



# Data Quality of the 2018 New Zealand Census

Barry Milne



**COMPASS  
RESEARCH CENTRE**

FACULTY OF ARTS  
**THE UNIVERSITY OF AUCKLAND**

Whare Wānanga o Tāmaki Makaurau

COMPASS Seminar  
Tuesday, 3 March 2020

- ▣ Background to the Census
- ▣ What happened with Census 2018?
  - ▣ Why did it happen?
- ▣ What fixes were undertaken?
- ▣ What are the data quality implications?
  1. Population counts
  2. Electoral implications
  3. Use of alternative data sources
  4. Poor/very poor quality variables
- ▣ Guidelines for users of the Census
- ▣ Some recommendations that (I think) should be taken on board

## ❑ New Zealand Census of Population and Dwellings

- Official count of how many people and dwellings there are in the country at a set point in time (by age, sex, ethnicity, region, community)
- Detailed social, cultural and socio-economic information about the total New Zealand population and key groups in the population
- Undertaken since 1851, and every five years since 1881, with exceptions
  - No census during the Great Depression (1931)
  - No census during the Second World War (1941)
  - The 1946 Census was brought forward to September 1945
  - The Christchurch earthquakes caused the 2011 Census to be re-run in 2013
- Since 1966, held on first Tuesday in March of Census year
- The most recent census was undertaken on March 6, 2018

<http://archive.stats.govt.nz/Census/2013-census/info-about-the-census/intro-to-nz-census/history/history-summary.aspx>

- ❑ Census is important for
  - Electorates and electoral boundaries
  - Central and local government policy making and monitoring
  - Allocating resources from central government to local areas
  - Academic and market research
  - Statistical benchmarks
  - A data frame to select samples for social surveys
  - Many other things beside...
  
- ❑ *“every dollar invested in the census generates a net benefit of five dollars in the economy” (Bakker, 2014, Valuing the census, p. 5)*

- ❑ Obligations under Te Tiriti o Waitangi relating to the production of official statistics
- ❑ Stats NZ identify responsibilities to support Māori well-being and development ‘on their own terms’ and ‘to have equity as citizens’
  
- ❑ Census 2018
  - ❑ ‘Digital first’ census – access codes mailed.
  - ❑ Paper questionnaires made available as a back-up upon request.

# What happened?



	2006			2013			2018 (Interim)		
	Individual form	Partial form	Total	Individual form	Partial form	Total	Individual form	Partial form	Total
National population response rates	94.5	0.6	95.1	92.2	1.0	93.2	83.3	4.2	87.5
Sub-group response rates <sup>a</sup>									
– Māori	93.1	0.6	93.7	88.5	1.2	89.7	68.2	6.1	74.3
– Asian	91.0	1.1	92.1	91.7	1.6	93.3	81.7	6.1	87.8
– Pacific	92.4	1.5	93.9	88.3	2.5	90.8	65.1	8.4	73.5
– 15–29-year-olds	91.9	0.9	92.8	88.5	1.8	90.3	75.0	6.1	81.1

Source: Modified from Jack and Graziadej, 2019

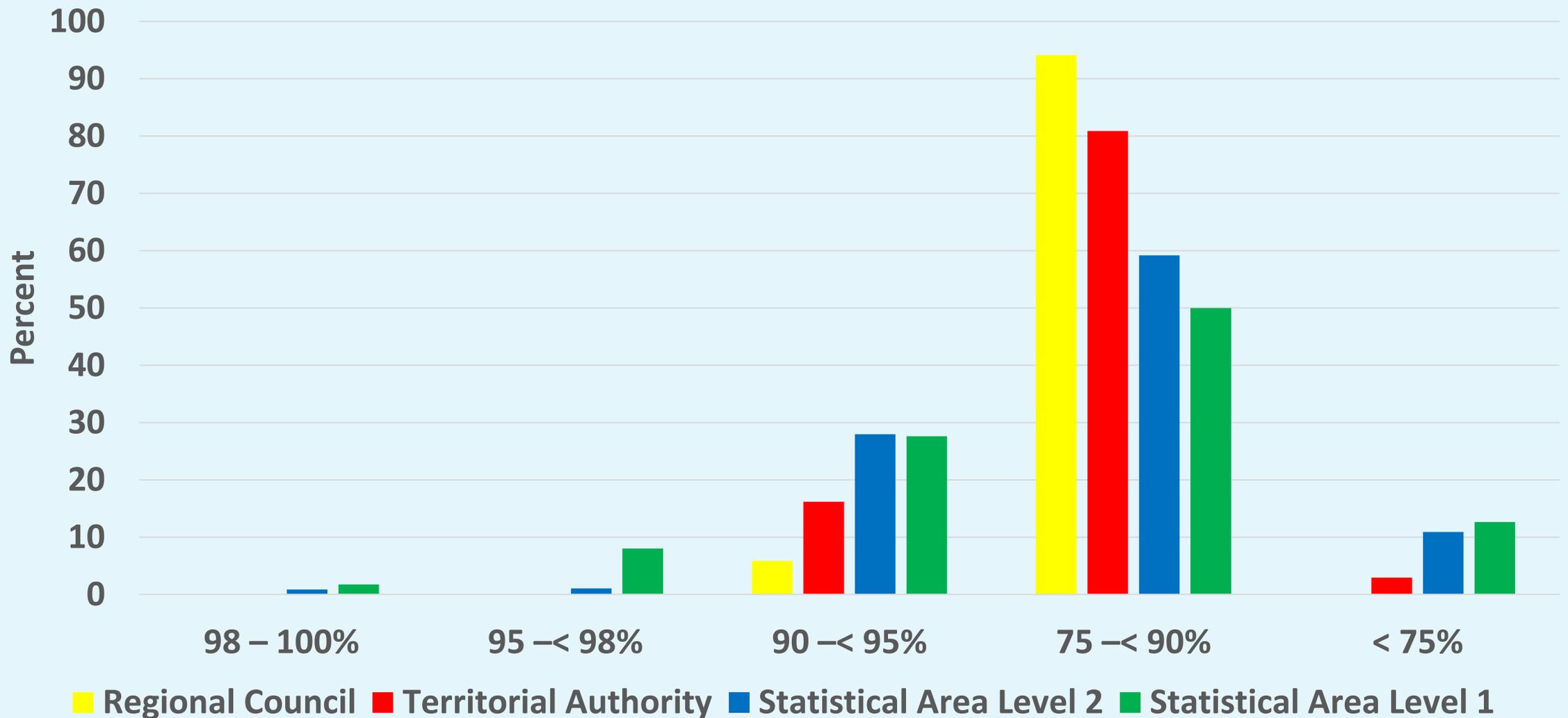
# What happened?



