



1           **National Evaluation System for health Technology Coordinating**  
2                           **Center (NESTcc) Data Quality Framework**

3  
4  
5    *Subcommittee Members: Lesley Curtis, PhD, MS (chair); Jeffrey Brown, PhD; John Laschinger,*  
6    *MD; Aaron Lottes, PhD; Keith Marsolo, PhD; Frederick Masoudi, MD, MSPH; Joseph Ross, MD,*  
7    *MHS; Art Sedrakyan, MD, PhD; Kara Southall, MS; James Tcheng, MD; Karen Ulisney, MS, CRNP;*  
8    *Charles Viviano, MD, PhD*

9  
10  
11    **Contents**

12    Introduction ..... 1  
13    Governance ..... 2  
14    Characteristics of Data ..... 5  
15    Data Capture and Transformation ..... 7  
16    Data Curation ..... 8  
17    NESTcc Data Quality Maturity Model ..... 11  
18    Conclusion ..... 14  
19    References ..... 14

20  
21    **Introduction**

22    In 2012, the National Evaluation System for health Technology (NEST) was born to “quickly identify  
23    problematic devices, accurately and transparently characterize and disseminate information about device  
24    performance in clinical practice, and efficiently generate data to support premarket clearance or approval  
25    of new devices and new uses of currently marketed devices.”<sup>1</sup>

26    In 2018, the Data Quality Subcommittee of the NEST Coordinating Center (NESTcc) was tasked with  
27    creating a Data Quality Framework for NESTcc Network Collaborators. The initial version of that  
28    framework, presented in this document, lays out the foundation for the capture and use of high-quality  
29    data for post-market evaluation of medical devices. Aligned with NESTcc’s pragmatic approach to device  
30    evaluation, this framework is grounded in the use of real-world data (RWD) gleaned from the clinical care  
31    setting instead of data collected specifically for research or evaluation purposes. This framework focuses  
32    on RWD from the electronic health record (EHR) rather than other clinically based data sources such as  
33    registries, which have been addressed elsewhere.<sup>2</sup>



34 This Data Quality Framework serves as a guide to Network Collaborators and organizations that wish to  
35 collaborate with NESTcc, to ensure the quality of their data related to medical devices. The overarching  
36 goal of this framework is to inform the capture and use of clinical information as high-quality data to  
37 support the generation of real-world evidence (RWE), which will ultimately, and most importantly,  
38 provide better care to patients.

39 This framework is composed of five sections that cover the topics most salient to achieving the highest  
40 data quality around medical devices:

- 41 1. **Governance:** Involving and engaging stakeholders is critical to good governance for RWD and  
42 RWE. Governance ensures stakeholder representation, limits the potential for bias or unethical  
43 behaviors, and results in trustworthy findings and conclusions.
- 44 2. **Characteristics of Data:** Choosing and using data appropriately first necessitates understanding  
45 and specifying the data needed, along with the context and limitations of potential sources of  
46 that data. Shortcomings of the data that potentially limit their application must also be  
47 identified.
- 48 3. **Data Capture and Transformation:** The use of EHR data for secondary analyses presents additional  
49 challenges in terms of data relevance and reliability. The processing and transformation of data  
50 into common data models provides a logical pathway for enabling analysis.
- 51 4. **Data Curation:** Curation turns raw data into information by organizing, assessing, and preparing  
52 the data for analysis. Data curation is an iterative process, with the goal to improve data quality  
53 over time.
- 54 5. **NESTcc Data Quality Maturity Model:** Maturity models are used by organizations to assess  
55 business capabilities, identify opportunities, and perform capacity planning. Maturity models also  
56 allow for benchmarking of relevant characteristics over time. The ability to capture data  
57 consistently and completely, to represent data via common data models, to validate the accuracy  
58 of data, and to then use the data through automated queries are examples of key processes that  
59 drive data quality. The five proposed stages of maturity reflect increasingly advanced and  
60 integrated levels of performance for health care systems to partner within the NESTcc ecosystem.  
61 The NESTcc Data Quality Maturity Model, by itself, does not ensure improvement but is rather an  
62 indicator of progress. The model can help researchers identify weaknesses, thereby enabling  
63 research teams to address them.

## 64 **Governance**

65 RWD are observational data that can be analyzed to produce RWE. RWD are defined by the U.S. Food  
66 and Drug Administration (FDA) as “data related to patient health status and/or the delivery of health care  
67 routinely collected from ... electronic health records (EHRs), claims and billing data, data from product  
68 and disease registries, patient-generated data including home-use settings, and data gathered from other  
69 sources that can inform on health status, such as mobile devices.”<sup>3</sup> To support the generation of RWE



70 from RWD, core principles must be agreed on to establish governance, including policies and processes  
71 for organizational transparency and integrity; data access, management, linkage and aggregation, and  
72 use; and submission, management, review, and acceptance of analytic requests.<sup>4,5</sup>

73 Stakeholder involvement and engagement is a critical component of good governance for RWD/RWE.  
74 The “Good Governance Standard for Public Services” has described stakeholder engagement as a core  
75 value of good governance.<sup>6</sup> As no individual party is free from bias or conflict of interest, governance  
76 provides a basis to balance stakeholder influences and provide equal representation, thereby limiting the  
77 potential for bias or unethical behaviors and allowing trustworthy research. The Patient-Centered  
78 Outcomes Research Institute (PCORI) has identified stakeholders to include patients, clinicians,  
79 researchers, purchasers, payors, industry, hospitals and health systems, policy makers, and training  
80 institutions,<sup>7</sup> and these same stakeholders remain relevant to the RWD and RWE domains. Additionally,  
81 engagement of stakeholders is necessary throughout the lifecycle of evaluation, from study and analysis  
82 planning and conduct through dissemination of results.

