imputation rates. An erroneous enumeration means a census person enumeration that should not have been counted for any of several reasons, such as, that the person (1) is a duplicate of a correct enumeration; (2) is inappropriate (e.g., the person died before Census Day); or (3) is enumerated in the wrong location for the relevant tabulation (<https://www.census.gov/coverage-measurement/definitions/>). A whole-person census imputation is a census microdata record for a person for which all characteristics are imputed.

Our analysis of the 2010 Census coverage errors (Census Coverage Measurement Estimation Report: Summary of Estimates of Coverage for Persons in the United States, Memo G-01) contains the relevant data. That study found that when the 2010 Census obtained a valid self-response (219 million persons),

the correct enumeration rate was 97.3%, erroneous enumerations were 2.5%, and whole-person census imputations were 0.3%. All erroneous enumeration and whole-person imputation rates are much greater for responses collected in NRFU. The vast majority of NRFU responses to the 2010 Census (59 million persons) were collected in May. During that month, the rate of correct enumerations was only 90.2%, the rate of incorrect enumeration was 4.8%, and the rate of whole-person census imputations was 5.0%. June NRFU accounted for 15 million persons, of whom only 84.6% were correctly enumerated, with erroneous enumerations of 5.7%, and whole-person census imputations of 9.6%. (See Table 19 of 2010 Census Memorandum G-01. That table does not provide statistics for all NRFU cases in aggregate.)

One reason that the erroneous enumeration and whole-person imputation rates are so much greater during NRFU is that the data are much more likely to be collected from a proxy rather than a household member, and, when they do come from a household member, that person has less accurate information than self-responders. The correct enumeration rate for NRFU household member interviews is 93.4% (see Table 21 of 2010 Census Memorandum G-01), compared to 97.3% for non-NRFU households (see Table 19). The information for 21.0% of the persons whose data were collected during NRFU is based on proxy responses. For these 16 million persons, the correct enumeration rate is only 70.1%. Among proxy responses, erroneous enumerations are 6.7% and whole-person census imputations are 23.1% (see Table 21).

Using these data, we can develop a cautious estimate of the data quality consequences of adding the citizenship question. We assume that citizens are unaffected by the change and that an additional 5.1% of households with at least one noncitizen go into NRFU because they do not self-respond. We expect about 126 million occupied households in the 2020 Census. From the 2016 ACS, we estimate that 9.8% of all households contain at least one noncitizen. Combining these assumptions implies an additional 630,000 households in NRFU. If the NRFU data for those households have the same quality as the average NRFU data in the 2010 Census, then the result would be 139,000 fewer correct enumerations, of which 46,000 are additional erroneous enumerations and 93,000 are additional whole-person census imputations. This analysis assumes that, during the NRFU operations, a cooperative member of the household supplies data 79.0% of the time and 21.0% receive proxy responses. If all of these new NRFU cases go to proxy responses instead, the result would be 432,000 fewer correct enumerations, of which 67,000 are erroneous enumerations and 365,000 are whole-person census imputations.

For Alternative B, our estimate of the incremental cost proceeds as follows. Using the analysis in the paragraph above, the estimated NRFU workload will increase by approximately 630,000 households, or approximately 0.5 percentage points. We currently estimate that for each percentage point increase in NRFU, the cost of the 2020 Census increases by approximately \$55 million. Accordingly, the addition of a question on citizenship could increase the cost of the 2020 Census by at least \$27.5 million. It is worth stressing that this cost estimate is a lower bound. Our estimate of \$55 million for each percentage point increase in NRFU is based on an average of three visits per household. We expect that many more of these noncitizen households would receive six NRFU visits.

We believe that \$27.5 million is a conservative estimate because the other evidence cited in this report suggests that the differences between citizen and noncitizen response rates and data quality will be amplified during the 2020 Census compared to historical levels. Hence, the decrease in self-response for citizen households in 2020 could be much greater than the 5.1 percentage points we observed during the 2010 Census.

Alternative C: Use administrative data on citizenship instead of add the question to the 2020 Census

Under this alternative, we would add the capability to link an accurate, edited citizenship variable from administrative records to the final 2020 Census microdata files. We would then produce block-level tables of citizen voting age population by race and ethnicity during the publication phase of the 2020 Census using the enhanced 2020 Census microdata.