COMPASS  
RESEARCH CENTRE

FACULTY OF ARTS  
THE UNIVERSITY OF AUCKLAND

Whare Wānanga o Tāmaki Makaurau



# What happened?



SA2	Percent	Territorial Authority/Local Board	Region
Wiri West	46.9	Manurewa	Auckland
Mount Eden North East	52.3	Albert-Eden/Waitemata	Auckland
Otara Central	54.9	Otara-Papatoetoe	Auckland
Ferguson	55.0	Otara-Papatoetoe	Auckland
Ngapuna	55.5	Rotorua	Bay of Plenty
Ngapuhi	55.6	Far North	Northland
Waima Forest	55.9	Far North	Northland
Otara West	56.5	Otara-Papatoetoe	Auckland
Flaxmere West	56.6	Hastings	Hawke's Bay
Panmure-Glen Innes Industrial	57.0	Orakei/Maungakiekie-Tamaki	Auckland
Otara South	57.1	Otara-Papatoetoe	Auckland
Harania North	57.3	Mangere-Otahuhu	Auckland
Burbank	58.0	Manurewa	Auckland
Fordlands	58.0	Rotorua	Bay of Plenty
Queenstown Central	58.1	Queenstown-Lakes	Otago
Otangarei	58.5	Whangarei	Northland
Mangere West	58.6	Mangere-Otahuhu	Auckland
Bridge Pa	58.7	Hastings	Hawke's Bay
Otara East	58.9	Otara-Papatoetoe	Auckland
Rowandale West	58.9	Manurewa	Auckland
Hokianga North	58.9	Far North	Northland
Grange	59.1	Otara-Papatoetoe	Auckland
Queen Street	59.2	Waitemata	Auckland
Clendon Park North	59.8	Manurewa	Auckland

- 1% (n=24) SA2 areas had <60% Census completion
- 15/24 (62.5%) in Auckland, which contains only 26% of all SA2s
- 10 from the South Auckland boards of Otara-Papatoetoe (6), Manurewa (4), and Mangere-Otahuhu (2)
- 4 from Northland (3 from Far North District)

# Why did it happen?



- ❑ Factors associated with low response rates (Independent Review of New Zealand's 2018 Census; Jack and Graziadei, 2019):
  - Not enough field staff employed in time.
  - The importance of paper forms in this model was underestimated.
  - Requests for paper forms often went unheeded, or took a long time to arrive
  - The same online access code was required for each individual within the household to complete their respective form
  - A form couldn't be saved – if not completed in a session the respondent had to start over again

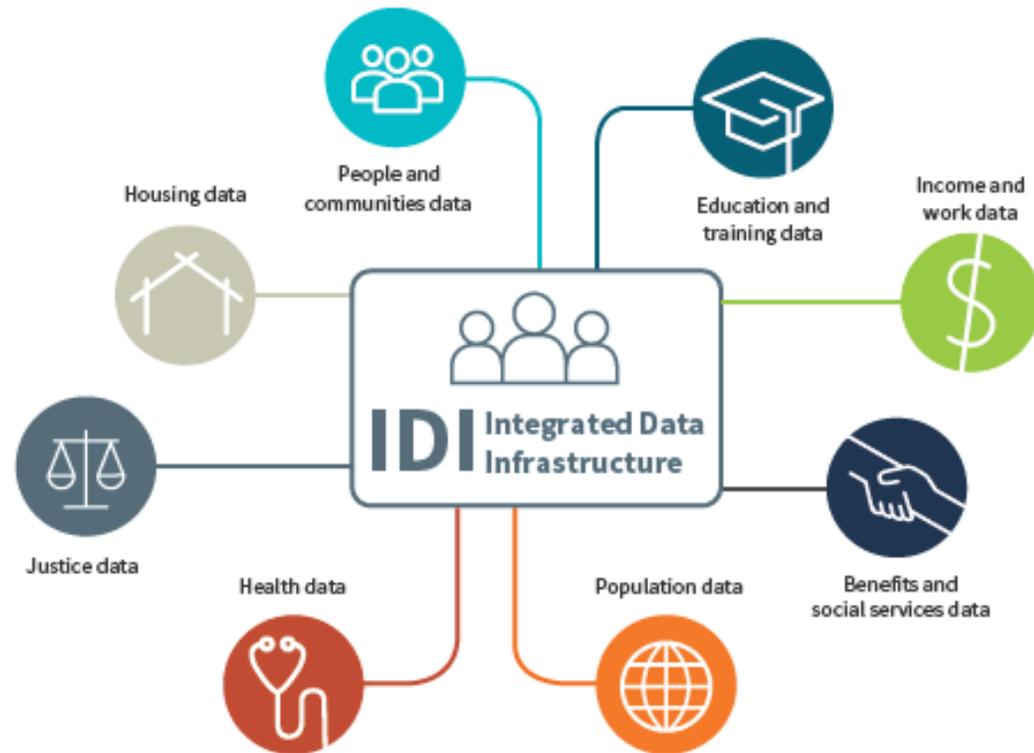
# Why did it happen?



