Section 2

Description of the Sample

This section describes the sample design and selection, the method of estimation, the sampling variability of the estimates, and the methodology of computing confidence intervals.

Domain of Study

The statistics in this report are estimates from a probability sample of unaudited Individual Income Tax Returns, Forms 1040, 1040A, 1040EZ, and 1040PC (including electronic returns) filed by U.S. citizens and residents during Calendar Year 1999.

All returns processed during 1999 were subjected to sampling except tentative and amended returns. Tentative returns were not subjected to sampling because the revised returns may have been sampled later, while amended returns were excluded because the original returns had already been subjected to sampling.

A small percentage of returns were not identified as tentative or amended until after sampling. These returns, along with those that contained no income information, were excluded in calculating estimates.

This resulted in a small difference between the population total (125,037,636 returns) reported in Table C and the estimated total of all returns (124,770,662) reported in other tables.

The estimates in this report are intended to represent all returns filed for Tax Year 1998. While about 98 percent of the returns processed during Calendar Year 1999 were for Tax Year 1998, the remaining returns were mostly for prior years, and a few for non-calendar years ending during 1998 and 1999. Returns for prior years were used in place of 1998 returns expected to be received and processed after December 31, 1999.

This was done based on the assumption that the characteristics of returns due, but not yet processed, can best be represented by the returns for previous income years that were processed in 1999.

Sample Design and Selection

The sample design is a stratified probability sample, in which the population of tax returns is classified into subpopulations, called strata, and a sample is randomly selected independently from each stratum. Strata are defined by:

- 1. Nontaxable with adjusted gross income or expanded income of \$200,000 or more and no alternative minimum tax.
- 2. High combined business and farm total receipts of \$50,000,000 or more.
- 3. Presence or absence of special Forms or Schedules (Form 2555, Form 1116, Form 1040 Schedule C, and Form 1040 Schedule F).

Bonnye Walker and Valerie Puckett designed the sample and prepared the text and tables in this section under the direction of Yahia Ahmed, Chief, Mathematical Statistics Section, Statistical Computing Branch.

- 4. Indexed positive or negative income. Sixty variables are used to derive positive and negative incomes. These positive and negative income classes are deflated using the Chain-Type Price Index for the Gross Domestic Product to represent a base year of 1991. (See footnote 1 for details.)
- 5. Potential usefulness of the return for tax policy modeling. Thirty-two variables are used to determine how useful the return is for tax modeling purposes.

Table C shows the population and sample count for each stratum after collapsing some strata with the same sampling rates. (See references 1 and 2 for details.) The sampling rates range from 0.05 percent to 100 percent.

Tax data processed to the IRS Individual Master File at the Martinsburg Computing Center during Calendar Year 1999 were used to assign each taxpayer's record to the appropriate stratum and to determine whether or not the record should be included in the sample. Records are selected for the sample either if they possess certain combinations of the four ending digits of the social security number, or if their ending five digits of an eleven-digit number generated by a mathematical transformation of the SSN is less than or equal to the stratum sampling rate times 100,000. (See reference 3 for details.)

Data Capture and Cleaning

Data capture for the SOI sample begins with the designation of a sample of administrative records. While the sample was being selected, the process was continually monitored for sample selection and data collection errors. In addition, a small subsample of returns was selected and independently reviewed, analyzed, and processed for a quality evaluation.

The administrative data and controlling information for each record designated for this sample was loaded onto an online database at the Cincinnati Service Center. Computer data for the selected administrative records were then used to identify inconsistencies, questionable values, and missing values as well as any additional variables that an editor needed to extract for each record. The editors use a hardcopy of the taxpayer's return to enter the required information onto the online system.

After the completion of service center review, data were further validated, tested, and balanced at the Detroit Computing Center. Adjustments and imputations for selected fields based on prior year data and other available information were used to make each record internally consistent. Finally, prior to publication, all statistics and tables were reviewed for accuracy and reasonableness in light of provisions of the tax law, taxpayer reporting variations and limitations, economic conditions, and comparability with other statistical series.

Some returns designated for the sample were not available for SOI processing because other areas of IRS needed the return at the same time. For Tax Year 1998, 0.08 percent of the sample returns were unavailable.

Method of Estimation

Weights were obtained by dividing the population count of returns in a stratum by the number of sample returns for that stratum. The weights were adjusted to correct for misclassified returns. These weights were applied to the sample data to produce all of the estimates in this report.

