CrossMark

COMMENTARY AND DISCUSSION ARTICLE

# The creation, management, and use of data quality information for life cycle assessment

Ashley Edelen[1] · Wesley W. Ingwersen[2]

## Abstract

*Purpose* Despite growing access to data, questions of "best fit" data and the appropriate use of results in supporting decision making still plague the life cycle assessment (LCA) community. This discussion paper addresses revisions to assessing data quality captured in a new US Environmental Protection Agency guidance document as well as additional recommendations on data quality creation, management, and use in LCA databases and studies.

*Approach* Existing data quality systems and approaches in LCA were reviewed and tested. The evaluations resulted in a revision to a commonly used pedigree matrix, for which flow and process level data quality indicators are described, more clarity for scoring criteria, and further guidance on interpretation are given.

*Discussion* Increased training for practitioners on data quality application and its limits are recommended. A multi-faceted approach to data quality assessment utilizing the pedigree method alongside uncertainty analysis in result interpretation is recommended. A method of data quality score aggregation is proposed and recommendations for usage of data quality scores in existing data are made to enable improved use of data quality scores in LCA results interpretation. Roles for data generators, data repositories, and data users are described in LCA data quality management. Guidance is provided on using data with data quality scores from other systems alongside data with scores from the new system. The new pedigree matrix and recommended data quality aggregation procedure can now be implemented in openLCA software.

*Future work* Additional ways in which data quality assessment might be improved and expanded are described. Interoperability efforts in LCA data should focus on descriptors to enable user scoring of data quality rather than translation of existing scores. Developing and using data quality indicators for additional dimensions of LCA data, and automation of data quality scoring through metadata extraction and comparison to goal and scope are needed.

✉ Wesley W. Ingwersen
  ingwersen.wesley@epa.gov

[1] Oak Ridge Institute for Science and Education (ORISE), Oak Ridge, TN, USA

[2] Life Cycle Assessment Center of Excellence, National Risk Management Research Laboratory, United States Environmental Protection Agency, Cincinnati, OH, USA

# 1 Introduction

The average life cycle assessment (LCA) model combines thousands of data points in order to describe a product system. LCA practitioners and generators are very familiar with the labor and time that accompanies data collection and processing. While the amount of life cycle inventory data is growing, and there are efforts to improve access to LCA data, questions of "best fit" data and the appropriate use of results in supporting decision making still plague the LCA community. Ultimately, these are questions of data quality.

A number of authors have called out data quality as an aspect limiting the power and reliability of LCA results (Björklund 2002; Coulon et al. 1997). In a survey of

🖆 Springer

unresolved problems in LCA, data quality was one of the problems identified to be only partially solved by existing methods (Reap et al. 2008). Key weaknesses of the current methodologies for data quality assessment (DQA) in LCA include limited coverage of data quality, limited aggregation, one-dimensional analysis of data quality (e.g., process or flow), and lack of reproducibility of results (Weidema 1998; Cooper and Kahn 2012). Another reason that data quality continues to be an elusive element in LCA is the lack of clarity surrounding the practice of data quality management. We distinguish between the terms data quality, data quality assessment, and data quality management to better clarify the multi-dimensional concept of data quality.

## 1.1 Data quality in LCA

Consideration of guidelines for data quality in LCA can be traced at least as far back as 1992 (Fava 1992). The international standards organization (ISO) maintains two standards, ISO 14040 and ISO 14044, that define data quality within LCA as

> "[the] characteristics of data that relate to their ability to satisfy stated requirements"
> (ISO 2006a, b)

ISO establishes ten key characteristics of data quality: time-related coverage, geographical coverage, technological coverage, precision, completeness, representativeness, consistency, reproducibility, sources of the data, and uncertainty of the information (ISO 2006a, b). In 2011, in an effort to provide further guidance, the Global Guidance Principles for LCA Databases (Shonan Guidance Principles) published further definition and suggested methods for addressing data quality (UNEP/SETAC 2011). The Global Guidance Principles similar to the ISO standards recommends ten data quality indicators and provides definitions for each of these indicators, but only states that a dataset developer is responsible for preparing a self-assessment of the data against these indicators. Neither ISO 14044, nor the Global Guidance Principles provide guidelines for data quality assessment. The US EPA defines data quality assessment as

> "The scientific and statistical evaluation of data to determine if data obtained from environmental operations are of the right type, quality, and quantity to support their intended use."
> (US EPA 2000)

Data quality assessment in the context of LCA has been defined as a comparative analysis of the data quality characteristics (DQCs) against the data quality goals, or qualitative statement that defines specifications for the adequacy of data

used in a life cycle inventory (LCI) or for certain LCI parameters (Bakst et al. 1995). This process of determining adequacy of data relevant to goals or purpose is related to the concept of "fitness for purpose." One of the objectives of the Global Life Cycle Data Access network (GLAD) metadata working group is to select indicators for assessing fitness for purpose (Canals et al. 2016). The GLAD metadata working group also distinguishes between types of indicators, either as intrinsic or contextual (Ciroth et al. 2017). Intrinsic data quality is an inherent data property that describes the data quality (Wang and Strong 1996). Contextual data quality is an aspect of the data that is situationally dependent (Wang and Strong 1996). Reliability is classified as an intrinsic indicator because users are always looking for believable, objective data with high accuracy and from a reputable source. Whereas, representativeness indicators are contextual because they are situationally dependent on the goal and scope of the project.

## 1.2 Approaches to data quality assessment in LCA databases

Of approximately 74 databases with life cycle data identified by the authors in a previous study that surveyed LCI, input-output and carbon footprinting databases, approximately 16 (22%) include data quality scores or other information (Edelen and Ingwersen 2015). Even fewer existing databases provide a description of their data quality assessment methods. Of those methods or data quality systems that are in current use, they can be classified as either qualitative or semi-quantitative. There are two dominant examples of semi-quantitative methods: the pedigree matrix approach utilized by ecoinvent, and the data quality ranking system utilized by the International Reference Life Cycle Data System (ILCD). A qualitative pass/fail method (Cooper and Kahn 2012) is used by the USDA LCA Commons. This method can be considered a binary qualitative pedigree matrix approach.

The two types of pedigree matrix methods use data quality indicators (DQI) linked with characteristic of the data quality, often related to the ISO 14044 DQCs. Table 1 shows the presence of the ISO 14044 DQCs in these three sources.