- ❑ Factors associated with low response rates (Independent Review of New Zealand's 2018 Census; Jack and Graziadei, 2019):
  - Communication and engagement strategies didn't engage enough communities
  - Strategies put in place for non-private dwellings didn't work
  - It was decided not to follow up partial responses, meaning there was substantially more of these than previous Censuses

- ❑ Census 2018 External Data Quality Panel set up to advise on
  - whether the methodologies used to produce quality information from the census are based on sound research and a strong evidence base
  - approaches to data processing and methodology, and increased use of administrative sources that affect the quality of the data
  - data issues that may affect the usefulness of the data for Māori and iwi as Treaty partners
  - any quality issues people need to consider when using 2018 Census
  
- Dick Bedford, Alison Reid, Len Cook, Ian Cope, Tahu Kukutai, Donna Cormack, Thomas Lumley, Barry Milne
- August 2018 – February 2020

Stats NZ's Integrated Data Infrastructure (IDI) is a large research database containing de-identified microdata about people and households.



- IDI: Collection of administrative data sets linked at the individual level, de-identified, and available for research
- IDI spine: list of people who are likely to have ever been a resident of NZ
- IDI ERP-Sure: List of people we can be pretty sure are currently resident in NZ (subset of IDI spine)
- Behind IDI (not available for research) is identifiable information for people in IDI spine (allows for datasets to be linked)
- FIX 1: Use the IDI ERP-Sure to get the people who didn't fill out the census.

- ❑ Fix 1: Link Census 2018 records linked to people in the IDI spine (using name, date of birth, meshblock)
  - ❑ 97.7% linked; 1.2% estimated to be missed; <1% estimated to be incorrect
  - ❑ Add people AND grab characteristics about those people
    - Adding to households; adding entirely new households
- ❑ Fix 2: Corrections gave Stats NZ unit-record data files for every prisoner; Ministry of Defence did the same for those in NZ Defence Force. Data for Census non-responders identified from IDI and placed in correct locations.
  - ❑ 4,700/9,700 prisoners; 800/3,200 of those in NZ Defence Force

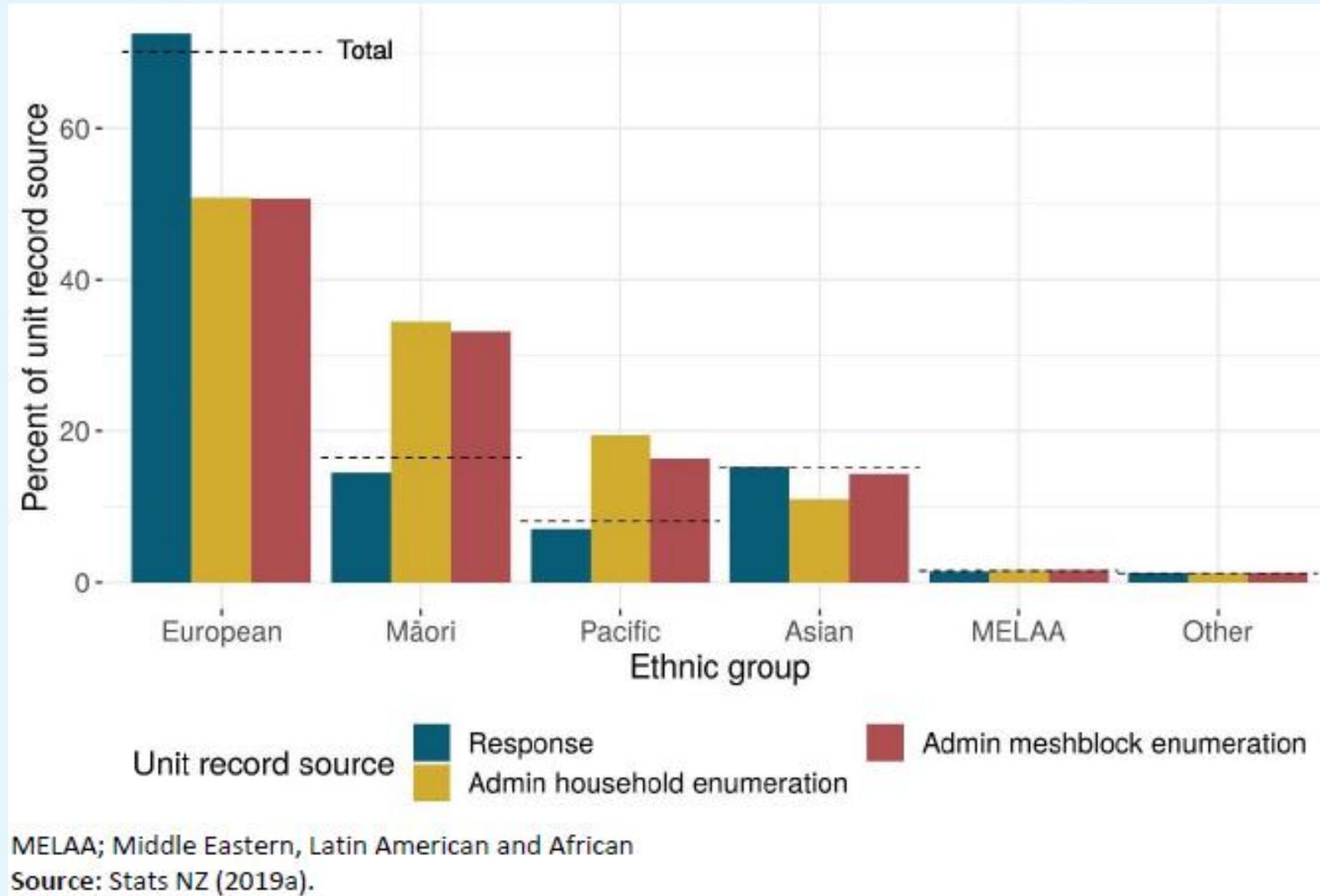
## 2018 Census usual resident population count

