



Codability of industry and occupation information from cancer registry records: Differences by patient demographics, casefinding source, payor, and cancer type

Sharon R. Silver MS¹  | Rebecca J. Tsai PhD¹ | Cyllene R. Morris DVM² |
James M. Boiano MS¹  | Jun Ju MS¹ | Marilyn S. Scocozza MS² |
Geoffrey M. Calvert MD¹

¹ National Institute for Occupational Safety and Health, Division of Surveillance, Hazard Evaluations, and Field Studies, Cincinnati, Ohio

² California Cancer Reporting and Epidemiologic Surveillance Program, Institute for Population Health Improvement, University of California Davis Health, Sacramento, California

Correspondence

Sharon R. Silver, National Institute for Occupational Safety and Health, Division of Surveillance, Hazard Evaluations, and Field Studies, MS R-17, 1150 Tusculum Ave, Cincinnati, OH.
Email: ssilver@cdc.gov

Funding information

No funding

Introduction: Industry and occupation (I&O) information collected by cancer registries is useful for assessing associations among jobs and malignancies. However, systematic differences in I&O availability can bias findings.

Methods: Codability by patient demographics, payor, identifying (casefinding) source, and cancer site was assessed using I&O text from first primaries diagnosed 2011-2012 and reported to California Cancer Registry. I&O were coded to a U.S. Census code or classified as blank/inadequate/unknown, retired, or not working for pay.

Results: Industry was codable for 37% of cases; 50% had “unknown” and 9% “retired” instead of usual industry. Cases initially reported by hospitals, covered by preferred providers, or with known occupational etiology had highest codable industry; cases from private pathology laboratories, with Medicaid, or diagnosed in outpatient settings had least. Occupation results were similar.

Conclusions: Recording usual I&O for retirees and improving linkages for reporting entities without patient access would improve I&O codability and research validity.

KEYWORDS

cancer, coding, industry, occupation, registry, surveillance

1 | INTRODUCTION

Population-based cancer registries are invaluable resources for research into multiple aspects of cancer and have often been used for assessing cross-sectional and longitudinal patterns of cancer incidence. Registry data have also facilitated preliminary assessment of associations among patient risk factors, including occupational exposures, and health outcomes. For example, industry and occupation data from the California Cancer Registry (CCR) have been used to evaluate the risk of cancer among firefighters,¹ occupations at

elevated risk for leukemia subtypes² and lung cancer among construction workers³; and to assess differences in risk of acute myeloid leukemia by industry and occupation.⁴ Given the importance of occupational exposures as risk factors for cancer, the 1992 Cancer Registries Amendment Act (1992 Act) included industrial and occupational (I&O) history among data items required to be collected, if available, for reported incident cases.⁵

If I&O information is not consistently available and of sufficient specificity to allow assignment to a standardized set of I&O codes, the validity, precision, and generalizability of the results of surveillance and epidemiologic analyses can be affected. CCR requires that “every effort be made to record the I&O in which the patient works or worked,” with the information ideally referring to the usual or longest

held job,⁶ a mandate more stringent than the 1992 Act, which requires only collection of I&O where available from the same record.⁵ CCR sources of I&O information include admission and discharge summaries, face sheets, patient history, oncology consultation reports, and health and social history questionnaires the patient has completed. Still, I&O data are often missing from state registry records, or are recorded in such a way that precludes assignment of standardized codes that are comparable and therefore usable for analyses.⁷

A study of I&O availability (categorizing “retired” and “non-working” as available) in the New Hampshire State Cancer Registry found differences by demographic characteristics and by broad groupings of malignancy type and data source.⁸ The current project extends this area of research by focusing on codability of I&O text to Census 2010 I&O codes by malignancy type and by demographic, source reporting, and payor characteristics. The aims of the current analyses are to identify areas for improvement of I&O data collection and to identify types of analyses likely to be most affected by missing I&O data.

2 | MATERIALS AND METHODS

CCR, California's statewide cancer surveillance system, has collected information on all cancers (except non-melanoma skin cancer and carcinoma in-situ of the cervix) diagnosed in California residents since 1988. CCR is housed within the California Department of Public Health, which receives funding from the Centers for Disease Control and Prevention's National Program of Cancer Registries. A system of regional registries, also affiliated with the National Cancer Institute's Surveillance, Epidemiology and End Results (SEER) program, collects and submits cancer data to CCR. Information collected includes patient demographics, diagnosis, tumor characteristics, types of treatment, and follow-up reports. The majority of cases are reported by hospital sources, but non-hospital sources, such as physicians and private pathology laboratories, provide cases as well. The registry collects information on primary and secondary payors (the entities financially responsible for covering the patient's costs) from each facility involved in the diagnosis or treatment of the tumor. Reports from each facility are retained (rather than overwritten); for this research, only the payor from the last facility treating the patient was analyzed. Inmates, children, and cases with the Veterans Benefits Administration as payor were excluded.

