

1 **CODIFICATION CHALLENGES FOR DATA SCIENCE IN CONSTRUCTION**

2
3 Ranjith K. Soman

4 PhD Student, Centre for Systems Engineering and Innovation, Department of Civil and Environmental
5 Engineering, Imperial College London, London, SW7 2AZ, UK

6 (corresponding author). Email: ranjithks17@imperial.ac.uk

7 <https://orcid.org/0000-0003-3967-9121>

8
9 Jennifer K. Whyte, PhD, FICE

10 Laing O'Rourke/Royal Academy of Engineering Professor in Systems Integration and Director, Centre for
11 Systems Engineering and Innovation, Department of Civil and Environmental Engineering, Imperial College
12 London, London, SW7 2AZ, UK. Email: j.k.whyte@imperial.ac.uk

13 <http://orcid.org/0000-0003-4640-2913>

14
15 **Abstract:** New forms of data science, including machine learning and data analytics, are enabled by machine-
16 readable information but are not widely deployed in construction. A qualitative study of information flow in three
17 projects using Building Information Modelling (BIM) in the late design and construction phase is used to identify
18 the challenges of codification which limit the application of data science. Despite substantial efforts to codify
19 information with 'Common Data Environment (CDE)' platforms to structure and transfer digital information
20 within and between teams, participants work across multiple media in both structured and unstructured ways.
21 Challenges of codification identified in this paper relate to *software usage* (interoperability, translation, modelling,
22 and file-based sharing), *information sharing* (unstructured information, document control, workarounds, process
23 change, and multiple CDEs), and *construction process information* (loss of constraints and low level of detail).
24 This paper contributes to the current understanding of data science in construction by articulating the codification
25 challenges and their implications for data quality dimensions, such as accuracy, completeness, accessibility,
26 consistency, timeliness, and provenance. It concludes with practical implications for developing and using
27 machine-readable information and directions for research to extract insight from data and support future
28 automation.

29 Keywords: Building information modelling (BIM), Codification; Artificial Intelligence (AI);
30 Automation; Data science; Machine readability; Construction.

31

32 **1. Introduction**

33 Machine-readable information is enabling new forms of data science, including machine learning and
34 data analytics. These methods offer value to the construction sector through resource and waste optimization,
35 data-driven design, prescriptive analytics for rule checking, visual analytics, performance predictions, operational
36 analytics, and more (Bilal et al., 2016). Yet the construction sector is not taking advantage of these developments
37 as data science is not widely deployed. Bilal et al. (2016) have identified poor data management as one of the
38 main factors which limit the application of data science in construction. The sector is actively trying to overcome
39 this limitation through Building Information Modelling (BIM), an approach to incrementally building structured
40 (and, hence, machine-readable) information throughout the project cycle (Eastman et al., 2008; Jordani, 2010),
41 and where it is machine readable, such structured information could support the use of data science. Recent ideas
42 such as the ‘digital twin’ are also predicated on the availability of such structured information (both geometries
43 and behaviours; see Bolton et al., 2018). Research has begun to develop approaches to, and document practices
44 of, codifying information (to convert it into a structured and machine-readable format) to improve delivery
45 practices (e.g. in relation to the construction phase (Goedert and Meadati, 2008), workspace planning (Choi et al.,
46 2014), construction defects (Kwon et al., 2014), and lessons learned in a project (Oti et al., 2018)). These studies
47 have advanced knowledge regarding frameworks and methods to codify construction information. In addition to
48 these studies, there is a need for work to extend understanding of the issues which prevent codification of
49 information to improve uptake of data science in construction.

50 By examining information sharing across projects which use BIM in the late design and construction
51 phase, this paper aims to identify the codification challenges which arise in practice by examining information
52 sharing across projects that use BIM. In this paper, codification is defined as the process of conversion of
53 information into a structured and machine-readable format to support the application of data science. Information
54 refers to the collection of data contextualized with relevant schema and semantics so that insights can be made
55 from the data. The codification challenges are the issues which reduce the machine readability and quality of
56 construction information and, in turn, limit the uptake of data science in the sector. Inspired by work which frames
57 BIM use as a complex social activity (Cao et al., 2014; Dossick and Neff, 2010), this paper builds on and
58 contributes to strands in the literature focused on data quality, machine readability, and BIM adoption and
59 implementation. Data analyses suggest codification challenges which are organizational as well as technical in
60 nature, relating to software use, information sharing, and construction process information.