83 NESTcc is fully committed to ensuring that the highest scientific and ethical standards are applied when  
84 using RWD to generate RWE. In doing so, evaluation activities (e.g., sharing patient data across various  
85 data sources) must incorporate patient protections such as patient privacy (e.g., HIPAA compliance) and  
86 comply with applicable local, state, and federal laws. Institutional review board review may be necessary.  
87 The best practices developed by the FDA Sentinel program offer a template for protecting patient privacy  
88 and institutional confidentiality when linking RWD across multiple health systems.<sup>8,9</sup>

89 The following are principles to guide health systems and other clinical organizations in forming policies  
90 and procedures for RWD/RWE:

### 91 **Organizational Transparency and Integrity**

- 92 • **Leadership:** Organization establishes executive leadership group for RWD/RWE
- 93 • **Data Stewardship:** Organization takes full responsibility for the organization’s RWD
- 94 • **Patient-centered:** Patients are engaged in the RWD/RWE process and provide consent when  
95 applicable; organization adheres to ethical standards for responsible conduct of research
- 96 • **Stakeholder Engagement:** Key stakeholders, including patients, clinicians, and other health system  
97 and organization staff, are engaged in RWD/RWE project development and execution
- 98 • **Transparency:** All involved individuals from the organization are made clear to the public,  
99 potential conflicts of interest are publicly disclosed/reported, and organization’s funding is  
100 publicly disclosed
- 101 • **Oversight:** Organization assembles independent advisory board with responsibility for  
102 organization’s local data warehouse and research portfolio, which may include legal counsel to  
103 manage liability risk

### 104 **Data Access, Management, Linkage and Aggregation, and Use**

- 105 • **Data Quality Assurance:** Data are accurate and complete
- 106 • **Data Storage:** Data are securely stored, minimizing risk of secondary use or distribution without
- 107 the appropriate permissions/agreements
- 108 • **Data Permission:** Appropriate agreements are in place for all data used for RWD/RWE, data are
- 109 de-identified to the greatest extent possible, and patient protections are in place, while still
- 110 allowing necessary analyses to be pursued; if identified data are used, analyses are conducted
- 111 within secure network areas from which only aggregated or de-identified data can be removed
- 112 • **Data Linkage:** Linkage of RWD within and across sources is performed with appropriate oversight
- 113 and processes in place, particularly patient privacy protection

114 **Submission, Management, Review, and Acceptance of RWD/RWE Requests**

- 115 • **Clear Criteria:** Criteria by which requests for RWD for RWE are considered are fair and publicly
- 116 disclosed, including preclusion of access for non-scientific purposes, such as in pursuit of
- 117 litigation, as well as qualifications for data security and storage
- 118 • **Transparent Submission and Review Process:** Requests for RWD for RWE are publicly disclosed and
- 119 considered by an independent approval panel (and ethics review as needed), whose
- 120 determinations are also publicly disclosed
- 121 • **Commitment to Responsible Analysis:** Requests for RWD for RWE include a description of
- 122 collaborators (including affiliations and conflicts of interest) and proposed use of the RWD/RWE,
- 123 including the research or evaluation question, data elements of interest, main outcome
- 124 measures, and statistical analysis plan, which is publicly disclosed; considerations may be made
- 125 for studies of as-yet-unapproved uses of medical products given commercial confidentiality
- 126 • **Efficiency:** Approved requests for RWD/RWE are managed expeditiously, within time frames that
- 127 are as rapid as possible, from initiation to analysis to dissemination
- 128 • **Data Use Agreements:** Contractual requirements for data protection and privacy must be
- 129 established for any approved RWD/RWE request in compliance with appropriate laws and
- 130 regulations
- 131 • **Commitment to Results Reporting:** All analyses pursued as part of RWD/RWE projects are publicly
- 132 reported (which could potentially include the project data dictionary and analytic code, as well as
- 133 all results), including both lay and scientific summaries, regardless of plans to publish in peer-
- 134 reviewed literature, and directly communicated to FDA when issues with medical product safety
- 135 are identified; considerations may be made for studies of as-yet-unapproved uses of medical
- 136 products given commercial confidentiality

137 Leveraging the use of RWD for RWE holds great promise for medical product evaluation. The principles  
138 described above should optimize the success of these efforts among health systems and other clinical

139 organizations, protect patient privacy, and guide the governance of policies and procedures for  
140 RWD/RWE.

## 141 **Characteristics of Data**

142 Generating evidence to inform and guide clinical and regulatory decisions requires data. For data to be  
143 useful, they must be both reliable (high quality) and relevant (fit to purpose) across a broad and  
144 representative population. A full understanding of the evaluation question(s) is a prerequisite for  
145 determining the assessments, outcomes, and endpoints needed for analysis, as well as the sources,  
146 settings, and methodologies needed for data accumulation or acquisition. Choice and use of data require  
147 understanding the limitations of the data source(s) and acknowledging that the shortcomings of the data  
148 may limit the questions that can be addressed. For example, retrospective observational data acquired  
149 from real-world sources, including EHR data, though typically more pragmatic and accurate for addressing  
150 real-world practice and outcomes of device use (i.e., questions of generalizability), lack the precision of  
151 data prospectively acquired in exploratory randomized clinical trials (i.e., questions of causality).  
152 However, rigorously designed prospective clinical trials that include assignment of therapy,  
153 randomization, and/or blinding can be embedded in existing RWD sources such as registries, supporting  
154 questions that address both causality and generalizability.<sup>10</sup>