### By unit record source

Unit record source	Count		Percent	
Individual form	3,971,892		84.5	
Individual from household listing	202,914		4.3	
Field enumerated rough sleeper	99		<0.1	
<i>Census responses</i>		<b>4,174,902</b>		<b>88.8</b>
Admin household enumerations – occupied	99,159		2.1	
Admin household enumerations – unoccupied	42,252		0.9	
<i>Admin household enumerations</i>		<b>141,411</b>		<b>3.0</b>
Admin enumeration at responding private dwelling	20,643		0.4	
Admin enumeration at prison or penal institution	4,707		0.1	
Admin enumeration at defence establishment	798		<0.1	
<i>Other dwelling-based admin enumeration</i>		<b>26,148</b>		<b>0.6</b>
<i>Admin meshblock enumerations</i>		<b>357,294</b>		<b>7.6</b>
<b>Total</b>		<b>4,699,755</b>		<b>100</b>

- ❑ The final Census usual resident population of 4,699,800 is estimated to cover 98.6 percent of the estimated New Zealand population at 6 March 2018 of 4,768,600 (using 'dual system estimation' based on Census & IDI-ERP-Sure).
- ❑ The under-count of 68,800 represents 1.4 percent of the estimated New Zealand population, compared to 2.4 percent in 2013 and 2.0 percent in 2006.
- ❑ However, the 2018 result is obtained only **after** 524,900 were added to the Census dataset from administrative data.

# Fixes



- ❑ EDQP endorsed the statistical approaches used to mitigate non-response
  - ❑ A census with 17% missing individual responses was not an option
  - ❑ Mitigation worked to get a census file that counts most New Zealanders
- ❑ Mitigations raise questions around social licence, cultural licence (collective mandate for the trusted use of Māori data), and Māori data sovereignty
  - ❑ No comprehensive and open public consultation with New Zealanders, including with the groups most affected by the use of alternative data, to gauge the acceptability of the revised census approach

- ❑ Was the data linkage legal?
  - Yes, according to Stats NZ’s legal advice
- ❑ Does the linking of admin data to census data enjoys social licence (i.e., tacit approval from the New Zealand public)?
  - Unclear...
  - SNZ “should ... provide clear notice to the public about ... the retention and use of names and addresses and integration with the IDI and explain that this is legitimate and adds value” (Simply Privacy, 2017, p. 13).
  - The individual and dwellings census forms did not contain this information
  - Retaining the trust of Māori is especially important, given that Māori have lower levels of institutional trust, but are among those most impacted by the extensive use of administrative data for census mitigation.

- ❑ Not clear that people in New Zealand understand the extent of data sets that are linked to the census, nor that it would not affect their willingness to provide data if they did understand
- ❑ Consent from prisoners and those in defence force was not obtained from the individuals concerned
- ❑ Also not clear whether there was cultural licence: collective mandate for the trusted use of Māori data, based on the trust that iwi and Māori Treaty partners have.

- ❑ Fix 3: Where data wasn't available from Census 2018, up to three other sources were used (depending on the variable)
  - ❑ 2013 census
  - ❑ Administrative data
  - ❑ Imputation
- ❑ Also used when Census was completed but response to a question was 'Not elsewhere included'
  - ❑ 'not stated', 'response outside scope', 'response unidentifiable', 'refused to answer', 'don't know'
- ❑ **A very different Census data file**
  - ❑ data from a mix of sources:

# Data quality implications

## 1. Population counts

### ❑ Did the final Census file provide an accurate count of the population?

Variable name	Stats NZ Quality rating	Q/A Panel Quality rating
Age	Very high	Very high – at the national and regional council levels of geography.
Census night address	Moderate	Moderate – at the national and regional council level. There is greater uncertainty at lower levels of geography.
Count of population – census night	Moderate	Moderate – The rating is mostly due to comparability with previous census estimates, particularly for overseas visitors.
Count of population – usually resident	Very high	Very high – at the national and regional council/territorial authority and Auckland Council local board areas (TALB) level  There are a small number of meshblocks where NPDs have been allocated to different meshblocks compared to 2013. Users should be careful if they come across such changes, but this will not impact on the quality of data at higher levels of geography.
Sex	Very high	Very high – down to the SA2 level of geography.
Usual residence address	High	High – at the national and regional council/TALB level.

- ❑ Yes. Post-enumeration survey results not available yet, but dual system estimation using IDI-ERP-Sure suggests only a small undercount, and accurate counts down to TALB area.
- ❑ Distributions by age and sex also appear to be accurate.