61 To develop the contribution, the rest of this paper is divided into four sections. Section 2 provides a brief
62 background on data quality and machine readability in projects and organizational issues associated with BIM
63 adoption and data quality. Section 3 describes the cases and method used in the study. Section 4 presents the
64 findings. These findings are discussed further in relation to the literature in section 5, and the conclusions are
65 presented in section 6.

66 **2. Background**

67 As the construction sector becomes increasingly digital, most information is stored digitally and is
68 accessible through servers or a common data environment (CDE) held in firms or projects (British Standards
69 Institution, 2018; Preidel et al., 2016). However, being digitally accessible does not mean that the information is
70 machine readable as the semantics may not be embedded in the data (Hendler and Pardo, 2012). Semantics could
71 be derived from the data using advanced machine learning techniques such as natural language processing and
72 deep learning (Carrillo et al., 2011; Wang, 2017). Nonetheless, this is a resource-intensive process which requires
73 training models for achieving satisfactory accuracy and has costs associated with it. The result of this process
74 may also not be of high-quality (Wang, 2017). Despite the existence of such data-cleaning algorithms across
75 sectors, poor data quality is costing \$3.1 trillion in the United States (Quintero et al., 2015). In addition, the poor
76 quality of data is increasing operational costs, decreasing revenue, and resulting in missed commercial
77 opportunities (Loshin, 2010). Within construction, Sacks et al. (2017) have described how the quality of input
78 information influences semantic enrichment of BIM when using machine learning. Moreover, Farias et al. (2018)
79 have shown the effects of poor quality data resulting in wrong inferences when they tried to extract building views
80 using a rule-based method. Whyte et al. (2016) articulated how managing change in large datasets becomes a
81 focus in an era of ‘big data’ in which project information is increasingly characterized by volume, velocity, and
82 variety. Recent work has further characterizes such data as also including characteristics of veracity and value
83 (e.g. Younas, 2019). These issues of data quality occur because construction data is heterogeneous, and its veracity
84 is not always known. Data cleaning related to the variety (heterogeneity) and veracity characteristics of big data
85 is especially difficult when compared to data cleaning related to other characteristics such as volume and velocity
86 (Fan, 2015; Janssen et al., 2017). Therefore, there is a need to keep the data of the highest quality and in a machine-
87 readable format (maintaining the data relationships) as far as possible to have the best inferencing.

88 **2.1 Quality and machine readability**

89 What constitutes a good quality dataset? According to Wang and Strong (1996), a good quality dataset
90 is the one that has enough information embedded in it for a particular use by the user. Researchers have set out

91 multiple dimensions to assess the data quality concerning big data analytics (Batini and Scannapieco, 2016; Cai
92 and Zhu, 2015; Delone and McLean, 2014; Naumann and Rolker, 2000; Wang and Strong, 1996). For this paper,
93 the focus is on the following data quality dimensions based on Batini and Scannapieco (2016) as they best reflect
94 the implications of codification challenges. *Accuracy* is the closeness of the measured/represented data and reality.
95 There are two kinds of accuracies: semantic and syntactic. Semantic accuracy relates to the closeness of the data
96 value to reality, whereas syntactic accuracy refers to the closeness of the data representation to the expected data
97 type/model. *Completeness* is the measure of information content present in the data compared to the extent of
98 information content required to be present in the data to perform a particular task. *Temporal* dimensions refer to
99 currency, volatility, and timeliness. *Currency* relates to the promptness of data updates. *Volatility* refers to the
100 frequency at which the data variance occurs. *Timeliness* refers to the suitability of the current data to perform a
101 task. *Consistency* refers to uniformity and constancy of data with respect to the semantic rules defined over
102 multiple data items. *Accessibility* refers to the ability of data to be accessed by a user (human user or computer
103 program) and generate information from it. *Data provenance* is the description of the origins of data and the
104 process by which it is manipulated. Jayawardene et al. (2015) have conducted a systematic literature review on
105 the extensive data quality dimensions and consolidated overlapping dimensions of quality. This means that some
106 of the data quality dimensions are interdependent. For example, the data quality dimension related to semantic
107 accuracy might depend on the timeliness dimension as the data may be accurate with respect to time. However,
108 the same value may be inaccurate at a different time. At the outset, a dataset is considered to be of good quality
109 when the measure of these dimensions is high, leading to better inferences.