The Census Bureau has conducted tests of its ability to link administrative data to supplement the decennial census and the ACS since the 1990s. Administrative record studies were performed for the 1990, 2000 and 2010 Censuses. We discuss some of the implications of the 2010 study below. We have used administrative data extensively in the production of the economic censuses for decades. Administrative business data from multiple sources are a key component of the production Business Register, which provides the frames for the economic censuses, annual, quarterly, and monthly business surveys. Administrative business data are also directly tabulated in many of our products.

In support of the 2020 Census, we moved the administrative data linking facility for households and individuals from research to production. This means that the ability to integrate administrative data at the record level is already part of the 2020 Census production environment. In addition, we began regularly ingesting and loading administrative data from the Social Security Administration, Internal Revenue Service and other federal and state sources into the 2020 Census data systems. In assessing the expected quality and cost of Alternative C, we assume the availability of these record linkage systems and the associated administrative data during the 2020 Census production cycle.

C.1. Quality of administrative record versus self-report citizenship status

We performed a detailed study of the responses to the citizenship question compared to the administrative record citizenship variable for the 2000 Census, 2010 ACS and 2016 ACS. These analyses confirm that the vast majority of citizens, as determined by reliable federal administrative records that require proof of citizenship, correctly report their status when asked a survey question. These analyses also demonstrate that when the administrative record source indicates an individual is not a citizen, the self-report is “citizen” for no less than 23.8% of the cases, and often more than 30%.

For all of these analyses, we linked the Census Bureau’s enhanced version of the SSA Numident data using the production individual record linkage system to append an administrative citizenship variable to the relevant census and ACS microdata. The Numident data contain information on every person who has ever been issued a Social Security Number or an Individual Taxpayer Identification Number. Since 1972, SSA has required proof of citizenship or legal resident alien status from applicants. We use this verified citizenship status as our administrative citizenship variable. Because noncitizens must interact with SSA if they become naturalized citizens, these data reflect current citizenship status albeit with a lag for some noncitizens.

For our analysis of the 2000 Census long-form data, we linked the 2002 version of the Census Numident data, which is the version closest to the April 1, 2000 Census date. For 92.3% of the 2000 Census long-form respondents, we successfully linked the administrative citizenship variable. The 7.7% of persons for whom the administrative data are missing is comparable to the item non-response for self-responders in the mail-in pre-ISR-option ACS. When the administrative data indicated that the 2000 Census respondent was a citizen, the self-response was citizen: 98.8%. For this same group, the long-form response was

noncitizen: 0.9% and missing: 0.3%. By contrast, when the administrative data indicated that the respondent was not a citizen, the self-report was citizen: 29.9%, noncitizen: 66.4%, and missing: 3.7%.

In the same analysis of 2000 Census data, we consider three categories of individuals: the reference person (the individual who completed the census form for the household), relatives of the reference person, and individuals unrelated to the reference person. When the administrative data show that the individual is a citizen, the reference person, relatives of the reference person, and nonrelatives of the reference person have self-reported citizenship status of 98.7%, 98.9% and 97.2%, respectively. On the other hand, when the administrative data report that the individual was a noncitizen, the long-form response was citizen for 32.9% of the reference persons; that is, reference persons who are not citizens according to the administrative data self-report that they are not citizens in only 63.3% of the long-form responses. When they are reporting for a relative who is not a citizen according to the administrative data, reference persons list that individual as a citizen in 28.6% of the long-form responses. When they are reporting for a nonrelative who is not a citizen according to the administrative data, reference persons list that individual as a citizen in 20.4% of the long-form responses.

