### 2020 Census: Data Reasonableness Review

National Academies of Sciences, Engineering, and Medicine July 21, 2021

Jason Devine
Assistant Division Chief for Census Programs – Population Division (POP)

Christine Borman 2020 Decennial Census Count Review – Population Division (POP)

The U.S. Census Bureau reviewed this data product for unauthorized disclosure of confidential information and approved the disclosure avoidance practices applied to this release. CBDRB-FY21-POP001-0183 Shape your future



# 2020 Census Presentation Roadmap

This presentation provides an overview of:

- 1. Objectives for reasonableness review in the 2020 Census
- 2. The focus of reviews at each stage of file processing
- 3. Review teams
- 4. Examples of findings from 2020 review



### 2020 Census Purpose of Data Review

Our goals for data review are:

- 1) To **identify data processing errors** and verify that edits and other processing steps have been properly applied.
- 2) To **assess data quality** by looking at item nonresponse/missing rates, population count only responses, proxy responses, and other early indicators of possible data quality issues.
- 3) To **evaluate demographic reasonableness** by looking at census responses and subsequent data files at multiple levels of geography compared to benchmarks, i.e., 2010 Census, American Community Survey data, and Population Estimates.

To do this, we take both a micro approach and a macro approach to the data review:

- At the micro level, we're looking to see if processing of individual records has been done correctly.
- At the macro level, we're looking to see if the aggregate results appear to be reasonable when compared to benchmark data.



# 2020 Census General Categories of Anomalies

There are two general categories of anomalies – data collection anomalies and data processing anomalies.

- 1) Data collection anomalies are less tangible and can lead to unexpected results. While we have developed many processes to address collection-based issues, some anomalies are just outliers or unexpected trends and <u>not</u> "problems" to be fixed. Example of common data collection anomalies that we have built procedures for are:
  - People responding multiple times or in multiple locations (i.e. residence rule confusion)
  - People not responding either at all or to specific demographic or housing questions
- 2) Data processing anomalies are more concrete and can be identified using our review programs where we double program many sections of the specification. Anomalies related to data processing could be because of:
  - A misinterpretation of the specification or a bug in the program itself
  - Situations that were not accounted for in processing specifications





#### 2020 Census Data Files Under Review

The response data are processed in multiple steps with reasonableness review conducted at each step:

#### DRF1

(Decennial Response File 1)

The DRF1 is the first file produced after data is collected.

It contains all response data including duplicate responses.

#### DRF2

(Decennial Response File 2)

The Primary Selection Algorithm selects which data from the DRF1 should represent a housing unit on the CUF.

#### CUF

(Census Unedited File)

Count imputation is applied to unresolved cases in the DRF2.

The CUF contains the final universe of addresses. enumeration status, and population count.

#### CEF

(Census Edited File)

Editing and imputation are applied to missing and erroneous values for all items.

The CEF then contains complete data for all items.

#### MDF

(Microdata Detail File)

Confidentiality protection and recodes are applied using the Disclosure Avoidance System to create the MDF.

After tabulation geography is added, the data are ready for tabulation and dissemination.

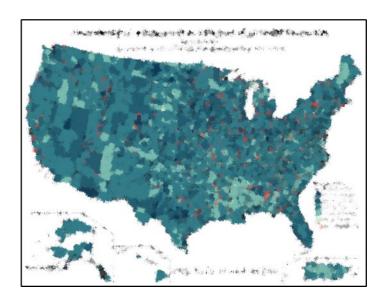
> your future **START HERE >**

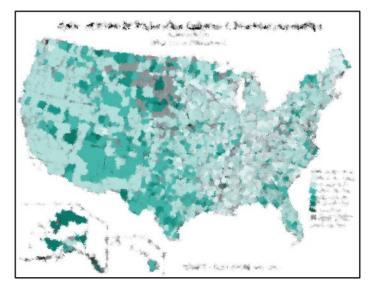


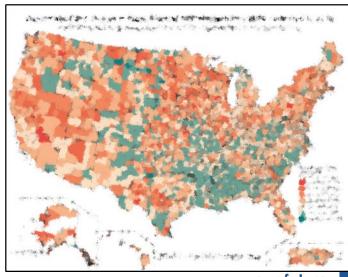
### 2020 Census Review of Early 2020 Census Results

Even before production DRF1 review began, experts in the Demographic Directorate (DEMO) reviewed the data as it was collected in order to identify and correct data processing errors and assess potential data quality and reasonableness concerns.