The qualitative assessment method used in the USDA LCA Commons includes seven indicators to address nine out of ten DQCs (excluding consistency) as defined by the ISO standards (Cooper and Kahn 2012). Indicators are applied at the flow level with no available methodology for aggregation. The system proposes a two-tiered scale for scoring data, as either pass or fail, using a score of A for pass or B for fail. The non-numeric system was developed to improve reproducibility of scores and allow for a broader adoption of DQA within the LCA community (Cooper and Kahn 2012). The ILCD DQA method includes six indicators, applied at the process level, scored based on compliance with a set of pre-defined standards for each indicator (EC JRC 2010). Indicator scores

**Table 1** Comparison of ISO 14044 data quality characteristics captured by data quality indicators from three LCI data providers

|  | ILCD | ecoinvent | USDA |
|---|---|---|---|
| Characteristics |  |  |  |
| Time-related coverage | x | x | x |
| Geographic coverage | x | x | x |
| Technological coverage | x | x | x |
| Precision | x[a] |  | x |
| Completeness | x | x | x |
| Representativeness | x | x | x |
| Consistency | x |  |  |
| Reproducibility |  |  | x[b] |
| Sources of the data |  | x | x[b] |
| Uncertainty of the information | x[a] | x[c] | x |

[a] Uncertainty is combined with precision

[b] Source of the data is combined with the reproducibility

[c] Uncertainty is calculated from the indicators, but is not a separate indicator

can vary from 1 (very good) to 5 (very poor or unknown), with an option of a "0" score for not applicable indicators. In the ILCD system, scores for each DQI are aggregated across indicators to formulate an overall data quality ranking score for the process. This process overall data quality score is then used in a three-tiered classification system: high quality, basic quality, and data estimate. The pedigree matrix DQA method used by ecoinvent was originally described by Weidema and Wesnaes (1996) and has evolved through the development of the ecoinvent database (Wernet et al. 2016; Weidema et al. 2013; Frischknecht et al. 2007). The DQIs in the pedigree matrix are applied at the flow level, or in other words, to each individual exchange in a unit process, and like the ILCD system, use a 1–5 scoring system, where 1 is the best and 5 is the poorest.

### 1.3 Reproducibility of data quality scores

As metrics that can be used by third parties to evaluate the quality of LCI datasets, data quality scores should be reproducible. Previous studies have shown problems in data quality score reproducibility, which has been attributed to lack of clarity in terminology and poor training in methodology (Weidema 1998). Following a similar test described by Weidema (1998), we conducted a data quality scoring reproducibility test among LCA practitioners at US EPA using two of the DQA methods described above. The details of the methodology and results of this reproducibility test are described in the Electronic Supplementary Material. The result of the reproducibility test suggests a lack of fundamental knowledge in

applying data quality to LCI datasets. The poor reproducibility in the test of both DQA assessment methods suggests that the problem of reproducibility still exists, and that efforts are still needed to improve it.

### 1.4 Aggregation of data quality scores and usage to perform uncertainty analysis

DQA methods have the potential to inform the practitioner about quality issues that are relevant to the LCA results in order to determine how well the underlying LCI data and model fulfill the goal and scope of the study. When models are composed of many (3 or more orders of magnitude for large background databases) processes, it becomes impractical to examine all processes directly with the intention of drawing conclusions on data quality. Even in smaller models, not all processes have equal quantitative influence (gravity) on the result. Therefore, there needs to be a method of examining data quality alongside results that aggregates data quality scores in a manner acceptable with LCA conventions. May and Brennan (2003) provide a practical comparison of some of the proposed approaches to data quality score aggregation across a LCI in an application to a comparative electricity LCA. Particularly, the authors describe and test methods by Wrisberg et al. (1997) and Rousseaux et al. (2001). Wrisberg et al. proposed calculating an average data quality score for a particular flow in the LCI giving equal weight to each exchange where the flow appears, as in Eq. (1):

$$\frac{\sum_e^n DQS_{i,f,e}}{n} = \text{LCIDQS}_{f,i} \tag{1}$$

where DQS is the flow data quality score for a given DQI, e.g., data reliability, in a given exchange $e$ in a process in the LCI (where an exchange is the use of a flow in a process), $n$ is the number of exchanges of the specified flow, $f$, in the LCI, and LCIDQS is the LCI data quality score for a given flow, $f$.

Rousseaux et al. describe a similar data quality aggregation procedure, but they weight each data quality score by the quantity of the associated exchange, as in Eq. (2):

$$\sum_e^n \frac{DQS_{i,f,e} \times FQ_e}{LCI_f} = \text{LCIDQS}_{f,i} \tag{2}$$

where DQS and LCIDQS are defined in Eq. 1, $FQ_e$ is the flow quantity (e.g., 100 kg) for the given exchange, $e$, and $LCI_f$ is the LCI life cycle flow quantity for the given flow, $f$.

Rousseaux et al. further use these scores in relation to a target data quality score to calculate "coefficients of acceptability" and "variability." This approach is similar to the approach used in the ILCD Handbook where acceptability of datasets is based on the data quality scores being better or equal to the target score (EC JRC 2010).

A very conservative approach, not known to have been described in the literature, would be to take the worst score for a given criterion of any of the flows in the LCI.

$$\max(\text{DQS}_i) = \text{LCIDQS}_{f,i} \tag{3}$$

Attention in the literature on methods of aggregating data quality scores has largely shifted in the last 5 years from aggregation methods for DQIs towards the use of data quality scores with uncertainty analysis. A method that appears to be subject to the most recent attention in the literature is to use scores from the pedigree matrix along with a flow amount to derive measures of spread that can be used in uncertainty analysis, referred to as a probability distribution function method (van den Berg et al. 1999), first proposed by Kennedy et al. (1996). Ecoinvent has adopted this type of approach by ascribing a total uncertainty value based on the uncertainty associated with the intrinsic variability and random error of the indicator and the uncertainty associated with the imperfect representativeness of the data (Ciroth et al. 2013; Muller et al. 2016). Gregory et al. (2016) and Noshadravan et al. (2013) propose a methodology for adapting uncertainty to evaluate comparative LCA scenarios with the inclusion of parameter uncertainty calculated in a similar method to ecoinvent. While this method provides a convenient and very direct way of aligning data quality with result interpretation, some argue that there is no sound justification for creating probability distributions from DQIs (van den Berg et al. 1999; May and Brennan 2003).

## 1.5 Data quality management

DQA is a method for assessing the characteristics and is distinct from data quality management. Data quality management is a rapidly developing field within Information Science, driven by the need to harness data for the purpose of gleaning information to better support decisions (Tayi and Ballou 1998). Data quality management in LCA requires a framework for defining, assessing, storing, and providing access to data quality information and requires a multi-dimensional perspective on data quality.

### 1.5.1 Levels

LCA data are complex and contain a number of levels at which data quality could theoretically be applied. Neither the ISO LCA standards nor Global Guidance Principles for LCA databases describe at which level (e.g., model, unit process or flow) data quality should be applied van den Berg et al. (1999), in a thorough study of data quality aspects of LCA, describe various levels of LCA data to which data quality assessment could be applied, including
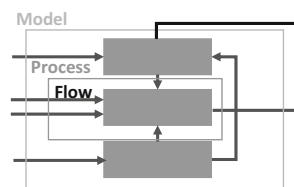


**Fig. 1** Levels in life cycle inventory data

flow, process, and system (or model). Those levels are synthesized in Fig. 1.