The current study is part of a larger project to code I&O for cancer registry cases from six states and then evaluate the use of job-exposure matrices to assign exposures to these cases. The project was approved by the Institutional Review Boards of the state collecting the data and the organizations that analyzed the data.

Industry and occupation data from 257 020 first primary cancers reported to CCR in either 2011 or 2012 were processed by the NIOSH Industry and Occupation Computerized Coding System (NIOCCS) software.⁹ NIOCCS processes industry and occupation text data using a two-step process: where possible, records are automatically assigned standardized (United States [U.S.] Census) codes for industry and

occupation; records that are not automatically coded due to I&O text ambiguity, incomplete I&O data, or other problems are presented to a human coder for computer-assisted coding.

For this analysis, I&O text from diagnosis years 2011–2012 CCR records was run through both steps of NIOCCS. The resulting I&O codes were then merged back into the CCR data. For each case, industry and occupation were separately classified as codable (through autocoding or computer-assisted coding) if a specific U.S. Census code could be assigned (2007 industry codes, 2010 occupation codes). When I&O did not match a Census code, the case was assigned to one of three categories of responses that could not be given a specific code: 1) retired; 2) not working for pay; or 3) unknown/missing/blank/information inadequate to code (referred to as “blank or uncodable” herein). “Blank or uncodable” includes blank fields, responses that were illegible or too vague to code (such as “family business” or responses that were somewhat more specific but could not be assigned 2010 Census codes). “Not a paid worker” comprises students, volunteers, homemakers, and unemployed persons and is a legitimate Census code for occupation but not for industry; because this designation provides no information about work-related exposures, these responses were treated as not having codable industry or occupation. Prevalences of these codability classifications were calculated for demographic (age, sex, race/ethnicity) and reporting (casefinding source, payor) categories.

For each case in CCR, information, including I&O text, may come from one or multiple health entity sources. While evaluation of sources of I&O information would optimally account for all contributing sources, that information was not available for the current study. CCR's electronic record identifies the source that first identified the tumor (casefinding source). In addition, CCR includes a variable that summarizes the overall best source for abstracting information about the tumor (henceforth called “best source”). The casefinding source variable offers specificity, separating sources into nine hospital groupings, five non-hospital groupings, and one source originating later with the CCR (quality control review). “Hospital” here indicates that the reporting sequence began in a hospital. Some hospital sources, such as pathology department review, also include reporting from non-hospital entities (ie external pathology laboratories) that process and report results for specimens sent by hospitals. Two hospital case identification sources, daily discharge review and disease index review, do not themselves have I&O information but lead registrars to seek case information from records of other hospital departments. Compared to casefinding source, the categories for the “best source” variable are broader, with only seven total categories: hospital inpatient and managed care plans with comprehensive, unified medical records; radiation treatment centers and medical oncology centers; laboratory; private medical practitioner; nursing home, convalescent hospital, or hospice; autopsy only; and death certificate only.

The source of specific data fields, including I&O, is not available in the registry's summary record. Thus the I&O text in the CCR might come from the initial casefinding source, the overall best source, from another source, or from more than one source. In the absence of information specifying the source of I&O text, our primary analysis

examined codability by initial casefinding source. Then, for records with non-hospital casefinding sources, we determined how codability differed when best source was also considered, as the latter can reflect information sources beyond the initial casefinding source. All analyses were conducted using SAS version 9.3 (SAS Institute).

3 | RESULTS

Results for industry and occupation were similar (almost all categories for all metrics differed by <3%), so only data for industry are presented here. Slightly more than 37% of cases diagnosed in 2011-2012 were codable to a specific industry. Uncodable industry fell into several groupings. A relatively small group (4% of cases) was classified as “not a paid worker.” “Retired” was the designation for 9% of cases. Most of the remaining 50% of cases included entries that were blank or uncodable.

Differences in codability by sex or race/ethnicity were relatively small compared to differences by age grouping (Table 1). Males were somewhat more likely (40%) to have codable industry than females (34%). This discrepancy may reflect that 5% more females than males were not in the paid workforce. Industry was codable for similar percentages of black and white non-Hispanic subjects. However, whites were 4% more likely to have industry listed as “retired,” while

for blacks, inadequate information was more likely to preclude coding (about 5%). Codability was slightly lower for Hispanics than for non-Hispanic blacks and whites.

Larger coding differences were seen by age category. In every age group, at least 45% of records had missing or inadequate industry information, and more than half of workers age 65 and above were in this category. Only about one-quarter of subjects younger than 25 had codable industry; this age group comprised by far the largest percentage of subjects classified “not a paid worker” (27.4%). Industry was codable for nearly half of subjects aged 25-55, but the codable percentage declined in older age groups as the percentage classified as retired increased, with the decline accelerating in the 65-69 year old category (data not shown).

