

1 **An Ontology-based Approach to Guide and Document Variable and**  
2 **Data Source Selection and Data Integration Process to Support**  
3 **Integrative Data Analysis in Cancer Outcomes Research**

4  
5 Hansi Zhang<sup>1</sup>, Yi Guo<sup>1</sup>, Jiang Bian<sup>1\*</sup>

6  
7 <sup>1</sup> Department of Health Outcomes and Biomedical Informatics, College of Medicine, University  
8 of Florida, Gainesville, Florida, USA

9  
10 Email: Hansi Zhang – [hansi.zhang@ufl.edu](mailto:hansi.zhang@ufl.edu); Yi Guo – [yiguo@ufl.edu](mailto:yiguo@ufl.edu); Jiang Bian –  
11 [bianjiang@ufl.edu](mailto:bianjiang@ufl.edu)

12  
13 \* Corresponding author: Jiang Bian, [bianjiang@ufl.edu](mailto:bianjiang@ufl.edu).

14  
15 Jiang Bian

16 Department of Health Outcomes and Biomedical Informatics, College of Medicine, University of  
17 Florida, 2197 Mowry Road, Suite 122, PO Box 100177, Gainesville, FL 32610-0177

18  
19  
20  
21  
22

## 23 **Abstract**

### 24 **Background**

25 To reduce cancer mortality and improve cancer outcomes, it is critical to understand the various  
26 cancer risk factors (RFs) across different domains (e.g., genetic, environmental, and behavioral  
27 risk factors) and levels (e.g., individual, interpersonal, and community levels). However, prior  
28 research on RFs of cancer outcomes, has primarily focused on individual level RFs due to the  
29 lack of integrated datasets that contain multi-level, multi-domain RFs. Further, the lack of a  
30 consensus and proper guidance on systematically identify RFs also increase the difficulty of RF  
31 selection from heterogenous data sources in a multi-level integrative data analysis (mIDA) study.  
32 More importantly, as mIDA studies require integrating heterogenous data sources, the data  
33 integration processes in the limited number of existing mIDA studies are inconsistently  
34 performed and poorly documented, and thus threatening transparency and reproducibility.

### 35 **Methods**

36 Informed by the National Institute on Minority Health and Health Disparities (NIMHD) research  
37 framework, we (1) reviewed existing reporting guidelines from the Enhancing the QUALity and  
38 Transparency Of health Research (EQUATOR) network and (2) developed a theory-driven  
39 reporting guideline to guide the RF variable selection, data source selection, and data integration  
40 process. Then, we developed an ontology to standardize the documentation of the RF selection  
41 and data integration process in mIDA studies.

### 42 **Results**

43 We summarized the review results and created a reporting guideline—ATTEST—for reporting  
44 the variable selection and data source selection and integration process. We provided an  
45 ATTEST check list to help researchers to annotate and clearly document each step of their mIDA

46 studies to ensure the transparency and reproducibility. We used the ATTEST to report two  
47 mIDA case studies and further transformed annotation results into semantic triples, so that the  
48 relationships among variables, data sources and integration processes are explicitly standardized  
49 and modeled using the classes and properties from OD-ATTEST.

## 50 **Conclusion**

51 Our ontology-based reporting guideline solves some key challenges in current mIDA studies for  
52 cancer outcomes research, through providing (1) a theory-driven guidance for multi-level and  
53 multi-domain RF variable and data source selection; and (2) a standardized documentation of the  
54 data selection and integration processes powered by an ontology, thus a way to enable sharing of  
55 mIDA study reports among researchers.

56

57

58 **Keywords:** Ontology, Integrative data analysis, Cancer outcomes research, Reporting guideline

59

60

61

62

63

64

65

66

67

## 68 **BACKGROUND**

69 Cancer is a major disease burden worldwide [1]. As the 2nd leading cause of death in the United  
70 States (US), about 1 in 4 deaths is due to various types of cancer [2]. In 2019, an estimation of  
71 1,762,450 new cancer cases diagnosed and 606,880 cancer deaths is reported by the American  
72 Cancer Society (ACS) in US [2]. The lifetime probabilities of being diagnosed with cancer are  
73 39.3% and 37.4% for male and female, respectively [3]. However, the risk factors for these high  
74 cancer incidence and mortality rates are still not fully understood.

75  
76 Nevertheless, to reduce cancer mortality rates and improve cancer outcomes (e.g., survival and  
77 prognosis), it is critical to understand the various risk factors of cancer. So far, evidence  
78 suggests that it is the interaction among many risk factors (RFs) together that affect the risk of  
79 cancer and cancer outcomes, rather than a single cause [4]. Further, the RFs involved are across  
80 different domains (e.g., genetic, environmental, and behavioral risk factors) and levels (e.g.,  
81 individual level, interpersonal level, and community level). However, there is not yet an  
82 agreement among the cancer research community regarding how these multi-level cancer RFs  
83 interact with each other. To do so, the first and most crucial step is to gain a comprehensive  
84 view of potential multi-level RFs associated with various cancer outcomes such as the stage of  
85 diagnosis (the most important prognostic factor) and survival.

86  
87 We surveyed existing research on RFs for late stage cancer diagnosis and poor survival, we  
88 found current studies about RFs for cancer outcomes are mostly from single-level analyses with  
89 mostly individual patient-level data. For instance, Andrew *et al.* assessed individual patient  
90 characteristics (e.g., age, gender, family history), and lifestyle factors (e.g., education, insurance

91 and socioeconomic status) to study their risks associated with colorectal cancer at late stage [5].  
92 These individual-level RFs have also been reported for other major types of cancers such as  
93 breast and cervical cancers [6–9]. Further, prior studies studying cancer RFs often only analyzed  
94 data from a single source, such as SEER [10], SEER-Medicare[11], or a state or hospital cancer  
95 registry [12]. Among these cancer risk factor studies, the complex interplay between difference  
96 levels RFs are often ignored (e.g., county-level smoking rate vs. individual smoking behavior).  
97 These single-level RF analyses (1) lead to biased effect estimates of RFs due to potential  
98 confounding from omitted factors, (2) omit critical cross-level RF interactions, such as race by  
99 residence, that could inform multi-level intervention design.

100

101 Nowadays, advances in technology created new ways for us to determine and measure disease  
102 risk factors across different levels (e.g., from advancements in genome sequencing for genetic  
103 markers to better sensors for producing more accurate estimates of environmental pollutants).  
104 The availability of such abundant data online in electronic formats enables researchers to pool  
105 data on an unprecedented scale and offers a great opportunity to do a thorough examination of  
106 multi-level RFs in a multi-level integrative data analysis (mIDA) so that confounding effects and  
107 across-level interactions can be studied. However, researchers face significant barriers to do so,  
108 especially because there is a lack of consensus and proper guidance to help researchers  
109 systematically think and discovery these variables from heterogenous sources. In 2017, National  
110 Institute on Minority Health and Health Disparities (NIMHD) of the National Institute of Health  
111 (NIH) proposed a Research Framework [13], an extension to the well-known social ecological  
112 model [14], to help investigators systematically study health disparities. Recognized by the  
113 NIMHD Framework, individuals are embedded within the larger social system and constrained

114 by the physical environment they live in. Within this framework, cancer outcomes are  
115 influenced by RFs from different levels (i.e., individual, interpersonal, community, and societal)  
116 and multiple domains (i.e., biological, behavioral, physical/built environment, sociocultural  
117 environment, and healthcare system). In this work, we adopted the NIMHD framework as the  
118 guiding theory for risk factor discovery and data source selection.

119

120 Further, mIDA for cancer outcomes research requires the integration of data from multiple  
121 sources. However, data integration processes in the very limited number of existing mIDA  
122 studies [15, 16] are inconsistently performed and poorly documented, and thus threatening  
123 transparency and reproducibility [17, 18]. The data integration processes are often time  
124 summarized in one or two sentences without explicitly documentation of the steps. For example,  
125 Guo *et al.* explored the impact of the relationships among socioeconomic status, individual  
126 smoking status, and community-level smoking rate on pharyngeal cancer survival [16]. The  
127 multi-level risk factors above were obtained and integrated from three different data sources (i.e.,  
128 Florida Cancer Data System [FCDS], U.S. Census, and Behavioral Risk Factor Surveillance  
129 System [BRFSS]) as mentioned in the abstract. However, for the rest of the paper, there is no  
130 description of how the individual-level records from FCDS are linked with county-level smoking  
131 rate from BRFSS and census tract-level poverty rate from U.S. Census. Even though the  
132 integration process might be as simple as integrating these multi-level variables through the  
133 geographic code (e.g., county code), it still needs to be standardized and explicitly documented  
134 to avoid ambiguity. For example, the paper discussed that “*regional smoking was measured as*  
135 *the average percentage of adult current smokers at the county level between 1996 and 2010*” and  
136 the readers might be able to make an educated guess that the regional smoking rates were more

137 likely to be generated using the BRFSS data rather than from the FCDS data; however, explicit  
138 documentation is needed as both BRFSS and FCDS data have individual smoking status.  
139 Keegan *et al.* explored and whether breast cancer survival patterns are influenced by factors such  
140 as nativity (individual level) and neighborhood socioeconomic status (community level).  
141 Similarly, they summarized integration process in one sentence by stating each patient was  
142 assigned a neighborhood socioeconomic status variable based the census block groups.  
143 However, the details such as variable names in each data sources, or whether the original  
144 geographic variables require pre-processing (e.g., derive census tract from zip codes) are not  
145 clearly documented [15]. The explicit documentation of these variable selection and data  
146 integration processes will help readers to better understand the study results, benefit other  
147 researchers who want to replicate the studies, but also more importantly, make it possible for  
148 machines to understand and replicate the steps (when these explicit documentations are encoded  
149 in a computable format such as with an ontology).

150

151 Further, even though these mIDA studies above did not emphasize the need for data integration  
152 or integrated datasets, the fact that they can only investigated a handful of variables at a time  
153 indicated the lack of but needed support on data integration. Even in studies on building  
154 frameworks or platforms to support or automate the data integration process (especially those  
155 related to creating integrated dataset to support cancer research), they often ignored the need for  
156 documenting the integration steps to guarantee the transparency and reproducibility of their  
157 approaches. For example, semantic data integration approach —connecting variables across  
158 different databases at the semantic level through mapping them to standardized concepts in a  
159 global schema (e.g., often time a global ontology) — has been proposed in data integration

160 studies in recent years to support generating integrated datasets for cancer research [19–21].  
161 However, none of these studies mentioned the need for standardizing and documenting their  
162 integration steps, for example, most of them did not even discuss the rationale for selecting the  
163 specific data sources to integrate. Nevertheless, when reporting mIDA studies, it is critical to  
164 document the steps that were followed to select, integrate, and process the data so that others can  
165 repeat the same steps and reproduce the findings.