### 1.5.2 Roles

Recently through the work of the GLAD metadata working group, the multi-role perspective of data quality has been described as a Janus face. The Janus face recognizes the different roles of the data developer and data user and the need for data quality management that encompasses a multi-role perspective (GLAD WG3 2016).

This paper builds on the Janus face of data quality by adopting a three-role data quality management perspective. These three roles are identified as follows (Wang et al. 2002):

1. Data generator: those who create, collect, or supply data.
2. Data repository/custodian: those who design/develop data guidelines and data system infrastructure; and maintain public access to the data.
3. Data user: those who utilize data, which may involve additional aggregation and integration.

### 1.5.3 Phases

Use and application of data quality in LCA may extend across the four phases of the LCA study (ISO 2006a, b), including the definition of data quality goals during the goal and scope phase; data quality documentation and assessment of inventory data during the inventory phase; data quality documentation and assessment of the life cycle impact assessment methods in the impact assessment phase; and interpretation of data quality scores in the interpretation phase. Although data quality applies to all phases of an LCA, most literature and methodologies have focused on data quality within only the goal and scope, inventory, and interpretation phases.

It is important to note that each level, role, and phase has different requirements, which is why a successful data quality assessment method must include a multi-dimensional management approach, so that the DQCs are transferable to end users. Therefore, this paper addresses improvements to data

quality assessment methods for determining the "fitness for use" within the context of a multi-dimensional data quality management perspective.

### 1.6 Translatability between data quality assessment methods

Translatability between assessment methods is important since no consensus exists within the LCA community on the "best" assessment method. Interoperability between methods is key when compiling a model composed of datasets from various sources. Not all DQCs required by the ISO standards appear as indicators for all assessment methods (see Table 1). The incompatibility of methods means direct translation from one method to another can cause difficulty because key characteristic information could be lost. The GLAD metadata working group has addressed this issue by proposing the storing of DQCs (Ciroth et al. 2017). Figure 2 highlights the difficulty in translatability when repositories do not store the original DQCs and only store the score and/or the data quality goal. Not storing the original DQC can also cause problems when users reassess data with their user defined goal. Users are required to locate the original data documentation in order to assess the data quality.

## 2 Discussion

This discussion addresses many aspects of data quality assessment and management in LCA. To adequately address data quality assessment, the authors recommend the following: an updated pedigree matrix, a method for data quality score aggregation for LCI and life cycle impact assessment (LCIA), and the need for DQA alongside, and not in place of or merged with uncertainty analysis. Data quality management improvements include clarification of generator, repository, and user roles in data quality management and recognition of the need for DQA training for LCA practitioners. We also explore the potential and pitfalls of translating data quality scores across assessment methods, and describe an implementation of the recommended data quality assessment and aggregation methods in openLCA software.

### 2.1 Updated data quality indicators and criteria

An updated version of the Weidema 2013 pedigree matrix is provided in Tables 2 and 3 and is described in detail in a US EPA data quality guidance report (Edelen and Ingwersen 2016b), and referred to here as the US EPA pedigree matrix. The updated pedigree matrix includes two levels of data quality, the flow and the process level. Flow level indicators are applied to any exchange (input or output) within a process except for the reference flow, with one score for each

exchange. Process indicators are designed to be applied to unit processes with one score for each process. Data quality systems currently implemented in major databases provides indicators only at either the flow (e.g., ecoinvent) or process (e.g., ILCD) level. Five indicators resembling those of Weidema et al. (2013) matrix are described at the flow level. These indicators were kept because they are the DQIs relevant for the flow level that could be scored in the ISO 14044 described data quality attributes. Two new indicators are provided at the process level. In previous matrices, all the data quality indicators were not necessarily orthogonal, in that the indicators were capturing overlapping information. In the updated table, all indicators are independent.

The scoring criteria for each indicator in the flow pedigree matrix were either updated or supplemented with more guidance to aid in assessment. Briefly, the following changes were made to the flow level pedigree matrix. For reliability, measurements, calculations, and estimates are defined as three distinct means of collecting data with descending reliability, and the process of verification is defined. For temporal representativeness, the criteria are unchanged but US EPA data quality guidance report clarifies that the date of the actual measurement should be assessed, and not the publication date. This choice was made because data sources often reuse or refer to previous publications, and actual data collected may be much older. For geographic representativeness, scores are defined based on the geographic scope of the data collection area using international guidelines for definitions, an improvement over the previous criteria requiring the user to determine if an area is "similar" or not. There was some debate over whether or not the geographic indicator was even relevant to data quality, as geographic differences are often truly technological differences. However, the US EPA data quality guidance report reinforces that any technology difference should be captured in the technology indicator, and the geographic indicator should just be used to capture differences in location or scope of the data collection areas. "Technology" in the technological correlation indicator was defined as consisting of four elements: process design, operating conditions, material quality, and process scale. The prior indicator "Completeness" was renamed "Data Collection Methods." For this indicator, additional guidance is provided on the "relevant market" and "adequate period of time," and the ranges were altered to require higher percentages than previously used.

Application of the pedigree matrix at the flow level has not always been clearly defined. In the case of technosphere flows, many practitioners are uncertain whether the flow amount or the linked process should be evaluated. For instance, in evaluating a flow of steel with the amount of 50 kg as an input into a process of making a wind turbine in the USA, what should be evaluated is the amount of 50 kg and its associated data quality, and not the flow from a background database that the user has chosen to model the upstream flow
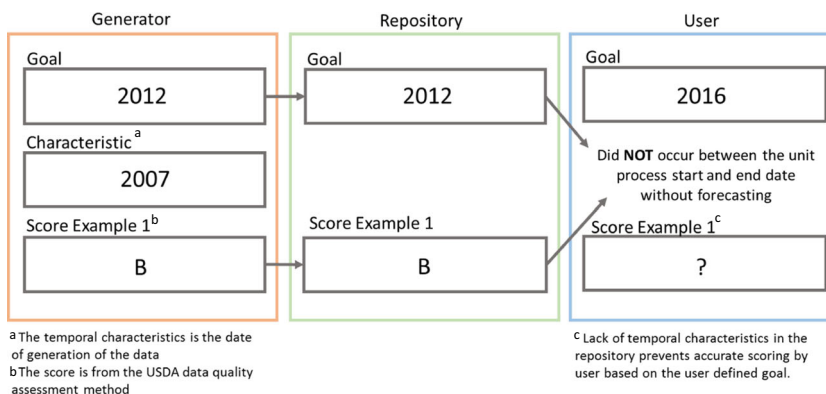
**Fig. 2** Data quality translatability



