Statistics Statistique
Canada Canada

# Methodology Branch

Business Survey Methods Division

# Direction de la méthodologie

Division des méthodes d'enquêtes-auprès
des entreprises

Canada

# The Methodology of the Quality Assurance Survey: 2002 and Beyond

**By**

Andrew Quigley

BSMD-2005-004E

# Abstract

The Quality Assurance Survey (QAS) is designed to determine the accuracy of classification of establishments on the Business Register (BR). During the most recent QAS in 2002, a review of the methodology was done in order to determine how the survey could be improved. In order to accommodate the continuing emphasis on revenue related estimates, one of the potential changes in the methodology is in how the stratification is carried out. Another change is in the collection process and it deals with who is performing the coding. This working paper examines both the current methodology used during the 2002 QAS and a template that can be used for a future QAS that incorporates the changes discussed in order to improve the quality of the results.

# Résumé

L'Enquête sur l'assurance de la qualité (EAQ) est conçue pour déterminer l'exactitude de la classification des établissements dans le Registre des entreprises (RE). Durant l'EAQ la plus récente, qui a eu lieu en 2002, on a examiné la méthodologie afin de déterminer comment on pourrait améliorer l'enquête. Étant donné l'importance qui continue d'être accordée aux estimations liées au revenu, l'une des modifications que l'on pourrait apporter à la méthodologie concerne la façon dont est réalisée la stratification. Un autre changement a trait au processus de collecte, plus précisément la question de savoir qui devrait effectuer le codage. Le présent document de travail porte sur la méthodologie courante utilisée durant l'EAQ de 2002 et sur un modèle qui pourrait être utilisé pour une EAQ future dans lequel sont intégrés les changements discutés en vue d'améliorer la qualité des résultats.

# Table of Contents

# 1. Introduction

The Quality Assurance Survey (QAS) is designed to determine the accuracy of the business activity coding on the Business Register (BR). It had not been performed since 1997 and the recently completed 2002 QAS can be used as a learning tool for future QAS's.

This report begins with a historical look at the QAS and why it was postponed after 1997. The next section focuses on the recently completed 2002 QAS. It includes a detailed description of the methodology used and the issues that arose during the analyses that may affect the execution of future QAS's. The final section introduces a "template" that may be used for the proposed 2005 QAS. It contains a description of the methodological issues for the elaboration of the sample design, collection method, edit and imputation, and estimation procedures. This section will also mention who may be using the results of the QAS and how they can be used.

# 2. History of the QAS

The QAS was carried out three times in the mid-nineties: 1993, 1995, and 1997. Using the Standard Industrial Classification (SIC80) system, the historical proportion of units on the BR that were 'consistently classified' was 90-92%. A unit was considered consistently classified if its business activity from the survey matched the BR at the 2-digit SIC80 (the letter groupings).

The QAS was postponed after 1997 due to the major changes on the BR. The first major change was the conversion from SIC80 to the North American Industrial Classification System (NAICS). NAICS was introduced to harmonize the classification systems used by Canada, The United States of America and Mexico. NAICS is a more difficult classification system to use when compared to SIC80. NAICS uses a 6-digit figure while SIC80 uses a 4-digit code, leading to a larger number of possible classifications using NAICS. Another reason why NAICS is more complicated than SIC80 is that NAICS is a 'production process' classification system while SIC80 uses an 'output' based system. NAICS considers the entire process or service of a business to determine its proper code while SIC80 only examines the end product. This implies that for NAICS, the end product or service of a business may not necessarily designate what its code will be and thus is harder to code properly. Because of these reasons, the conversion from SIC80 to NAICS in 1997 was not always a 1:1 relationship and may have led to an increase in misclassification errors.

The second change in 1997 was the introduction of the non-employers to the BR. Prior to 1997, the BR was made up almost entirely of employer businesses. Only the few large non-employers were found on the BR before 1997. This addition of the remainder of the non-employers nearly doubled the size of the base and added unpredictability to the BR as these new units had never been classified before.