We analyzed the 2010 and 2016 ACS citizenship responses using the same methodology. The 2010 ACS respondents were linked to the 2010 version of the Census Numident. The 2016 ACS respondents were linked to the 2016 Census Numident. In 2010, 8.5% of the respondents could not be linked, or had missing citizenship status on the administrative data. In 2016, 10.9% could not be linked or had missing administrative data. We reached the same conclusions using 2010 and 2016 ACS data with the following exceptions. When the administrative data report that the individual is a citizen, the self-response is citizen on 96.9% of the 2010 ACS questionnaires and 93.8% of the 2016 questionnaires. These lower self-reported citizenship rates are due to missing responses on the ACS, not misclassification. As we noted above, the item nonresponse rate for the citizenship question has been increasing. These item nonresponse data show that some citizens are not reporting their status on the ACS at all. In 2010 and 2016, individuals for whom the administrative data indicate noncitizen respond citizen in 32.7% and 34.7% of the ACS questionnaires, respectively. The rates of missing ACS citizenship response are also greater for individuals who are noncitizens in the administrative data (2010: 4.1%, 2016: 7.7%). The analysis of reference persons, relatives, and nonrelatives is qualitatively identical to the 2000 Census analysis.

In all three analyses, the results for racial and ethnic groups and for voting age individuals are similar to the results for the whole population with one important exception. If the administrative data indicate that the person is a citizen, the self-report is citizen at a very high rate with the remainder being predominately missing self-reports for all groups. If the administrative data indicate noncitizen, the self-report is citizen at a very high rate (never less than 23.8% for any racial, ethnic or voting age group in any year we studied). The exception is the missing data rate for Hispanics, who are missing administrative data about twice as often as non-Hispanic blacks and three times as often as non-Hispanic whites.

C.2. Analysis of coverage differences between administrative and survey citizenship data

Our analysis suggests that the ACS and 2000 long form survey data have more complete coverage of citizenship than administrative record data, but the relative advantage of the survey data is diminishing. Citizenship status is missing for 10.9 percent of persons in the 2016 administrative records, and it is missing for 6.3 percent of persons in the 2016 ACS. This 4.6 percentage point gap between administrative and survey missing data rates is smaller than the gap in 2000 (6.9 percentage points) and 2010 (5.6

percentage points). Incomplete (through November) pre-production ACS data indicate that citizenship item nonresponse has again increased in 2017.

There is an important caveat to the conclusion that survey-based citizenship data are more complete than administrative records, albeit less so now than in 2000. The methods used to adjust the ACS weights for survey nonresponse and to allocate citizenship status for item nonresponse assume that the predicted answers of the sampled non-respondents are statistically the same as those of respondents. Our analysis casts serious doubt on this assumption, suggesting that those who do not respond to either the entire ACS or the citizenship question on the ACS are not statistically similar to those who do; in particular, their responses to the citizenship question would not be well-predicted by the answers of those who did respond.

The consequences of missing citizenship data in the administrative records are asymmetric. In the Census Numident, citizenship data may be missing for older citizens who obtained SSNs before the 1972 requirement to verify citizenship, naturalized citizens who have not confirmed their naturalization to SSA, and noncitizens who do not have an SSN or ITIN. All three of these shortcomings are addressed by adding data from the United States Citizen and Immigration Services (USCIS). Those data would complement the Census Numident data for older citizens and update those data for naturalized citizens. A less obvious, but equally important benefit, is that they would permit record linkage for legal resident aliens by allowing the construction of a supplementary record linkage master list for such people, who are only in scope for the Numident if they apply for and receive an SSN or ITIN. Consequently, the administrative records citizenship data would most likely have both more accurate citizen status and fewer missing individuals than would be the case for any survey-based collection method. Finally, having two sources of administrative citizenship data permits a detailed verification of the accuracy of those sources as well.

C.3. Cost of administrative record data production

For Alternative C, we estimate that the incremental cost, except for new MOUs, is \$450,000. This cost estimate includes the time to develop an MOU with USCIS, estimated ingestion and curation costs for USCIS data, incremental costs of other administrative data already in use in the 2020 Census but for which continued acquisition is now a requirement, and staff time to do the required statistical work for integration of the administrative-data citizenship status onto the 2020 Census microdata. This cost estimate is necessarily incomplete because we have not had adequate time to develop a draft MOU with USCIS, which is a requirement for getting a firm delivery cost estimate from the agency. Acquisition costs for other administrative data acquired or proposed for the 2020 Census varied from zero to \$1.5M. Thus the realistic range of cost estimates, including the cost of USCIS data, is between \$500,000 and \$2.0M