110 What constitutes machine-readable data? Data in a structured format that could easily be processed by
111 computers are considered as ‘machine-readable’. Berners-Lee (2006) has stated a set of ‘rules’ for creating
112 structured data so that it can be connected and interpreted easily by machines. The first one is to index the data to
113 make it digitally accessible by storing it on online servers so that it can be easily accessed by computers as well
114 as people. This relates to the accessibility dimension of data quality. Indexing the data and storing it online
115 increases the accessibility as it is easy to find. The second rule is to structure the data with relevant schemas for
116 easy interpretation by machines. This step makes the data structured such that semantic relations are embedded
117 within it, resulting in better inferencing and, thus, improving the syntactic accuracy of the data. The third rule is
118 to make the schemas public and machine readable by using open-source schemas to describe the data model so
119 that interpretations can be made by computers without any proprietary data interfaces. Proprietary data formats
120 limit data inferencing as the schema by which data is modelled is only accessible to few applications. Therefore,

121 using open-source schema would increase the measure of accessibility dimension as more applications can use
122 the open schema to derive the context of data for inferencing. The last rule is to link the data with other datasets
123 so that better inferences can be made by deriving the context information. This improves the consistency
124 dimensions of data quality as the same data is linked to multiple datasets. Linking would ensure that there are no
125 conflicts in the data about a concept stored in multiple databases. Based on these rules, structured information can
126 be classified into five types in the increasing order of machine readability, as shown in Table 1. Increasing the
127 machine readability of the data, in turn, increases the data quality as the dimensions relating to accessibility,
128 accuracy, and consistency are improved in the process.

129 <<Insert Table 1>>

130 **2.2 Limits of existing research**

131 What is limiting the generation of good quality machine-readable information within the sector? It does
132 not appear to be technical development as novel technical solutions are being developed by construction
133 informatics researchers with a focus on the integration of data in the sector, for example, through the use of data
134 standards (Krijnen and Beetz, 2017; Pazlar and Turk, 2008), cloud-based BIM (Beetz et al., 2010; Singh et al.,
135 2011; Zhang et al., 2017), and linked-data technologies (Kim et al., 2018; Pauwels et al., 2015; Pedro et al., 2017,
136 Zhang and Beetz, 2015). It does not appear to be policy interventions either. Standards and public mandates have
137 placed BIM at the heart of the information management required to coordinate processes in project delivery and
138 operation of infrastructure, making BIM central to digital tools and workflows in projects (British Standards
139 Institution, 2018; Sacks et al., 2018). Instead, the literature suggests that the issues may be both organizational
140 and technical in nature.

141 Prior research has assessed BIM adoption across different markets, using a model of diffusion area,
142 macro-maturity components, macro-diffusion dynamics, and so forth, and validated this model by applying it to
143 assess BIM adoption amongst 21 countries (Kassem and Succar, 2017). This work has determined the BIM project
144 objectives, critical success factors, and operative critical success factors for effective implementation of BIM
145 (Chegu, Badrinath, and Hsieh, 2019). It has identified the success factors for adoption of BIM in a company,
146 selection of projects within the company to implement BIM, and selection of BIM services and software (Won et
147 al., 2013). It has surveyed the degree of implementation of BIM statistically by evaluating the level of BIM
148 implementation and quality of collaboration and communication in BIM-enabled projects, and linking discussing
149 its influence on uptake of integrated delivery systems (Chang et al., 2017), and developing strategies for using
150 BIM to reduce rework in construction (Hwang, et al., 2019) and improve collaboration through the development

151 of BIM-based platforms by analyzing requirements and details of elements needed for a collaborative work model
152 (Zhang et al., 2017). In addition, Gu and London (2010) have created a collaborative BIM decision framework to
153 facilitate BIM adoption through a four-part method. The framework first defines the scope, purpose, roles,
154 relationships, and project phases, followed by developing a work process roadmap, identifying the technical
155 capabilities and the limitations of tools, and finally customizing these to suit the capabilities and skillsets of the
156 project team. Building on this study, Singh et al. (2011) have determined technical requirements for a BIM server
157 to serve as a collaboration platform. These studies have extended the existing knowledge base, and they give a
158 deeper understanding of the problems associated with BIM implementation, and suggested steps for the effective
159 implementation of BIM. However, data quality issues emerging from the problems associated with BIM
160 implementation is relatively less studied.