155 The characteristics of satisfactory data are predicated upon a detailed understanding of the question  
156 which allows the investigator to prospectively define:

- 157 • The appropriate study population;
- 158 • The specific data elements required to measure device or medical product utilization;
- 159 • The specific data elements required to assess performance and outcomes (including adverse  
160 events and their timing) in the course of the disease or treatment;
- 161 • The appropriate settings and sources for data acquisition;
- 162 • The development or revision of standardized data sets (i.e., development of a common data  
163 dictionary for common data elements and key outcomes and endpoints);
- 164 • The experimental methods required (e.g., causal inference from prospective randomized  
165 controlled trial vs. informed decision making from available or collected observational data)

166 To generate information and evidence of sufficient quality for generating actionable insights and  
167 informing clinical or regulatory decisions, data must satisfy four characteristics:<sup>3</sup>

- 168 1. High quality;
- 169 2. Relevant to purpose and context;
- 170 3. Amendable to the application of appropriate analytic methods (i.e., convertible to evidence);
- 171 4. Interpretable using clinical and scientific judgment



172 High-quality data are (to the greatest extent possible) complete, accurate, and free from errors that  
173 matter. The quality of the raw data increases when common definitional and temporal frameworks exist  
174 for disparate sources accessible for analysis. To obtain key endpoints or outcomes, adjudication, use of  
175 modular datasets with defined data elements, outcome verification from multiple sources, or other  
176 additional mechanisms might be needed to provide additional assurances of the accuracy of available  
177 data.

178 Data must be relevant or fit to purpose. This means the data are reliable and have the scope and content  
179 needed to answer the question(s) at hand. Pre-study planning and assessment of the various available  
180 data sources must be sufficient to determine whether existing data are contextually appropriate and  
181 complete, and whether additional data need to be acquired. Linkages of multiple high-quality data sets  
182 for either retrospective or prospective data generation may be used as needed to ensure all needed data  
183 are available. Accurate assessments of the totality of the data that will be available for the prespecified  
184 analyses are essential.

185 The combination and analysis of data and information is the final step in the production of evidence.  
186 Effective analysis requires the application of appropriate analytic and statistical tools. Prespecified  
187 statistical analysis plans are essential to minimize bias. Rigorous analysis makes information  
188 interpretable, transforming it into evidence. Objective evaluation of the totality of the evidence coupled  
189 with clinical and/or regulatory judgment leads to insights that can be used to inform clinical and  
190 regulatory decisions based on the question (see figure below).

191 Data and information should be viewed as a continuum, capable of developing evidence over the total  
192 product life cycle of a device or procedure.<sup>1</sup> Accessibility of the evidence as it evolves requires  
193 continuous data access coupled with seamless curation, analysis, and interpretation. Integrated data  
194 solutions that allow permanent linkages between previously isolated sources of data and development of  
195 open standards will foster a cooperative environment where duplication and costs are minimized and the  
196 value of evidence and the underlying infrastructure is maximized.<sup>11,12</sup>

197

198

199

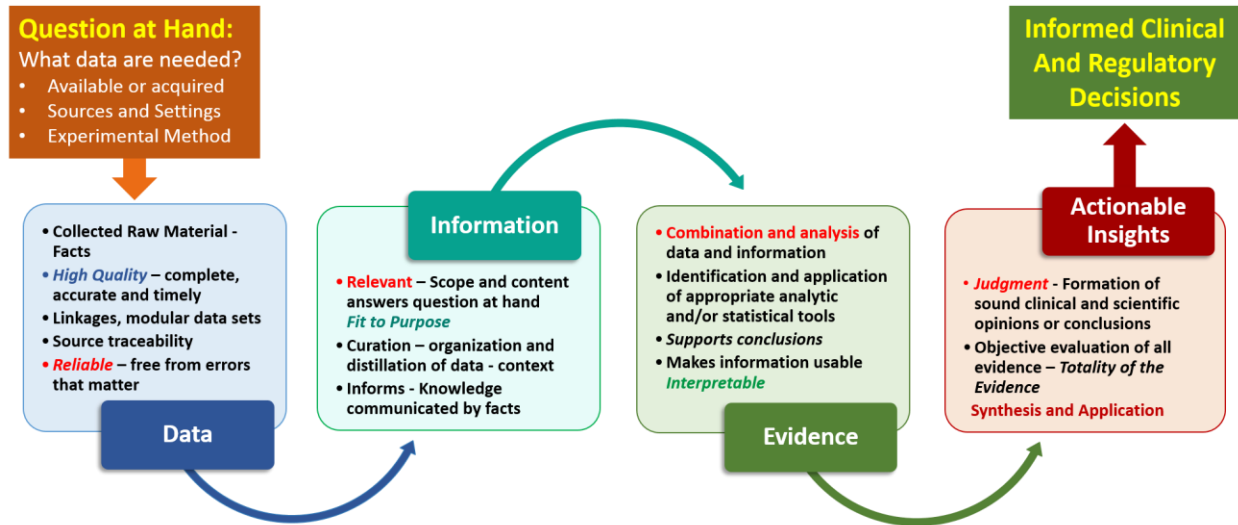
200

201

202

203 Figure 1. Evidence generation and evaluation: Actionable insights for informed clinical and regulatory  
204 decisions (adapted from Califf RM, Sherman R, What we mean when we talk about data. MassDevice.  
205 December 11, 2015. <https://www.massdevice.com/44947-2/>)