The third change to the BR was the introduction of the Business Number (BN). The BN was introduced by Canada Revenue Agency (CRA) so that each business would have a unique identifier which could be used to link each business to various other files and agencies. Previous to 1997, the employer businesses used the payroll deduction account (PAYDAC)

number as their unique identifier. But as the non-employers do not have a PAYDAC, the BN was introduced for all businesses to use.

## 3. The Methodology of the 2002 Quality Assurance Survey

The 2002 QAS was initially designed as a one-phase survey. Based on the results of the one-phase (Phase 1) survey, a second phase was added (Phase 2). This section describes the methodology of the two phases, including sample design, sample size determination, allocation, selection, collection, edit and imputation and estimation.

### 3.1 Phase 1

### 3.1.1 Sample design

The target population of Phase 1 of the 2002 QAS was all of the establishments considered to be active on the BR. The population consisted of 2.18 million establishments, based on the Unified Enterprise Survey (UES) survey universe file (SUF) from the fall of 2002. The UES's SUF contains all of the establishments in the population and not just the establishments that are in-sample in the UES. The sample design was a stratified/systematic sample, using industrial sector (2-digit NAICS) and method of tax remittance as stratification variables. There were 20 industrial sectors and 3 methods of tax remittance (T1, T2, and other), creating a total of 60 strata in the original sample design. Two strata, the 'others' in the management sector and the 'T1s' in public administration, were not surveyed because their stratum population sizes were negligible when compared to their industrial sector populations. The final Phase 1 sample consisted of 58 strata. The breakdown of the target population can be found in Table A1 in the Appendix I (taken from Latendresse (1), 2003).

### 3.1.2 Sample size determination, allocation and selection

For sample size determination, a target CV of 2.5% was desired at the national level and 5% at the sector or tax remittance level. Using these requirements, a sample size of approximately 5000 establishments was determined.

The sample size was allocated using proportional allocation to the square root of N, the population total. The allocation was performed twice. The first allocation was at the industrial sector level. For two industrial sectors, utilities and public administration, the sample was increased to reach a minimum sample size of 100. Once the industrial sector sample sizes were determined, proportional allocation was repeated within each industrial sector using the tax remittance stratification variable. No additional constraints were set for this allocation. The Phase 1 final sample size (n) was 5093 establishments. The breakdown of the Phase 1 sample size can be found in Table A2 in the Appendix I (taken from Latendresse (1), 2003). To select the sample, units were sorted by descending gross business income (GBI) within each stratum and systematic sampling was used to select the sample.

### 3.1.3 Collection

The Phase 1 collection was mainly performed by the Regional Offices (RO). All of the simple units (n=4875) were contacted by the RO and all of the complex units (n=218) were contacted by the Business Register Division (BRD). The complex units are the establishments that belong to an enterprise which meet at least one of the complexity flags as defined by the BR (multi-province, multi-industry, multi-legal, multi-establishment, or multi-location). All other units are considered simple. The phone call during the collection was blind, meaning that the coder did not have any prior knowledge regarding the unit they were calling. Based on the respondent's description, the coder generated a 6-digit NAICS code. A 4-digit SIC80 code was also generated but was only used to make comparisons between NAICS and SIC80. Unless otherwise stated, a unit was deemed to be 'consistently classified' if the business activity obtained from the survey matched the activity description on the BR at the 2-digit NAICS level.

### 3.1.4 Edit and Imputation

Some minor edits were performed on the responses before and after the collection of the 2002 QAS. They are summarized in a document by Latendresse (2) (2003). After the sample was initially selected, some of the selected establishments were listed as being dead on the BR (n=19) or did not have a BN (n=13). These establishments were assigned a frame response code of '60' (inactive) and were not contacted.

To deal with the issue of non-response, the 2002 QAS used the techniques of removing units and re-weighting. The refusal rate was extremely low (3 out of 5093), and these units were dropped from the analysis. The units that were active (n=3870) were re-weighted to account for those units that the coders were unable to contact (n=152). The units that were inactive (n=893) were re-weighted to those units that the coders were unable to locate (n=175). All of the re-weighting was done within each stratum and the formulas can be found in Figure 1.