161 Previous research on data quality within the construction sector has studied semantic and syntactic
162 accuracy of BIM, BIM quality assurance processes for the design stage and data quality issues in the design model,
163 and completeness of information in BIM for facility management. Solihin et al. (2015) have identified
164 requirements for good quality for BIM in an Industry Foundation Class (IFC) format. Building on this study, Lee
165 et al. (2018) have presented a semantic rule-checking process to ensure data quality pertaining to semantic and
166 syntactic accuracy is maintained whilst BIM in an IFC format is exchanged. In another study, the quality of the
167 information in the design phase was assessed using a structured and quantifiable process based on a BIM quality
168 assurance by Donato et al. (2017). Mirarchi and Pavan (2019) have analyzed the data quality issues concerning
169 accuracy, consistency, and completeness dimensions of the BIM models created during the design. For facility
170 management, Zadeh et al. (2017) have proposed a framework to assess the quality of BIM, focusing on the
171 completeness dimension. These studies have advanced the knowledge on data quality issues associated with BIM
172 models. A limitation of these studies is that they do not address the data quality issues emerging due to the wider
173 practice of using document-based and model-based information sharing in the sector.

174 To address the data quality issues emerging due to such document-based and model-based information-
175 sharing practices, it is necessary to study information flow across project teams in detail. By better characterizing
176 the practice, such empirical work can then inform future technical developments (e.g. Hartmann, 2008) and
177 address challenges raised in prior work in areas such as automated scheduling (e.g. Han and Golparvar-Fard,
178 2017). Previous research on the use of BIM in organizations articulates antecedents to BIM uptake (Taylor, 2007)
179 and identifies organizational issues which affect BIM implementation, such as organizational divisions (Dossick
180 and Neff, 2010). It describes how practices are always 'hybrid', overlaying a range of old and new media and

181 processes (Harty and Whyte, 2010; Whyte, 2011), with the roles of construction professionals also evolving
182 (Akintola et al., 2017; Jaradat et al., 2013; Sebastian, 2011). Such work draws attention to the organizational
183 factors associated with information use, whereby technological integration cannot be assumed to foster closer
184 collaboration across companies (Dossick and Neff, 2010). This literature, which understands BIM use as a
185 complex socialized activity (Cao et al., 2014), provides an approach that can be used to study the codification
186 challenges and design quality issues which emerge in leading practice.

187

188 **3. Research method**

189 To understand the challenges of codification in construction, three construction projects are studied
190 qualitatively to investigate the digital tools and workflows used in the projects, structured and unstructured
191 information flows in these projects, and the problems associated with the information flow. These three leading
192 projects are a multi-storey residential student apartment block in the United Kingdom (Case 1), a metro rail
193 infrastructure project in India (Case 2), and a major water infrastructure megaproject in the United Kingdom (Case
194 3). The multi-storey residential apartment project (Case 1) is an exemplary project exhibiting the use of BIM in
195 the United Kingdom, constructed by a leading contractor and using state-of-the-art offsite manufacturing
196 approaches in construction. The metro rail project in India (Case2) is pioneering BIM implementation amongst
197 the metro projects in India, incorporating learning on digital implementation from global megaprojects. The water
198 infrastructure project (Case 3) is one of the biggest construction projects in the United Kingdom, using innovative
199 technological solutions to futureproof construction and deliver a physical asset as well as a digital asset for
200 operation. Early in the study, the first author, who collected the data, also visited another infrastructure project
201 and a commercial retrofit project, and these three projects were then chosen due to their significance and because
202 they use a level of digital collaboration categorized as BIM level 2. These projects follow the BIM level 2
203 recommendations set out by the mandate (it is mandated in Case 3 and seen as best practice in the other cases).
204 Although the metro rail project (Case 2 is in a country which does not have a regulatory framework for BIM level
205 2 recommendation, the owner required the adoption of BIM level 2 as international best practice, which justifies
206 our choice for selecting the case. Information sharing across project teams in the late design and construction
207 phases of the three projects is studied qualitatively by visiting the projects, analyzing internal and publicly
208 available documents, observing meetings, and conducting informal and formal interviews on the use of product
209 and process information during the construction stage (refer Appendix 1).

210 Within each project, there was an initial setup meeting to present the study and identify interviewees.
211 The interview protocol covered questions of communication, software tools used, BIM, collaboration, and
212 information flow. The data analysis phase overlapped with the data collection phase. The taped interviews were
213 transcribed, and field notes were typed up. These were read and reread between the project meetings. Summaries
214 of interviews were sent back to the interviewees for member checks. All the data was organized into cases and
215 stored into the qualitative analysis software. These methods draw on a qualitative case study approach (Eisenhardt
216 and Graebner, 2007), building insights across the three cases from multiple sources (site visits, documents, field
217 notes, and interview transcripts).