206  
207



208  
209

## Data Capture and Transformation

210 The use of EHR data for research purposes poses additional challenges to data relevance and reliability.  
211 Standardizing definitions for identifying patient cohorts and study endpoints or outcomes, increasing  
212 health care systems’ interoperability to capture longitudinal patient data, and universally implementing  
213 the unique device identifier (UDI) capture will improve the relevance of EHR-based medical device  
214 research. Considerations for use of EHR data to conduct research also include understanding  
215 provenance, completeness, accuracy and consistency of the data, as well as awareness of what internal  
216 and external validation checks have been performed to evaluate the quality of data entry.<sup>13</sup>

217 Regardless of improvements in data collection systems to accommodate EHR-based research, researchers  
218 have little control over data recording and collection processes in clinical care facilities. Individuals using  
219 EHR data to derive RWE should understand how and why the data of interest were originally obtained as  
220 well as data provenance including subsequent data processing and other nuances that might affect  
221 reliability of the data. This understanding will help the researcher determine whether the EHR data are of  
222 suitable quality for a particular evaluation.

223 Improving data quality at the point of care and point of data entry should be the ultimate goal. Wherever  
224 possible, the key stakeholder communities should agree regarding the clinical concepts that need to be  
225 captured as data for use within the medical device evaluation ecosystem. This might be accomplished at  
226 the point of data entry (e.g., when tied to reimbursement or as required in clinical decision support). In  
227 these contexts, clinical concepts must be specified and defined as domain-specific common data  
228 elements (CDEs), which ideally use standardized definitions and are harmonized with common data  
229 models (CDMs) for optimal utility. Currently, most clinical information in an EHR is conveyed as free text



230 versus the ideal state where EHR data capture would be predominantly in discrete, structured fields.  
231 Clinical workflows and documentation systems will likely require modifications to ensure capture of  
232 structured data at the point of care.

233 Once data are captured as discrete elements, extraction, transformation, and loading (ETL) are more  
234 amenable to standardization. Discrete data elements, semantic interoperability, compatibility of data  
235 capture, and appropriate specification of CDM conventions allow for application of a CDM that  
236 subsequently permits execution of standardized analyses or queries by data partners.

237 Current capabilities may only allow a hybrid approach that combines auto-populating certain discrete  
238 structured data elements (e.g., demographics, numerical values, ICD codes) complemented by manual  
239 abstraction of other data into a CDM. Alternatively, or in combination, a process such as natural language  
240 processing could be used to obtain the data of interest from unstructured text. Of note, natural language  
241 processing cannot synthesize data elements or derive inferential conclusions. Moving toward greater  
242 agreement and use of computable phenotypes to assist with population identification and perhaps  
243 endpoint or outcome identification might address this issue.<sup>14</sup> Such an approach would need to be as  
244 comprehensive as possible and include input from all members of the health care ecosystem.

245 The capture of quality data is one component that determines the quality of the ETL process, which  
246 describes how data are extracted and transformed to conform to data standards and CDM specifications.  
247 These data are then loaded into a defined location and available for queries (e.g., via a distributed  
248 research network).<sup>15</sup> Additional consideration must also be given to applicable patient privacy  
249 requirements and agreements or contracts.

250 Designing the ETL process should follow established best practices, such as seeking input from CDM and  
251 data experts to design the ETL process, clinical experts to create coding maps for the process, technical  
252 experts to implement the process, and all stakeholders to design and implement quality control  
253 procedures.<sup>16</sup>

254 Data assurance and quality control are essential to the reliability of the RWD for RWE generation. Quality  
255 control processes should be integrated throughout, including a review of the ETL design documentation  
256 and verification and validation of each step of the ETL process.<sup>16</sup>

257 Consistency in data element definitions on the data capture side, along with the use of standards to  
258 support consistency in the ETL process, will allow researchers to have confidence in the quality of the  
259 data extracted from the EHR. Data aggregation is relatively straightforward when data are captured and  
260 transformed consistently and reproducibly.

261

262

## 263 **Data Curation**



264 Data curation is one of the steps used to turn raw data into information. Through the curation process,  
 265 data are organized, assessed, and prepared for analysis. Many frameworks exist to guide this translation  
 266 of RWD into fit-for-purpose data, but one approach is to consider a two-stage process. The first,  
 267 foundational stage takes the raw data, applies a series of transformations and quality checks to make the  
 268 dataset “research ready.” The second, study-specific stage applies another series of transformations and  
 269 quality checks to ensure that the dataset is “fit-for-purpose” for the specific question at hand (some  
 270 networks/projects may combine both stages into a single process). Foundational data curation examines  
 271 the data repository or datamart in the context of broad research concepts (e.g., are laboratory results  
 272 mapped to an appropriate coding scheme?), whereas study-specific data curation also considers a  
 273 specific-study context (e.g., are outcomes complete for the study population?). Surveys or metadata  
 274 about data elements, the workflows that give rise to them, and source system provenance further inform  
 275 the process of data curation and, when combined with information about data latency and extraction and  
 276 transformation processes, help ensure that fitness-for-use can be assessed as needed. Examples of data  
 277 curation processes developed by distributed research networks are shown in the table below.