# Data quality implications

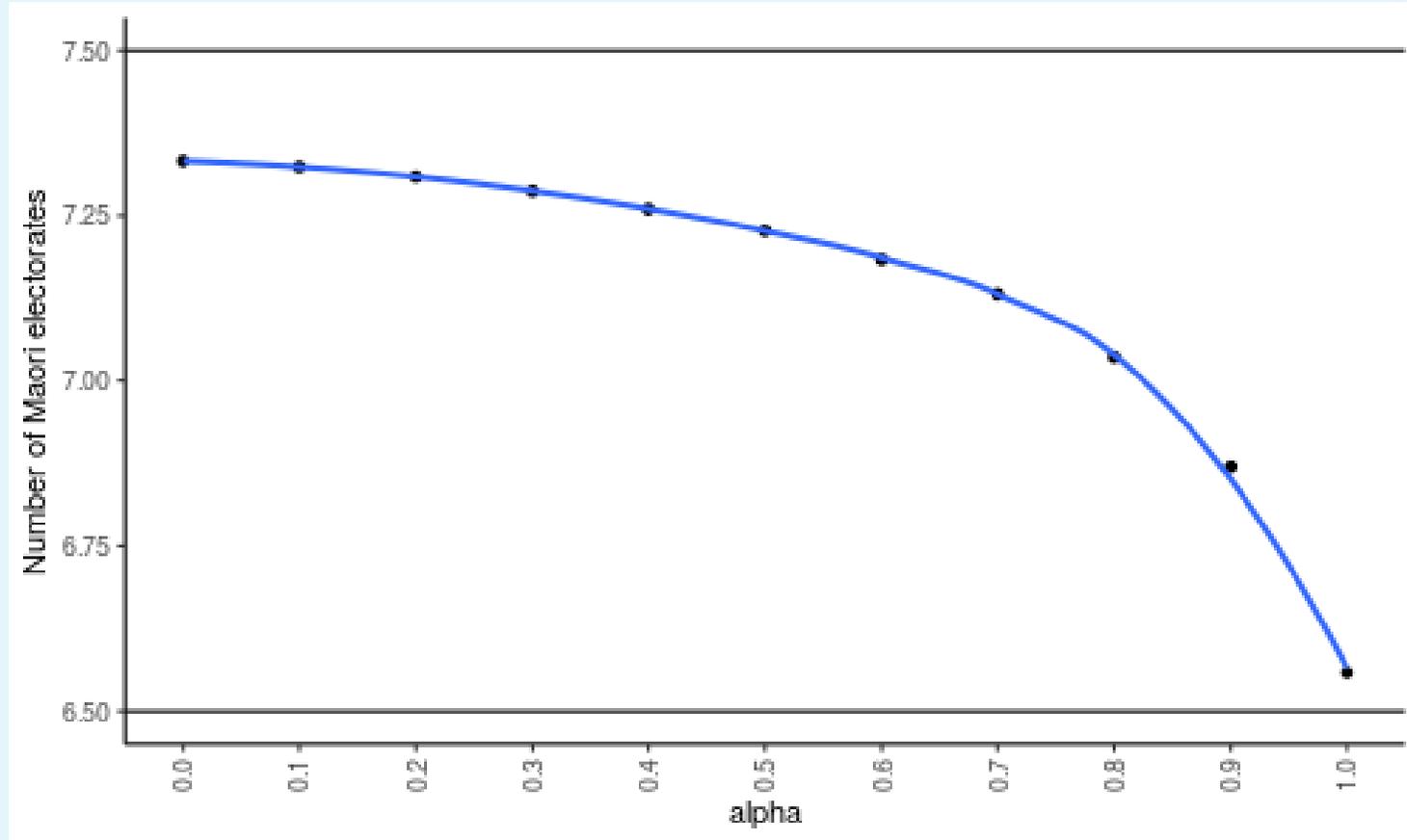
## 2. Electoral allocations

- ❑ Did the 2018 Census file allow for the number of electorates to be accurately determined
  - ❑ Yes
- ❑ Some background...
  - ❑ Māori electoral population (MEP) = electoral Māori descent usually resident population count multiplied by the percent of enrolled Māori voters choosing the Māori roll (52%).
  - ❑ The General electoral population (GEP) = the census usually resident population count minus the MEP.
  - ❑ The number of South Island general electorates is fixed at 16 (Electoral Act, 1993), so South Island GEP/16 = South Island quota
  - ❑ MEP/South Island quota = Number of Māori electorates
  - ❑ North Island GEP/South Island quota = Number of General electorates in the North Island
  - ❑ All electorates must have roughly the same population,  $\pm 5\%$

# Data quality implications

## 2. Electoral allocations

- Stats NZ used a threshold (alpha) for determining whether knowledge of a person's location was accurate enough to add that person to the Census file (1.0 = absolutely certain; 0.0 = a guess). Stats NZ use 0.5 for Census 2018.
- **REGARDLESS OF THE THRESHOLD CHOSEN, THE RESULT IS ALWAYS 7 MĀORI ELECTORATES**
- Unrealistic assumptions about population change would be needed for the number of Māori electorates to not be 7.