218 • Multi-storey residential student apartment block in the United Kingdom (Case 1): To study this case, the first
219 author visited the construction site and offices of the projects, had informal conversations with the digital and
220 planning engineers, examined construction documents and models, and studied the software tools used to
221 understand the information embedded in the BIMs, construction schedules, and other reports such as design
222 calculation, method statements, and so on. These documents were centrally hosted on a CDE, and the first
223 author had access to it whilst being there at the office.

224 • Metro rail infrastructure project in India (Case 2): To examine this case, documents such as the BIM execution
225 plan, presentation documents for training, and press releases were studied to understand digital information
226 management practices. The project manager, chief site engineer, casting yard engineer, and BIM consultants,
227 who form a cohort of the key decision makers during the construction stage, were interviewed informally to
228 get an insight into the extent of codification in the information flow during their daily work practices. Field
229 notes were taken during the interview. In addition to these interviews, the casting yard, viaduct construction
230 site, and a station site were visited to understand the on-ground practice of various activities. Further insight
231 into this case was obtained through a workshop, co-organized by the authors, with 40 participants, including
232 client representatives of six major Indian metro-rail projects along with technology providers and delivery
233 teams. The workshop provided a perspective on the digitization of this project in the broader landscape of
234 Indian metro rail construction.

235 • Water infrastructure project in the United Kingdom (Case 3): To understand the codification challenges in
236 this case, eight semi-structured interviews were conducted. All eight interviewees had more than ten years of
237 experience in the construction sector and had worked with different major projects in the United Kingdom
238 and abroad. The interviewees' areas of expertise covered design, planning, project engineering, digital
239 engineering, prefabricated construction, and information management. All interviewees had teams working

240 with them on their areas of specialization and also interacted with the other stakeholders in the project. These
241 characteristics make them ideal for case-based research. Following the semi-structured approach ensured
242 participants would talk broadly on their experiences with information flow using digital collaboration tools.
243 Seven out of the eight interviews were recorded, and transcripts were made from the recordings. In addition
244 to the interviews, the first author conducted multiple visits over two weeks to the project office, observing
245 meetings and the work practice. The first author also had access to CDE and documents such as a construction
246 programme, look-ahead schedule, and method statements.

247 Data analysis took place in three steps. First, each of the cases was separately analyzed. Second, the cases
248 were compared and contrasted. The initial analyses were conducted during data collection, so early analyses
249 focused and directed later data collection. The ‘within case’ analyses and the ‘cross-case comparisons’ also led to
250 an iteration between these steps of data analysis as the comparison across cases was instructive in directing
251 analytic attention within cases. Finally, the third step was a more in-depth analysis of the third case study (for
252 which there was more detailed information).

253 A starting assumption of this current study is that there are data quality issues caused by the way digital
254 tools and workflows are used in late design and construction stages. This research thus addresses the questions:
255 How do codification challenges arise because of the different digital workflows and working practices across
256 projects? How do these lead to data quality issues? ‘Within case’ analysis of the multi-storey student apartment
257 raised issues of data interoperability, information loss, use of 2D CAD, and lack of detail in the schedule. In the
258 metro rail project, the issues of data interoperability, use of 2D CAD, and lack of detail in the schedule were also
259 present, but there were also issues of unstructured communication channels, document control, and lack of skills
260 to adopt digital technologies. In the water infrastructure project, additional issues were identified concerning
261 problems with CDE, lack of process codification, and long processing times. In the water infrastructure project,
262 the design workflow, work package plan, construction programme, BIM models, and drawings in the CDE were
263 studied to understand the level of detail and machine readability of the documents. Coding was done on the field
264 notes and interview transcripts to identify different issues related to codification and information sharing. The
265 software was used to track the patterns emerging from these data. These codes were organized to find themes.
266 The identified themes were then analyzed based on the data quality dimensions to understand their implications
267 about data quality.

268

269 **4. Codification challenges in construction**

270 Table 2 summarizes the codification challenges observed from studying the projects. Low machine
271 readability of data is a significant challenge for codification, which was observed across the projects. Product
272 information is well codified through BIM, CAD drawings, analysis models, and so on in all the projects. However,
273 the codified information is distributed amongst different formats and databases, limiting the application of
274 analytics. In addition, multiple modes of communication, multiple CDEs, and lack of process change also limit
275 the codification of information in the projects studied. Different codification challenges observed in the cases have
276 been mapped in Table 2. These topics are discussed in detail in this section.