Table 2 Updated data quality pedigree matrix—flow indicators (Edelen and Ingwersen 2016b)

| Indicator | | 1 | 2 | 3 | 4 | 5 (default) |
|---|---|---|---|---|---|---|
| Flow reliability | | Verified[a] data based on measurements | Verified data based on a calculation *or* non-verified data based on measurements | Non-verified data based on a calculation | Documented estimate | Undocumented estimate |
| Flow representative-ness | Temporal correlation | Less than 3 years of difference[b] | Less than 6 years of difference | Less than 10 years of difference | Less than 15 years of difference | Age of data unknown or more than 15 years |
| | Geographical correlation | Data from same resolution *and* same area of study | Within one level of resolution *and* a related area of study[c] | Within two levels of resolution *and* a related area of study | Outside of two levels of resolution *but* a related area of study | From a different or unknown area of study |
| | Technological correlation | All technology categories[d] are equivalent | *Three* of the technology categories are equivalent | *Two* of the technology categories are equivalent | *One* of the technology categories is equivalent | *None* of the technology categories are equivalent |
| | Data collection methods | Representative data from >80% of the relevant market[e], over an adequate period[f] | Representative data from 60 to 79% of the relevant market, over an adequate period *or* representative data from >80% of the relevant market, over a shorter period of time | Representative data from 40 to 59% of the relevant market, over an adequate period *or* representative data from 60 to 79% of the relevant market, over a shorter period of time | Representative data from <40% of the relevant market, over an adequate period of time *or* representative data from 40 to 59% of the relevant market, over a shorter period of time | Unknown *or* data from a small number of sites *and* from shorter periods |

[a] Verification may take place in several ways, e.g., by on-site checking, by recalculation, through mass balances or cross-checks with other sources. For values calculated from a mass-balance or another calculation method, an independent verification method must be used in order to qualify the value as verified

[b] Temporal difference refers to the difference between date of data generation and the date of representativeness as defined by the goal of the project

[c] A related area of study is defined by the user and should be documented in the geographical metadata. The relationship established in the metadata of the unit process should be consistently applied to all flows within the unit process. Default relationship is established as within the same hierarchy of political boundaries (e.g., Denver is within Colorado, is within the USA, is within North America)

[d] Technology categories are process design, operating conditions, material quality, and process scale

[e] The relevant market should be documented in the DQG. The default relevant market is measured in production units. If the relevant market is determined using other units, this should be documented in the DQG. The relevant market established in the metadata should be consistently applied to all flows within the unit process

[f] Adequate time period can be evaluated as a time period long enough to even out normal fluctuations. The default time period is 1 year, except for emerging technologies (2–6 months) or agricultural projects >3 years

**Table 3** Updated data quality pedigree matrix—process indicators (Edelen and Ingwersen 2016b)

| Indicator | 1 | 2 | 3 | 4 | 5 (default) |
|---|---|---|---|---|---|
| Process review | Documented reviews by a minimum of two types[a] of third party reviewers | Documented reviews by a minimum of two types of reviewers, with one being a third party | Documented review by a third party reviewer | Documented review by an internal reviewer | No documented review |
| Process completeness | >80% of determined flows have been evaluated and given a value | 60–79% of determined flows have been evaluated and given a value | 40–59% of determined flows have been evaluated and given a value | <40% of determined flows have been evaluated and given a value | Process completeness not scored |

[a] Types are defined as either industry or LCA experts

with, for instance "steel plating/EU." The rationale is that the user is modeling the amount and type of a steel material needed to make the wind turbine. Using best practices, the user should specify exactly what type of steel and where it comes from. However, in practice, users make connections to background data and are forced to represent steel with the best available process in his/her database. However, the choice of what technosphere flow to use from a background database as an input is an aspect of model level data quality and not flow level data quality, and we recommend that the flow indicators score should not be convoluted with this aspect of model data quality. The updated pedigree matrix flow level indicators are visually designed to separate intrinsic data quality indicators from contextual indicators.

At the process level, two indicators are defined: process completeness and process review. Process completeness is a new indicator representing the extent to which all expected technosphere and elementary flows are included in a dataset, and a separate scoring table is provided just for this indicator. The process completeness indicator is derived from the completeness DQC in the ISO 14044 standard. Process review was inspired by a recent guidance document establishing review criteria for LCI datasets which further elaborates the basic review requirement specific in the Global Guidance principles for LCI databases (Ciroth et al. 2015; UNEP/SETAC 2011).

In general, data quality within LCA lacks clear definitions for relevant terminology, specifically around indicators and attributes of data quality. Vague and unclear definitions create uncertainty and confusion around the proper evaluation of data quality dimensions and undermines the validity and reproducibility of data quality assessment results (Chen et al. 2014). Terminology such as "similar," "partially," "qualified," "unqualified," "some," and "many" have historically been used within LCA pedigree matrices scoring criteria. This language requires individual evaluators to make subjective

judgments to distinguish data quality scores. The updated pedigree matrix attempts to either eliminate vague and confusing language or clearly define terms to improve reproducibility of score results across multiple users.

The contextual aspects of data quality often require practitioners to make subjective decisions on the relevance of certain data quality indicators. This has led some to selectively apply different data quality indicators in their evaluations. For consistency, the US EPA data quality guidance specifies that all indicators should always be evaluated. Since the importance of indicators is situationally dependent, a practitioner must exercise judgment when using the indicator to interpret results. Practitioners should apply their knowledge of the system and review all indicators together before making judgments on the best use of data. Practitioners should document the decisions about the significance of indicators in their interpretation of LCA results based on the data.

A comprehensive data quality methodology should use more than just the pedigree matrix. It is important that practitioners be aware of the limitations of the pedigree matrix. Not all important data quality characteristics are addressed using a pedigree matrix. Some areas can only be addressed qualitatively through a methodology description of the LCA study (e.g., consistency and reproducibility). The pedigree matrix is not designed to capture all areas of data quality, but to semi-quantitatively address certain key area to improve communication of data quality results.

## 2.2 Training

Updates to the clarity of the pedigree matrix are not enough to solve the reproducibility problem identified by the study described in the Electronic Supplementary Material. The low reproducibility of results and scoring mistakes show a lack of training in the proper application of data quality for a dataset by LCA practitioners. Training in the appropriate

evaluation methodologies associated with data quality evaluation is necessary to ensure consistent application of a data quality system. In LCA, current data quality systems documentation is focused on describing the system. Limited resources from data quality system creators are available to practitioners for training in the appropriate application of the data quality system. The US EPA Guidance on Data Quality Assessment for Life Cycle Inventory Data is stylistically unique in that it not only describes the updated matrix, but guides practitioners through a step-by-step example of the application of the matrix using a test dataset (Edelen and Ingwersen 2016b). A workshop on the application of the updated pedigree matrix was offered at the American Center for LCA XVI conference in Charleston, SC (Edelen and Ingwersen 2016a), and the US EPA is developing an online training module. With the adaptation of the framework that contextual data quality continuously changes and must be re-applied situationally, it is important that individual practitioners and data collectors be offered training on the application of a data quality system to ensure representational consistency of data quality results.