166  
167 To address challenges above, in this paper, we first developed a reporting guideline to guide and  
168 document the RF variable selection, data source selection, and data integration process. The  
169 guideline is informed by (1) the NIMHD research framework that provides guidance and  
170 promotes structural thinking on identifying multi-level cancer RFs; and (2) reviewing existing  
171 reporting guidelines from the Enhancing the QUALity and Transparency Of health Research  
172 (EQUATOR) network [22]. Then, we proposed an ontology-based approach to annotate and  
173 document the RF selection and data integration process in mIDA studies based on the reporting  
174 guideline we developed. To do so, we developed the **O**ntology for the **D**ocumentation of  
175 **v**ariable selec**T**ion and **da**Ta sour**E** **S**election and in**T**egration proess (OD-ATTEST) so that the  
176 RF selection and data integration report can be (1) explicitly modeled with a shared, controlled  
177 vocabulary, (2) understandable to humans and computable to computers, and (3) adaptive to  
178 changes when the reporting process is refined.

179  
180 In our prior work [23], we proposed a preliminary reporting guideline for RF variable and data  
181 source selection based on our own experience of pooling multi-level RFs from different data  
182 sources to support mIDAs of cancer survival [24, 25]. In this extended journal paper, we



183 significantly expanded our ontology-based reporting guideline—ATTEST (vAriable selecTion  
184 and daTa sourcE Selection and inTegration):

- 185 • We conducted a systematic search of existing reporting guidelines from the EQUATOR  
186 network to extract reporting elements relevant to variable selection and data integration.
- 187 • We updated our reporting guideline based on the result of the systematic review to  
188 include new items regarding data integration (e.g., data processing, data integration  
189 strategy, data validation, etc.) as well as variable and data source selection.
- 190 • We completed building the OD-ATTEST following the best practice in ontology  
191 development to provide a formal presentation for the reporting guideline with  
192 standardized and controlled vocabularies.
- 193 • We provided an ontology (OD-ATTEST) annotated report generated based on a prior  
194 mIDA study to represent the annotated items and their relationships in reporting  
195 guideline.

196

## 197 **METHODS**

### 198 **Development of a reporting guideline for risk factor selection, data source selection, and** 199 **data integration**

200 To develop the reporting guideline, we started with summarizing our previous studies where we  
201 assessed the effect of data integration on predictive ability of cancer survival models [24] and  
202 created a semantic data integration framework to pool multi-level RFs from heterogenous data  
203 sources to support mIDA [25]. In the above studies, we went through the process of RF  
204 selection, data source selection, and data integration. To be able to ensure the reproducibility of  
205 these studies, a number of middle steps need to be documented as detailed in our previous paper

206 [23]. For example, both rural-urban commuting area (RUCA) codes [26] and the National  
207 Center for Health Statistics (NCHS) urban-rural classification scheme [27] are often used to  
208 represent an geographic area's rurality status. The difference between the two resides in the  
209 classification granularity, where RUCA focuses on classifying U.S. census tracts (i.e., tens levels  
210 from rural to metropolitan) while the NCHS urban-rural classification scheme focuses on  
211 classifying U.S. counties (i.e., a hierarchal definition with six levels). Thus, we need to clearly  
212 document which rural definition we used in the data analysis since different representations of  
213 the same variable (i.e., rurality in this case) have different impacts on model results, as shown in  
214 our prior work [24]. Further, before integration RFs from various data sources at different levels  
215 (e.g., census tract level vs. county level) and covered different time periods, we made an  
216 assumption assumed that area-level characteristics (e.g., social vulnerability index) derived from  
217 2000 U.S. Census data were applicable across different time periods (as our individual level data  
218 from FCDS covered 1996 and 2010). Above experiences suggest that we must document these  
219 data integration nuances so that other researchers can repeat our data integration and data  
220 processing pipeline and reproduce the same results (e.g., integrated dataset). In sum, we  
221 summarized 3 key items that need to be documented: (1) RF selection (e.g., individual vs county-  
222 level variables), (2) data source selection (e.g., individual-level data from FCDS and contextual-  
223 level data from US Census), and (3) data integration and data preprocessing strategies.

224

225 Through discussions with expert biostatisticians, data analysts, and cancer outcomes researchers,  
226 we summarized the typical mIDA process and found there is little structured thinking when  
227 investigators selecting and identifying risk factors and their data sources. We thus propose to use  
228 the NIMHD research framework to provide a theory-driven guidance for multi-level and multi-

229 domain RF and data source selections. The NIMHD framework is originally designed to depicts  
230 a wide range of health determinants (i.e., RFs from different levels and domains) relevant to  
231 understanding and addressing minority health and health disparities. The goal of using the  
232 NIMHD framework is to help investigators to structurally and comprehensively think and  
233 identify relevant RFs and corresponding data sources in their IDA studies.

234  
235 To build upon existing established reporting guidelines, we searched and identified relevant  
236 reporting guidelines from the Enhancing the QUALity and Transparency Of health Research  
237 (EQUATOR) network—a comprehensive searchable database of guidelines for health research  
238 reporting. The EQUATOR network categorizes health researches into 13 study types (e.g.,  
239 quantitative studies, experimental studies, and observational studies), where reporting guidelines  
240 for observational studies are most relevant to our mIDA use case. To further identify relevant  
241 reporting guidelines in EQUATOR, we developed a set of screening criteria to determine  
242 whether a reporting guideline in EQUATOR contains the information that can be used to  
243 improve our ATTEST reporting guideline as shown below:

- 244 • The reporting guideline is designed for secondary data analysis studies.
- 245 • The reporting guideline contains at least one of the following sections: data, outcomes  
246 (variables), and methods, as these sections will contain information related to variable  
247 selection, data source selection, and data integration methods.
- 248 • The reported data within the guideline must be health related.
- 249 • The use of the guideline (at least part of the guideline) can be extended to the cancer  
250 outcomes research, especially those related to variable selection, data source  
251 selection, and data integration.

252 We reviewed all reporting guidelines designed for observational studies and eliminated  
253 guidelines that do not involve the tasks of RF and data source selection and integration. We then  
254 identified all reporting guidelines that contain the following sections: data, outcomes (variables),  
255 and methods. For those that do not have sections clearly marked, we manually reviewed the  
256 entire reporting guideline to identify whether they discussed one of the three aspects. We then  
257 extracted reporting items in the selected reporting guidelines that are relevant to RF selection,  
258 data source selection, and data integration. Two reviewers (HZ and JB) independently extracted  
259 these reporting items of interest and resolved conflicts with a third reviewer (YG). We further  
260 analyzed these extracted reporting items and discussed with experts (i.e., biostatisticians, data  
261 analysts and cancer outcomes researchers) to summarize items needed in our reporting guideline,  
262 especially those related to the data integration process.

263

### 264 **Construction of an ontology for the documentation of variable and data source selection** 265 **and integration process (OD-ATTEST)**

266 The ATTEST reporting guideline we developed is used to guide the variable and data source  
267 selection and integration process in cancer outcomes research. We propose to use an ontology-  
268 based approach to annotate and document the items in the reporting guideline. The goal of the  
269 OD-ATTEST ontology is to standardize the terminology used in documenting the selection and  
270 integration steps of RF variables and data sources to support mIDA.

271

272 The OD- ATTEST is developed using Protégé 5. We used Basic Formal Ontology (BFO) [28]  
273 as the upper-level ontology. We first adopted a top down approach to enumerate important  
274 entities (classes and relations) based on the reporting guideline we developed. Following the

275 best practice, we reviewed existing widely accepted ontologies using the National Center for  
276 Biomedical Ontology (NCBO) BioPortal [29] to find the entities can be reused in OD-ATTEST.  
277 Then, we started with the definitions of the most general concepts in the domain and subsequent  
278 specialization of the concepts to develop the class hierarchy. We also took a bottom-up process,  
279 where we started with the definitions of the most specific classes, and then subsequent grouped  
280 similar classes into more general concepts. We also examined how these reporting items are  
281 associated with each other (e.g., “*sample size*” is determined by “*primary outcome*”) and  
282 determined what additional classes and relations were needed to fully represent these entities in  
283 OD-ATTEST.

284

### 285 **An OD-ATTEST-annotated report generated based on a mIDA case study following the** 286 **reporting guideline**

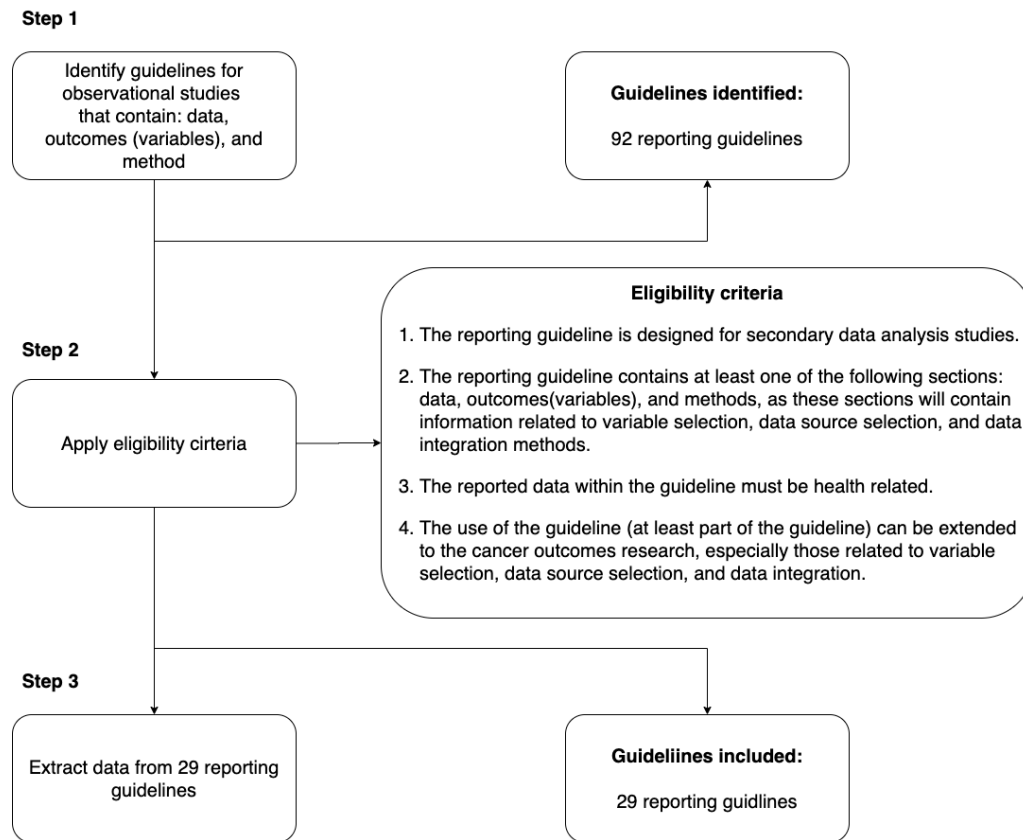
287 To test the developed ATTEST reporting guideline and the OD-ATTEST ontology, we first  
288 created a ATTEST report based on our previous mIDA case study, where we explored the impact  
289 of the relationships among socioeconomic status, individual smoking status, and community-  
290 level smoking rate on pharyngeal cancer survival [16]. To annotate the ATTEST report using  
291 OD-ATTEST, we used the following annotation process: 1) identify information related to the  
292 reporting items in ATTEST through reviewing the original publication and supplementary  
293 materials; 2) annotate the information using the entities in OD-ATTEST; and 3) transform  
294 annotation results into semantic triples in Resource Description Framework (RDF) format using  
295 Turtle syntax [30].