For active units : $w_{ijk}^{*} = w_{ijk} \times \left( \dfrac{\text{\# units active} + \text{\# units unable to contact}}{\text{\# units active}} \right)$,

For inactive units : $w_{ijk}^{*} = w_{ijk} \times \left( \dfrac{\text{\# units inactive} + \text{\# units unable to locate}}{\text{\# units inactive}} \right)$,

Where :

$w_{ijk}^{*}$ = Adjusted Phase 1 sampling weight for unit $k$ in tax remittance $j$ and industrial sector $i$,

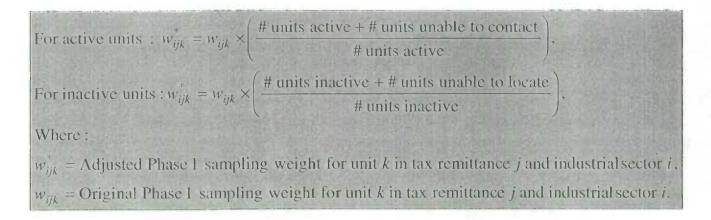$w_{ijk}$ = Original Phase 1 sampling weight for unit $k$ in tax remittance $j$ and industrial sector $i$.

Figure 1: Formulas for new weights to account for non-response during Phase 1

## 3.2 Phase 2

### 3.2.1 Rationale and sample design

The proportion of active units that were consistently classified at the 2-digits NAICS was found to be 72%. A second phase was added to the 2002 QAS after this result was obtained because it was far below the historical level of 90-92% and further validation was required. The rationale of Phase 2 was to substantiate the results of Phase 1 and to explore reasons for the misclassification. Note that public administration was not sampled in Phase 2 because its coding was concurrently being addressed elsewhere in the BRD.

The Phase 2 sample was created by 2 different methods. After evaluating the Phase 1 results, it was decided to contact all of the units from Phase 1 that were active and inconsistently coded at the 6-digit NAICS level (n=1962). This complete call back initially went out to three industrial sectors: professional services, manufacturing and accommodation and food services. Only the inconsistent units were contacted. To improve on the time requirements and cost of Phase 2, a sample design was devised for the remaining 16 industrial sectors. A stratified sample with simple random sampling was used, stratified only by industrial sector.

### 3.2.2 Sample size determination, allocation and selection

The sample size of the three complete call back industrial sectors was 366 establishments. The sample design for the 16 remaining industrial sectors was done considering four factors: a target CV of 10%, the proportion of Phase 1 units that were consistently classified per industrial sector, an estimated response rate of 50% for Phase 2 units and an estimated 50% misclassification rate of Phase 2 units (that is they would remain inconsistent after the re-contact). No pre-determined sample size was set and a simple random sample was selected within each of the strata. The sample size of the 16 sampled industrial sectors was 461 establishments, generating a total Phase 2 sample size of 827 establishments.

### 3.2.3 Collection

For both of the Phase 2 sample methods, the contacting and coding was performed by a single individual working in the BRD. This person was considered to be the BRD's most experienced coder and he processed the entire Phase 2 sample to ensure that the 6-digit NAICS code that he assigned would be the most accurate code to represent each business. For each Phase 2 unit, the coder was supplied with both the original NAICS code found on the BR and the NAICS code given during Phase 1. He was then able to use the information from the respondents to determine if the code on the BR was the most correct, if the Phase 2 code was the most correct or if neither was correct and a new, third NAICS, was more appropriate.

### 3.2.4 Additional Phase 2 question: Volatility

Partway through the collection of Phase 2, an additional question was asked to those units that had remained misclassified after the Phase 2 coding. The question was: "Were you at one time what the BR had you classified as?" The purpose of this question was to determine, for those units that were still misclassified after Phase 2, whether the unit had simply been miscoded. A unit may have once been coded correctly, but a change in business activity not updated on the BR may have led to it being misclassified. As this question was introduced partway through Phase 2, units from the three industrial sectors that had been collected first (manufacturing, professional services and accommodation and food services) were not asked this question.