### 2.3 Data quality management roles

A data quality management approach that encompasses a multi-role perspective is needed in order to improve the interoperability of data quality within LCA. Suggested responsibilities in data quality assessment, management and use are described for data generators, data repositories, and data users in Table 4.

Data generators can also be repositories and even users, however the blending of these roles can lead to inadequate documentation and confusion when improper or inadequate documentation is provided, especially for third party users of data. Table 4 assigns the generator the role of documentation of DQCs, both the generator and the user the responsibility of data quality assessment, and the repository the responsibility of maintaining and disseminating guidance measures to ensure the assessment and utilization of the DQCs by the users is supported by appropriate documentation from the generators. Data repositories are also unique in design because they should store data linked with its DQCs.

Identifying and fulfilling the needs of all roles is important in order to adequately address data quality from a multi-perspective approach. Three major data generators/repositories (e.g. ILCD, ecoinvent and USDA) were assessed for implementation of the data quality management responsibilities identified in Table 4. ILCD, ecoinvent and USDA were identified as data generators and custodians. The assessment in Table 5 highlights potential improvements in data quality management for all three generators/repositories. In general, DQCs are not documented and stored separately from DQ scores. Although datasets are documented, the DQC of the

original data is either missing or partially stored in the background documentation. The lack of clear documentation of the DQC by generators is a hindrance to the interoperability of the data, since users must search through background documentation and or find original documentation of data in order to perform an evaluation on the contextual indicators. The GLAD Metadata Working Group draft report proposes the first guidance on storing DQC and is developing a list of key characteristics that should be documented (Ciroth et al. 2017).

### 2.4 Usage alongside quantitative uncertainty

DQCs, as captured by the updated pedigree matrix, do not capture uncertainty and variability in LCA data. When evaluating changes to final LCA results based on model or data uncertainty (Lloyd and Ries 2007), a more traditional uncertainty or variability analysis methodology should be used. We side with many previous scholars that data quality scores are best used independent of uncertainty analysis (van den Berg et al. 1999; Cooper and Kahn 2012; May and Brennan 2003). Data quality assessment can complement uncertainty evaluation and vice versa, providing a more complete understanding of strengths and weaknesses of underlying LCI data and the LCA study results.

### 2.5 Aggregation of data quality indicators in an LCA model

We recommend using a flow-weighted average approach for aggregating flow level data quality scores for use in interpretation, as was first described by Rousseaux et al. (2001), in Eq. (2) previously shown. This approach can be extended to impact assessment so that data quality scores can be viewed for LCIA scores for a given impact category, without using any subjective weighting factors, using the following equation:

$$\sum_{f}^{n} \frac{\text{LCI}_f \times \text{CF}_{f,c} \times \text{LCIDQS}_{f,i}}{\text{LCIAS}_c} = \text{LCIADQS}_{i,c} \qquad (4)$$

where LCIADQS is the life cycle impact assessment data quality score for a given data quality indicator, $i$, and impact category, $c$; LCIDQS is the calculated data quality score for a given flow, $f$, and data quality indicator, $i$, calculated in Eq. (2); LCI is the life cycle inventory total for flow $f$; CF is the characterization factor for flow, $f$, for a given impact category, $c$; and LCIAS is the life cycle impact assessment score for category, $c$.

A simple example is summarized in Fig. 3 for an LCA model with two processes, a truck transport process and a process for diesel fuel. This example comes from US EPA

**Table 4.** Roles in LCA data quality assessment, management and use

| Responsibility | Generator | Repository/ custodian | User |
|---|---|---|---|
| Documentation of inventory DQCs | ✓ | | |
| Documentation of DQGs | ✓ | | ✓ |
| Documentation of DQ scores for data | ✓ | | ✓ |
| Data quality scores linked with interpretation of LCA results | | | ✓ |
| Developing/Adopting DQC documentation guidelines | | ✓ | |
| Developing/Adopting DQA guidelines | | ✓ | |
| Developing DQA training | | ✓ | |
| Storage of DQCs | | ✓ | |
| Storage of data quality scores | | ✓ | |

(2016) and it is used with LCIA characterization factors for the human health effects of criteria air pollutants from TRACI 2.1 (US EPA 2012). The data quality scores were not provided in the original dataset, so hypothetical scores are assigned here for demonstration purposes. Flow level data quality scores are assigned for $PM_{2.5}$ and $PM_{10}$ emissions to air that are common to both processes. In this system, the emissions from the truck transport process contributes >30 and >40 ($PM_{2.5}$ and $PM_{10}$, respectively) times more to the life cycle inventory than the "diesel, dispensed at pump" process and therefore the aggregate data quality scores are more influenced by the truck transport process. If these scores were unacceptable for the data quality goals of the study or the practitioners, the user would be more directed to improve the data quality for the truck transport process than the diesel process, at least to improve the data quality for these particular results. The practitioner may be able to use the aggregate scores to determine confidence in aspects of the LCI results. In this example, the $PM_{2.5}$ emissions are derived from mostly reliable sources, however the data point with the greatest flow quantity is based on an undocumented estimate, and the age of the underlying data is on average older than 10 years, suggesting lower confidence in the results. We do not recommend aggregating flow data quality scores across the categories (e.g., taking the average of

the five aggregate scores) in an ISO14044-compliant LCA. Just like in regard to scores from different LCIA indicators that represent independent categories, combining them is a value-based decision.

Process level data quality scores will generally not be able to be aggregated over the life cycle, because processes have different units that cannot be aggregated, as in the example above, where the first process has units of mass and the second a unit of mass × distance. There may be an exception to this rule for certain LCA models, such as environmentally-extended input-output models where the units are all in a currency, where a weighted-average aggregate process score may be calculated.

The updated data quality pedigree matrix and the aggregation method for data quality scores have been demonstrated with application to a gate-to-gate LCI for acetic acid production (Cashman et al. 2016), and for a new environmentally-extended input-output model of the US economy (Yang et al. 2017).