- This early review focused on housing unit responses from self-response and nonresponse followup.
- We used a lot of mapping to allow us to see trends in the data. We know different groups respond at different times, which we saw as the data came and the maps "filled in."
- For 2030, we have a lot of ideas about how this early review could be expanded to help us identify potential anomalies even earlier including adding Group Quarters to the mix.







# 2020 Census DRF1 Focus of Review

The DRF1 universe is larger than the final universe because it contains multiple responses for some addresses.

- For population totals, we focus our reviews on areas with suspected undercounts or extremely large overcounts. This includes looking at the household population and the group quarters population – where potential data collection issues may be more easily evident.
- **For characteristics**, we focus our review on identifying potential errors from data collection and response processing steps that are applied to the DRF1. We also take an early look at data quality and reasonableness.

We get this file for states on a flow basis, so we can't see the full picture for the nation until all states are delivered.





## 2020 Census DRF2 Focus of Review

The DRF2 universe is *smaller* than the DRF1 because all duplicate and multiple responses have been resolved.

The DRF2 population counts should be much closer to benchmarks. Accordingly, it is easier to identify areas with populations that are outliers when compared to benchmark data.

**For characteristics**, the smaller universe allows us to better assess data quality and reasonableness concerns, such as item nonresponse.

We get this file for all states in the nation at the same time – allowing us to do both a national and state-level analysis as soon as we get the data.



# 2020 Census **CUF** Focus of Review

The CUF universe is slightly larger than the DRF2 universe because count imputation has been applied.

Because the CUF is the basis for apportionment, we do a critical final look at the population totals compared to benchmarks.

**For characteristics,** we continue our review for data quality and reasonableness. We also assess the Hispanic origin and race codes from residual coding of write-in responses to ensure there were no coding or processing errors.

We get this file for states on a flow basis, so we can't see the full picture for the nation until all states are delivered.



# 2020 Census CEF Focus of Review

The population counts are locked in with the CUF and do not change in the CEF.

At this point, the data in the CEF should closely represent the demographics of the nation.

The focus of this review is on characteristics, including sex, age/date of birth, Hispanic origin, race, relationship, tenure, detailed vacancy status, and group quarters type. All characteristics should now have a valid and consistent response.

We get this file for states on a flow basis, so we can't see the full picture for the nation until all states are delivered.



## 2020 Census MDF Focus of Review

#### At this stage:

- State-level population totals do <u>not</u> change.
- The number of housing units and group quarters facilities by type at the block level do not change.
- Below-state population counts, person characteristics, and housing unit characteristics
   <u>do</u> change at the lower levels of geography because of the application of disclosure
   avoidance.

**For both population totals and characteristics**, teams review the MDF with the goal of ensuring that the MDF data are consistent with the expected effects of disclosure avoidance.

We get this file for all states in the nation all at the same time – allowing us to do both a national and state-level analysis as soon as we get the data.

#### 2020 Census Review Teams

There are three DEMO review teams that actively review the census data files – Subject Matter Experts (SMEs), General Experts (GEs), and the Special Total Population team (STP).

DEMO Subject Matter Experts (SME)	DEMO General Experts (GE)	DEMO Special Total Population (STP)
Focus on reasonableness and quality of characteristics using comparisons to benchmarks	Focus on aggregate population, group quarters, and housing units totals	Focus on reasonableness of state population totals for apportionment
Identify data collection and response processing errors	Identify deviations from benchmarks for lower levels of geography	Identify deviations from benchmarks as well as possible demographic trends



#### 2020 Census Subject Matter Experts (SMEs) - Overview

SMEs come from content-specific branches within DEMO. These analysts focus on specific topics – such as age, sex, race, or Hispanic origin – and have worked on the specifications for both data collections and response processing throughout the decade.