### 3.2.5 Edit and imputation

Only minor edits were required for the Phase 2 sample. There were a few cases where a unit was given a response code of '3: new 6-digit NAICS' when they should have been given a response code of '1: NAICS code on BR was correct'. This was discovered when the new 6-digit NAICS supplied during Phase 2 actually matched the original code on the BR. Due to a very high response rate, no imputation or re-weighting was performed on the Phase 2 sample.

## 3.3 Estimation and analysis of Phase 1 and Phase 2

The information gained from Phase 2 was integrated with the results of Phase 1 using domain estimation. Estimates of interest included the proportion of units classified consistently at the global level and broken down by industrial sector and/or type of tax remittance. The majority of the estimates focused only on the units that were determined to be active during the survey. The general form of the formula for the global proportion of units that were consistent at the 2-digit NAICS can be seen in Figure 2. The first term in the numerator represents the estimated number of units that were consistent at the 6-digit NAICS after Phase 1. Because these units are consistent at the 6-digit NAICS, they are also consistent at the 2-digit NAICS. The second term in the numerator represents the estimated number of units that were found to be consistent at the 2-digit NAICS after Phase 2. The denominator represents the total number of units.

The issue of volatility was also integrated in a similar fashion to the Phase 2 information. For volatility, the formula in Figure 2 was modified so that $C_{ijk_2}$ would also be equal to 1 if the Phase 2 respondents had answered 'yes' to the volatility question.

Further analyses were also performed using the data. A 'hot-spots' table was created by making a cross tabulation chart based on the industrial sector that each unit was originally coded as on the BR versus the industrial sector determined during Phase 1 of the 2002 QAS. Areas were highlighted to indicate occurrences where common coding errors were made.

$$P = \frac{\displaystyle\sum_{i=1}^{20}\sum_{j=1}^{3}\sum_{k=1}^{n_{ij}} w_{ijk}^{*} \cdot I_{ijk} + \sum_{i=1}^{20}\sum_{j=1}^{3}\sum_{k=1}^{n_{ij_2}} w_{ijk}^{*} \cdot w_{ijk_2} \cdot C_{ijk_2}}{\displaystyle\sum_{i=1}^{20}\sum_{j=1}^{3}\sum_{k=1}^{n_{ij}} w_{ijk}^{*}}$$

Where :

$P$ = Proportion of units consistent at the 2 - digit NAICS,

$k$ = The $k^{th}$ unit in the sample,

$n_{ij}$ = Phase 1 sample size in tax remittance $j$ and industrial sector $i$,

$w_{ijk}^{*}$ = Adjusted Phase 1 sampling weight,

$$I_{ijk} = \begin{cases} 1 \text{ if unit } k \text{ is correct at the 6 - digit NAICS after Phase 1 in} \\ \qquad \text{tax remittance } j \text{ and industrial sector } i, \\ 0 \text{ otherwise,} \end{cases}$$

$n_{ij_2}$ = Phase 2 sample size in tax remittance $j$ and industrial sector $i$,

$w_{ijk_2}$ = Phase 2 post - stratification sampling weight ( = 0 if not in Phase 2),

$$C_{ijk_2} = \begin{cases} 1 \text{ if unit } k \text{ is now correct at the 2 - digit NAICS after Phase 2 in} \\ \qquad \text{tax remittance } j \text{ and industrial sector } i, \\ 0 \text{ otherwise.} \end{cases}$$

Figure 2: Formulas for proportional estimates integrating Phase 1 and 2

The results of the 2002 QAS were also mapped onto the Unified Enterprise Survey (UES). The UES is composed of approximately 26 surveys that use a similar methodology and the same sampling frame. The mapping from the QAS to the UES was based on survey and sampling stratum (i.e., take all, take some large, take some small, and take none). The purpose was to identify problem areas in coding with regards to units that should be eligible for a specific UES but are considered out of scope to the UES based on their code on the BR and units that should be out of scope but are coded as being eligible. See Macfarlane (2004) for a complete look at the hot spots, UES analysis and all other results of the 2002 QAS.