**2.6 Use of data quality scores in existing datasets**

Practitioners typically assemble LCA models using background database LCI or other LCI prepared for a different

**Table 5.** Data quality management assessment of three data generators/repositories

| Responsibility | ILCD | Ecoinvent | USDA |
|---|---|---|---|
| Documentation of inventory DQCs | | | |
| Documentation of DQGs | ✓ | ✓ | ✓ |
| Documentation of DQ scores for data | ✓[a] | ✓ | |
| Developing/Adopting DQC documentation guidelines | | | |
| Developing/Adopting DQA guidelines | ✓ | ✓ | b |
| Developing DQA training | | | |
| Storage of DQCs | | | |
| Storage of data quality scores | ✓[a] | ✓ | |

[a] Only compliant datasets

[b] No official guidance; Cooper and Kahn (2012) discusses methodology

**Fig. 3** Example aggregate flow level data quality scoring for a simplified two process model

Diesel, dispensed at pump → Truck transport

*Columns (diagonal headers): Flow Reliability, Temporal correlation, Geographical correlation, Technological correlation, Data Collection Method*

**Process 1 - Diesel, dispensed at pump**

| Direction | Flow | Amount | Unit | Flow Reliability | Temporal correlation | Geographical correlation | Technological correlation | Data Collection Method |
|---|---|---|---|---|---|---|---|---|
| Output | Diesel, dispensed at pump | 1 | kg | | | | | |
| Output | PM2.5/air | 2.08E-05 | kg | 1 | 2 | 3 | 4 | 5 |
| Output | PM10/air | 2.09E-05 | kg | 3 | 3 | 3 | 3 | 3 |

**Process 2 - Truck transport**

| Direction | Flow | Amount | Unit | Flow Reliability | Temporal correlation | Geographical correlation | Technological correlation | Data Collection Method |
|---|---|---|---|---|---|---|---|---|
| Input | Diesel, dispensed at pump | 3.02E-02 | kg | | | | | |
| Output | Truck transport | 1 | tkm | | | | | |
| Output | PM2.5/air | 2.00E-05 | kg | 5 | 4 | 3 | 2 | 1 |
| Output | PM10/air | 2.79E-05 | kg | 1 | 1 | 1 | 1 | 1 |

**Life cycle inventory for 1 tkm Truck Transport**

| Direction | Flow | Contributing Process | Amount | Unit | Flow Reliability | Temporal correlation | Geographical correlation | Technological correlation | Data Collection Method | |
|---|---|---|---|---|---|---|---|---|---|---|
| Output | PM2.5/air | Truck transport | 2.00E-05 | kg | 5 | 4 | 3 | 2 | 1 | |
| Output | PM2.5/air | Diesel, dispensed at pump | 6.27E-07 | kg | 1 | 2 | 3 | 4 | 5 | |
| Output | PM10/air | Truck transport | 2.79E-05 | kg | 1 | 1 | 1 | 1 | 1 | |
| Output | PM10/air | Diesel, dispensed at pump | 6.31E-07 | kg | 3 | 3 | 3 | 3 | 3 | |
| Output | PM2.5/air | TOTAL | 2.06E-05 | kg | 4.9 | 3.9 | 3.0 | 2.1 | 1.1 | LCIDQS |
| Output | PM10/air | TOTAL | 2.86E-05 | kg | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | LCIDQS |

**Life cycle impact assessment method (TRACI 2.1 Human Health Criteria category)**

| Direction | Flow | Characterization Factor | Unit |
|---|---|---|---|
| Output | PM2.5/air | 1 | kg PM2.5-eq/kg |
| Output | PM10/air | 0.23 | kg PM2.5-eq/kg |

**Life cycle impact assessment score**

| | LCIA Score | Unit | Flow Reliability | Temporal correlation | Geographical correlation | Technological correlation | Data Collection Method | |
|---|---|---|---|---|---|---|---|---|
| TOTAL | 2.71E-05 | kg PM2.5-eq | 4.0 | 3.2 | 2.5 | 1.8 | 1.1 | LCIADQS |

end purpose than the practitioners' study goal. The existing datasets might have data quality scores associated with them. Practitioners may use these scores to evaluate the datasets, or may intend to use these scores along with their own scores for their primary data to calculate aggregated scores as in Eqs. (2) through (4). Previous papers have proposed, described, or implemented data quality aggregation methods, the implications for the use of existing data with data quality scores for aggregation was not addressed. These scores may come from the same or a different pedigree matrix system and review of the scores themselves may or may not have been conducted with any dataset or study review. For these reasons alone, practitioners should not use the existing scores as they are for either evaluation or for data quality aggregation, with few exceptions. Table 6 provides a key that might be cautiously used for data quality score translation to the US EPA data quality system from datasets scored using the same system or some of the alternative systems reviewed above.

The next question that arises is how many datasets from a background system would need to be scored, and how to select those datasets. With newer background datasets consisting of 10,000 or more processes, it would be infeasible to re-evaluate data quality for all background datasets. If the flow quantity-weighted scoring algorithms in Eqs. (2) and (4) are used, the practitioners can prioritize the scoring based on the influence of the flows in the total LCI or impact score calculations. For instance, building on the example summarized in Fig. 3, if the process "Truck transport" were from a background database, the flow-weighted aggregation method would suggest that this process is very influential on the LCI and LCIA results related to PM emissions and human health impacts and that data quality of this process should be rescored before calculating the aggregated data quality scores.

### 2.7 Application in LCA software

Working with GreenDelta, the handling of data quality assessment in openLCA has been modified in a number of ways to support these recommendations starting with version 1.6.1. The user can now define one or more data quality systems in

**Table 6** Key for adaptation of existing data quality scores for use in a new model that uses the USEPA DQA method

| System | Type | Level | Data quality indicator | Equivalency | USEPA DQI | Recommendation |
|---|---|---|---|---|---|---|
| USEPA | Intrinsic | Flow | Flow reliability | = | Flow reliability | Use as is |
| | | Process | Process review | = | Process review | Use as is |
| | Contextual | Flow | Temporal correlation | = | Temporal correlation | Adjust scores to reflect time variation between "old" and "new" study |
| | | Flow | Geographical correlation | = | Geographical correlation | Rescore |
| | | Flow | Technological correlation | = | Technological correlation | Rescore |
| | | Flow | Data collection methods | = | Data collection methods | Rescore |
| | | Process | Process completeness | = | Process completeness | Adjust scores for any new expected flows |
| Weidema et al. (2013) | Intrinsic | Flow | Reliability | ≈ | Flow reliability | Use as is for score of 1, rescore for others if metadata available |
| | | Flow | Uncertainty correlation | X | | No equivalent |
| | | Flow | Precision | X | | No equivalent |
| | Contextual | Flow | Completeness | ≈ | Data collection methods | Rescore if metadata are available |
| | | Flow | Temporal correlation | = | Temporal correlation | Adjust scores to reflect time variation between "old" and "new" study |
| | | Flow | Geographical correlation | ≈ | Geographical correlation | Rescore if metadata are available |
| | | Flow | Further Technological correlation | ≈ | Technological correlation | Rescore if metadata are available |
| LCA Digital Commons (USDA) | Intrinsic | Flow | Reliability and reproducibility | ≈ | Flow reliability | Rescore if metadata are available |
| | | Flow | Uncertainty | X | | No equivalent |
| | | Flow | Precision | X | | No equivalent |
| | Contextual | Flow | Flow data completeness | ≈ | Data collection methods | Rescore if metadata are available |
| | | Flow | Temporal correlation | ≈ | Temporal correlation | Rescore if metadata are available |
| | | Flow | Geographical coverage | ≈ | Geographical correlation | Rescore if metadata are available |
| | | Flow | Technological coverage | ≈ | Technological correlation | Rescore if metadata are available |
| ILCD | Intrinsic | Process | Precision/uncertainty | X | | No equivalent |
| | Contextual | Process/model | Methodological appropriateness and consistency | X | | No equivalent |
| | | Process/model | Completeness | ≈ | Process completeness | Rescore if metadata are available |
| | | Process | Time-related representativeness | ≈ | Temporal correlation | Rescore if metadata are available |
| | | Process | Geographical representativeness | ≈ | Geographical correlation | Rescore if metadata are available |
| | | Process | Technological representativeness | ≈ | Technological correlation | Rescore if metadata are available |
| | | Process | Technological coverage | ≈ | Technological correlation | Rescore if metadata are available |