Most cancer cases (91.7%) were initially reported to the registry by hospital casefinding sources (Table 2). Hospital pathology department review was by far the largest source of casefinding, comprising 60% of all cases; industry was coded for 40% of records identified by this source. Of non-hospital sources, private pathology laboratories reported the largest percentage of all cases (4.6%), but provided codable industry information for fewer than 15% of those cases. Some non-hospital sources comprised less than one percent of reports, but provided either very good (death certificate follow-up, >60% coded) or very poor (physician as casefinding source, <20% coded) industry information.

TABLE 1 Codability of industry by sex, race, and age category

| Category | % of category (n) ^a | Industry codability (% by category) | | | |
|------------------------|--------------------------------|-------------------------------------|-------------------|--------------------|----------------|
| | | Retired | Not a paid worker | Blank or uncodable | Coded industry |
| Sex ^b | | | | | |
| Male | 47.6 (122 223) | 9.0 | 1.4 | 49.2 | 40.4 |
| Female | 52.4 (134 757) | 9.4 | 6.4 | 49.9 | 34.3 |
| Race | | | | | |
| Non-Hispanic white | 59.9 (153 991) | 10.1 | 3.4 | 47.4 | 39.1 |
| Non-Hispanic black | 6.3 (16 175) | 6.0 | 3.2 | 52.2 | 38.5 |
| Hispanic | 19.3 (49 696) | 8.0 | 6.0 | 51.9 | 34.2 |
| Asian/Pacific Islander | 11.5 (29 564) | 9.6 | 4.8 | 48.5 | 37.0 |
| American Indian | 0.5 (1217) | 6.2 | 5.3 | 52.3 | 36.2 |
| Other/unknown | 2.5 (6377) | 3.5 | 1.1 | 82.1 | 13.3 |
| Age category (years) | | | | | |
| <25 | 1.0 (2435) | 0.04 | 27.4 | 46.0 | 26.5 |
| 25-<35 | 2.8 (7175) | 0.06 | 6.4 | 45.9 | 47.6 |
| 35-<45 | 6.5 (16 763) | 0.04 | 6.1 | 45.6 | 48.3 |
| 45-<55 | 16.0 (41 061) | 0.3 | 4.8 | 46.6 | 48.3 |
| 55-<65 | 25.7 (65 998) | 2.5 | 4.0 | 48.9 | 44.6 |
| 65-<75 | 25.4 (65 331) | 15.2 | 2.6 | 50.2 | 32.0 |
| >=75 | 22.7 (58 257) | 20.4 | 3.1 | 53.6 | 22.9 |
| Total | 100 (257 020) | 9.2 | 4.0 | 49.6 | 37.2 |

^aPercentages may not sum to 100 due to rounding.

^bOther/transsexual/transgender not reported separately (because *n* < 50) but included in total.

TABLE 2 Codability of industry by case identification (casefinding) source

| Casefinding source ^a | % of all sources ^b (n) | Industry codability (% by category) | | | |
|---|-----------------------------------|-------------------------------------|-------------------|--------------------|----------------|
| | | Retired | Not a paid worker | Blank or uncodable | Coded industry |
| Hospital sources | 91.7 (235 736) | 9.4 | 4.1 | 48.2 | 38.4 |
| Reporting hospital, NOS | 19.9 (51 070) | 6.6 | 4.1 | 55.6 | 33.7 |
| Hospital pathology department review | 60.4 (155 146) | 9.2 | 4.1 | 46.5 | 40.2 |
| Daily discharge review | 0.25 (644) | 2.2 | 8.4 | 61.5 | 28.0 |
| Disease index review | 8.5 (21 883) | 16.8 | 4.5 | 44.4 | 34.2 |
| Radiation therapy department/center | 1.7 (4426) | 9.0 | 2.7 | 41.3 | 47.0 |
| Outpatient chemotherapy | 0.23 (590) | 6.3 | 5.9 | 38.1 | 49.7 |
| Diagnostic imaging/radiology | 0.43 (1097) | 14.0 | 3.8 | 36.1 | 46.0 |
| Tumor board | 0.17 (432) | 10.9 | 5.8 | 47.4 | 35.9 |
| Other hospital reporting source, including clinic | 0.17 (446) | 5.2 | 2.5 | 44.4 | 48.0 |
| Non-hospital sources | 7.6 (19 539) | 7.1 | 2.7 | 67.1 | 23.1 |
| Physician report | 0.28 (709) | 12.8 | 1.6 | 66.2 | 19.5 |
| Consultation-only or pathology-only report | 1.7 (4412) | 7.2 | 1.0 | 66.7 | 25.1 |
| Private pathology laboratory report | 4.6 (11 881) | 7.0 | 1.3 | 77.2 | 14.5 |
| Death certificate follow-back | 0.83 (2140) | 3.4 | 13.5 | 20.2 | 62.9 |
| Other non-hospital reporting source | 0.15 (397) | 20.3 | 8.7 | 31.8 | 39.2 |
| Quality control review | 0.33 (842) | 5.2 | 3.3 | 56.4 | 35.0 |
| Missing/invalid source Information | 0.33 (836) | 11.6 | 5.1 | 42.0 | 41.3 |

NOS = not otherwise specified.