## 4. Template for the 2005 QAS

### 4.1 Goals of the 2005 QAS

The 2002 QAS was able to produce reliable estimates of the accuracy of NAICS coding on the BR. Based on suggestions and requests from methodologists and the BRD, potential changes that could improve the QAS were proposed. These changes could meet the needs of the BRD and business survey programs and increase the effectiveness of the QAS. The goals of the 2005 QAS will determine what potential changes will be implemented. The primary goal will be to get an accurate understanding of how well the business activity descriptions on the BR match the actual activity of the businesses. Estimates for this goal will be desired globally and for all industrial sectors based on both units and revenue. Besides this primary goal, other areas may be considered as potential focus points. The secondary goals could be to further examine the issue of volatility that was introduced during Phase 2 of the 2002 QAS and to examine the quality of the codes supplied by the coders. A tertiary goal could be to use the information regarding active/inactive businesses obtained from the QAS to make adjustments to other surveys which are using tax data replacement. Units that are being replaced by tax data may be inactive or out of scope, and adjustments to estimates could be applied. This goal would require a listing of which surveys the adjustments are aimed at to ensure that adequate sample sizes would be allocated in these areas.

### 4.2 Frame

The SUF from the fall 2004 UES can be used as the survey frame for the 2005 QAS. It contains all known-to-be-alive units on the BR and requires no additional creation of a survey frame. The sampling unit will be at the statistical establishment level, similar to the 2002 QAS.

### 4.3 Sample design

The 2002 QAS used a stratified sample design with systematic sampling (Phase 1) and simple random sampling (Phase 2) based on population counts. This method was suitable for the desired outcome: determining the proportion of the population on the BR that was active and consistently classified. Throughout our meetings with the BRD, the importance of revenue-related estimates was frequently discussed. There would be less concern with a 20% misclassification rate if these units represented less than 5% of the total revenue of the population. We were able to produce revenue-weighted estimates using the 2002 QAS, but the sample design was inefficient because it did not contain any criterion based on revenue. This led to estimates with large variances. To get more accurate estimates with smaller variances, such as the proportion of total revenue misclassified and the proportion of consistent units within a defined revenue stratum, future QAS's should use revenue as a stratification variable. A prospective table is shown in Table 1, using industrial sector and revenue size as stratification variables. The smallest revenue range contains units that are below the Royce-Maranda

threshold (revenue $\leq$ \$30,000). To establish the other revenue ranges the values of X and Y would be determined by the BRD and BSMD.

Table 1: Proposed breakdown using industrial sector and revenue size for stratification

| Industrial Sector | | Revenue Range | | | | Global |
|---|---|---|---|---|---|---|
| | | Rev≤\$30K | \$30K<Rev≤X | X<Rev≤Y | Rev>Y | |
| Agriculture | Proportion CV | | | | | |
| Mining and Oil | Proportion CV | | | | | |
| Utilities | Proportion CV | | | | | |
| Construction | Proportion CV | | | | | |

To meet the goals of the 2005 QAS, the sample design used should be a stratified sample, stratified on industrial sector (2-digit NAICS) and revenue. The industrial sector stratification variable is important in creating a sample which can be used to make inferences regarding the wide range of businesses in Canada. The 20 industrial sectors will be used. The revenue stratification variable should be used to ensure that the sample contains units that will reflect the diversity of revenue in Canadian businesses. This will allow valid estimation of revenue-related estimates. To stratify by revenue, the population can be split into 4 categories:

1) Revenue $\leq$ \$30K
2) \$30K < Revenue $\leq$ \$500K
3) \$500K < Revenue $\leq$ \$300M
4) Revenue > \$300M

During the 2002 QAS, some analyses were performed using revenue as the variable of interest. Initially, units were post-stratified based on their establishment revenue. When presented to the BRD, there was concern in the lack of units in the 'greater than \$300 million' category. This was rectified by post-stratifying based on a unit's enterprise revenue. For the 2005 QAS, it will need to be decided which revenue figure will be used for stratification. Stratifying by the establishment's revenue will require adjustments to the proposed size categories as there are few statistical establishments that have a revenue greater than \$300M.