Equivalency: "=" information captured in this indicator directly corresponds to a USEPA DQI, "≈" information captured in this indicator is similar to a USEPA DQI, "**X**" no equivalent is found within the USEPA system

openLCA in the pedigree matrix format. For each system defined, a user can determine whether or not these systems should be used with uncertainty analysis. Within any process in the database, a data quality system can be chosen for use at the flow or process level. Only one flow and one process level data quality system can be defined for each model. Analysis options have been added that will now perform aggregation of data quality scores using the algorithm in Eq. (2) and Eq. (4) by default, but also using other algorithms, such the use of a maximum score as in Eq. (3) or a variation of Eq. (4), a weighted squared average. Options are available as well for how data quality scores are rounded (normal rounding by default) and how missing scores are handled (omitted by default) in the aggregation. The new data quality results provide statistics on the percentage of data quality scores provided in the LCI for each DQI, show process data quality for all processes used in a product system along with the total requirement for each process, and show aggregated data quality scores for elementary flows in the LCI totals and LCIA scores along with the LCI amounts and LCIA results. The openLCA JSON-LD format was also expanded to capture the data quality system information, so when datasets are transferred they will include a definition of the data quality system used in the datasets. The US EPA pedigree matrix is available in JSON-LD format as two systems (one each for flow and process levels) for users to import into openLCA through the Environmental Dataset Gateway (http://edg.epa.gov) and openLCA (http://www.openlca.com) websites.

# 3 Future work

## 3.1 Interoperability and limitations

Data quality scores are, as a general rule, not translatable from one methodology to another. There is a current effort by the GLAD Metadata Working Group to improve data quality interoperability (Ciroth et al. 2017). This effort focuses on the need for standard documentation instead of a standard data quality assessment methodology, because as Table 6 shows, direct translation between methodologies is not always possible. It is important that users of data quality methods understand the situational dependency of data quality and not use data quality indicator scores from previous assessments for any of the contextual indicators, as this practice will be misleading in aggregate data quality calculations.

## 3.2 Additional data quality levels

Although the updated pedigree matrix includes a new multidimensional component to capture flow and process data quality, it fails to capture all dimensions of LCA data and models. LCIA data quality presents several unique differences from

LCI data quality. Unique challenges to LCIA data quality include, but are not limited to the following: lack of scientific knowledge of environmental impacts, inability to connect aggregated LCI data with environmental impacts, high level of input data into assessment models (e.g., toxicity, persistence, bioaccumulation, and equivalency factors), and quality of assumptions made by assessment models (US EPA 1995; Bare et al. 1999; ISO 2000).

## 3.3 Automating data quality judgments

A certain amount of variance is to be expected with data quality since it relies heavily on the user's personal judgment, especially since data quality is largely contextual. The high amount of subjectivity in the production of data quality judgments has an impact on the trustworthiness and added value of data quality scores. The rapid increase in data generation coupled with the inability of humans to manually assess data at the same rate has driven efforts to automate or partially automate data quality (Isaac and Lynes 2003). Current LCA data quality methods require data generators to document data quality and then individually translate documentation into data quality scoring. Wang et al. (2002) propose the novel concept of using an automated data quality reasoner. The data quality reasoner is a knowledge-based approach for dealing with the subjective, decision-analytic nature of data quality judgment. The ability to use objective terminology to delineate differences in scoring of criteria is a key step towards semi- or fully automating data quality judgments.

# 4 Conclusions

LCA practitioners continue to struggle to reproduce data quality scores and to fully evaluate and use these scores in interpretation of LCA results. This paper summarizes improvements to data quality indicators and associated guidance issued in a US EPA data quality guidance report (Edelen and Ingwersen 2016b). The guidance does not recommend a major departure from the common pedigree approach described by Weidema et al. (2013), but intends to increase clarity and consistency in application to make scoring more objective and useful. Indicators are defined at the flow and process levels and the contextual nature of data quality assessment is emphasized. More guidance is generally given on performing data quality assessment, and the need for more practitioner training is expressed. Roles in management of data quality information are outlined for data generators, data repositories, and data users. A method for data quality aggregation is proposed that extends earlier work (Rousseaux et al. 2001) to provide aggregate data quality scores for LCIA results. Initial recommendations are made on how to use existing data quality scores from a background database along with new scores,

including mappings of US EPA-based scores and other systems. The use of data quality is recommended alongside, and not mixed with, quantitative uncertainty assessment. The new data quality system and aggregation methods can now be used in openLCA software. Future work is needed to better address interoperability of data quality scores between systems, to automate data quality scores, and to extend their scope to cover other components of LCA data including the model level and to LCIA characterization factors.

**Compliance with ethical standards**

**Disclaimer** This article does not reflect the endorsement or opinion of any of the institutions that employ the authors.

# References

Bakst J, Lacke C, Weitz K, Warren J (1995) Guidelines for assessing the quality of life-cycle inventory data. Research Triangle Institute and the US Environmental Protection Agency, Washington

Bare J, Pennington D, Udo de Haes H (1999) Life cycle impact assessment sophistication. Int J Life Cycle Assess 4:299–306

van den Berg N, Huppes G, Lindeijer E, van der Ven B, Wrisberg N (1999) Quality assessment for LCA. CML Center for Environmental Science, Netherlands

Björklund A (2002) Survey of approaches to improve reliability in LCA. Int J Life Cycle Assess. doi:10.1007/bf02978849

Canals L, Boureima F, Brago T, Ciroth A, Czaga P, Fazio S, Goedkoop M, Ingwersen W, Krewer C, Leite C, Schally H, Suh S, Tahara K, Tonda E, Vigon B, Wang F, Wernet G (2016) Towards an interoperable network of LCA databases. Paper presented at the SETAC Europe, Nantes, 22–26 May

Cashman S, Meyer D, Edelen A, Ingwersen W, Abraham J, Barrett W, Gonzalez M, Randall P, Ruiz-Mercado G, Smith R (2016) Mining available data from the United States Environmental Protection Agency to support rapid life cycle inventory modeling of chemical manufacturing. Environ Sci Technol 50(17):9013–9025