^aSources with < 50 reported cases (laboratory reports, hospital rehabilitation service or clinic, nursing home initiated case, Coroner's Office records review, managed care organization or insurance records, out-of-state case sharing) are not shown separately but are included in italicized category totals.

^bPercentages may not sum to 100 due to rounding.

Industry was most frequently blank or uncodable when obtained from one of three non-hospital sources: private pathology laboratories, consultation only/pathology only reports, and physician reports. However, cross-tabulation of these sources with the "best source" variable showed marked differences between the three sources (Figure 1): when the "best source" was a hospital or managed care source, industry codability was higher than when private physicians or laboratories were designated as "best source."

Industry codability differed by payor as well. HMOs, PPOs, and private insurance under managed care plans collectively were the payment source for 51.4% of cases (Table 3). Industry was reported for almost half of these cases, with cases covered by PPOs having the best codability (data not shown). In contrast, industry was codable for only 17.6% of cases involving dual-eligibility (patients eligible for Medicare and Medicaid). The dual-eligibility group had blank or uncodable industry information for 53.9% of cases and had 20.7% of cases designated as "retired." When Medicaid was the provider or when the subject was uninsured, prevalence of codable industry was low (approximately 25% of reported cases) and prevalence of missing/inadequate industry information was high.

Industry codability for cancer sites comprising at least 257 cases (1% of malignancies diagnosed 2011-2012 and reported to CCR) is shown in Table 4; codability for cancers with at least 100 cases reported to the CCR can be found in online appendix A. Three malignancies had

industry coding rates above 50%: cancer of the tonsil (54.0%, appendix A); pleural cancer (57.1%, data not shown); and mesothelioma (54.0%, appendix A). Kaposi sarcoma was least likely to have codable industry (27.1%, appendix A) and most likely to have missing or inadequate industry (65.9%). Codability was at or below 33% (appendix A) for every cancer accessible on colonoscopy except cancers of the anus, rectum, and rectosigmoid junction. For leukemia, codability varied by subtype, with acute lymphocytic leukemia having the highest (45.1%, appendix A) and chronic lymphocytic leukemia the lowest (35.9%, appendix A).

Differences in coding rates by malignancy type were not independent of distributions of reporting sources and/or payors (data not shown). For example, malignant melanoma and other non-epithelial skin cancers were six times as likely as all other cancers combined to be diagnosed by private pathology labs and ten times as likely as all other cancers to be physician-initiated cases. Lung cancer, a common malignancy, had low codability, reflecting in part relatively large proportions of patients with "retired" recorded as industry and a higher than average percentage of cases paid by Medicare.

4 | DISCUSSION

Results of this study, which found overall codability of I&O data to be less than 40%, point to the need for improvement in the recording and

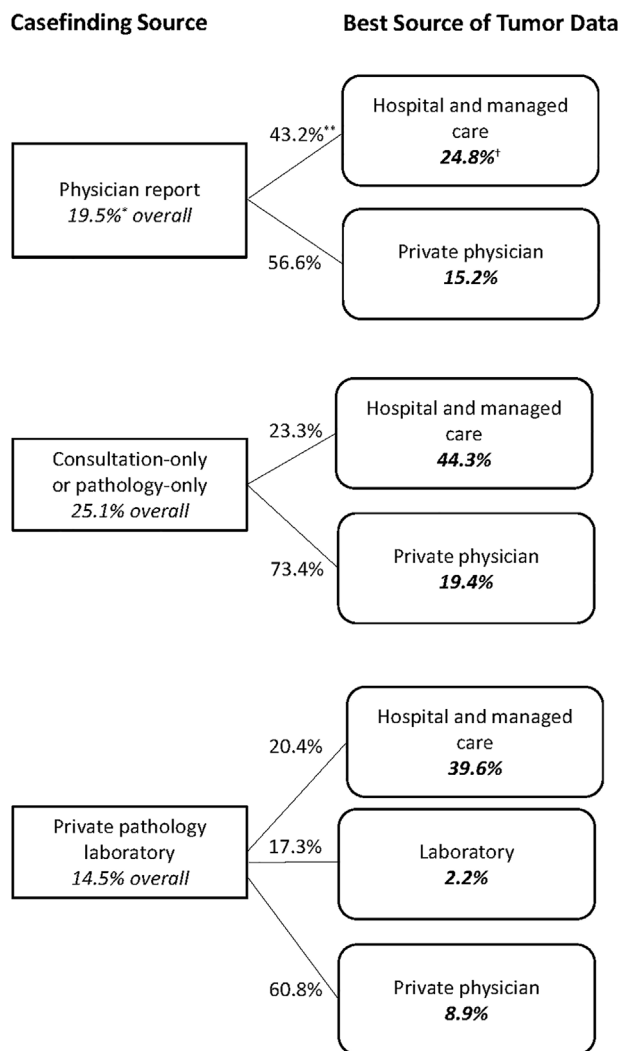


FIGURE 1 Supplementary sources can augment industry information for cases initially reported by entities with poor codability