Sampling Variability and Confidence Intervals

The sample used in this study is one of a large number of samples that could have been selected using the same sample design. The estimates calculated from these different samples would vary. The standard error (SE) of an estimate is a measure of the variation among the estimates from the possible samples and, thus, is a measure of the precision with which an estimate from a particular sample approximates the average of the estimates calculated from all possible samples.

The standard error may be expressed as a percentage of the value being estimated. This ratio is called the coefficient of variation (CV). Table 1.4 CV contains estimated CV's for the estimates included in Table 1.4 of this report.

The sample estimate and an estimate of its standard error permit the construction of interval estimates with

20

Description of the Sample

prescribed confidence that the interval includes the population value. If all possible samples were selected under essentially the same conditions and an estimate and its estimated standard error were calculated from each sample, then:

- 1. About 68 percent of the intervals from one standard error below the estimate to one standard error above the estimate would include the population value. This is a 68 percent confidence interval.
- 2. About 95 percent of the intervals from two standard errors below the estimate to two standard errors above the estimate would include the population value. This is a 95 percent confidence interval.

For example, from Table 1.4, the amount estimate for State Income Tax Refunds, X, is 14.708 billion, and its related coefficient of variation, CV(X), is 0.98 percent. The standard error of the estimate, SE(X), needed to construct the confidence interval estimate, is:

SE (X) = X • CV(X) = $(\$14.708 \times 10^{9}) • (0.0098)$ = \$0.144 billion

The p percent confidence interval is calculated using the formula:

 $X \pm z \bullet SE(X)$

where z takes the value 1, 2, or 3 when p is 68, 95, or 99, respectively. Based on these data, the 68 percent confidence interval is from \$14.564 billion to \$14.852 billion, and the 95 percent confidence interval is from \$14.420 billion to \$14.996 billion.

Table Presentation

Whenever a weighted frequency is less than 3, the estimate and its corresponding amount are combined or deleted in order to avoid disclosure of information for specific taxpayers. (The combined or deleted data, if any, are included in the corresponding column totals.)

These combinations and deletions are indicated by a double asterisk (**). Estimates based on less than 10 sampled returns are considered to be unreliable. These

estimates are noted by a single asterisk (*) to the left of the data unless all of the sampled returns are selected with certainty (at the 100 percent rate).

In the tables, a dash (- or --) in place of a frequency or an amount indicates that either no returns in the population had the characteristic or the characteristic was so rare that it did not appear on any of the sampled returns.

Footnote

 Indexing of positive and negative income is done by dividing each by the ratio of the Chain-Type Price Index for the Gross Domestic Product for the fourth quarter of 1997 to the fourth quarter of the base year of 1991. The indices can be found in U. S. Department of Commerce, Bureau of Economic Analysis, *Survey of Current Business* (November 1998) Vol. 78, number 11.

References

- [1] Hostetter, S., Czajka, J. L., Schirm, A. L., and O'Conor, K. (1990), "Choosing the Appropriate Income Classifier for Economic Tax Modeling," in *Proceedings of the Section on Survey Research Methods*, American Statistical Association, 419-424.
- [2] Schirm, A. L., and Czajka, J. L. (1991), "Alternative Designs for a Cross-Sectional Sample of Individual Tax Returns: the Old and the New," *Proceedings of the Section on Survey Research Methods*, American Statistical Association, 163-168.
- [3] Harte, J.M. (1986), "Some Mathematical and Statistical Aspects of the transformed Taxpayer Identification Number: A Sample Selection Tool Used at IRS," *Proceedings of the Section on Survey Research Methods*, American Statistical Association, 603-608.