296

## 297 **RESULTS**

298 **The ATTEST reporting guideline for RF variable and data source selection and data**  
299 **integration**

300 We extended our preliminary reporting guideline [23] through a review of existing relevant  
301 reporting guidelines published in the EQUATOR network. *Figure 1* shows our review process.  
302 We reviewed 94 reporting guidelines designed for from observational studies in the EQUATOR  
303 network. Out of the 94 reporting guidelines, 30 contain the required data, outcomes (variables),  
304 and method sections, which we retained for data extraction. In the data extraction step, for each  
305 reporting guideline, we extracted items relevant to RF and data source selection and integration,  
306 where the data and outcomes (variables) sections often contain information regarding how RF  
307 variables and data sources are selected, while the method section contains information about how  
308 data are processed and integrated.



309

310 *Figure 1. Review of relevant reporting guidelines in the EQUATOR network.*

311 We categorized these reporting guidelines (*Table 1*) based on the domains and levels of the data  
 312 sources reported in the guidelines and mapped them to the NIMHD framework. As shown in  
 313 *Table 1*, these 29 reporting guidelines cover data sources from all domains and levels of  
 314 influences. Among them, 9 guidelines focused on providing a general reporting guideline for  
 315 observational studies without specifying a specific domain of influence; while the rest of the  
 316 guidelines are designed for different domains. For example, the Genetic Risk Prediction Studies  
 317 (GRIPS) statement [31] is designed for risk prediction studies using genetic data. Furthermore,  
 318 most guidelines only considered the data sources from individual level, while 2 of them  
 319 considered data sources from multi-levels. For example, the Checklist for One Health  
 320 Epidemiological Reporting of Evidence (COHERE) [32] considered both individual and  
 321 environmental risk factors when studying a disease.

322  
 323 *Table 1. Summary of reporting guidelines based on the data source domains and levels guided*  
 324 *by the NIMHD framework.*

Domain of influences		Level of influences	Guidelines
Not specified*		Individual level	[33], [34], [35], [36], [37], [38], [39], [40], [41]
		Societal level	[42]
Biological data	Genetics data	Individual level	[31], [43]
	Immunogenomic data		[44]
	Molecular epidemiological data		[45], [46]
	Drug safety data from biologics registers		[47], [48]
Behavioral data	Crime, violence data	Individual level	[49]
	Dietary or nutritional data		[50]
	Medication adherence		[51]
Sociocultural environment	Environmental data	Individual/ Community/ Societal/Interpersonal	[32]
Physical environment			

Healthcare system	Administrative data, Electronic health records, Claim data, Patient or disease registries, Quality or safety surveillance databases	Individual level	[52], [53], [54], [55], [56], [57], [58], [59]
*When reporting data sources or RF variables, these studies did not specify a specific data domain.			

325

326 In our preliminary reporting guideline [23], we focused only on reporting items relevant to RF  
327 variables and data sources selection. In this review, we extracted items that can be used to  
328 improve our initial reporting guideline but with a focus on documenting the data integration  
329 process. In total, we found 3 reporting guidelines [53–55] contain information about data  
330 integration processes. However, items included in these 3 guidelines focus on data linkage and  
331 do not contain enough details about how to solve the heterogeneities of data from different  
332 sources. For example, when integrating variables across different levels (e.g., combine  
333 individual-level patient data and county-level smoking rate), none of the 3 guidelines have items  
334 on documenting the cross-level integration choices (e.g., layering the county-level smoking rate  
335 to individual based on residence of the individuals and county code), while this type of choices is  
336 frequently encountered in mIDA studies. Further, data processing steps such as the choices and  
337 algorithms used for creating new data elements (e.g., compute a body mass index variable from  
338 two separate variables, weight and height) are not documented in existing reporting guidelines.  
339 Therefore, we further extended the ATTEST to include these important data integration and data  
340 processing procedures based on our previous research experience on building data integration  
341 framework [25].

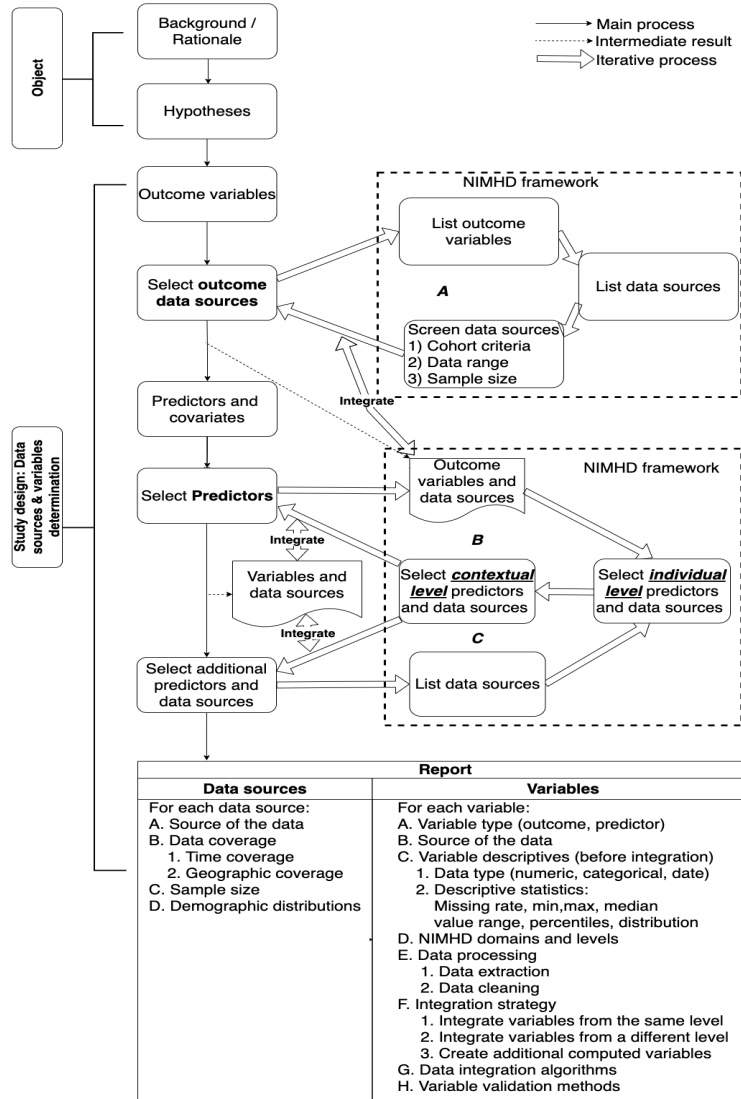
342

343 Informed by the NIMHD research framework and consistent with our prior work, the ATTEST  
344 reporting guideline consists of two main parts as shown in *Figure 2*, reporting (1) the objective



345 of the study including explaining the background and rationale for designing the study in one or  
346 two sentences and describing the hypothesis of the study; and (2) the study design for variable  
347 and data source selection processes and describing the data along with the data integration and  
348 processing strategies. The variable and data source selection process consist of five key steps:  
349 (1) define the outcome variables for primary and (if necessary) secondary outcomes; (2) for each  
350 outcome variable, follow an iterative process (see *Figure 2.A*) to determine the data sources  
351 according to NIMHD framework. After selecting each outcome variable and data sources,  
352 investigators need to think about how to select or consolidate similar outcome variables from the  
353 different selected data sources. For example, if the outcome of interest is an individual's lung  
354 cancer risk, we shall first identify potential data sources (e.g., cancer registries or electronic  
355 health records [EHRs]) that contain individual-level patient data where lung cancer incidence  
356 data are available. Then, based on the cohort criteria and other information such as required  
357 sample size and data range (e.g., time coverage and geographic information) of the potential data  
358 sources, the investigator could determine the qualified data sources and choose an adequate one  
359 based on the objective and design of the study. For example, if 2 data sources, cancer registry  
360 and EHRs, are both available and contain individual-level lung cancer incidence data, the  
361 investigator has the choices to (1) choose one data source over the other, or (2) link the two data  
362 sources and integrate variables from the two data sources. If the investigator chooses to link and  
363 integrate the two data sources, she needs to explicitly document the linkage and integration  
364 processes for each of variables as shown in *Figure 2* (Report – Variables – E, F, G, H) so that  
365 others can repeat the processes to generate the same analytical dataset; (3) determine the  
366 individual-level predictors and covariates of the study; (4) for each individual-level predictor or  
367 covariate, follow loop B in *Figure 2* to identify the different levels/domains of predictors or

368 covariates according to NIMHD framework. Similar to the outcome variables, different data  
369 sources could potentially contain the same predictor or covariate variable, thus, it is important to  
370 contrast and consolidate a new predictor or covariate with the existing selected predictors and  
371 covariates to resolve duplicates. If an investigator chooses to integrate the “duplicate” variables  
372 (e.g., choosing smoking status from cancer registry data over EHRs because cancer registries  
373 data are manually abstracted and typically have better data quality than raw EHRs), these data  
374 integration choices also need to be explicitly documented. Nevertheless, it is often a difficult  
375 choice and these “duplicate” variables might all need to be tested in models before a selection  
376 can be made. Regardless, these decisions and data processing steps need to be clearly  
377 documented; and (5) after selecting individual-level predictors and covariates, one can use a  
378 similar process, following loop C in **Figure 2** to identify additional contextual-level predictors  
379 and covariates and data sources of interest. In the end, a report of the selected data and data  
380 sources as well as the data integration processes shall be generated as shown in **Figure 2**. The  
381 corresponding ATTEST reporting guideline checklist is shown in **Table 2**.



382

383 **Figure 2.** An overview of the reporting guideline for RF variable and data source selection and

384 data integration.