collection of I&O data. In addition, the finding that I&O data availability varies by demographic, reporting, and outcome characteristics raises questions about the potential impact of missing I&O data on epidemiologic analyses of registry data. Research has shown that when data are not missing at random and missingness depends on both covariate values and disease status, biases may be introduced.^{10,11} The complete case approach, in which cases with incomplete data are excluded, has been shown to produce biased estimates with even 25% of data missing.¹² Multiple imputation, a common approach to missing data, is not feasible for filling in missing industry and occupation codes, with their large numbers of categories. No gold-standard population-based dataset that could be used to determine the actual distribution of cancer cases across industries and occupations exists. However, the current findings raise the possibility of selection bias on industry and occupation. If, for example, lung cancer incidence is 50% higher in industry A than in industry B, but industry information is 50% more likely to be missing for industry A (eg if most workers in industry A are

uninsured and most workers in industry B have private/HMO/PPO insurance), lung cancer incidence would appear to be equal in the two groups, rather than double for industry A.

Increasing I&O codability, and therefore its utility for public health research, will require enhanced efforts in eliciting, recording, abstracting, and coding I&O. While this task appears daunting, recording usual I&O instead of “retired” would increase codability by nearly 10% in the CCR, and potentially to a greater extent in states like Texas, where “retired” has been reported to comprise 15% of I&O listed for cases from the Texas Cancer Registry.¹³

Codability varies by casefinding source. The finding in this assessment that private pathology labs and physicians reporting to the CCR have particularly low I&O codability echoes results from a study of New York State Cancer Registry data in which private physician offices and laboratories (not further specified) had the highest percentages of unknown values for race and Hispanic ethnicity, as well as tumor staging information.¹⁴ The impact of these deficiencies is particularly strong for melanoma, a malignancy for which diagnosis frequently involves submission of specimens by a dermatologist to an outside pathologist or pathology laboratory.¹⁵ Electronic linkages between pathology labs and cancer registries are increasing,¹⁶ but these pathology laboratories do not generally have any contact with patients, so linkage to the provider submitting the sample would likely be necessary in order to obtain I&O data. The administrative burden of seeking approval for such linkages and ensuring that they are executed is likely to be substantial, although increased use of electronic health records should facilitate these efforts. Targeted encouragement of physicians to report I&O with melanoma diagnoses has been suggested.¹⁵ Such approaches (electronic linkages and educational efforts targeted to specific types of providers) are likely to be needed to encourage inclusion of I&O with reporting for melanoma and for other malignancies increasingly diagnosed in outpatient settings, such as colorectal malignancies found on colonoscopy. With this move to diagnosis in the outpatient setting, connecting to ambulatory care providers has become more important to ensure case completeness;¹⁷ the same approach will be needed to enhance reporting of I&O by these sources. Electronic linkages will only improve I&O availability if these fields are consistently included and collected in electronic health records.

In the interim, hospitals might be a good starting point for efforts to improve I&O reporting, as they comprise the large majority of reports, are relatively centralized, and have some departments (those involved with the diagnosis and treatment of malignancies) reporting I&O at higher frequencies. Some potential solutions to the lack of I&O availability include: information sharing within the hospital (including from departments such as registration, which are not linked directly to the medical records but collect job information for billing purposes); clearer instructions for eliciting I&O in registration systems and related paperwork; and training for providers who collect the information to ensure that they consistently elicit and report usual specific I&O. Specific barriers to reporting I&O in a 2005 study of Connecticut hospitals were lack of awareness of reporting requirements, lack of hospital reporting requirements, and insufficient time to report.¹⁸

TABLE 3 Codability of industry by payor

| Payor | % of all sources (n) ^a | Mean age at diagnosis | Industry codability (% by category) | | | |
|--|-----------------------------------|-----------------------|-------------------------------------|-------------------|--------------------|----------------|
| | | | Retired | Not a paid worker | Blank or uncodable | Coded industry |
| Not insured /not insured self-pay | 2.0 (5033) | 54.8 | 1.8 | 3.1 | 67.4 | 27.8 |
| Insured: HMO, PPO, private insurance, or NOS | 51.4 (132 354) | 58.6 | 5.1 | 3.4 | 44.6 | 46.9 |
| Medicaid or covered by county | 9.5 (24 467) | 53.5 | 2.8 | 7.7 | 63.8 | 25.8 |
| Medicare without supplemental coverage or NOS | 10.9 (27 937) | 73.4 | 17.6 | 4.0 | 52.9 | 25.5 |
| Medicare: with supplemental coverage or via managed care | 16.7 (42 880) | 74.4 | 18.7 | 3.0 | 46.5 | 31.8 |
| Medicare with medicaid eligibility | 4.9 (12 523) | 71.8 | 20.7 | 7.8 | 53.9 | 17.6 |
| Miscellaneous ^b | 4.6 (11 826) | 63.6 | 4.3 | 3.5 | 68.5 | 24.1 |

NOS = not otherwise specified.