278 Table 1. Data curation processes for specific distributed research networks.

Network	Collaborators		Approach to Data Characterization
	Health systems	Payors	
HCSRN	X	X	Detailed checks look at ranges, cross-field agreement, implausible data patterns, and cross-site comparisons. Partners execute data characterization package each time data are refreshed. Results are returned to the HCSRN Coordinating Center. Potential quality issues are flagged and mitigated at the partner level. <sup>18</sup>
Sentinel	X	X	Detailed checks look at ranges, cross-field agreement, implausible data patterns, and cross-site comparisons. Partners execute data characterization package each time data are refreshed. Results are returned to the Sentinel Coordinating Center. Potential quality issues are flagged and mitigated at the partner level. <sup>19</sup>
PCORnet	X	X	Includes <i>foundational data curation</i> process, which establishes a baseline level of research readiness for all network partners to support prep-to-research queries, and <i>study-specific data curation</i> , which includes assessments of outcomes/variables or other derived concepts for the cohort under study. <sup>20</sup>
OHDSI	X	X	Optional – each datamart can generate a standardized data profile that is viewable through a web-based tool (Achilles). Institutions can choose whether to share these profiles or retain them locally. <sup>21</sup>
ACT	X		Under development.

279  
 280 ACT = Accrual for Clinical Trials; HCSRN = Health Care Systems Research Network; OHDSI = Observational  
 281 Health Data Sciences and Informatics; PCORnet = National Patient-Centered Clinical Research Network.



282 Key to the curation process are *data characterization routines*, which run against a collaborator’s data  
283 repository or CDM and describe their performance against a series of *data quality checks* through  
284 descriptive statistics such as summaries of missing values, outliers, and frequency distributions. Many  
285 data checks rely on concepts analogous to conformance (“does the format of the data adhere to the  
286 underlying model?”), completeness (“are there values where we expect to see data populated?”), and  
287 plausibility (“do the values that appear make sense?”), as well as comparisons across collaborators.<sup>17</sup> As  
288 an example, the most recent PCORnet data characterization process consists of a set of SAS procedures  
289 that execute against the tables of the PCORnet CDM.<sup>22</sup> There are 31 unique data checks,<sup>23</sup> many of which  
290 apply to multiple fields or tables within the CDM (e.g., required fields are present, tables do not have  
291 orphan patient identifiers), for a total of 1,144 individual quality queries. These routines also generate  
292 additional tables of descriptive statistics, including the frequencies of specific data elements, crosstabs of  
293 data (e.g., procedure and procedure type), and counts of missing, non-missing, and distinct records. FDA  
294 Sentinel follows a similar process<sup>24</sup> and, as described below, NESTcc expects collaborators to utilize an  
295 approach that is suitable for the dataset and the question(s) being asked. While the data characterization  
296 routines are necessarily designed to assess quality within a collaborator’s data repository, the summary  
297 results are aggregated and analyzed across a network to establish baseline trends and identify outliers or  
298 other anomalies.

### 299 **Metadata About Data Provenance**

300 The results of data characterization alone are not always enough to determine whether a given data set is  
301 fit to purpose. Information on provenance also plays a role, as there is widespread variability in how data  
302 are entered into EHRs or processed as claims, as well as how health systems and health plans extract  
303 those data to populate a given table within their repository or CDM. Knowledge about data collection  
304 practices and the decisions made to translate the source material into the target CDM can help provide  
305 additional context.<sup>25</sup> Many networks ask their collaborators to complete surveys that describe the  
306 provenance of their data sources, providing additional insight into the characteristics of their clinical  
307 workflows and/or source systems.<sup>26,27</sup> In some cases, provenance can also be derived automatically as  
308 part of the data capture or data transformation process (e.g., did the record originate from a billing  
309 system, or was it entered by a clinician?). This is important, because in studies on inpatient medication  
310 usage, for instance, one must know whether a datamart has included records only for medications that  
311 were administered to patients or all medications that were ordered, including prescriptions written  
312 prophylactically (or both), as they will generate markedly different characterization profiles.

### 313 **Documentation of the Iterative Process of Data Curation**

314 Data curation is an iterative process, with the expectation that characterization activities will help quality  
315 improve over time. Therefore, the operational definition of a given data check should stay consistent to  
316 allow comparisons over time. Networks may have data checks that are required or investigative. Given  
317 the variability in health system data, networks often limit required checks to those related to  
318 conformance. Investigative data checks may be remediable by a health system (e.g., >80% of laboratory  
319 results have a Logical Observation Identifiers Names and Codes [LOINC] code), or not be remediable due  
320 to source system limitations (e.g., <10% of medication orders include an end date). Investigative data



321 checks that are broadly remediable across the network are good candidates for having thresholds that  
322 are raised or lowered to reflect improvements in data quality (e.g., requiring that >50% of laboratory  
323 results be mapped to LOINC initially, gradually raising the minimum threshold to >80% as collaborators  
324 develop their mappings). Collaborators should track their efforts to address failed investigative data  
325 checks and networks should ensure that they perform purpose-specific curation for the  
326 population/question in these areas in order to determine whether the data support the study of interest.  
327 All of these steps should be documented and included as part of any analysis plan. As the base of RWE  
328 studies grows and the FDA releases more guidance, we expect to see best practices and standards  
329 emerge as to how to convey this information.

### 330 **Data Curation Should Be Fit-for-purpose**

331 The minimum requirements for data curation will vary depending on the dataset and the study, but there  
332 should be sufficient evidence that the data can answer the question of interest within the context of the  
333 intended use. For example, studies of overall utilization patterns for exploratory analyses will require a  
334 different level of certainty than a comparative study intended for policy or regulatory decision-making.  
335 Studies that use data from emerging domains (e.g., patient-generated data, information derived from  
336 natural language processing) may require a higher level of interrogation than a prep-to-research query  
337 using a well-known data source. Collaborators that participate in distributed research networks with  
338 formalized curation processes may be “pre-cleared” to support a range of activities if their data pass all  
339 relevant checks. Collaborators that are not part of any existing network will need to decide how much to  
340 invest in data curation. Ensuring that the resulting dataset can be used to answer operational questions  
341 that are of value to the health system/health plan is one way to justify the potential expense.  
342 Collaborators with data that have only been subjected to a cursory level of curation may still be able to  
343 participate but may find themselves restricted to high-level or preliminary exercises.