385

386

387

388

389

390

391 **Table 2. ATTEST reporting guideline checklist**

	<b>Item No</b>	<b>Recommendation</b>	<b>Page No</b>
<b>Objectives</b>			
Background/rationale	1	Explain the scientific background and rationale for the study being reported in one or two sentences	
Prespecified hypotheses:	2	State prespecified hypotheses in on or two sentences	
<b>Study design: data sources selection &amp; variables selection &amp; data integration</b>			
Data sources	3a	Describe the time coverage	
	3b	Describe the geographic coverage	
	3c	Describe the sample size	
	3d	Describe the demographic distribution	
	3e	Describe the cohort criteria	
	3f	Describe the sources of biases (e.g., sample bias)	
	3g	Describe the data collection approach	
Dependent variables	4a	State the variable definition and variable type (e.g., primary outcome variable, secondary outcome variable)	
	4b	State the data source of dependent variable	
	4c	State the data type (e.g., numerical, categorical, date-time) of dependent variable	
	4d	State descriptive statistics (e.g., min, max, median, value range, percentile) of dependent variable	
	4e	State the NIMHD domains and levels of dependent variable	
Independent variables	5a	State the variable definition and variable type (e.g., primary predictor, secondary predictor)	
	5b	State the data source of dependent variable	
	5c	State the data type (e.g., numerical, categorical, date-time) of dependent variable	
	5b	State descriptive statistics (e.g., min, max, median, value range, percentile) of independent variable	
	5e	State the NIMHD domains and levels of independent variable	
Controlled variables	6a	State the variables type (e.g., numerical, categorical) of controlled variable	
	6b	State the data source of controlled variable	
	6c	State descriptive statistics (e.g., min, max, median, value range, percentile) of controlled variable	
	6d	State the NIMHD domains and levels of controlled variable	
Missing data	7a	For each data source, describe whether required or expected variable that is not present	
	7b	For each variable, describe method of how to handle missing data	
	7c	For each variable, describe the missing rate	
<b>Data integration</b>			
Data processing	8a	Data extraction: for each variable, describe how to process the raw data source to extract the variable	
	8b	Data cleaning: for each variable, describe the method used to detect and correct (or remove) the incorrect records, missing values or outliers	
Integration strategy	9	Describe the integration strategy for each variable:1) Integrate with variables from same level, 2) Integrate with variables from different levels, and 3) Creation of additional computed elements	

Integration algorithm	10	For each variable, describe the algorithm used to integrate it with variables from other data sources	
Variable validation	11	For each variable, describe data validation rule for the selected variable. Rule should identify both the variable and the validation algorithms	
Integrated variable	12	Describe the variable after integration and basic descriptive statistics (e.g., min, max, median, value range, percentile)	

\*National Institute on Minority Health and Health Disparities (NIMHD)

Note: Please document the items for each data source and variable separately.

392

### 393 **Development of the OD-ATTEST ontology**

394 Based on the ATTEST reporting protocol above, we identified that 48 classes and 25 properties  
395 are needed in OD-ATTEST to represent the ATTEST reporting guideline. *Figure 3* shows the  
396 class hierarchy of OD-ATTEST. We reused classes from the following existing well-known  
397 ontologies: Ontology for Biomedical Investigations (OBI), Information Artifact Ontology (IAO),  
398 National Cancer Institute Thesaurus (NCIt), Statistics Ontology (STATO) and Semanticscience  
399 Integrated Ontology (SIO) as shown in *Table 3*. Note that there are very few existing ontologies  
400 designed for the purpose of documenting the variable and data source selection and data  
401 integration process. The limited number of properties in these existing ontologies are not  
402 informative to represent the elements in the reporting guideline and their relationships, requiring  
403 us to create a large number of new properties in **OD-ATTEST**.



404

405 **Figure 3.** The class hierarchy of OD-ATTEST.

406 **Table 3.** The classes and properties reused or created for OD-ATTEST.

	Label	Internationalized Resource Identifiers (IRIs)*	Reference ontology
Classes	objective	iao:0000005	IAO <sup>1</sup>
	data source	iao:0000100	
	measurement datum	iao:0000109	
Classes	dependent variable	obi:0000751	OBI <sup>2</sup>
	independent variable	obi:0000750	
	controlled variable	obi:0000785	
	data processing	obi:0200000	
Classes	study	ncit:C63536	NCIt <sup>3</sup>
	hypothesis	ncit:C28362	
	rationale	ncit:C80263	
	primary outcome	ncit:C142644	
	secondary outcome	ncit:C142680	
	sample size	ncit:C53190	
	missing data	ncit:C142610	
	data validation	ncit:C142500	
	data type	ncit:C42645	

	data collection method	ncit:C103159	
	data analysis	sio:001051	SIO <sup>4</sup>
	minimum value	stato:0000150	STATO <sup>5</sup>
	maximum value	stato:0000151	
	median	stato:0000574	
	mean	stato:0000573	
	value range	stato:0000035	
	percentile	stato:0000293	
	data distribution	stato:0000161	
	statistical sampling	stato:0000502	
	outlier	stato:0000036	
	primary predictor	od-attest:000015	
	secondary predictor	od-attest:000016	
	demographic distribution	od-attest:000093	ATTEST <sup>6</sup>
	outcome variable data source	od-attest:000019	
	predictor data source	od-attest:000094	
	cohort criteria	od-attest:000008	
	descriptive statistic	od-attest:000012	
	missing rate	od-attest:000068	
	data source time coverage	od-attest:000023	
	data source geographic coverage	od-attest:000024	
	sources of bias	od-attest:000051	
	data integration	od-attest:000052	
	data extraction	od-attest:000054	
	data cleaning	od-attest:000055	
	integration strategy	od-attest:000056	
	integrate variables from same level	od-attest:000057	
	integrate variables from different levels	od-attest:000058	
	creation of additional elements	od-attest:000059	
	integration algorithm	od-attest:000060	
	validation strategy	od-attest:000068	
	integrated variable	od-attest:000096	
Properties	is determined by	od-attest:000097	OD-
	has rationale	od-attest:000098	
	has objective	od-attest:000099	ATTEST
	has data source	od-attest:000100	
	has cohort criteria	od-attest:000101	
	has demographic distribution	od-attest:000102	
	has sources of bias	od-attest:000103	
	has controlled variable	od-attest:000104	
	has independent variable	od-attest:000105	
	has dependent variable	od-attest:000106	
	has data type	od-attest:000107	
	has descriptive statistics	od-attest:000108	
	has NIMHD level	od-attest:000109	

	has NIMHD domain	od-attest:000110	
	has data collection approach	od-attest:000111	
	has sample size	od-attest:000112	
	has missing data	od-attest:000113	
	has data integration	od-attest:000114	
	has data processing	od-attest:000115	
	has data validation	od-attest:000116	
	has integration strategy	od-attest:000117	
	extracted from	od-attest:000118	
	has description	od-attest:000119	
	has time coverage	od-attest:000120	
	has geographic coverage	od-attest:000121	
<sup>4</sup> Prefix: iao: < <a href="http://purl.obolibrary.org/obo/IAO_">http://purl.obolibrary.org/obo/IAO_</a> > obi: < <a href="http://purl.obolibrary.org/obo/OBI_">http://purl.obolibrary.org/obo/OBI_</a> > ncit: < <a href="http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#">http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#</a> > sio: < <a href="http://purl.obolibrary.org/obo/SIO_">http://purl.obolibrary.org/obo/SIO_</a> > stato: < <a href="http://purl.obolibrary.org/obo/STATO_">http://purl.obolibrary.org/obo/STATO_</a> > od-attest: < <a href="http://purl.obolibrary.org/obo/OD-ATTEST/">http://purl.obolibrary.org/obo/OD-ATTEST/</a> > <sup>1</sup> Information Artifact Ontology <sup>2</sup> Ontology for Biomedical Investigations <sup>3</sup> National Cancer Institute Thesaurus <sup>4</sup> Statistics Ontology <sup>5</sup> Semanticscience Integrated Ontology <sup>6</sup> Ontology for the Documentation of Variable and Data Source Selection and Integration Process			

407

408 **An OD-ATTEST-annotated report generated based on a mIDA case study following the**  
 409 **reporting guideline**

410 We annotated two of our previously published mIDA case studies: (1) one study that explored  
 411 the impact of the relationships among socioeconomic status, individual smoking status, and  
 412 community-level smoking rate on pharyngeal cancer survival [16], and (2) another study that  
 413 created a semantic data integration framework to pool multi-level RFs from heterogenous data  
 414 sources to support mIDA [25]. **Table 4** is the filled ATTEST checklist for the two studies.  
 415 **Figure 4** shows a snippet of the ontology annotated variable and data source selection and  
 416 integration process for the second study [25], while the corresponding semantic triples in RDF  
 417 format using Turtle syntax is shown in **Table 5**. The items: RF variables, data sources, and data  
 418 integration steps and their relationships are explicitly standardized and modeled using the classes  
 419 and properties from OD-ATTEST.



420 **Table 4.** An example of two previous mIDA case studies annotated using ATTEST checklist.

Item No	Recommendation	Page No Study (1) [16]	Page No Study (2) [25]
<b>Objectives</b>			
Background/rationale	1 Explain the scientific background and rationale for the study being reported in one or two sentences	Page 1, section “ <i>Abstract</i> ”, paragraph 1, line 1-7	Page 1, section “ <i>Abstract</i> ”, paragraph 1, line 1-4
Prespecified hypotheses	2 State prespecified hypotheses in on or two sentences	Page 2, section “ <i>Introduction</i> ”, paragraph 3, line 1-2	N/A
<b>Study design: data sources selection &amp; variables selection &amp; data integration</b>			
Data source	3a Describe the time coverage	FCDS: Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 2	FCDS: Page 4, section “Data sources”, paragraph 1, line 11
		BRFSS: Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 6	BRFSS: N/A
		2000 U.S. census data: Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 7	United States Census Bureau: Page 4, section “Data sources”, paragraph 1, line 23
	3b Describe the geographic coverage	FCDS: Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 4-5”	FCDS: Page 4, section “Data sources”, paragraph 1, line 12-14
		BRFSS: N/A	BRFSS: Page 10, section “Result”, paragraph 2, line 7-8
		2000 U.S. census data: N/A	United States Census Bureau: N/A
			ATSDR: N/A County Health Ranking & Roadmaps: N/A
	3c Describe the sample size	FCDS: Page 2, section “ <i>Data source and case selection</i> ”, paragraph 2, line 7	FCDS: Page 4, section “Data sources”, paragraph 2, line 6-7
		BRFSS: N/A	BRFSS: N/A
		2000 U.S. census data: N/A	United States Census Bureau: N/A
	3d Describe the demographic distribution	FCDS: Page 2, Table 1	N/A
		BRFSS: N/A 2000 U.S. census data: N/A	
3e Describe the Cohort criteria	FCDS: Page 2, section “ <i>Data source and case selection</i> ”, paragraph 2, line 1-5	FCDS: Page 4, section “Data sources”, paragraph 2, line 1-6	
	BRFSS: N/A	BRFSS: N/A	
	2000 U.S. census data: N/A	United States Census Bureau: N/A	
3f Describe the sources of bias	N/A	N/A	
3g Describe the data collection approach	N/A	FCDS: N/A	FCDS: N/A
		BRFSS: Page 4, section “Data sources”, paragraph 2, line 6-7	BRFSS: Page 4, section “Data sources”, paragraph 2, line 6-7
		United States Census Bureau: N/A	United States Census Bureau: N/A
		ATSDR: N/A County Health Ranking & Roadmaps: N/A	ATSDR: N/A County Health Ranking & Roadmaps: N/A
Dependent variable	4a State the variable definition and variable type (e.g., primary outcome)	Survival time: Page 2, section “ <i>Variable definitions</i> ”, line 1-3	Cancer survival: Page 4, section “Data integration use case: The multi-level integrative data analysis of Cancer survival”, paragraph 1, line 1-2