^aPercentages may not sum to 100 due to rounding.

^bUnknown, military, or Indian Health Service.

While some of these barriers may have shifted with the advent of electronic reporting, the lack of a standardized requirement for inclusion of I&O information in the medical record is still perhaps the main challenge for the collection of meaningful I&O data.

Training of cancer registrars is also important. Cancer registrars may need to look at records from multiple sources to determine I&O, particularly when cases are initially identified by sources such as pathology labs that have no patient contact. Results of the current

TABLE 4 Industry codability by type of malignancy, descending order by codability

| Malignancy ^a | % of all cases (n) ^b | Average age at diagnosis | Industry codability (% by category) | | | |
|--|---------------------------------|--------------------------|-------------------------------------|-------------------|--------------------|----------------|
| | | | Retired | Not a paid worker | Blank or uncodable | Coded industry |
| Oral cavity and pharynx | 2.3 (5887) | 61.2 | 6.5 | 2.9 | 46.6 | 44.0 |
| Esophagus | 0.8 (2064) | 67.7 | 11.0 | 2.0 | 43.6 | 43.4 |
| Thyroid and other endocrine | 3.2 (8265) | 49.2 | 4.0 | 7.1 | 45.8 | 43.2 |
| Multiple myeloma | 1.3 (3390) | 66.5 | 9.9 | 3.4 | 43.9 | 42.8 |
| Breast | 18.8 (48 438) | 60.0 | 7.5 | 6.7 | 44.4 | 41.4 |
| Prostate | 13.3 (34 116) | 66.0 | 8.2 | 0.5 | 50.9 | 40.3 |
| Rectum | 2.1 (5299) | 61.9 | 9.2 | 3.6 | 47.2 | 40.0 |
| Leukemia | 2.2 (5648) | 62.4 | 8.8 | 4.5 | 47.0 | 40.0 |
| Non-hodgkin lymphoma | 3.9 (10 080) | 63.6 | 9.2 | 3.7 | 47.8 | 39.4 |
| Pancreas | 2.5 (6331) | 68.6 | 12.7 | 4.1 | 46.3 | 36.9 |
| Lung, bronchus | 8.6 (22 131) | 69.7 | 13.4 | 3.8 | 46.9 | 35.8 |
| Kidney and renal pelvis | 3.0 (7831) | 62.2 | 8.9 | 3.5 | 52.3 | 35.3 |
| Liver and intrahepatic bile duct | 2.0 (5218) | 64.1 | 8.8 | 4.4 | 51.8 | 35.0 |
| Ovary | 1.5 (3859) | 60.2 | 9.2 | 5.9 | 49.9 | 34.9 |
| Stomach | 1.8 (4508) | 65.8 | 13.2 | 4.1 | 48.1 | 34.7 |
| Corpus uteri, uterus unspecified | 3.4 (8864) | 61.0 | 8.8 | 5.5 | 51.1 | 34.6 |
| Cervix uteri | 1.0 (2561) | 50.4 | 3.4 | 7.9 | 55.1 | 33.7 |
| Colon and appendix | 6.0 (15 346) | 67.1 | 12.7 | 3.6 | 53.0 | 30.8 |
| Urinary bladder | 3.5 (8918) | 70.7 | 14.6 | 2.0 | 54.3 | 29.1 |
| Melanoma of the skin/other non-epithelial skin cancers | 7.4 (19 104) | 61.9 | 7.1 | 1.9 | 64.5 | 26.6 |

^aLimited to sites comprising at least 1% of all malignancies reported to California Cancer Registry, 2011- 2012.

^bPercentages may not sum to 100 due to rounding.

study show that such cases have markedly greater I&O codability when hospital records are also available and accessed for case information. A New Hampshire study showed a decrease in the number of records judged to have no I&O data from 74% to 14% after detailed records review followed by targeted registrar training in I/O collection.⁸ The percentage of records judged to have complete, codable I&O data was only 48% even after review training (in part because 20% of records belonged to individuals not in the paid workforce), but this was nearly triple the level before these additional steps. However, the New Hampshire study noted that the location of I&O data within a medical record varied by type of facility and record system and that centralized training efforts might be more difficult to implement in larger states with more and varied healthcare facilities. This finding points to the advantage of standardizing where and how I&O are recorded across health record systems. If registration and billing information is made available to registrars in the future, expansion of registrar trainings to include abstraction of this information would be beneficial.

In the current study, records with death certificate follow-back as the casefinding source had the highest I&O codability (though only 0.83% of cases had this source). However, the mixing of I&O data from death certificates with other sources of these data can be problematic for epidemiologic analysis. Death certificate data is, of necessity, obtained by proxy report (ie next-of-kin). A study comparing occupation from death certificates to occupation self-reported at midlife¹⁹ found that while agreement for broad occupational categories was reasonable (67%), agreement for job titles was poor (32%). More importantly, death certificate data are only available for deceased cases, so their use would introduce additional potential for bias.