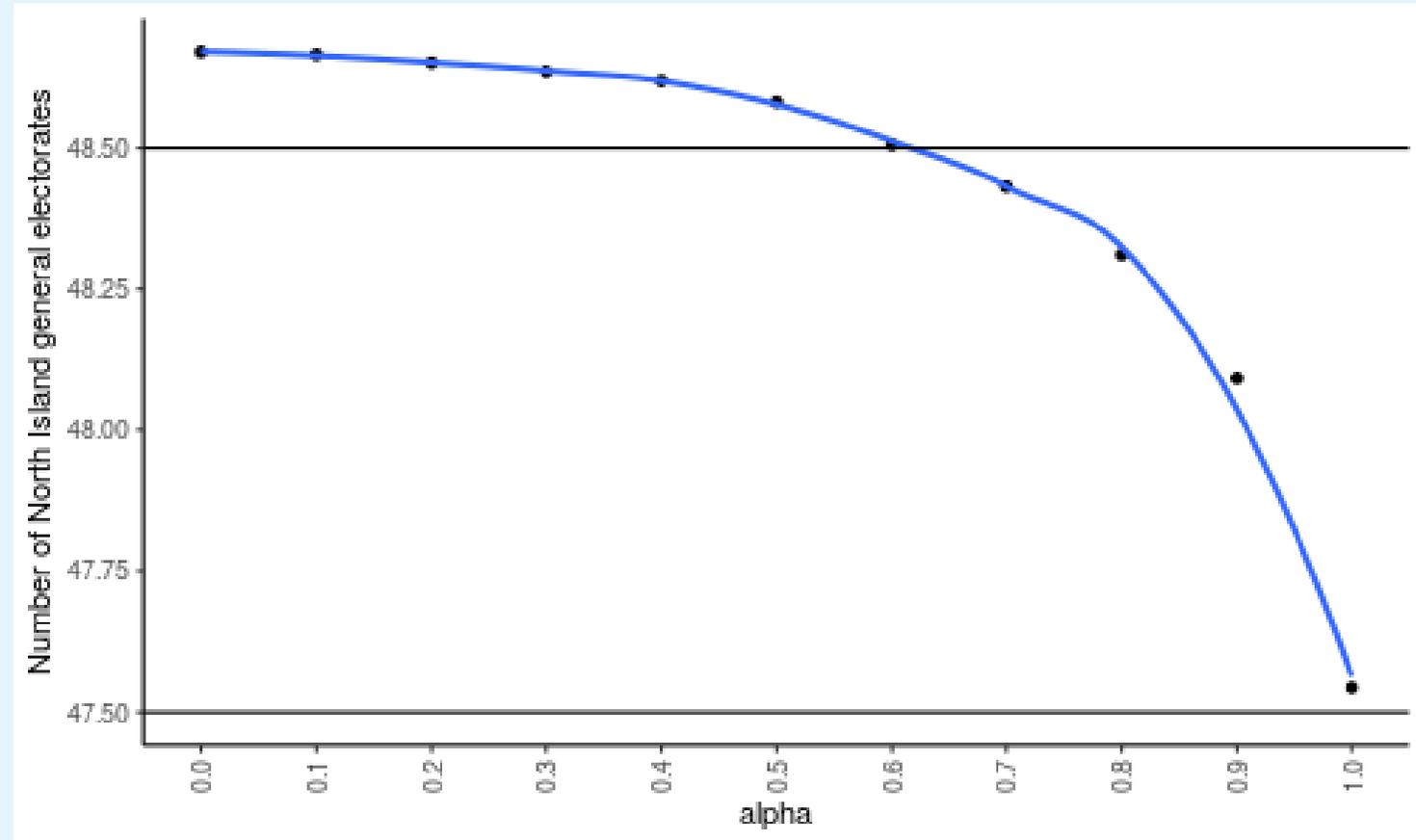


*Dot Loves Data (2019). Sensitivity analysis of 2018 Census for electoral boundaries. Unpublished report provided to Statistics NZ.*

# Data quality implications

## 2. Electoral allocations

- What about North Island general electorates?
- Here the threshold matters a little. Most thresholds ( $\leq 0.6$ ) suggest 49 electorates. Only strict thresholds suggest 48 electorates (as there was at the 2017 election).
- The Electoral Act 1993 enables the Government Statistician to exercise a degree of discretion; this would include the selection of alpha (and 0.5 seems reasonable).



*Dot Loves Data (2019). Sensitivity analysis of 2018 Census for electoral boundaries. Unpublished report provided to Statistics NZ.*

# Data quality implications

## 3. Alternative data sources

- ❑ The use of administrative data, Census 2013, and imputation data has improved the quality of the Census results
  - ❑ Census undercount reduced (a good thing)
  - ❑ Use of alternative data sources better than doing nothing, but not as good as if census as more complete
    - 2015 Cabinet Paper *Census Transformation - Promising Future*: “a census based on administrative data is not yet possible.”
  - ❑ But there are issues...

# Data quality implications

## 3. Alternative data sources: Admin data



- ❑ Admin data may not be contemporaneous with Census (6/3/2018)
  - A value of ethnicity from education data might have been supplied to the IDI in December 2017, from an enrolment in February 2017, which might itself have defaulted to the value given the first time that student enrolled.
- ❑ Admin data may not measure exactly the same thing
  - Taxable income from IRD is not the same as personal income reported at the Census
- ❑ >10% Admin data: Sector of Ownership, Industry, Workplace Address, Income, Sector of Landlord, Usual Residence Address, Weekly Rent Paid by Households, Age, Sex

# Data quality implications

## 3. Alternative data sources: Census 2013

- ❑ 2013 Census used for variables which do not change or change very little over time
  - Degree of change for variables will be underestimated for variable that do change over time
  - Sometime analysis of those that do change (ethnicity, smoking, religion) are of interest to researchers.
  
- ❑ >7% use of 2013 census: Usual residence 5 years ago, birth place, Māori descent, religion, languages spoken, ethnicity, smoking, years since arrival in NZ, highest secondary school qualification

# Data quality implications

## 3. Alternative data sources: Imputation

- ❑ CANCEIS (CANadian Census Edit and Imputation System) imputation system searches records that are near neighbours to find potential donors who are good matches on a set of matching variables. Closest match is chosen.
  - ❑ Unbiased, so should produce accurate counts
  - ❑ Accuracy may be low at the individual level and this will affect estimates bivariate associations
    - May increase estimates of association with variable included in the imputation model
    - May decrease estimates of association with variable not included in the imputation model
- ❑ >15% imputation: occupation, work and labour force status, main means of travel to work, main means of travel to education

# Data quality implications

## 4. Poor/very poor quality variables



- ❑ Stats NZ assessed the quality of variables using a five point scale (very high, high, moderate, poor, very poor), based on:
  - ❑ Metric 1 – data sources and coverage
    - A score (0–1) is given based on the contribution of each data source, weighted by a quality rating (0–1) give to each data source.
  - ❑ Metric 2 – consistency and coherence
    - comparability with the expected trends
    - comparability with other sources
  - ❑ Metric 3 – data quality
    - Including aspects such as coding, level of detail/classification, accuracy of responses

<https://www.stats.govt.nz/methods/data-quality-assurance-for-2018-census>

<http://datainfoplus.stats.govt.nz/Item/nz.govt.stats/ca28210f-3fd6-415c-a162-ecc07b4a28b0>