### 344 **NESTcc Data Quality Maturity Model**

345 Organizational maturity can be described as an expression of the capabilities of an organization in a  
346 specific domain, with the intent to foster continuous improvement across those capabilities. Maturity  
347 models organize levels of maturity into a framework, typically assessing culture, process, and/or  
348 technology.<sup>28</sup> Maturity models are typically self-administered by organizations to assess current state,  
349 model business capabilities, identify opportunities, and perform capacity planning. A key benefit is the  
350 benchmarking of relevant characteristics over time. In health care, the Healthcare Information and  
351 Management Systems Society (HIMSS) has published several maturity models, including a Health IT  
352 Usability Maturity Model ([www.himss.org/himss-usability-maturity-model](http://www.himss.org/himss-usability-maturity-model)), EHR Adoption Model  
353 ([www.himssanalytics.org/emram](http://www.himssanalytics.org/emram)), and Adoption Model for Analytics Maturity  
354 ([www.himssanalytics.org/amam](http://www.himssanalytics.org/amam)). Specific to models developed for enterprise data governance, a  
355 detailed descriptive model from Stanford addresses the axes of people, policies, and capabilities across  
356 the dimensions of awareness, formalization, metadata, stewardship, data quality, and master data.<sup>29</sup>

357 To articulate a high level of expectations at different levels of organizational maturity with respect to  
358 RWD quality, we have developed the NESTcc Data Quality Maturity Model. The model is based on the



359 expectations of health care systems regarding source systems for RWD capture and management,  
 360 principally via EHR and other clinical documentation systems.

361 We propose five stages of maturity of increasingly advanced and integrated levels of performance for  
 362 health care systems to partner within the NESTcc ecosystem. The stages are at least partially aligned with  
 363 previous maturity models, of which the HIMSS Usability Model is most informative:

364 Table 2. Comparability of Stages of NESTcc and Other Maturity Models

NESTcc Stage	HIMSS Usability Model	Capability Maturity Model Integration (CMMI) Model	Stanford Model
1. Conceptual	Unrecognized	Initial	Awareness
2. Reactive	Preliminary	Managed	Formalization
3. Structured	Implemented	Defined	Stewardship
4. Complete	Integrated	Quantitatively managed	Data quality
5. Advanced	Strategic	Optimizing	Master data

365  
 366 **Stage 1** – Clinical processes capture information primarily in verbose, unstructured documents, not as  
 367 discrete data; lack of organizational awareness of data utility, no effort to systematically manage health  
 368 care data, lack of consistent or centralized governance, policies, and/or resources, data not organized  
 369 centrally; data not available for organizational use and analysis; individual data units are project oriented  
 370 or focused on immediate profits.

371 **Stage 2** – Able to react to requests for analysis, respond to research requests – but mostly accomplished  
 372 by manual chart review and abstraction; data management inefficient and expensive, with only sporadic  
 373 recognition of data utility beyond immediate use; tacit support from leadership regarding need for  
 374 centralized data governance and management, but only limited allocation of resources; data not available  
 375 for organizational use and analysis beyond individual requests; individual data units are project-oriented  
 376 or focused on immediate profits.

377 **Stage 3** – Clinical systems manage transactional data types (e.g., orders, transactions, laboratory results,  
 378 medication prescriptions) as discrete data; support from leadership (with resources provided) for  
 379 centralized data governance and management of these data types at the enterprise level (e.g., support  
 380 for ETL among internal systems); commitment to centralized enterprise data governance, management,  
 381 and curation via managed processes, people, and technologies (e.g., enterprise data warehouse [EDW]);  
 382 non-administrative queries (clinical questions, research) conducted mostly as one-offs via individual  
 383 queries, still moderate-to-high cost to extract data for analysis; able to support a CDM but not done  
 384 routinely and automatically; data transmission to registries still largely accomplished by manual chart  
 385 review and abstraction.

386 **Stage 4** – Granular and complete clinical data based on standardized clinical CDEs captured in the  
 387 processes of care, integrated into those care processes; UDI captured in the processes of care and  
 388 available in EHR and in the EDW; EDW routinely and systematically represents data externally via various



389 CDMs, including efficient queries, support for large number of research projects; leadership provides  
 390 centralized data governance, management, and curation at the enterprise level, ensuring performance  
 391 and data quality of local units and achieving financial sustainability.

392 **Stage 5** – Data linkage and aggregation across systems enabled and open to external queries;  
 393 interoperability of clinical data enabled; multiple sources of sustainable funding support for research;  
 394 engagement of regulatory and industry enterprises with enterprise data; leadership responsible for  
 395 centralized data governance, management, and curation at the enterprise level, business benefit well  
 396 understood, with financial sustainability, and recognition and participation in initiatives external to the  
 397 organization.