	variable, secondary outcome variable)		
	4b State the data source of dependent variable	<b>Survival time:</b> Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 2	<b>Cancer survival:</b> Page 4, section “Data sources”, paragraph 1, line 9-14
	4c State the data type (e.g., numerical, categorical, date-time) of dependent variable	<b>Survival time:</b> Page 2, section “ <i>Variable definitions</i> ”, paragraph 1, line 1	<b>Cancer survival:</b> N/A
	4d State descriptive statistics (e.g., min, max, median, value range, percentile) of dependent variable	<b>Survival time:</b> Page 4, Table 1	Cancer survival: N/A
	4e State the NIMHD domain and levels of dependent variable	<b>Survival time:</b> Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 1-2	<b>Cancer survival:</b> Page 4, section “Data sources”, paragraph 2, line 15
Independent variable	5a State the variable definition and variable type (e.g., primary predictor, secondary predictor)	<b>Socioeconomic status:</b> Page 2, section “ <i>Variable definitions</i> ”, paragraph 3, line 1-2	<b>Demographic variables:</b> Page 5, table 1
		<b>Individual smoking:</b> Page 2, section “ <i>Data source and case selection</i> ”, paragraph 2, line 1-2	<b>Smoking status:</b> Page 10, section “The ontology for Cancer research variables (OCRv)”, paragraph 2, line 13- 27
		<b>Regional smoking:</b> Page 3, section “ <i>Data source and case selection</i> ”, paragraph 2, line 4-6	<b>Marital status:</b> Page 14, section “Type 4: Queries that generate results based on the knowledge encoded in ontology”, paragraph 2, line 7- 10
			<b>Insurance payer:</b> Page 5, table 1
			<b>Residency:</b> Page 5, table 1
			<b>Age at diagnosis:</b> Page 5, table 1
			<b>Year of diagnosis:</b> Page 5, table 1
			<b>Tumor stage:</b> Page 5, table 1
			<b>Tumor type:</b> Page 5, table 1
			<b>Treatment procedure:</b> Page 5, table 1
			<b>Census Tract SVI:</b> Page 14, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 5-16
			<b>Census tract high school completion rates:</b> Page 5, table 1
			<b>Census tract family poverty rates:</b> Page 5, table 1
			<b>Census tract rurality status:</b> Page 4, section “Data integration use case: The multi-level integrative data analysis of Cancer survival”, paragraph 1, line 8-11
			<b>County adult mental and physical health status:</b> Page 5, table 1
			<b>County density of primary care physicians:</b> Page 5, table 1
			<b>County smoking rate:</b> Page 10, section “The ontology for Cancer research variables (OCRv)”, paragraph 2
			<b>County alcohol consumption rate:</b> Page 5, table 1
	5b State the data type (e.g., numerical, categorical) of independent variable	<b>Socioeconomic status:</b> Page 2, section “ <i>Variable definitions</i> ”, paragraph 3, line 9-10	<b>Demographic variables:</b> N/A
		<b>Individual smoking:</b> Page 2, section “ <i>Data source and case selection</i> ”, paragraph 2, line 2-3	<b>Smoking status:</b> Page 13, table 3
		<b>Regional smoking:</b> Page 3, section “ <i>Data source and case selection</i> ”, paragraph 2, line 4-6	<b>Marital status:</b> Page 14, section “Type 4: Queries that generate results based on the knowledge encoded in ontology”, paragraph 2, line 7- 10

			<b>Insurance payer:</b> N/A <b>Residency:</b> N/A <b>Age at diagnosis:</b> Page 16, figure 6 <b>Year of diagnosis:</b> Page 16, figure 6 <b>Tumor stage:</b> N/A <b>Tumor type:</b> Page 4, section “Data sources”, paragraph 2, line 1-6 <b>Treatment procedure:</b> Page 5, table 1 <b>Census Tract SVI:</b> Page 14, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 5-16 <b>Census tract high school completion rates:</b> N/A <b>Census tract family poverty rates:</b> N/A <b>Census tract rurality status:</b> N/A <b>County adult mental and physical health status:</b> N/A <b>County density of primary care physicians:</b> N/A <b>County smoking rate:</b> Page 10, section “The ontology for Cancer research variables (OCRv)”, paragraph 2 <b>County alcohol consumption rate:</b> N/A	
	5c	State the data source of independent variable	<b>Socioeconomic status:</b> Page 2, section “Data source and case selection”, paragraph 1, line 6-7 <b>Individual smoking:</b> Page 2, section “Data source and case selection”, paragraph 1, line 1-2 <b>Regional smoking:</b> Page 2, section “Data source and case selection”, paragraph 1, line 7-10	Page 5, table 1
	5d	State descriptive statistics (e.g., min, max, median, value range, percentile) of independent variable	Page 4, table 1	N/A
	5e	State the NIMHD domain and levels of independent variable	<b>Socioeconomic status:</b> Page 2, section “Data source and case selection”, paragraph 1, line 6 <b>Individual smoking:</b> Page 2, section “Data source and case selection”, paragraph 2, line 1 <b>Regional smoking:</b> Page 3, section “Data source and case selection”, paragraph 2, line 4-6	Page 5, table 1
Controlled variable	6a	State the controlled variable and variable type (e.g., numerical, categorical) of controlled variable	<b>Age of diagnosis:</b> Page 2, section “Variable definitions”, paragraph 1, line 10-13 <b>Anatomic site:</b> Page 2, section “Variable definitions”, paragraph 1, line 2-9 <b>Race-ethnicity:</b> Page 4, table 1 <b>Marital status:</b> Page 4, table 1 <b>Insurance:</b> Page 4, table 1 <b>Year of diagnosis:</b> Page 4, table 1 <b>Gender:</b> Page 4, table 1 <b>Stage of diagnosis:</b> Page 4, table 1 <b>Treatment:</b> Page 4, table 1	N/A
	6b	State the data source of controlled variable	Page 2, section “Data source and case selection”, paragraph 1, line 2 <sup>7</sup>	N/A

	6c	State descriptive statistics (e.g., min, max, median, value range, percentile) of controlled variable	Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 2”	N/A
	6d	State the NIMHD domain and levels of controlled variable	Page 2, section “ <i>Data source and case selection</i> ”, paragraph 1, line 1-5”	N/A
Missing data	7a	For each data source, describe whether required or expected variable that is not present	N/A	N/A
	7b	For each variable, describe method of how to handle missing data	N/A	N/A
	7c	For each variable, describe the missing rate	N/A	N/A
Data processing	9a	Data extraction: for each variable, describe how to process the raw data source to extract the variable	N/A	<b>Demographic variables:</b> Page 15, figure 5 <b>Age at diagnosis:</b> Page 16, figure 6 <b>Census Tract SVI:</b> Page 16, figure 7 <b>County smoking rate:</b> Page 17, figure 8 <b>Marital status:</b> Page 18, figure 9
	9b	Data cleaning: for each variable, describe the method used to detect and correct (or remove) the incorrect records, missing values or outliers	N/A	N/A
Integration strategy	10	Describe the integration strategy for each variable: 1) Integrate with variables from same level, 2) Integrate with variables from different levels, and 3) Creation of additional computed elements	<b>Socioeconomic status:</b> Page 2, section “ <i>Variable definitions</i> ”, paragraph 3, line 6-7. <b>Regional smoking:</b> Page 2, section “ <i>Variable definitions</i> ”, paragraph 2, line 4-5.	<b>Demographic variables:</b> Page 15, figure 5 <b>Age at diagnosis:</b> Page 16, figure 6 <b>Census Tract SVI:</b> Page 16, figure 7 <b>County smoking rate:</b> Page 17, figure 8 <b>Marital status:</b> Page 18, figure 9 <b>Census tract high school completion rates:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3 <b>Census tract family poverty rates:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3 <b>Census tract rurality status:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3 <b>County adult mental and physical health status:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3

				<p><b>County density of primary care physicians:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p> <p><b>County alcohol consumption rate:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p>
Integration algorithms	11	For each variable, describe the algorithm used to integrate it with variables from other data sources	N/A	<p><b>Demographic variables:</b> Page 15, figure 5</p> <p><b>Age at diagnosis:</b> Page 16, figure 6</p> <p><b>Census Tract SVI:</b> Page 16, figure 7</p> <p><b>County smoking rate:</b> Page 17, figure 8</p> <p><b>Marital status:</b> Page 18, figure 9</p> <p><b>Census tract high school completion rates:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p> <p><b>Census tract family poverty rates:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p> <p><b>Census tract rurality status:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p> <p><b>County adult mental and physical health status:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p> <p><b>County density of primary care physicians:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p> <p><b>County alcohol consumption rate:</b> Page 15, section “Type 3: Queries that are used to link a patient to contextual factors through geographic variables”, paragraph 1, line 1- 3</p>
Variable validation	12	For each variable, describe data validation rule for the selected variable. Rule should identify both the variable and the validation algorithms	N/A	<p><b>Demographic variables:</b> Page 19, section “Data quality and consistency checks of the source data using the ontology”</p>
Integrated variable	13	Describe the variable after integration and basic descriptive statistics (e.g., min, max, median, value range, percentile)	N/A	Page 18, Table 4

\*Note: if the reported items for all variables or data sources are described at the same place, you can list the page/section/table information at once. For the integration related items, we only presented variables that have the information (N/A will not be showed in the table).