The smallest of the proposed revenue size categories represents those units which are below the Royce-Maranda threshold and would therefore not be eligible for survey sampling. They are still required in the QAS because the information from these units is used for making macro-adjustments on estimates and allows for a greater comparison between the current and previous QAS's. The units in this stratum would be selected using simple random sampling. The largest revenue size category represents those units which are members of the large businesses (LBUS) and are contacted on an ongoing basis. This category would contain a small fraction of the total number of businesses in the population and could therefore be collected as a

take-all stratum. The 2 intermediate revenue size categories are used to split up the remaining units into units with small revenues and units with large revenues. These strata would also be sampled using simple random sampling. An alternative method to deriving the stratum boundaries for the two intermediate size categories is to use the Lavallée-Hiridoglou algorithm, which can define the boundaries to minimize the sample size of each stratum for a pre-determined coefficient of variation (CV). Once the optimal stratum boundaries are determined, sampling can be performed in a similar fashion to that mentioned above.

There are 2 important questions that need to be considered related to the sample design: How many size categories are needed and should these categories be based on a unit's establishment revenue or the revenue from its enterprise? With the previously mentioned revenue size categories, the sample design contains 80 strata. Will there be a large enough sample size to ensure proper coverage within each of these strata or should the number of revenue size categories be decreased to increase the number of units per stratum? Answers to these questions will need to be established before sampling can begin.

## 4.4 Collection method and data capture

The majority of the units in Phase 1 were contacted by the RO and all of the units in Phase 2 were contacted by the BRD's most experienced coder (MEC). Table 2 displays the breakdown of how the 6-digit NAICS codes given during Phase 1 by the RO compared with the 6-digit NAICS codes given by the BRD's MEC during Phase 2. The RO coder and the MEC gave the exact 6-digit NAICS code 59% of the time. The remaining 41% of the time they differed. This is a large rate of discrepancy between the coders. Even more of a concern is that of the entire Phase 2 sample, 26% (n=164) of the units differed at the industrial sector level (2-digit NAICS) between the 2 coders. This indicates that the person performing the coding can have a large effect on the quality of the coding.

Table 2: Comparison of coding performed by RO and MEC after Phase 2

| Comparison between Regional Office (RO) coder and BR's Most Experienced Coder (MEC) | Number | Percent | Different at NAICS 2 |
|---|---|---|---|
| Different at NAICS 6, with MEC matching BR base | 169 | 27% | 105 |
| Same at NAICS 6, with neither matching BR base | 373 | 59% | |
| Different at NAICS 6, neither matching BR base | 88 | 14% | 59 |
| Total | 630 | 100% | 164 |

To improve the quality and accuracy of the coding for the 2005 QAS, all of the coding should be performed by members of the BRD who have the most experience in using and understanding NAICS. This ensures that all NAICS codes supplied by the coders accurately reflect what each unit does. If it is felt that there may be a 'coder effect' on the quality of the codes supplied, the identity of the coders could be recorded for each unit. This information

could then be used to determine if this is a valid hypothesis. If it is, the 'coder effect' could be accounted for in any potential analyses using an appropriate sample design.

The data collection will involve a blind phone call, ensuring that the coder has no previous knowledge of business activities of the unit they are calling. The coder will ask each unit "What is your business activity?" and will derive a 6-digit NAICS based on the response given. The blind phone call is used to limit any bias that may be introduced if the coder were to know what the unit was coded as on the BR beforehand. Having a 1-phase survey may eliminate the ability to explore the areas of volatility and the quality of the coders as the coder would be unable to ask the question 'Were you at one time 'X', but have since changed your business activity?" because the coder would not know what 'X' is. An alternative to the question could be to ask if the establishment had changed their business activity since 'Y', where 'Y' is the date of the most recent BR update. If the respondent answers yes, additional pertinent information could then be obtained.