SOURCE: IRS, Statistics of Income – 1998, Individual Income Tax Returns, Publication 1304, Revised 4-01

Table C.—Number of Individual Income Tax Returns in the Population and Sample by Sampling Strata for 1998

										Number of returns	
Description of the sample strata										Population counts	Sample counts
Grand total Form 1040 returns only with adjusted gross income or expanded income of \$200,000 and over, with no income tax after credits and no additional tax for tax preferences, total Form 1040 returns only with combined Schedule C (business or profession) total receipts of \$50,000,000 and over, total Other Returns, total										125,037,636 2,907 86 125,034,643	164,340 2,907 86 161,347
		Number of Returns by type of form attached									
		Form 1040. with Form 1116 or Form 2555		Form 1040, with Schedule C but without Form 1116 or Form 2555		Form 1040, with Schedule F but without Schedule C, Form 1116 or Form 2555		All other forms			
	Degree of	Population	Sample	Population	Sample	Population	Sample	Population	Sample		
Description of the sample strata	(1)	(4)	counts (5)	counts (6)	(7)	(8)	counts (9)	(10)	(11)		
Total		2,465,704	31,283	17,118,658	35,196	1,546,182	4,403	103,904,099	90,465		
Indexed Negative Income ⁴ \$10,000,000 or more \$5,000,000 under \$10,000,000 \$2,000,000 under \$5,000,000 \$1,000,000 under \$2,000,000 \$250,000 under \$20,000 \$120,000 under \$250,000 \$60,000 under \$120,000 Under \$60,000 Indexed Positive Income ⁴ Under \$30,000	Ali Ali Ali Ali Ali Ali Ali Ali Ali 1	78 96 671 1,444 2,905 	78 96 101 113 48 39 ** *	484 622 2,247 4,932 13,819 34,515 80,691 117,403 329,368	484 622 732 760 456 334 356 299 431	73 121 534 1,313 4,010 10,370 19,244 20,446 38,812	73 121 182 220 123 102 82 47 60	577 699 2,536 5,073 11,835 26,680 57,692 87,359 315,763 27,033,158	577 699 8200 833 404 247 265 228 407 13,427	1,212 1,538 5,649 11,989 31,108 74,470 157,627 225,208 683,943 27,033,158	1,212 1,538 1,835 1,926 1,031 722 703 574 898
Under \$30,000 Under \$30,000 \$30,000 under \$60,000 \$60,000 under \$60,000 \$60,000 under \$120,000 \$60,000 under \$120,000 \$120,000 under \$250,000 \$120,000 under \$250,000	2 3-4 1-2 3-4 1-3 4 1-3 4	74,599 251,336 167,242 307,440 389,111 322,834 227,463 284,199	44 272 89 349 181 330 362 810	1,870,503 3,531,202 1,718,727 3,365,480 1,817,425 2,201,565 436,583 1,018,129	970 3,645 827 3,611 967 2,258 639 2,925	116,952 185,113 194,278 286,784 230,539 182,835 106,446 69,351	55 200 91 299 112 172 161 180	29,186,440 6,939,555 20,134,368 5,517,314 9,414,789 2,182,254 1,399,085 883,073	14,582 7,322 9,956 6,023 4,698 2,207 2,009 2,551	31,248,494 10,907,206 22,214,615 9,477,018 11,851,864 4,889,488 2,169,577 2,254,752	15,651 11,439 10,963 10,282 5,958 4,967 3,171 6,466
\$250,000 under \$500,000 \$500,000 under \$1,000,000 \$1,000,000 under \$2,000,000 \$2,000,000 under \$5,000,000 \$5,000,000 under \$10,000,000 \$10,000,000 or more	AII AII AII AII AII AII	243,266 112,063 46,392 23,483 6,595 4,155	1,636 2,676 5,677 7,632 6,595 4,155	421,230 113,905 27,867 9,191 1,879 891	2,850 2,834 3,462 2,964 1,879 891	57,665 15,408 4,018 1,457 277 136	358 388 513 451 277 136	493,608 144,754 45,033 17,220 3,472 1,762	3,298 3,594 5,562 5,522 3,472 1,762	1,215,769 386,130 123,310 51,351 12,223 6,944	8,142 9,492 15,214 16,569 12,223 6,944

¹This population includes an estimated 266,974 returns that were excluded from other tables in this report because they contained no income information or represented amended or tentative returns identified after sampling.

² This population includes 167 Form 1040 returns that were misclassified because of bad data collected during revenue processing.

³ Each population member is assigned a degree of interest based on how useful it is for tax modeling purposes. Degree of interest ranges from one (1) to four (4), with a one being assigned to returns that are the least

interesting, and a four being assigned to those that are the most interesting. 'All' refers to income classes for which returns with all four degrees of interest are assigned.

⁴ Positive and Negative Income classes are divided by a Chain-Type Price Index for the Gross Domestic Product of 1.1403 to represent a base year of 1991.

** Data combined.