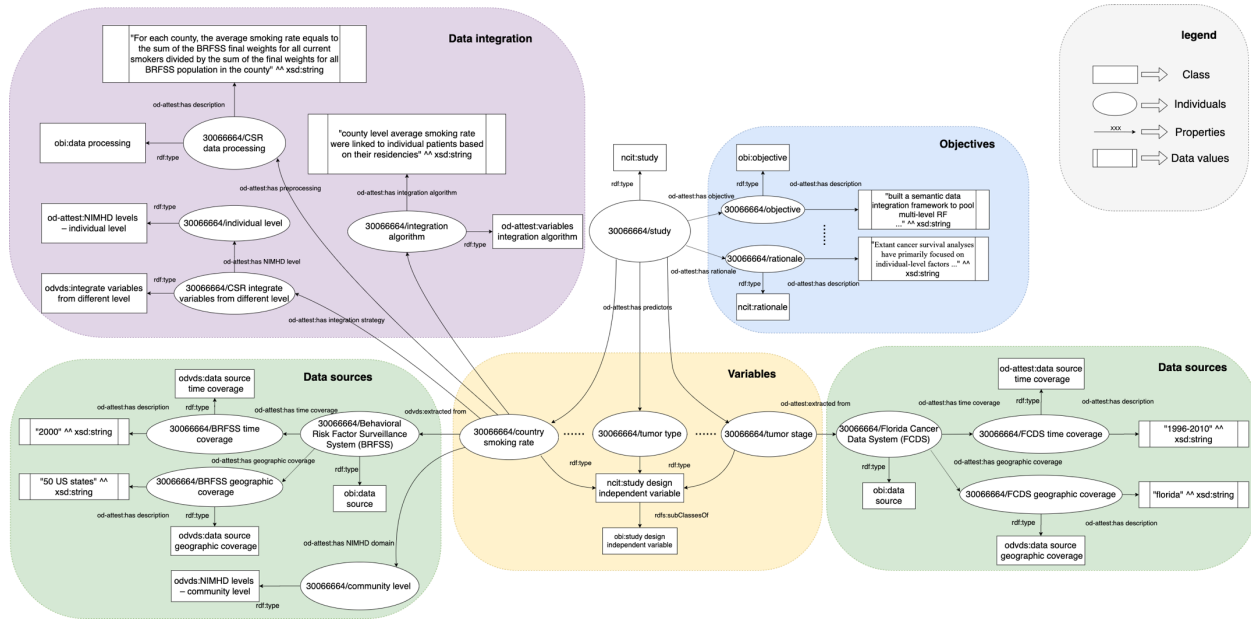
FCDS: Florida Cancer Data System

ATSDR: Agency for Toxic Substances & Disease Registry

BRFSS: Behavioral Risk Factor Surveillance System

421  
422

423 **Figure 4.** An OD-ATTEST-annotated report generated based on a mIDA case study.



424

425 **Table 5.** An example of annotated semantic triples represented in RDF format using Turtle

426 syntax.

<b>Prefix</b>	<p>@prefix od-attest: &lt;<a href="http://www.semanticweb.org/od-attest/">http://www.semanticweb.org/od-attest/</a>&gt;.                  @prefix ncit: &lt;<a href="http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl/">http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl/</a>&gt;.                  @prefix rdfs: &lt;<a href="http://www.w3.org/2000/01/rdf-schema/">http://www.w3.org/2000/01/rdf-schema/</a>&gt;.                  @prefix xsd: &lt;<a href="http://www.w3.org/2001/XMLSchema/">http://www.w3.org/2001/XMLSchema/</a>&gt;.</p>	
<b>RDF<sup>1</sup> triples</b>	<p>od-attest:30066664                  rdf:type ncit:study;                  od-attest:has rationale od-attest:30066664/rationale;                  od-attest:has objective od-attest:30066664/objective.</p> <p>od-attest:30066664/rationale                  rdf:type ncit:rationale;                  od-attest:has description "Extant cancer survival analyses have..." ^^ xsd:string.</p> <p>od-attest:30066664/objective                  rdf:type ncit:objective;                  od-attest:has description "built a semantic data integration ..." ^^ xsd:string.</p>	
	<p><sup>1</sup> Resource Description Framework</p>	

427

## 428 **DISCUSSION**

429 In this study, we first developed a reporting guideline, ATTEST, to provide a theory-driven  
430 approach to guide the RF variable and data source selection and integration process in cancer  
431 outcomes research. We then proposed an ontology-based approach to annotate the items in our  
432 reporting guideline so that information relevant to variables, data sources and data integration in  
433 mIDA studies can be explicitly documented. To develop the reporting guideline, we conducted a  
434 systematic search to identify useful reporting items to improve our selection and data integration  
435 process. We categorized these reporting guidelines based on their reported data source domains  
436 and levels according to NIMHD framework, so that we can identify items need to be reported  
437 when selecting variables or data sources from different domains and levels. For example, when  
438 report population-level estimates (variables) [42], the information regarding the sources of bias  
439 (e.g., selection bias) need to be documented. Therefore, we updated our previous reporting  
440 guideline and added “*sources of bias*” as a reporting item when documenting data sources. This  
441 is important, because subsequent data processing steps might be needed to correct the bias.

442 Further,

443  
444 The use of NIMHD framework can also help researchers to systematically think and structure the  
445 variable and data source selection process when considering multi-level RF variables from  
446 heterogenous data sources. For example, if an investigator is considering smoking related risk  
447 factors in cancer outcomes research, following the NIMHD framework, one can start with  
448 variables in the behavioral domain and then list potential smoking related variables for each level  
449 of influences step by step, such as individual smoking status at the individual level, second hand  
450 smoke exposure at the interpersonal level, county level smoking rate at the community level, and

451 smoking policies or laws (e.g., federal minimum age to purchase tobacco products) at the  
452 societal level. The same process can be applied to select other smoking related variables from  
453 other domains of influences. In this way, investigators can systematically think and evaluate the  
454 confounding effects and cross-level interactions among those selected variables which are  
455 usually ignored in previous cancer outcome studies using a single data source.

456

457 We provided a ATTEST checklist (1) to help researchers clearly document each step of their RF  
458 and data source selection and integration process, and (2) to improve the completeness and  
459 transparency of their mIDA studies. As shown in **Table 4**, we used the ATTEST checklist to  
460 report two previous mIDA studies. Based on the checklist, we can easily (1) check whether  
461 these mIDA studies document required items that can help other researchers replicate their  
462 studies, and (2) compare their variables, data sources and data integration processes. As shown  
463 in **Table 4**, we found that there are 3 items never discussed in either of the two studies including  
464 “*sources of bias*”, “*missing data*” for selected variables, and “*data cleaning*” (i.e., method used  
465 to detect and correct or remove the incorrect records, missing values or outliers). All three items  
466 are relevant to data quality issues, where rarely being discussed or documented in these mIDA  
467 studies or even more broadly in cancer outcomes research. Nevertheless, data quality issues such  
468 as missing data can dramatically affect the results of the cancer outcomes research (e.g., in  
469 cancer survival prediction) [60]. Comparing the two case mIDA studies, the data integration  
470 process was not well-documented in the first study [16], where most of the items relevant to data  
471 integration are blank; while, in the other study [25], the processes about data processing, data  
472 integration, and data validation were all clearly documented according to the ATTEST checklist.



473 Therefore, using this checklist, one can improve the completeness of their documentation on the  
474 selection and integration process as shown in *Table 4*.

475  
476 The OD-ATTEST ontology provides a way to standardize the documentation of the mIDA study  
477 process from variable and data source selection to data integration. Also, the ontology-based  
478 annotations of the report is beneficial because it provides an initial step towards a report that is  
479 not only readable and understandable by human but also potentially executable by machines.  
480 After transforming these annotations into semantic triples, the report can be stored into a  
481 knowledge base and represented as knowledge graphs (*Figure 4*) to facilitate examination and  
482 analysis of these mIDA reports, enabling robust sharing and comparison of different mIDA  
483 studies.

484

## 485 **Limitations and future work**

486 Most of the reporting guidelines we reviewed from the EQUATOR network have limited  
487 information on how to document the data integration process, indicating a significant gap in  
488 existing practice. Nevertheless, we were able to summarize the key elements need to be reported  
489 for the integration process based on 3 existing guidelines and our own previous experience on  
490 semantic data integration case studies. As a future study, one shall conduct a systematic review  
491 on data integration literatures to summarize relevant reporting items to improve the reporting  
492 guideline. Meanwhile, we will conduct a yearly review of existing reporting guidelines  
493 following the reviewing process discussed in *Figure 1* to identify new reporting items of interest  
494 and keep our framework up to date. Further, beyond standardized reporting, our ultimate goal is  
495 to let computers understand the ontology-annotated report (in RDF triples) regarding (1) how

496 different variables are defined and represented and (2) how different variables are selected and  
497 integrated, so that machines can automatically repeat these processes and generate integrated  
498 dataset based on an executable ontology-annotated report. For variable definition and  
499 representation, it is important to recognize and being interoperable with existing data standards  
500 and common data models (CDM) such as those that standardized exchanging of EHRs data  
501 including the national Patient-Centered Clinical Research Network (PCORnet) CDM, the  
502 Observational Medical Outcomes Partnership (OMOP) from the Observational Health Data  
503 Sciences and Informatics (OHDSI) network, and the uprising Fast Healthcare Interoperability  
504 Resources (FHIR) protocol adopted by major EHR system vendors. Developing the ontology  
505 against these CDMs that have already standardized existing data resources would be critical to  
506 assure the generalizability of our framework. Nevertheless, for modeling the variable selection  
507 and integration processes as shown in **Figure 4**, more fine-grained information regarding the  
508 variables, data sources and the integration process are currently documented as free-text  
509 descriptions. We face challenges in transforming these “free-text” information into executable  
510 algorithms (e.g., a data processing step that calculates BMI using weight and height). Such  
511 information is related to the concept of data provenance—“*a type of metadata, concerned with*  
512 *the history of data, its origin and changes made to it*” [61]. The importance of data provenance  
513 is widely recognized, especially for study reproducibility and replicability. More than one-half  
514 of the systematic efforts to reproduce computational results across different fields have failed,  
515 mainly due to insufficient detail on digital artifacts, such as data, code, and computational  
516 workflow [62]. However, descriptions of data provenance are often neglected or inadequate in  
517 scientific literature due to the lack of a tractable, easily operated approach with supporting tools.  
518 Future studies that focus on the development of easy-to-use tools with a standardized framework

519 to persist end-to-end data provenance with high integrity including intermediate processes and  
520 data products are urgently needed. Further, future developments of tools and platforms to  
521 automate the documentation process, where the data elements and associated information (e.g.,  
522 levels and domains) are also automatically annotated with the standardized ontology are  
523 warranted.

524

## 525 **CONCLUSIONS**

526 In this paper, we have proposed and developed an ontology-based reporting guideline solving  
527 some key challenges in current mIDA studies for cancer outcomes research, through providing  
528 (1) a theory-driven guidance for multi-level and multi-domain RF variable and data source  
529 selection; and (2) a standardized documentation of the data selection and integration processes  
530 powered by an ontology, thus a way to enable sharing of mIDA study reports among researchers.