If the issues of volatility and quality of the coders are to be examined, either a change in the Phase 1 calling procedure or the introduction of a second phase will be required. A change in the calling procedure could be that the coder would initially record the unit's activity description and only then find out what the BR had the unit coded as. The coder would then be able to ask the volatility question to the respondent. The major concern with this change is that knowing what the BR has the unit coded as may introduce bias to how the coder interprets the activity description that was just obtained. The most appropriate method to examine volatility and quality of the coders is to perform a second phase. To carry out a second phase, a sub-sample of units from all strata, both consistent and inconsistent after the first phase, would be selected and called back by a different coder. This method allows for an examination of the quality of the initial code and can also allow for the examination of volatility by focusing only on the units that are still inconsistent after the second phase.

## 4.5 Editing

For the 2002 QAS, the active and inactive units were re-weighted within each stratum to incorporate the unable to contact and unable to locate units, respectively. This was an effective method for handling the non-response as they made up only a small proportion of the total sample units (6.5%) and thus did not significantly increase the sampling weights. This method can be used for the 2005 QAS. Minor edits similar to those used during the 2002 QAS will be performed to ensure logical responses.

## 4.6 Estimation

The 2002 QAS used the sample design weights and domain estimation to incorporate the results of Phase 1 and Phase 2. This method worked well when the estimate to be calculated was based on population counts but did not perform well when the estimate was based on population revenue. With the introduction of a revenue stratification variable, along with the retention of the industrial sector variable, the sample design weights and domain estimation will provide valid estimates for both population counts and revenue.

# 5 Conclusions

The QAS is used to determine the quality of classification of businesses on the BR. The 2002 QAS was the first instance of this survey in 5 years. A general comparison of the methodology of the 2002 QAS to the potential methodology of the 2005 QAS is summarized in Table 3. The 2002 QAS was performed in 2 phases due to the uncertain results of the initial survey, therefore adding added extra cost and time to the survey. To ensure higher levels of confidence during the collection stage it is recommended that the coding of businesses only be performed by coders with high levels of experience and knowledge of NAICS. This eliminates the requirement for a second phase, though one may be desired for other reasons.

During the course of the 2002 QAS, revenue emerged as a variable of major interest. We were able to produce revenue-weighted estimates, but we were not as confident with these estimates due to the fact that the sample design did not include any revenue-weighted stratification variables. It is now being proposed to use revenue as a stratification variable for the 2005 QAS to ensure that revenue-weighted estimates can be validly established.

Another emerging factor from the 2002 QAS was the idea that some businesses had changed their business activities and that these changes had not been updated on the BR, thus leading to an incorrect NAICS code. This idea, deemed volatility in this report, could explain nearly half of the estimated errors in coding from the 2002 QAS. It will need to be decided if this issue is important enough to focus on for the 2005 QAS. Doing so will most likely require a second phase, and will therefore lead to an increase in cost, time, and complexity of the analyses.

The analysis of the 2002 QAS was very extensive and time-consuming as many different issues were examined. With some pre-planning, the 2005 QAS can focus immediately on the central issues of importance and therefore be a more effective survey.

Table 3: Comparison of the 2002 and 2005 Quality Assurance Surveys

| Area of focus | 2002 QAS | 2005 QAS |
|---|---|---|
| Goal(s) | Determine the proportion of units on the BR that are consistently classified using NAICS | 1) Same as 2002 QAS<br>2) Provide revenue-related estimates<br>3) Potentially to examine volatility and quality control<br>4) Potentially to use results to make macro adjustments to other surveys |
| Frame | Fall 2002 UES SUF | Fall 2004 UES SUF |
| Sampling design | Stratified sample with some systematic sampling (Phase 1) and SRS (Phase 2) | Stratified sample with SRS and some census taking (i.e. take all strata) |
| Stratification variables | 1) Industrial sector (20)<br>2) Method of tax remittance (3) | 1) Industrial sector (20)<br>2) Revenue size (3 or 4) |
| Coders | Regional offices (Phase 1) and BR (Phase 1 and 2) | Coders with high levels of experience with NAICS (mainly BR) |
| Method of collection | Phase 1: Blind phone call<br>Phase 2: non-blind phone call | Same as 2002 QAS for Phase 1 and for Phase 2 (if one is needed) |
| Edit and Imputation | Respondents re-weighted to include unable to locate or unable to contact units (the non-respondants) | Same as 2002 QAS |
| Estimation | Sample design weights with some domain estimation for analytical | Same as 2002 QAS |

# 6 Resources

LATENDRESSE, EDITH (1) (2003). Stratégie d'échantillonnage. Internal document.