Pending efforts to improve collection of codable data, awareness of which study populations and outcomes are most affected by non-codability is important for planning epidemiologic studies. Non-codable data fell into several categories, each impacting certain malignancies more strongly. Recording "retired" instead of usual industry generally had greater impact on malignancies occurring more frequently in older adults, such as bladder cancer. In contrast, young adults may be students or may not yet have established a "usual" occupation, leading to higher levels of categorization as "not a paid worker" for malignancies like testicular cancer that are more commonly diagnosed in younger adults. Low codability due to blank or uncodable I&O fields was associated with lack of insurance or with public coverage (Medicaid or county coverage); the high level of missing I&O among patients with Medicaid and county coverage could be associated with the reporting source (differences in information collection or reporting systems) and/or with characteristics of the covered case population (lower employment stability reducing the likelihood of a specific "usual" occupation and increasing likelihood of being unemployed at time of diagnosis; as well as disincentives to reporting employment). Designing studies to compare groups that are similar with respect to factors such as insurance coverage should decrease the potential for selection bias.

Several limitations pertain to this study. Collection of "usual" I&O is preferred to current I&O for research purposes (and a full work

history would be better still), but the prevalence of "retired" in I&O fields suggests that directions on healthcare intake forms may not request "usual" I&O and that there may be suboptimal probing by providers, who have many competing demands. In addition, information collection forms may simply elicit "occupation" or "industry" or even "job" rather than specifying usual or longest-held I&O. Correlations between current and usual I&O have been found to be good for high-level I&O groupings, but the concordance decreased as more detailed I&O groupings were used.²⁰ Recommendations and/or incentives from healthcare accreditation organizations such as The Joint Commission, or by the American College of Surgeons Cancer Programs could increase awareness of the value of I&O data for cancer research and improve the regular collection of this information by healthcare providers.

Despite these limitations, cancer registry I&O data present the potential for meaningful assessment of associations between different types of work and malignancies. Findings of an examination of cancer outcomes among firefighters using CCR data¹ were generally consistent with those of a cohort study of firefighters.²¹ While traditional cohort and case-control studies usually have access to more detailed work history information, primary advantages of analyzing registry data are the much lower cost and larger study populations.

Enhancing the utility of cancer registries for public health research by increasing I&O reporting will involve prioritizing data collection at the source, training of providers and registrars, and development of feedback loops to ensure continuous improvement. In the interim, consideration of these limitations while planning epidemiologic analyses of registry data is important to limit the potential for bias in the results.

AUTHORS' CONTRIBUTIONS

Ms Silver planned and implemented the analyses and wrote the manuscript. Dr Tsai performed data management and provided extensive comments and revisions for the manuscript. Dr Morris and Ms Scocozza provided key information about the procedures of the California Cancer Registry and provided revisions to the manuscript. Mr Boiano provided an industrial hygiene perspective during project development and provided comments and revisions for the manuscript. Ms Ju led programming and data management for the project and provided review of the manuscript. Dr Calvert planned and secured funding for the National Occupational Research Agenda project under which the data were received and provided extensive comments and revisions for the manuscript. All authors have approved this version of the manuscript for submission to AJIM and agree to be accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved.

ACKNOWLEDGMENTS

The authors would like to thank Karla Armenti, MaryBeth Freeman, and Marie Haring Sweeney for the valuable comments on the

manuscript; Pam Schumacher for managing industry and occupation (I&O) coding for the project; Matthew Hirst for QA of I&O coding; Sue Burton, Ashley Fite, and Surprese Watts for I&O coding; and the California Cancer Registry personnel for providing the data used in these analyses. The collection of cancer incidence data used in this study was supported by California Department of Public Health as part of the statewide cancer reporting program mandated by the California Health and Safety Code Section 103885; the National Cancer Institute's Surveillance, Epidemiology and End Results Program under contracts awarded to the Cancer Prevention Institute of California, the University of Southern California, and the Public Health Institute; and the Centers for Disease Control and Prevention's National Program of Cancer Registries, under agreement awarded to the California Department of Public Health.

FUNDING

The authors report that there was no funding source for the work that resulted in the article or the preparation of the article. The article is based on previously collected data; funding sources for that data collection are listed in the acknowledgements section.

ETHICS APPROVAL AND INFORMED CONSENT

The project was approved by the NIOSH Institutional Review Board (IRB) and was approved by the California Committee for the Protection of Human Subjects. Per the IRBs, informed consent was not needed as these were analyses on existing data with no patient contact or follow-back.

DISCLOSURE (AUTHORS)

The authors declare no conflicts of interest.

DISCLOSURE BY AJIM EDITOR OF RECORD

Steven B. Markowitz declares that he has no conflict of interest in the review and publication decision regarding this article.