531

## 532 **List of abbreviations**

533 ACS American Cancer Society  
534 BRFSS Behavioral Risk Factor Surveillance System  
535 EQUATOR Enhancing the QUALity and Transparency Of health Research  
536 FCDS Florida Cancer Data System  
537 mIDA Multi-level Integrative Data Analysis  
538 NIH National Institute of Health  
539 NIMHD Minority Health and Health Disparities  
540 OD-ATTEST Ontology for the Documentation of Variable and Data Source Selection and  
541 Integration Process

542	RF	Risk Factor
543	US	United States
544	RUCA	Rural-Urban Commuting Area
545	NCHS	National Center for Health Statistics
546	BFO	Basic Formal Ontology
547	NCBO	National Center for Biomedical Ontology
548	RDF	Resource Description Framework
549	GRIPS	Genetic Risk Prediction Studies
550	COHERE	Checklist for One Health Epidemiological Reporting of Evidence
551	EHR	Electronic Health Records
552	OBI	Ontology for Biomedical Investigations
553	IAO	Information Artifact Ontology
554	NCIt	National Cancer Institute Thesaurus
555	STATO	Statistics Ontology
556	SIO	Semanticscience Integrated Ontology
557	CDM	Common Data Model
558	PCORnet	The national Patient-Centered Clinical Research Network

559

## 560 **Declarations**

### 561 **Ethics approval and consent to participate**

562 Not applicable

563

### 564 **Consent to publish**

565 Not applicable

566

567 **Availability of data and materials**

568 The reviewed reporting guidelines are available in the public: Enhancing the QUALity and

569 Transparency Of health Research (EQUATOR) network (<https://www.equator->

570 [network.org/reporting-guidelines/](https://www.equator-network.org/reporting-guidelines/)).

571

572 **Competing interests**

573 The authors declare that they have no competing interests.

574

575 **Funding**

576 This study was supported in part by the National Institute of Health (NIH) awards

577 UL1TR001427 and R01CA246418 and Patient-Centered Outcomes Research Institute (PCORI)

578 award ME-2018C3-14754. The content is solely the responsibility of the authors and does not

579 necessarily represent the official views of the NIH or PCORI.

580

581 **Authors' contributions**

582 The work presented here was carried out in collaboration among all authors. YG and JB

583 designed the study. YG, QL and HZ were involved in acquisition of the data and review of

584 existing reporting guidelines. HZ wrote the initial draft of the manuscript with substantial

585 support from YG and JB. YG and JB provided expert opinion during the ontology curation

586 process and guided the design of the ontology. All authors provided critical feedback on the

587 study design, reviewed and edited the manuscript. All authors have read and approved the final  
588 manuscript.

589

## 590 **Acknowledgements**

591 None

592

## 593 **References**

594 1. World Health Organization. Cancer - key facts. 2018. [https://www.who.int/news-room/fact-](https://www.who.int/news-room/fact-sheets/detail/cancer)  
595 [sheets/detail/cancer](https://www.who.int/news-room/fact-sheets/detail/cancer). Accessed 2 Jan 2020.

596 2. Atlanta: American Cancer Society. Cancer Facts & Figures 2019. 2019.  
597 [https://www.cancer.org/research/cancer-facts-statistics/all-cancer-facts-figures/cancer-facts-](https://www.cancer.org/research/cancer-facts-statistics/all-cancer-facts-figures/cancer-facts-figures-2019.html)  
598 [figures-2019.html](https://www.cancer.org/research/cancer-facts-statistics/all-cancer-facts-figures/cancer-facts-figures-2019.html). Accessed 2 Jan 2020.

599 3. Siegel RL, Miller KD, Jemal A. Cancer statistics, 2019. *CA Cancer J Clin.* 2019;69:7–34.  
600 doi:10.3322/caac.21551.

601 4. National Cancer Institute. Cancer Risk Factors.  
602 <https://training.seer.cancer.gov/disease/cancer/risk.html>. Accessed 2 Jan 2020.

603 5. Andrew AS, Parker S, Anderson JC, Rees JR, Robinson C, Riddle B, et al. Risk Factors for  
604 Diagnosis of Colorectal Cancer at a Late Stage: a Population-Based Study. *J Gen Intern Med.*  
605 2018;33:2100–5. doi:10.1007/s11606-018-4648-7.

606 6. Mobley LR, Kuo T-M. Demographic Disparities in Late-Stage Diagnosis of Breast and  
607 Colorectal Cancers Across the USA. *J Racial Ethn Health Disparities.* 2017;4:201–12.  
608 doi:10.1007/s40615-016-0219-y.

609 7. Markossian TW, Hines RB. Disparities in late stage diagnosis, treatment, and breast cancer-  
610 related death by race, age, and rural residence among women in Georgia. *Women Health.*  
611 2012;52:317–35.

612 8. Chatterjee NA, He Y, Keating NL. Racial differences in breast cancer stage at diagnosis in the  
613 mammography era. *Am J Public Health.* 2013;103:170–6.

614 9. Montealegre JR, Zhou R, Amirian ES, Follen M, Scheurer ME. Nativity disparities in late-stage  
615 diagnosis and cause-specific survival among Hispanic women with invasive cervical cancer: an  
616 analysis of Surveillance, Epidemiology, and End Results data. *Cancer Causes Control.*  
617 2013;24:1985–94. doi:10.1007/s10552-013-0274-1.

- 618 10. Baquet CR, Mishra SI, Commiskey P, Ellison GL, DeShields M. Breast cancer epidemiology in  
619 blacks and whites: disparities in incidence, mortality, survival rates and histology. *J Natl Med*  
620 *Assoc.* 2008;100:480–8.
- 621 11. Yasmeen S, Xing G, Morris C, Chlebowski RT, Romano PS. Comorbidities and mammography  
622 use interact to explain racial/ethnic disparities in breast cancer stage at diagnosis. *Cancer.*  
623 2011;117:3252–61.
- 624 12. Echeverría SE, Borrell LN, Brown D, Rhoads G. A local area analysis of racial, ethnic, and  
625 neighborhood disparities in breast cancer staging. *Cancer Epidemiol Biomark Prev Publ Am*  
626 *Assoc Cancer Res Cosponsored Am Soc Prev Oncol.* 2009;18:3024–9.
- 627 13. NIMHD. NIMHD Research Framework. [https://www.nimhd.nih.gov/about/  
628 overview/research-framework.html](https://www.nimhd.nih.gov/about/overview/research-framework.html). Accessed 28 Jun 2019.
- 629 14. Dahlberg LL, Krug EG. Violence a global public health problem. *Ciênc Saúde Coletiva.*  
630 2006;11:277–92. doi:10.1590/S1413-81232006000200007.
- 631 15. Keegan TH, Quach T, Shema S, Glaser SL, Gomez SL. The influence of nativity and  
632 neighborhoods on breast cancer stage at diagnosis and survival among California Hispanic  
633 women. *BMC Cancer.* 2010;10:603. doi:10.1186/1471-2407-10-603.
- 634 16. Guo Y, Logan HL, Marks JG, Shenkman EA. The relationships among individual and regional  
635 smoking, socioeconomic status, and oral and pharyngeal cancer survival: a mediation analysis.  
636 *Cancer Med.* 2015;4:1612–9.
- 637 17. Giordano A. Data integration blueprint and modeling: techniques for a scalable and  
638 sustainable architecture. Upper Saddle River, NJ: IBM Press Pearson; 2011.
- 639 18. Schloss PD. Identifying and Overcoming Threats to Reproducibility, Replicability,  
640 Robustness, and Generalizability in Microbiome Research. *mBio.* 2018;9:e00525-18,  
641 /mbio/9/3/mBio.00525-18.atom. doi:10.1128/mBio.00525-18.
- 642 19. Alonso-Calvo R, Paraiso-Medina S, Perez-Rey D, Alonso-Oset E, van Stiphout R, Yu S, et al. A  
643 semantic interoperability approach to support integration of gene expression and clinical data  
644 in breast cancer. *Comput Biol Med.* 2017;87:179–86. doi:10.1016/j.combiomed.2017.06.005.
- 645 20. Kondylakis H, Claerhout B, Keyur M, Koumakis L, van Leeuwen J, Marias K, et al. The  
646 INTEGRATE project: Delivering solutions for efficient multi-centric clinical research and trials. *J*  
647 *Biomed Inform.* 2016;62:32–47. doi:10.1016/j.jbi.2016.05.006.
- 648 21. METABRIC Group, Papatheodorou I, Crichton C, Morris L, Maccallum P, Davies J, et al. A  
649 metadata approach for clinical data management in translational genomics studies in breast  
650 cancer. *BMC Med Genomics.* 2009;2. doi:10.1186/1755-8794-2-66.

- 651 22. Centre for Statistics in Medicine, NDORMS, University of Oxford. Enhancing the QUALity and  
652 Transparency Of health Research. <https://www.equator-network.org/reporting-guidelines/>.  
653 Accessed 28 Jan 2020.
- 654 23. Zhang H, Guo Y, Bian J. Ontology for Documentation of Variable and Data Source Selection  
655 Process to Support Integrative Data Analysis in Cancer Outcomes Research. In: SEPDA@ISWC.  
656 2019.
- 657 24. Guo Y, Bian J, Modave F, Li Q, George TJ, Prosperi M, et al. Assessing the effect of data  
658 integration on predictive ability of cancer survival models. *Health Informatics J*.  
659 2019;:1460458218824692.
- 660 25. Zhang H, Guo Y, Li Q, George TJ, Shenkman E, Modave F, et al. An ontology-guided semantic  
661 data integration framework to support integrative data analysis of cancer survival. *BMC Med*  
662 *Inform Decis Mak*. 2018;18. doi:10.1186/s12911-018-0636-4.
- 663 26. Rural-Urban Commuting Area Codes. 2019. [https://www.ers.usda.gov/data-products/rural-](https://www.ers.usda.gov/data-products/rural-urban-commuting-area-codes.aspx)  
664 [urban-commuting-area-codes.aspx](https://www.ers.usda.gov/data-products/rural-urban-commuting-area-codes.aspx). Accessed 28 Jan 2020.
- 665 27. National Center for Health Statistics, Office of Analysis and Epidemiology. NCHS Urban-Rural  
666 Classification Scheme for Counties. 2017.  
667 [https://www.cdc.gov/nchs/data\\_access/urban\\_rural.htm#2013\\_Urban-](https://www.cdc.gov/nchs/data_access/urban_rural.htm#2013_Urban-Rural_Classification_Scheme_for_Counties)  
668 [Rural\\_Classification\\_Scheme\\_for\\_Counties](https://www.cdc.gov/nchs/data_access/urban_rural.htm#2013_Urban-Rural_Classification_Scheme_for_Counties). Accessed 28 Jan 2017.
- 669 28. Arp R, Smith B, Spear AD. *Building Ontologies With Basic Formal Ontology*. The MIT Press;  
670 2015. doi:10.7551/mitpress/9780262527811.001.0001.
- 671 29. Whetzel PL, Noy NF, Shah NH, Alexander PR, Nyulas C, Tudorache T, et al. BioPortal:  
672 enhanced functionality via new Web services from the National Center for Biomedical Ontology  
673 to access and use ontologies in software applications. *Nucleic Acids Res*. 2011;39 Web Server  
674 issue:W541-545.
- 675 30. David Beckett, Tim Berners-Lee, Eric Prud'hommeaux, Gavin Carothers, Lex Machina. RDF  
676 1.1 Turtle. 2014. <https://www.w3.org/TR/2014/REC-turtle-20140225/Overview.html>. Accessed  
677 31 Jan 2020.
- 678 31. Janssens ACJW, Ioannidis JPA, van Duijn CM, Little J, Khoury MJ, GRIPS Group. Strengthening  
679 the reporting of Genetic Risk Prediction Studies: the GRIPS Statement. *PLoS Med*.  
680 2011;8:e1000420.
- 681 32. Davis MF, Rankin SC, Schurer JM, Cole S, Conti L, Rabinowitz P, et al. Checklist for One  
682 Health Epidemiological Reporting of Evidence (COHERE). *One Health*. 2017;4:14–21.  
683 doi:10.1016/j.onehlt.2017.07.001.