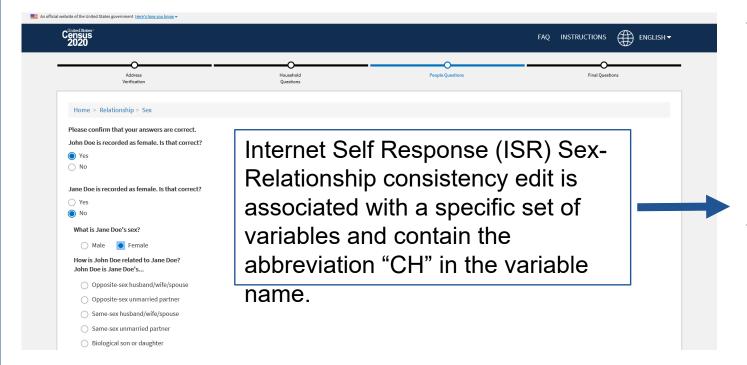
Their primary focuses are the reasonableness of specific characteristics of the population and characteristic-level data quality – particularly looking at item nonresponse, imputation rates, and for specific processing errors that may be impacting the quality of the data. Example of review activities are:

- Ensure all expected variables have values and verify that data collection variables are output as expected
- Examine data reasonableness and item nonresponse rates by data collection mode to check for issues with collection and capture of the data
- Verify specific processing steps
- Review percent changes and raw differences in characteristic distributions at multiple levels
  of geography compared to benchmarks





### 2020 Census Examples of SME Review Steps – Mode-Based Variable Checks



The sex-relationship consistency edit-specific variables should **never** be filled by responses coming from the paper mode. If it was, then this might indicate a data-transfer or processing error.

The purpose of this review step is to ensure the data itself has not been corrupted or transferred incorrectly during file creation and upstream data collection processes.

#### 2020 Census Examples of SME Review Steps – Mode-Based Reasonableness Checks

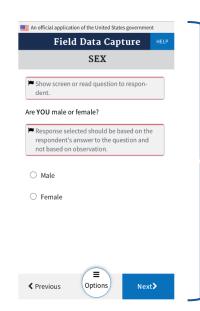
SMEs examine data reasonableness and item nonresponse rates (INR) by data collection mode to check for issues with collection and capture of the data.

Data Collection Mode	N		Missing Sex	Missing Age and Year of Birth	Missing Hispanic Origin	Missing Race	Missing Relationship
2010 Census Total Pop INR from the 2010 CUF 2020 DRF2 Total Population Internet Self Response (ISR) Paper Self Response NRFU Production Census Questionnaire Assistance (CQA) Coverage Improvement (CI) Update Leave (UL) GQ eResponse GQ Facility Self-enumeration GQ Paper Listing		<ul> <li>Whi</li> <li>Whi</li> <li>How to the How</li> </ul>	ch mode(s Is the INR r considera ch person v does the ne 2010 Ce v does the	tions? characteris INR for the ensus? INR compo	highest INF based on stics have t total popu are betwee		INR? pare ehold
NRFU Administrative Records Enumeration							
2020 DRF2 Household Population							
2020 DRF2 Group Quarters Population							

### 2020 Census Examples of SME Review Steps – DRF1 Processing Check

An official application of the United States government

SMEs verify specific processing steps by double programming portions of the DRF1 specification. In one step, DRF1 processing chooses between the responses based on a set of rules created by POP and SEHSD to determine the respondent's final response.



Field Data Capture The REVIEW enumerator I am going to read you a summary of the information I have recorded. Please let me selects the know if anything is incorrect. Select each box that contains incorrect infor "male" mation, and needs updating checkbox for the Date of birth: 1/1/1991 respondent Age (on 04/01/2020): 29 on the initial Sex screen.

The enumerator goes over the selected responses on the Review Screen and the respondent indicates that their sex is wrong.



The respondent instructs the enumerator the change their response to "female" on the review Sex screen.

In this example, the respondent's final response to the sex question should be "Female" in the DRF1.

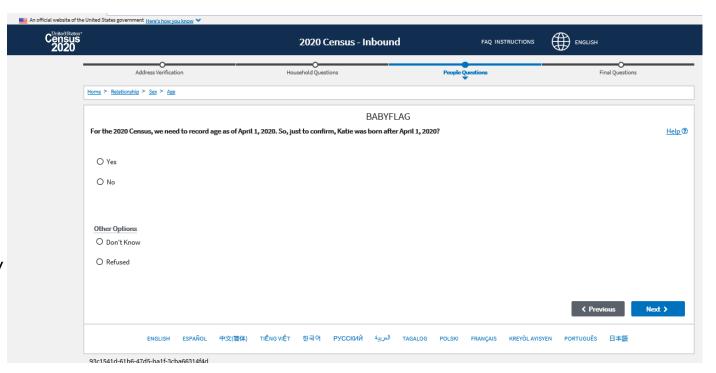
# 2020 Census Examples of SME Review Steps – DRF1 Processing Check (Baby Flag)

When people are confirmed to be born after Census day in the Census Questionnaire Assistance (CQA) instrument, they are flagged as confirmed babies born after Census day in DRF1 processing and removed from the Census.