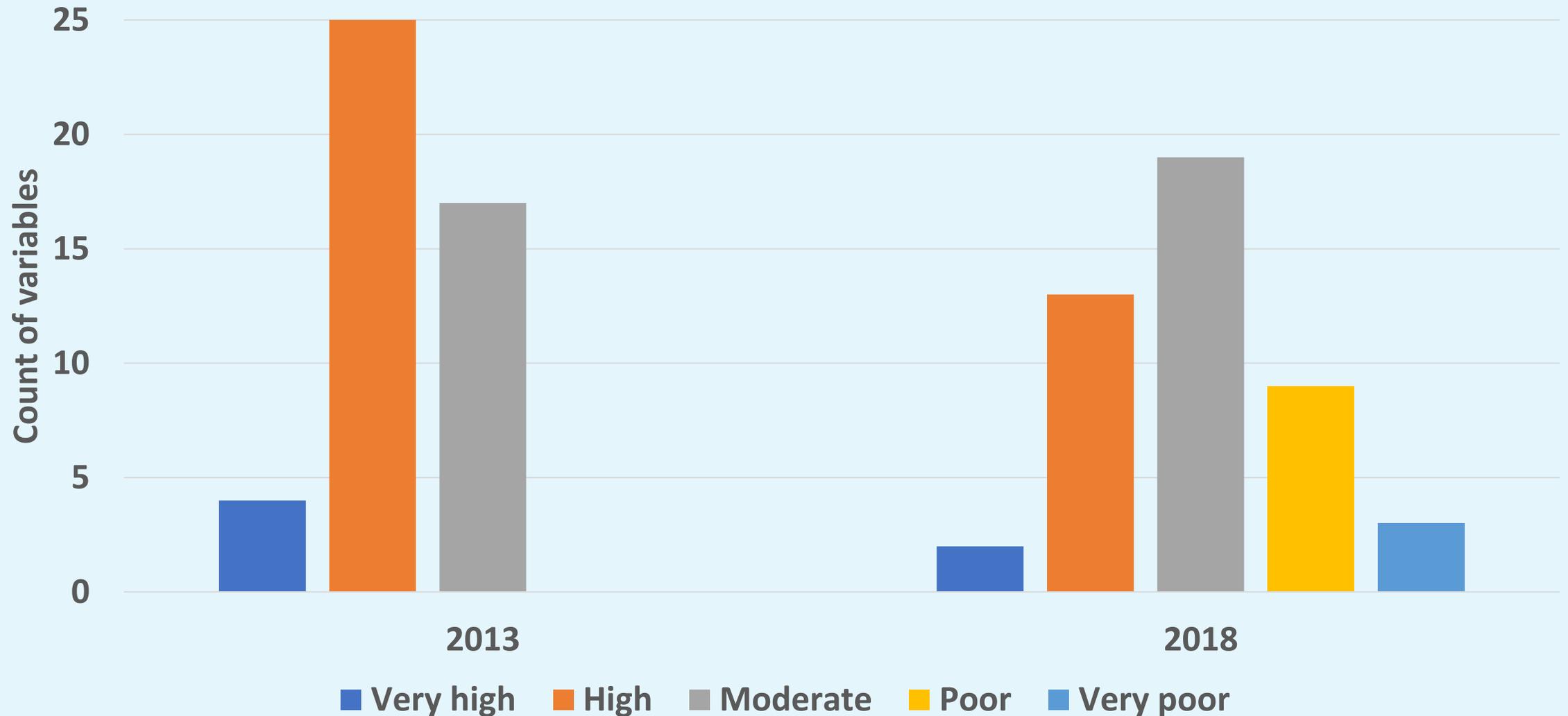
# Data quality implications

## 4. Poor/very poor quality variables

- ❑ EDQP adapted the quality framework used by Stats Canada to assess Census variables according to:
  - Coverage
    - *For the overall population, and by ethnic group (individual variables only) and regions*
  - Consistency
    - *Was a consistent classification used, and was data collection consistent across online and paper data collection methods?*
  - Comparability
    - *How does census 2018 compare to recent Censuses and other measures of the same variable?*
  - Contemporaneity
    - *Were all data sources used for the variable obtained at the same time?*
- ❑ EDQP tended to rate variables as lower quality than Stats NZ for subgroups and at lower levels of classifications

# Data quality implications

## 4. Poor/very poor quality variables



# Data quality implications

## 4. Poor/very poor quality variables

- ❑ There are only two individual/personal variables that have been rated by Stats NZ as having data of overall **very poor** quality. These are iwi affiliation and absentees from the household.
- ❑ There does not appear to be a robust or reliable way to address missing iwi data in Census 2018
  - ❑ Iwi administrative data are sparse. Data that do exist are of poorer quality than the census (Ministry of Education, Corrections, NZ Police)
  - ❑ At the aggregate level, there is very significant inter-censal change in iwi identification. It would be difficult to justify the use of an individual's 2013 census response to replace their missing 2018 response.
  - ❑ Significant changes to the iwi classification in 2017 classified a number of iwi for the first time. For these, no prior census data exists.