- 684 33. Leech NL, Onwuegbuzie AJ. Guidelines for Conducting and Reporting Mixed Research in the  
685 Field of Counseling and Beyond. *J Couns Dev.* 2010;88:61–9. doi:10.1002/j.1556-  
686 6678.2010.tb00151.x.
- 687 34. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent Reporting of a multivariable  
688 prediction model for Individual Prognosis Or Diagnosis (TRIPOD): The TRIPOD Statement. *Ann*  
689 *Intern Med.* 2015;162:55. doi:10.7326/M14-0697.
- 690 35. Kerr KF, Meisner A, Thiessen-Philbrook H, Coca SG, Parikh CR. RiGoR: reporting guidelines to  
691 address common sources of bias in risk model development. *Biomark Res.* 2015;3:2.
- 692 36. Jason LA, Unger ER, Dimitrakoff JD, Fagin AP, Houghton M, Cook DB, et al. Minimum data  
693 elements for research reports on CFS. *Brain Behav Immun.* 2012;26:401–6.
- 694 37. Fitchett EJA, Seale AC, Vergnano S, Sharland M, Heath PT, Saha SK, et al. Strengthening the  
695 Reporting of Observational Studies in Epidemiology for Newborn Infection (STROBE-NI): an  
696 extension of the STROBE statement for neonatal infection research. *Lancet Infect Dis.*  
697 2016;16:e202–13.
- 698 38. White RG, Hakim AJ, Salganik MJ, Spiller MW, Johnston LG, Kerr L, et al. Strengthening the  
699 Reporting of Observational Studies in Epidemiology for respondent-driven sampling studies:  
700 “STROBE-RDS” statement. *J Clin Epidemiol.* 2015;68:1463–71.
- 701 39. von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP, et al. The  
702 Strengthening the Reporting of Observational Studies in Epidemiology (STROBE) statement:  
703 guidelines for reporting observational studies. *Ann Intern Med.* 2007;147:573–7.
- 704 40. Jackson DL. Reporting results of latent growth modeling and multilevel modeling analyses:  
705 some recommendations for rehabilitation psychology. *Rehabil Psychol.* 2010;55:272–85.
- 706 41. Wolfe F, Lassere M, van der Heijde D, Stucki G, Suarez-Almazor M, Pincus T, et al.  
707 Preliminary core set of domains and reporting requirements for longitudinal observational  
708 studies in rheumatology. *J Rheumatol.* 1999;26:484–9.
- 709 42. Stevens GA, Alkema L, Black RE, Boerma JT, Collins GS, Ezzati M, et al. Guidelines for  
710 Accurate and Transparent Health Estimates Reporting: the GATHER statement. *The Lancet.*  
711 2016;388:e19–23. doi:10.1016/S0140-6736(16)30388-9.
- 712 43. Little J, Higgins JPT, Ioannidis JPA, Moher D, Gagnon F, von Elm E, et al. STrengthening the  
713 REporting of Genetic Association Studies (STREGA): an extension of the STROBE statement.  
714 *PLoS Med.* 2009;6:e22.
- 715 44. Hollenbach JA, Mack SJ, Gourraud P-A, Single RM, Maiers M, Middleton D, et al. A  
716 community standard for immunogenomic data reporting and analysis: proposal for a  
717 STrengthening the REporting of Immunogenomic Studies statement. *Tissue Antigens.*  
718 2011;78:333–44.

- 719 45. Field N, Cohen T, Struelens MJ, Palm D, Cookson B, Glynn JR, et al. Strengthening the  
720 Reporting of Molecular Epidemiology for Infectious Diseases (STROME-ID): an extension of the  
721 STROBE statement. *Lancet Infect Dis*. 2014;14:341–52.
- 722 46. Gallo V, Egger M, McCormack V, Farmer PB, Ioannidis JPA, Kirsch-Volders M, et al.  
723 STrengthening the Reporting of OBservational studies in Epidemiology - Molecular  
724 Epidemiology (STROBE-ME): an extension of the STROBE statement. *Eur J Clin Invest*.  
725 2012;42:1–16.
- 726 47. Dixon WG, Carmona L, Finckh A, Hetland ML, Kvien TK, Landewe R, et al. EULAR points to  
727 consider when establishing, analysing and reporting safety data of biologics registers in  
728 rheumatology. *Ann Rheum Dis*. 2010;69:1596–602.
- 729 48. Zavada J, Dixon WG, Askling J, EULAR Study group on Longitudinal Observational Registers  
730 and Drug Studies. Launch of a checklist for reporting longitudinal observational drug studies in  
731 rheumatology: a EULAR extension of STROBE guidelines based on experience from biologics  
732 registries. *Ann Rheum Dis*. 2014;73:628.
- 733 49. Singh JP, Yang S, Mulvey EP, RAGEE Group. Reporting guidance for violence risk assessment  
734 predictive validity studies: the RAGEE Statement. *Law Hum Behav*. 2015;39:15–22.
- 735 50. Lachat C, Hawwash D, Ocké MC, Berg C, Forsum E, Hörnell A, et al. Strengthening the  
736 Reporting of Observational Studies in Epidemiology—Nutritional Epidemiology (STROBE-nut):  
737 An Extension of the STROBE Statement. *PLOS Med*. 2016;13:e1002036.  
738 doi:10.1371/journal.pmed.1002036.
- 739 51. De Geest S, Zullig LL, Dunbar-Jacob J, Helmy R, Hughes DA, Wilson IB, et al. ESPACOMP  
740 Medication Adherence Reporting Guideline (EMERGE). *Ann Intern Med*. 2018;169:30–5.
- 741 52. Wang SV, Schneeweiss S, Berger ML, Brown J, de Vries F, Douglas I, et al. Reporting to  
742 Improve Reproducibility and Facilitate Validity Assessment for Healthcare Database Studies  
743 V1.0. *Value Health J Int Soc Pharmacoeconomics Outcomes Res*. 2017;20:1009–22.
- 744 53. Kahn MG, Brown JS, Chun AT, Davidson BN, Meeker D, Ryan PB, et al. Transparent reporting  
745 of data quality in distributed data networks. *EGEMS Wash DC*. 2015;3:1052.
- 746 54. Langan SM, Schmidt SA, Wing K, Ehrenstein V, Nicholls SG, Filion KB, et al. The reporting of  
747 studies conducted using observational routinely collected health data statement for  
748 pharmacoepidemiology (RECORD-PE). *BMJ*. 2018;k3532. doi:10.1136/bmj.k3532.
- 749 55. Benchimol EI, Smeeth L, Guttman A, Harron K, Moher D, Petersen I, et al. The REporting of  
750 studies Conducted using Observational Routinely-collected health Data (RECORD) Statement.  
751 *PLOS Med*. 2015;12:e1001885. doi:10.1371/journal.pmed.1001885.
- 752 56. Bennett DA, Brayne C, Feigin VL, Barker-Collo S, Brainin M, Davis D, et al. Development of  
753 the Standards of Reporting of Neurological Disorders (STROND) checklist: A guideline for the

- 754 reporting of incidence and prevalence studies in neuroepidemiology. *Neurology*. 2015;85:821–  
755 8.
- 756 57. Berger ML, Mamdani M, Atkins D, Johnson ML. Good research practices for comparative  
757 effectiveness research: defining, reporting and interpreting nonrandomized studies of  
758 treatment effects using secondary data sources: the ISPOR Good Research Practices for  
759 Retrospective Database Analysis Task Force Report--Part I. *Value Health J Int Soc*  
760 *Pharmacoeconomics Outcomes Res*. 2009;12:1044–52.
- 761 58. Holtfreter B, Albandar JM, Dietrich T, Dye BA, Eaton KA, Eke PI, et al. Standards for reporting  
762 chronic periodontitis prevalence and severity in epidemiologic studies: Proposed standards  
763 from the Joint EU/USA Periodontal Epidemiology Working Group. *J Clin Periodontol*.  
764 2015;42:407–12.
- 765 59. Tacconelli E, Cataldo MA, Paul M, Leibovici L, Kluytmans J, Schröder W, et al. STROBE-AMS:  
766 recommendations to optimise reporting of epidemiological studies on antimicrobial resistance  
767 and informing improvement in antimicrobial stewardship. *BMJ Open*. 2016;6:e010134.  
768 doi:10.1136/bmjopen-2015-010134.
- 769 60. Barakat MS, Field M, Ghose A, Stirling D, Holloway L, Vinod S, et al. The effect of imputing  
770 missing clinical attribute values on training lung cancer survival prediction model performance.  
771 *Health Inf Sci Syst*. 2017;5:16.
- 772 61. Glavic B, Dittrich KR. Data Provenance: A Categorization of Existing Approaches. In:  
773 *Datenbanksysteme in Business, Technologie und Web (BTW)*. Aachen: Ges. für Informatik;  
774 2007. p. 227–41.
- 775 62. Committee on Reproducibility and Replicability in Science, Board on Behavioral, Cognitive,  
776 and Sensory Sciences, Committee on National Statistics, Division of Behavioral and Social  
777 Sciences and Education, Nuclear and Radiation Studies Board, Division on Earth and Life  
778 Studies, et al. *Reproducibility and Replicability in Science*. Washington, D.C.: National  
779 Academies Press; 2019. doi:10.17226/25303.
- 780
- 781