398 **Key Data Process Domains that Drive Data Quality**

399 Optimally, use of health care system RWD requires competency across several data process domains,  
 400 including data consistency, completeness, and automation.<sup>13</sup> Building on those data process domains,  
 401 the table below describes expectations at each NESTcc maturity stage. A foundational requirement is  
 402 **consistent** clinical data based on standardized data dictionaries and/or applicable data standards. While  
 403 data consistency can be most easily understood within the confines of an individual health care  
 404 organization, ideally the data are semantically interoperable (i.e., have the same clinical and  
 405 computational meaning) across organizations. Once standards have been implemented, the ability to  
 406 capture complete data sets (including interpretation and accounting of the absence of data) characterizes  
 407 the data **completeness** domain. The ability to represent data via **CDMs**, to validate the **accuracy** of data,  
 408 and to then use the data through **automation** of queries are additional domains that describe business  
 409 capabilities related to data quality.

410 Table 3. Organizational Operational Characteristics Typical of NESTcc Maturity Model Stages

	NESTcc Data Quality Domain				
	Consistency <sup>a</sup>	Completeness <sup>b</sup>	CDM <sup>c</sup>	Accuracy <sup>d</sup>	Automation <sup>e</sup>
1. Conceptual					
2. Reactive	+	+	+/-		
3. Structured	+	+	+	+/-	
4. Complete	+	+	+	+	+
5. Advanced	+	+	+	+	+

411  
 412 <sup>a</sup>Data Consistency: Relevant uniformity in data: Across all hospitals, providers, and outpatients (e.g.,  
 413 population/cohort identification, clinical documentation practices/policies between entities, workflow  
 414 descriptions)

415 <sup>b</sup>Data Completeness: Presence of the necessary data elements for outcome assessment, CDEs used, all  
 416 data are electronically available and either complete or with little missing data

417 <sup>c</sup>Data Models: CDMs include all data needed for decision making (e.g., clinical data elements, UDI)



418 <sup>d</sup>Data Accuracy: Validation: EHR data are validated systematically, with comparison to the source,  
419 independent measurement, upstream data source, and known standard or valid values (e.g., audits from  
420 charts)

421 <sup>e</sup>Data Automation: Queries able to be run automatically against CDMs

## 422 **Conclusion**

423 High-quality data are essential to support the post-market evaluation of medical devices and to inform  
424 regulatory decision-making. In this initial version of the NESTcc Data Quality Framework, we discuss the  
425 most salient topics associated with achieving high-quality data including data governance, characteristics  
426 of data, approaches to data capture and transformation, and best practices in data curation. We  
427 synthesize these topics in the NESTcc Data Quality Maturity Model, which enables collaborators to  
428 indicate their progress toward achieving the highest quality data. The next iteration of this framework  
429 will include the NESTcc Data Quality Self-Evaluation, a checklist that charts the specific actions  
430 organizations can take to move between stages of the maturity model. We welcome further discussion  
431 about how the framework can be operationalized by health systems, given the variability in maturity  
432 among individual clinics that compose a health system.

## 433 **References**

- 434 1. Shuren J, Califf RM. Need for a national evaluation system for health technology. *JAMA*.  
435 2016;316(11):1153-4.
- 436 2. Agency for Healthcare Research and Quality. Registries for Evaluating Patient Outcomes: A User's  
437 Guide: 3rd Edition. Research Report. April 30, 2014.  
438 <https://effectivehealthcare.ahrq.gov/topics/registries-guide-3rd-edition/research>.
- 439 3. U.S. Food and Drug Administration. Use of Real-World Evidence to Support Regulatory Decision-Making  
440 for Medical Devices: Guidance for Industry and Food and Drug Administration Staff. August 31,  
441 2017.  
442 <https://www.fda.gov/downloads/medicaldevices/deviceregulationandguidance/guidancedocument/s/ucm513027.pdf>.
- 443
- 444 4. Cole A, Garrison L, Mestre-Ferrandiz J, Towse A. Data Governance Arrangements for Real-World  
445 Evidence. Office of Health Economics Consulting. November 2015.  
446 <https://www.ohe.org/system/files/private/publications/420%20-%20Data%20Governance%20for%20RWE.pdf>.
- 447
- 448 5. International Medical Device Regulators Forum, IMDRF Registry Working Group. Patient Registry:  
449 Essential Principles. October 2, 2015. <http://www.imdrf.org/docs/imdrf/final/consultations/imdrf-cons-essential-principles-151124.pdf>.
- 450
- 451 6. The Independent Commission on Good Governance in Public Services. The Good Governance Standard  
452 for Public Services. 2004. Office for Public Management, Chartered Institute of Public Finance and  
453 Accountancy, and Joseph Rowntree Foundation. <https://www.jrf.org.uk/report/good-governance-standard-public-services>.
- 454