# Data quality implications

## 4. Poor/very poor quality variables



❑ 29 family and household variables are currently rated **very poor** quality, though Stats NZ are reviewing these ratings:

- Number of People in Family
- Number of Children in Family
- Number of Usual Residents in Household
- Number of Usual Residents Aged 15 and Over in Household
- Number of Usual Residents Aged Under 15 in Household
- Identification of Individual's Family Nucleus
- Individual's Role in Family Nucleus
- Dependent Child Under 18
- Dependent Young Person Indicator
- Number of Dependent Children in Family
- Number of Adult Children in Family
- Age of Youngest Child in Family
- Age of Youngest Dependent Child in Family
- Family Type
- Family Type with Type of Couple
- Family Type by Number of Children
- Extended Family Type
- Family Type by Child Dependency Status
- Household Composition
- Number of Dependent Children in Household
- Age of Youngest Child in Household
- Age of Youngest Dependent Child in Household
- Household Composition by Child Dependency Status
- Type of Couple
- Age of Male Partner in Opposite-Sex Couple
- Age of Female Partner in Opposite-Sex Couple
- Age of Older Partner in Same-Sex Couple
- Age of Younger Partner in Same-Sex Couple
- Sex of Sole Parent

# Data quality implications

## 4. Poor/very poor quality variables



- ❑ Lack of coverage of families in admin data means the potential for producing census-type information on families is currently minimal
  - ❑ ~357,000 people (from admin data) not able to be placed into a dwelling
  - ❑ a disproportionate number of these are for meshblocks in areas where Māori and Pacific populations are high

# Data quality implications

## 4. Poor/very poor quality variables

- ❑ Problems in how Stats NZ's new processing system handled the complex processing and coding of household and family data
  - A large decrease in one-parent families
  - Potential undercount of children under 5 years old
  - Underage partners in opposite sex couples
  - Some very old “children”, very young “parents”, and very young people living alone
  - There is an overcount in same-sex couples
    - Implausibly large increases in the age of the older partner in same-sex couples
  - Major increase in the number of households comprising a couple and other person(s)
- ❑ Chose not to dedicate staff to family coding issues
  - Fewer households to manual coding (3% vs 18% previously)
- ❑ Too hard to fix given time constraints

# Data quality implications

## 4. Poor/very poor quality variables



- ❑ The following variables have been rated by Stats NZ as having overall **poor** quality data:
  - Activity limitations;
  - Individual home ownership;
  - Number of rooms;
  - Qualifications: Post-school qualification field of study;
  - Relationship status: Legally registered relationship status, and partnership status in current relationship;
  - Unpaid activities;
  - Usual residence one year ago;
  - Usual residence five years ago;
  - Years at usual residence.

# Data quality implications

## 4. Poor/very poor quality variables

- ❑ In addition, EDQP assess data to be poor or very poor for some levels of some classifications, and for some ethnic groups
  - ❑ Very poor
    - Level 4 of the ethnicity classification for 45 “Middle Eastern, Latin American and African (MELAA)” ethnicities
  - ❑ Poor
    - “Te reo” under the language classification (and perhaps other Level 4 languages)
    - Smoking for Māori, Pacific and MELAA (through over-reliance on Census 2013 data)
    - Hours worked in employment for Pacific (nearly 40% imputation)
    - Occupation (overall)

# Data quality implications

## 4. Poor/very poor quality variables



### EDQP believe

- **Very poor** quality variables should not be released.
- Data rated overall as being of **poor** quality overall has the potential to mislead and that such data should not be released as official statistics.
- Access to data rated as **poor** quality overall should be restricted to accredited individuals working in controlled environments who are able to work closely with Stats NZ to understand the quality of the data.

## ▣ Guidelines for use of 2018 Census data

- ▣ Data quality is differential by ethnicity (and region and other factors) so caution is advised when undertaking comparisons.
- ▣ Read the EDQP's [assessments](#) and the relevant [Stats NZ DataInfo+ page](#)
- ▣ Check the use of [alternative data sources](#), overall, by subpopulation, and for small areas.
- ▣ Analyses may be affected by high levels of imputation.
  - Sensitivity analyses should test if imputation impacts results
  - Sensitivity analyses using missing-data techniques (e.g. multiple imputation) can be considered.
- ▣ Less 'no information' in Census 2018 needs to be accounted for when comparing across censuses.

# Recommendations for Stats NZ (selected)

- R 1.** **Stats NZ should** ensure data collection in future censuses is comprehensive enough to accurately measure iwi affiliation, and should take responsibility, in partnership with iwi, for investigating alternative ways to measure iwi affiliation so that the census is not the only source.
- R 2a.** **Stats NZ should** ensure there is genuine partnership with Māori communities, organisations and iwi to develop and implement decision-making and governance mechanisms, to ensure meaningful involvement of Māori in future censuses. This includes Stats NZ actively addressing the acceptability of the extensive use of administrative data in future censuses and issues of social license and Māori data sovereignty specifically for the 2023 Census.
- R 2b** **Stats NZ should** ensure there is a real voice for members of all communities, especially Pacific peoples and new migrants, in decision-making on data about them, including the use of admin data in the census.
- R 3.** **Stats NZ should** ensure individual census responses from prisoners are obtained in the 2023 Census.
- R 6.** **Stats NZ should** review the extent to which the way the online forms were administered contributed to missing responses in 2018, with a focus on the differential impacts for different population groups, and consider whether changes are needed for the 2023 Census.

# Recommendations for Stats NZ (selected)

- R 12. Stats NZ should** systematically investigate the impact of the use of alternative data sources (previous census data, data from a range of admin sources, imputed data) on the quality of data across variables. Analyses should focus ... on estimates of inter-censal change, the impact on the sizes of ethnic groups and small areas (e.g. SA2s), and the impact on bivariate associations between variables.
- R 17. Stats NZ should** support a dedicated team for the 2023 Census to undertake post-processing for families and households data, and other complex variables, and not divert this team to other tasks.
- R 19. Stats NZ should** only make data rated as being of **poor** or **very poor** quality overall available where project proposals are considered by Stats NZ on a case-by-case basis
- R 21. Stats NZ should** have an organisational commitment to, and focus on, achieving effective partnership with Māori to develop a census delivery model that will achieve a very high response (>94 percent) from Māori in the 2023 Census.
- R 22. Stats NZ should** set response rate targets for particular Territorial Authority and Auckland Local Board areas and ethnic groups that had low response rates in 2018.



# QUESTIONS?