Eligible date of births include:

- any birthdays starting with April 2, 2020
- incomplete date of births, such as May 2020, where the day is missing but the month is after April

SMEs verify that this flag is set correctly through their double programming.





### 2020 Census Subject Matter Experts (SMEs) – Reasonableness Review

SMEs review percent changes and raw differences in characteristic distributions at multiple levels of geography compared to benchmarks.

Race	CUF 2010 Number	ACS 2019 Number	DRF 2020 Number	Number Difference CUF 2010 vs DRF 2020	Number Difference ACS 2019 vs DRF 2020	CUF 2010 Percent	ACS 2019 Percent	DRF 2020 Percent	Percent Difference CUF 2010 vs DRF 2020	Percent Difference ACS 2019 vs DRF 2020
White alone										
Black alone										
AIAN alone				March Land	erachering		eastween my	Annual College Balley	Annah at Marian 18 70	Section of the sectio
Asian alone			25/13/400	- 20 0 - 20 - 2 + 6		- A.		A. 2000.0	had being being	
NHPI alone						- 5.* ·	The same to the state of the same			
Some other race alone										
Two or more races						- 3 7				4
Non-response				200						
Hispanic Origin	CUF 2010 Number	Estimates 2019 Number	DR Nu			erce				
Not Hispanic			Section 2	and the second	The second second	34 467		and the second	n <sup>2</sup>	į.
Hispanic					100					
No response					T.	A 100 1				*
Age (Computed or Reported)			PRF2 20 Numbe			2010 cent	* * * * *	10 of 10 mg		
18-64			gert ge		out the depotent of the second					2020
65+										ZUZU

#### 2020 Census General Experts (GEs) - Overview

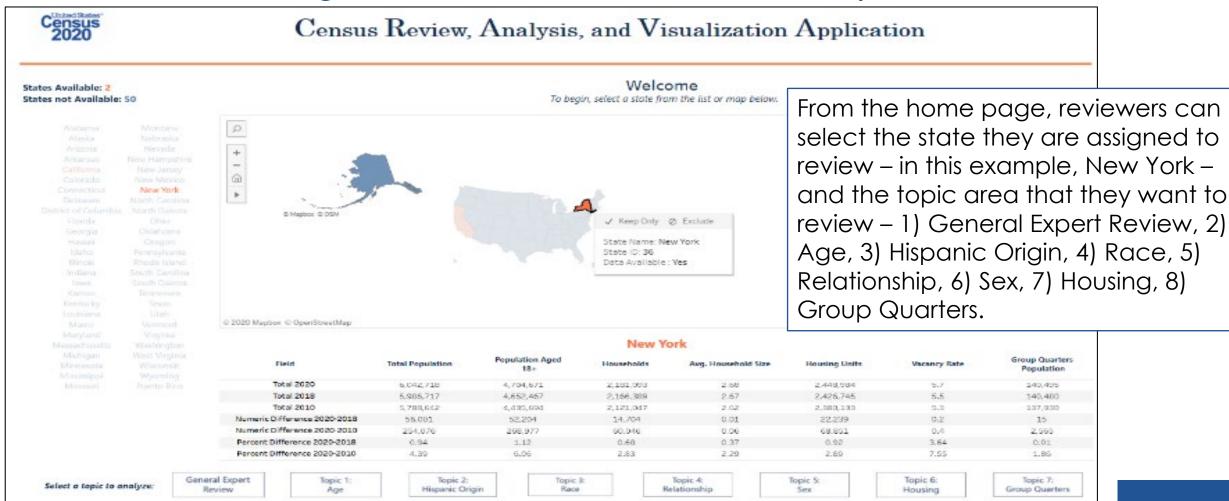
GEs are volunteers from within DEMO who have experience with reviewing and analyzing data for reasonableness from their work on other surveys.