DISCLAIMER

The findings and conclusions presented in this article are those of the authors and do not necessarily represent the views of the National Institute for Occupational Safety and Health/Centers for Disease Control and Prevention; the State of California, Department of Public Health; the National Cancer Institute; or their Contractors and Subcontractors.

ORCID

Sharon R. Silver  <http://orcid.org/0000-0002-7679-5028>

James M. Boiano  <http://orcid.org/0000-0003-2738-4588>

REFERENCES

1. Tsai RJ, Luckhaupt SE, Schumacher P, Cress RD, Deapen DM, Calvert GM. Risk of cancer among firefighters in California, 1988–2007. *Am J Ind Med*. 2015;58:715–729.
2. Luckhaupt SE, Deapen D, Cress R, Schumacher P, Shen R, Calvert GM. Leukemia among male construction workers in California, 1988–2007. *Leuk Lymphoma*. 2012;53:2228–2236.
3. Calvert GM, Luckhaupt S, Lee SJ, et al. Lung cancer risk among construction workers in California, 1988–2007. *Am J Ind Med*. 2012; 55:412–422.
4. Tsai RJ, Luckhaupt SE, Schumacher P, Cress RD, Deapen DM, Calvert GM. Acute myeloid leukemia risk by industry and occupation. *Leuk Lymphoma*. 2014;55:2584–2591.
5. Cancer Registries Amendment Act, 3372 (Public Law 102–515).
6. Kairon RS. California cancer reporting system standards, Volume I: Abstracting and coding procedures. 2017; http://www.ccrca.org/paqa-pubs/v1_2017_online_manual/index.htm. Accessed 5/26/2017, 2017.
7. Freeman MB, Pollack LA, Rees JR, et al. Capture and coding of industry and occupation measures: findings from eight national program of cancer registries states. *Am J Ind Med*. 2017;60:689–695.
8. Armenti KR, Celaya MO, Cherala S, Riddle B, Schumacher PK, Rees JR. Improving the quality of industry and occupation data at a central cancer registry. *Am J Ind Med*. 2010;53:995–1001.
9. Schmitz M, Forst L. Industry and occupation in the electronic health record: an investigation of the national institute for occupational safety and health industry and occupation computerized coding system. *JMIR Med Inform*. 2016;4:e5.
10. Vach W, Blettner M. 2005. *Missing Data in Epidemiologic Studies*. *Encyclopedia of Biostatistics*. 5. <https://doi.org/10.1002/0470011815.b2a03085>
11. Heitjan DF. Incomplete data: what you don't know might hurt you. *Cancer Epidemiol Biomarkers Prev*. 2011;20:1567–1570.
12. Marshall A, Altman DG, Holder RL. Comparison of imputation methods for handling missing covariate data when fitting a Cox proportional hazards model: a resampling study. *BMC Med Res Methodol*. 2010;10:112.
13. Weiss NS, Cooper SP, Socias C, Weiss RA, Chen VW. Coding of central cancer registry industry and occupation information: the Texas and Louisiana experiences. *J Registry Manag*. 2015;42:103–110.
14. Kahn AR, Schymura MJ, Juster TM. Completeness of source-level data items based on type of source. NAACCR Annual Conference. St. Louis, MO 2016.
15. Cockburn M, Swetter SM, Peng D, Keegan TH, Deapen D, Clarke CA. Melanoma underreporting: why does it happen, how big is the problem, and how do we fix it? *J Am Acad Dermatol*. 2008;59: 1081–1085.
16. Gal TS, Durbin EB. Successes and challenges in population-based electronic pathology reporting. NAACCR Annual Conference. Quebec City, QC Canada 2010.
17. Ayers CM, Celaya MO, Rees JR. Increasing non-hospital reporting: the NH experience. North American Association of Central Cancer Registries (NAACCR) Annual Conference. St. Louis, MO 2016.
18. Polednak AP. Obtaining occupation as an indicator of patients' socioeconomic status in a population-based cancer registry. *J Registry Manag*. 2005;32:176–181.
19. Bidulescu A, Rose KM, Wolf SH, Rosamond WD. Occupation recorded on certificates of death compared with self-report: the atherosclerosis risk in communities (ARIC) study. *BMC Public Health*. 2007;7:229.
20. Luckhaupt SE, Cohen MA, Calvert GM. Concordance between current job and usual job in occupational and industry groupings: assessment of the 2010 national health interview survey. *J Occup Environ Med*. 2013;55:1074–1090.

21. Daniels RD, Kubale TL, Yiin JH, et al. Mortality and cancer incidence in a pooled cohort of US firefighters from San Francisco, Chicago and Philadelphia (1950–2009). *Occup Environ Med*. 2014;71:388–397.

SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

How to cite this article: Silver SR, Tsai RJ, Morris CR, et al. Codability of industry and occupation information from cancer registry records: Differences by patient demographics, casefinding source, payor, and cancer type. *Am J Ind Med*. 2018;61:524–532. <https://doi.org/10.1002/ajim.22840>