Chen H, Hailey D, Wang N, Yu P (2014) A review of data quality assessment methods for public health information systems. Int J Environ Res Public Health 11(5):5170–5207

Ciroth A, Muller S, Weidema B, Lesage P (2013) Empirically based uncertainty factors for the pedigree matrix in ecoinvent. Int J Life Cycle Assess 21(9):1338–1348

Ciroth A, Hildenbrand J, Zamagni A, Foster C (2015) Life cycle inventory dataset review criteria development. version 2. UNEP/SETAC Life Cycle Initiative, Paris

Ciroth A, Arbuckle P, Cherubini E, Ugaya C, Edelen A (2017) Task 3: Core meta-data descriptors and guidance on populating descriptors. WG3 of the Global Life Cycle Data Access Network (GLAD)

Cooper J, Kahn E (2012) Commentary on issues in data quality analysis in life cycle assessment. Int J Life Cycle Assess 17:499–503. doi:10.1007/s11367-011-0371-x

Coulon R, Camobreco V, Teulon H, Besnainou J (1997) Data quality and uncertainty in LCI. Int J Life Cycle Assess 2(3):178–182

EC JRC (2010) International reference life cycle data system (ILCD) handbook—specific guide for life cycle inventory data sets. European Commission–Joint Research Centre–Institute for Environment and Sustainability, Luxembourg

Edelen A, Ingwersen W (2015) Addressing data quality indicators to support interoperability: recommendations for further developments in life cycle assessment data quality systems. In: LCA XV, Vancouver, October 6–8

Edelen A, Ingwersen W (2016a) Data quality assessment in LCA. In: LCA XVI, Charleston, September 26-29

Edelen A, Ingwersen W (2016b) Guidance on data quality assessment for life cycle inventory data. EPA/600/R-16/096. United States Environmental Protection Agency, Cincinnati, OH

Fava J (1992) Life cycle assessment data quality: a conceptual framework. Society of Toxicological and Chemistry (SETAC), Wintergreen

Frischknecht R, Jungbluth N, Althaus H, Doka G, Dones R, Heck T, Hellweg S, Hischier R, Nemecek T, Rebitzer G, Spielmann M, Wernet G (2007) Overview and methodology. Data v2.0. ecoinveint centre, Dubendorf

GLAD WG3 (2016) Bruce Vigon, Andreas Ciroth (co-chairs) with input from WG members: Meta-Data Needs Assessment – Element 1, v. 3, March 2016

Gregory J, Noshadravan A, Olivetti E, Kirchain R (2016) A proposal for extending the pedigree matrix approach to quantify uncertainty in the application of intermediate flows. In: LCA XVI, Charleston, September 26-29

Isaac D, Lynes C (2003) Automated data quality assessment in the intelligent archive. Technical Report: Intelligent Data Understaing. NASA, Goddard Space Flight Center, Washington

ISO (2000) ISO 14043: Environmental management-Life cycle assessment-Life cycle interpretation. vol ISO 14043:2000(E). International Standards Organization, Switzerland

ISO (2006a) ISO 14040: Environmental management—life cycle assessment—principles and framework. International Organization for Standardization, Switzerland

ISO (2006b) ISO 14044: Environmental management—life cycle assessment—requirements and guidelines. International Organization for Standardization, Switzerland

Kennedy D, Montgomery D, Quay B (1996) Data quality. Int J Life Cycle Assess 1(4):199–207

Lloyd S, Ries R (2007) Characterizing, propagating, and analyzing uncertainty in life-cycle assessment: a survey of quantitative approaches. J Ind Ecol 11(1):161–179

May J, Brennan D (2003) Application of data quality assessment methods to an LCA of electricity generation. Int J Life Cycle Assess 8(4):215

Muller S, Lesage P, Samson R (2016) Giving a scientific basis for uncertainty factors used in global life cycle inventory databases: an algorithm to update factors using new information. Int J Life Cycle Assess 21(8):1185

Noshadravan A, Wildnauer M, Gregor J, Kirchain R (2013) Comparative pavement life cycle assessment with parameter uncertainty. Transport Res D-Tr E 25:131–138

Reap J, Roman F, Duncan S, Bras B (2008) A survey of unresolved problems in life cycle assessment. Int J Life Cycle Assess 13(5):374

Rousseaux P, Labouze E, Suh Y, Blanc I, Gaveglia V, Navarro A (2001) An overall assessment of life cycle inventory quality. Int J Life Cycle Assess 6(5):299–306

Tayi G, Ballou D (1998) Examining data quality. Commun ACM 41(2): 54–57

UNEP/SETAC (2011) Global guidance principles for life cycle assessment databases: a basis for greener processes and products. Glossary. United Nations Environment Programme, Paris

US EPA (1995) Life-cycle impact assessment: a conceptual framework, key issues, and summary of existing methods. United States Environmental Agency, Research Triangle park

US EPA (2000) Guidance for data quality assessment: practical methods for data analysis. QA00 update edn. United States Environmental Protection Agency, Washington

US EPA (2012) Tool for the reduction and assessment of chemical and other environmental impacts: TRACI version 2.1—user's manual. US EPA National Risk Management Laboratory, Cincinnati

US EPA (2016) Life cycle inventory data—chemicals, construction, transportation other processes for use in environmental footprint analyses. US EPA National Risk Management Research Laboratory, Cincinnati

Wang R, Strong D (1996) Beyond accuracy: what data quality means to data consumers. J Manage Infom Syst 12(4):5–34

Wang R, Ziad M, Lee Y (2002) Data quality. Advances in database systems. Kluwer Academic Publishers, New York

Weidema B (1998) Multi-user test of the data quality matrix for product life cycle inventory data. Int J Life Cycle Assess 3(5):259–265

Weidema B, Wesnaes M (1996) Data quality management for life cycle inventories—an example for using data quality indicators. J Clean Prod 4(3–4):167–174

Weidema B, Bauer C, Hischier R, Mutel C, Nemecek T, Reinhard J, Wernet G (2013) Overview and methodology: data quality guideline for the ecoinvent database version 3. The ecoinvent Centre, St. Gallen

Wernet G, Bauer C, Steubing B, Reinhard J, Moreno-Ruiz E, Weidema B (2016) The ecoinvent database version 3 (part I): overview and methodology. Int J Life Cycle Assess 21(9):1218–1230

Wrisberg M, Lindeijer E, Mulders P, Ram A, Van der Ven B, Van der Wel H (1997) A semi-quantitative approach for assessing data quality in LCA. In: 7th Annual Meeting of SETAC-Europe, Amsterdam, April 6–10 1997

Yang Y, Ingwersen W, Hawkins T, Srocka M, Meyer D (2017) USEEIO: a new and transparent United States environmentally-extended input output model. J Clean Prod 158:308–318