- The focus of GE review is the same across files to identify any systematic issues or anomalies between the census files and benchmark data for the total population, group quarters population, and housing units at various levels of geography.
- GEs use a review tool specifically designed for their use, called the Census Review, Analysis, and Visualization Application (CRAVA).
- In the context of the GE review, an outlier is considered a geography with a 2020 tabulation that appears unreasonable or unexpected when compared to its benchmark tabulation.
- Outliers may exist because of a processing or data collection error, or because of explainable population shifts, such as:
  - changes to the geography's boundaries,
  - natural disasters that require evacuations from a particular area,
  - a new prison or college that increases the group quarters population,
  - economic forces that pull or push people away from a geography.

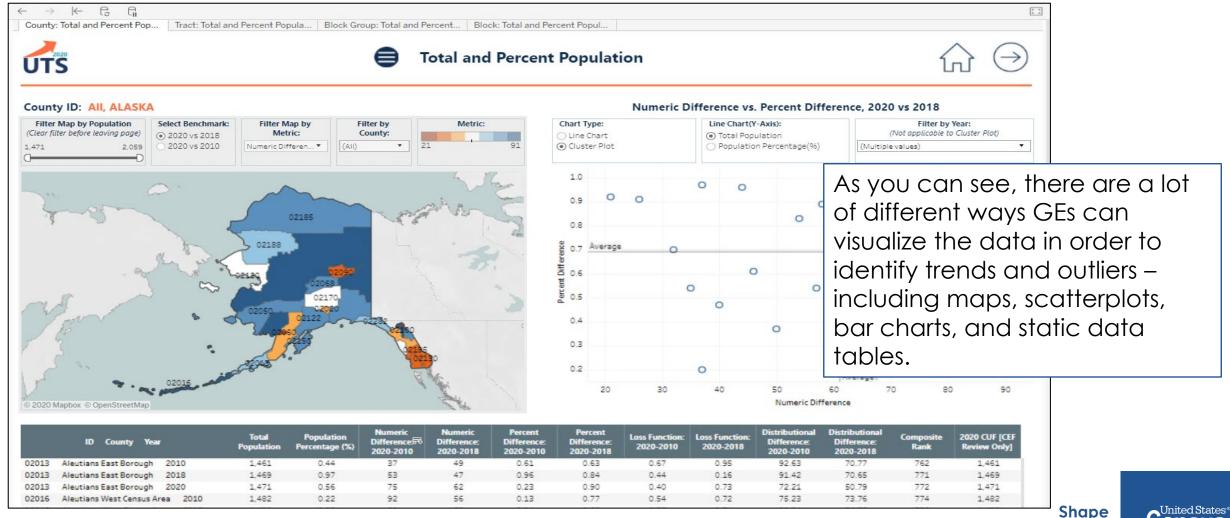




### 2020 Census CRAVA Home Page

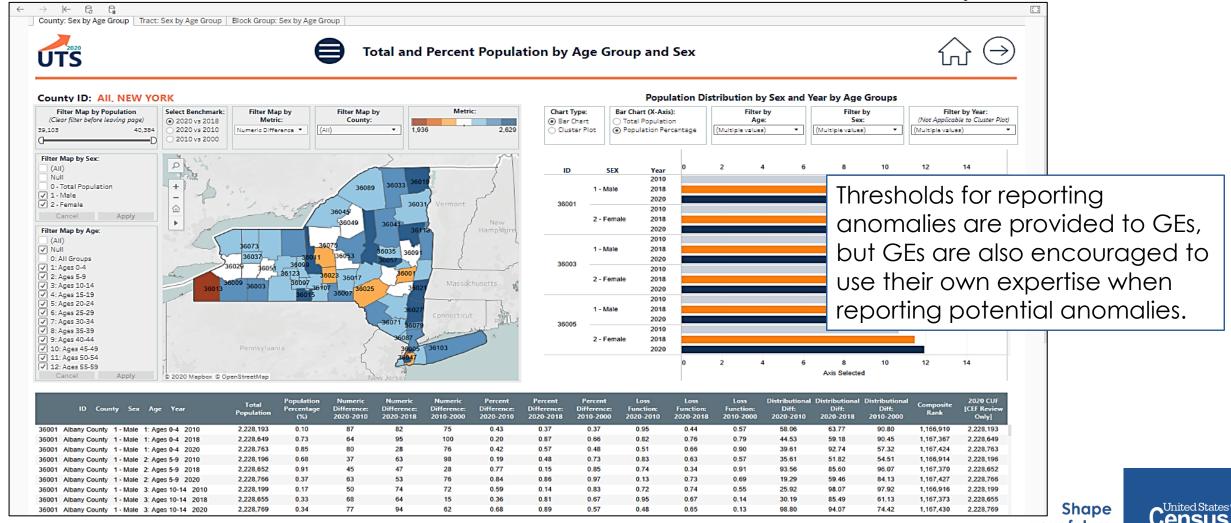


### 2020 Census CRAVA Review Screen Example 1

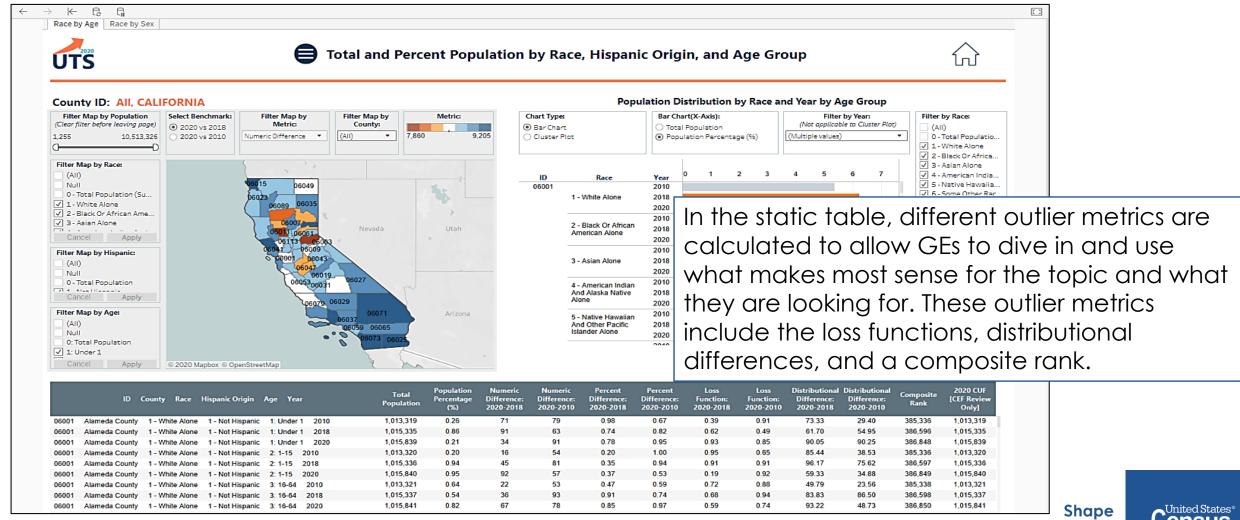




### 2020 Census CRAVA Review Screen Example 2



### 2020 Census CRAVA Review Screen Example 3



### 2020 Census Special Total Population Team (STP) – Overview

The STP is made up of experienced analysts and demographers with a variety of backgrounds.

- The STP uses characteristics data to assess the total population counts looking for potential over and under counts as well as duplication at the state level and below.
- The team uses a combination of CRAVA and independent tabulations and mapping to identify deviations from benchmarks as well as demographic trends, possible impacts of COVID-19, and data collection operational changes on the census results.
- Beyond nation, state, and county totals, the STP also looks at incorporated places and other levels of geography to identify potential over and undercounts.
- For each stage of review, the STP compares the results to benchmarks based on how we expect the numbers to look at that point in processing.





#### 2020 Census Examples of Findings from 2020 Review

- Processing errors both where the code itself was incorrect and where the data led to
  processing interacting in unexpected ways that needed to be corrected. See Michael
  Thieme's blog: <u>Finding 'Anomalies' Illustrates 2020 Census Quality Checks Are Working</u>
- Undercounts and overcounts of group quarters populations. See Deborah Stempowski's blog: 2020 Census Group Quarters
- Anomalies where the population shifted and did not require a fix. See Jason Devine, et al's blog: 2020 Census Data Review



### Questions





# END