- 455 7. Patient-Centered Outcomes Research Institute. The Value of Engagement. Posted: October 30, 2018.  
456 <https://www.pcori.org/engagement/what-we-mean-engagement>.
- 457 8. Sentinel. Background. <https://www.sentinelinitiative.org/background>.
- 458 9. Evans BJ. Panel 3: Appropriate Human-Subject Protections for Research Use of Sentinel System Data.  
459 Engelberg Center for Health Care Reform at Brookings. FDA Sentinel Initiative Meeting Series: Issue  
460 Brief. March 2010. [https://www.brookings.edu/wp-content/uploads/2012/04/Panel-3-Issue-](https://www.brookings.edu/wp-content/uploads/2012/04/Panel-3-Issue-Brief.pdf)  
461 [Brief.pdf](https://www.brookings.edu/wp-content/uploads/2012/04/Panel-3-Issue-Brief.pdf).
- 462 10. Frobert O, Lagerqvist B, Olivecrona GK, et al. Thrombus aspiration during ST-segment elevation  
463 myocardial infarction. *N Engl J Med* 2013;369:1587-97.
- 464 11. Califf RM. Benefit-risk assessments at the US Food and Drug Administration: finding the balance.  
465 *JAMA*. 2017;317(7):693-4.
- 466 12. Faris O, Shuren J. An FDA viewpoint on unique considerations for medical-device clinical trials. *N Engl J*  
467 *Med*. 2017;376:1350-7.
- 468 13. Zozus MN, Hammond WE, Green BB, et al. Assessing Data Quality for Healthcare Systems Data Used  
469 in Clinical Research. Version: 1.0, last updated July 28, 2014. NIH Collaboratory.  
470 [https://www.nihcollaboratory.org/Products/Assessing-data-quality\\_V1%200.pdf](https://www.nihcollaboratory.org/Products/Assessing-data-quality_V1%200.pdf).
- 471 14. Richesson R, Smerek M, Rusincovitch S, et al. Electronic Health Records-Based Phenotyping. NIH  
472 Collaboratory Living Textbook of Pragmatic Clinical Trials. June 27, 2014.  
473 <http://rethinkingclinicaltrials.org/resources/ehr-phenotyping/>.
- 474 15. PCORnet: The National Patient-Centered Clinical Research Network. PCORnet Glossary. December  
475 2015. [https://pcornet.org/wp-content/uploads/2014/07/PCORnet\\_CDM\\_Glossary.pdf](https://pcornet.org/wp-content/uploads/2014/07/PCORnet_CDM_Glossary.pdf).
- 476 16. Observational Health Data Sciences and Informatics. ETL Creation Best Practices. Last modified June  
477 28, 2017. [http://www.ohdsi.org/web/wiki/doku.php?id=documentation:etl\\_best\\_practices](http://www.ohdsi.org/web/wiki/doku.php?id=documentation:etl_best_practices).
- 478 17. Kahn MG, Callahan TJ, Barnard J, et al. A harmonized data quality assessment terminology and  
479 framework for the secondary use of electronic health record data. *EGEMS (Wash DC)*.  
480 2016;4(1):1244.
- 481 18. Health Care Systems Research Network. Data Resources. <http://www.hcsrn.org/en/About/Data/>.  
482 Accessed January 26, 2018.
- 483 19. Sentinel. Sentinel Operations Center. Sentinel Common Data Model - Data Quality Review and  
484 Characterization Process and Programs. Program Package version: 3.3.4. 2017.  
485 <https://www.sentinelinitiative.org/sentinel/data-quality-review-and-characterization>. Accessed  
486 January 26, 2018.
- 487 20. Qualls LG, Phillips TA, Hammill BG, et al. Evaluating foundational data quality in the National Patient-  
488 Centered Clinical Research Network (PCORnet®). *EGEMS (Wash DC)*. 2018;6(1):3.
- 489 21. Observational Health Data Sciences and Informatics. ACHILLES for Data Characterization.  
490 <https://www.ohdsi.org/analytic-tools/achilles-for-data-characterization/>. Accessed January 26,  
491 2018.
- 492 22. PCORnet Distributed Research Network Operations Center. PCORnet Data Curation Query Package.  
493 <https://github.com/PCORnet-DRN-OC/PCORnet-Data-Curation>. Accessed January 26, 2018.





- 494 23. PCORnet: The National Patient-Centered Clinical Research Network. PCORnet Data Checks v5.  
495 <https://pcorner.org/wp-content/uploads/2018/05/PCORnet-Data-Checks-v5.pdf>.
- 496 24. Sentinel. Sentinel Data Quality Assurance Practices.  
497 [https://www.sentinelinitiative.org/sentinel/data/distributed-database-common-data-](https://www.sentinelinitiative.org/sentinel/data/distributed-database-common-data-model/sentinel-data-quality-assurance-practices)  
498 [model/sentinel-data-quality-assurance-practices](https://www.sentinelinitiative.org/sentinel/data/distributed-database-common-data-model/sentinel-data-quality-assurance-practices).
- 499 25. Johnson KE, Kamineni A, Fuller S, Olmstead D, Wernli KJ. How the provenance of electronic health  
500 record data matters for research: a case example using system mapping. *EGEMS (Wash DC)*.  
501 2014;2(1):1058.
- 502 26. Curtis LH, Weiner MG, Boudreau DM, et al. Design considerations, architecture, and use of the Mini-  
503 Sentinel distributed data system. *Pharmacoepidemiol Drug Saf*. 2012;21(Suppl 1):23-31.
- 504 27. Brown JS, Kahn M, Toh S. Data quality assessment for comparative effectiveness research in  
505 distributed data networks. *Med Care*. 2013;51(8 Suppl 3):S22-29.
- 506 28. Mettler T. Maturity assessment models: a design science research approach. *Int J Soc Syst Sci*.  
507 2011;3(1/2):213–222.
- 508 29. Stanford University. Data Governance Maturity Model: Guiding Questions for Each Component-  
509 Dimension. [http://web.stanford.edu/dept/pres-provost/cgi-bin/dg/wordpress/wp-](http://web.stanford.edu/dept/pres-provost/cgi-bin/dg/wordpress/wp-content/uploads/2011/11/StanfordDataGovernanceMaturityModel.pdf)  
510 [content/uploads/2011/11/StanfordDataGovernanceMaturityModel.pdf](http://web.stanford.edu/dept/pres-provost/cgi-bin/dg/wordpress/wp-content/uploads/2011/11/StanfordDataGovernanceMaturityModel.pdf).