LATENDRESSE, EDITH (2) (2003). Correction et imputation des simples. Internal document.

MACFARLANE, ALISTAIR (2004). QA Survey 2002: Concluding report. Internal document.

QUIGLEY, ANDREW and MACFARLANE, ALISTAIR (2004). The 2002 Quality Assurance Survey: An analysis of the accuracy of NAICS coding on the business register. Presentation, Statistics Canada.
http://method/BiblioStat/Seminars/BSMD/2003-2004E/re040623.htm

# 7 Appendix I

Table A1: Distribution of the 2002 SUF by industrial sector and T1T2 flag

| Industrial Sector | Flag T1T2 | | | Total |
|---|---|---|---|---|
| | Others | T1 | T2 | |
| Agriculture | 13496 | 126345 | 56068 | 195909 |
| Mining and Oil | 96 | 1481 | 13840 | 15417 |
| Utilities | 55 | 219 | 1641 | 1915 |
| Construction | 8744 | 77528 | 155927 | 242199 |
| Manufacturing | 2289 | 19825 | 81996 | 104110 |
| Wholesale Trade | 1912 | 20473 | 99083 | 121468 |
| Retail Trade | 8046 | 53467 | 154547 | 216060 |
| Transportation | 4189 | 41180 | 59711 | 105080 |
| Information | 512 | 4412 | 24578 | 29502 |
| Finance and Insurance | 1154 | 4456 | 99190 | 104800 |
| Real Estate | 3638 | 50918 | 105132 | 159688 |
| Professional Services | 9114 | 96241 | 176157 | 281512 |
| Management | 147 | 1110 | 72661 | 73918 |
| Administrative and Support | 4324 | 32446 | 62366 | 99136 |
| Educational Services | 1162 | 5534 | 12511 | 19207 |
| Health Care | 10137 | 37747 | 44128 | 92012 |
| Arts | 2365 | 12848 | 24207 | 39420 |
| Accommodation and Food Services | 4987 | 25085 | 80316 | 110388 |
| Other Services | 15822 | 57624 | 90268 | 163714 |
| Public Administration | 764 | 29 | 6785 | 7578 |
| Total | 92953 | 668968 | 1421112 | 2183033 |

Table A2: Distribution of the 2002 QAS Phase 1 sample by industrial sector and T1T2 flag

| Industrial Sector | Flag T1T2 | | | Total |
|---|---|---|---|---|
| | Others | T1 | T2 | |
| Agriculture | 60 | 184 | 122 | 366 |
| Mining and Oil | 6 | 24 | 73 | 103 |
| Utilities | 12 | 24 | 65 | 101 |
| Construction | 50 | 148 | 210 | 408 |
| Manufacturing | 27 | 79 | 161 | 267 |
| Wholesale Trade | 25 | 82 | 181 | 288 |
| Retail Trade | 48 | 125 | 212 | 385 |
| Transportation | 34 | 106 | 128 | 268 |
| Information | 13 | 38 | 91 | 142 |
| Finance and Insurance | 22 | 43 | 203 | 268 |
| Real Estate | 33 | 122 | 176 | 331 |
| Professional Services | 51 | 165 | 223 | 439 |
| Management | | 25 | 200 | 225 |
| Administrative and Support | 35 | 95 | 132 | 262 |
| Educational Services | 18 | 39 | 58 | 115 |
| Health Care | 50 | 97 | 104 | 251 |
| Arts | 25 | 59 | 80 | 164 |
| Accommodation and Food Services | 38 | 85 | 152 | 275 |
| Other Services | 63 | 121 | 151 | 335 |
| Public Administration | 25 | | 75 | 100 |
| Total | 635 | 1661 | 2797 | 5093 |