277 <<Insert Table 2 here>>

278 **4.1 Software usage**

279 This section presents the codification challenges related to software usages such as interoperability,
280 information loss during conversions, and multiple modelling techniques during the late design and construction
281 phase. The implication of these challenges on the data quality with respect to dimension accuracy, completeness,
282 accessibility, and data provenance is explained in this section.

283 **4.1.1 Interoperability**

284 Interoperability was raised as a central problem by the interviewees, especially when working with
285 multiple CAD tools in projects, resulting in data loss during format conversions. Even within the same software
286 environment, there are problems related to data compatibility whilst working between different software versions.

287 *“Sometimes the drawings where I am using this MicroStation, but sometimes they were drafted from*
288 *the client, let’s say, in Autodesk. And transferring things from Autodesk to MicroStation, you lose*
289 *data [...] the 2019 will open the 2018, 2017, and 2016. But when the 2020 comes into play, then you*
290 *cannot open it anymore with the 2019 files that are generated with the 2020” (Technical manager,*
291 *C315)*

292 Here, the drawings are made using different CAD tools such as AutoCAD and MicroStation. However, the data
293 may not be opened (or edited) using the same tools with which they were created. This creates problems of
294 interoperability and loss of information when data in one format is converted to another. Similar problems also
295 occur when using multiple versions of the same tool.

296 This issue was observed in all three cases but at different scales. Uneven distribution of tools would result
297 in issues of data interoperability. The student apartment (Case 1) and the metro project (Case 2) had predominantly
298 used tools from a single software vendor (Autodesk for the student apartment and Bentley for the metro project)

299 for executing most of their tasks, resulting in better interoperability when compared to the water project (Case 3),
300 which uses tools from different software vendors. The scale of the project has an influence on this diverse
301 distribution of software use. The scale of the student apartment was smaller than that of the metro and water
302 projects, with a leading firm involved in both design and construction, resulting in evenness of software usage
303 (contractor office visit, C1S1). Even though the metro project had different firms engaged in design and
304 construction, the information management was handled by an owner support organization, resulting in evenness
305 in the data (BIM consultant 1, C2I4). The water project, on the other hand, had multiple firms working on different
306 phases of the projects with their own sets of tools, resulting in issues of data interoperability (digital engineer,
307 C3I3).

308 Data interoperability is a significant problem when it comes to data quality and machine readability. If
309 the data is locked to a proprietary format, it limits the application of data science. Software vendors provide a
310 proprietary Application Program Interface (API) to access the data. However, access to the data through APIs is
311 limited, and information access is limited to proprietary domain-specific programs. This limits data science as
312 different systems cannot talk to each other and derive the context from the information. Although there are open-
313 data formats available for the exchange of information, these are not often used in the construction phase and are
314 only submitted at specific data drops. Evidence also shows that there is a loss of information when converting
315 between formats because of inefficient exporters and importers. This problem reflects data quality related to the
316 accessibility dimension. As long as this data remains inaccessible, algorithms make inferences with limited data,
317 resulting in incomplete inferences.

318 ***4.1.2 Information loss during the conversion***

319 For structured information flow through a CDE, the files are converted to a PDF format. The original file
320 may be uploaded as a supporting document, but this is not a necessary requirement.

321 *“If we’re conveying CAD information, it’s being uploaded as a supporting file. There’s a facility in*
322 *[CDE1] that when you upload a PDF, you can also upload a secondary file.” (Information manager,*
323 *C3I4)*

324 The student apartment and the water project followed a workflow which required documents to be
325 uploaded as a PDF to the CDE, resulting in the loss of information during conversions (access to CDE C1D4,
326 access to CDE1, C3D2). However, the metro project used a workflow without this requirement, resulting in
327 retaining the information (digital project management, C2D2).

328 The conversion from native formats to a PDF format results in the loss of semantic relationships
329 embedded in the file. The loss of the semantic relationship between datasets results in data silos and limits machine
330 readability. This aspect of the loss of information results in incomplete information and lowers the data quality
331 related to the completeness dimension. Furthermore, during the conversions, the metadata related to the original
332 file is lost. This lowers the data quality dimension associated with data provenance.

333 **4.1.3 Multiple modelling techniques**

334 Software tools allow different methods for creation of information. However, not all methods lead to the
335 information being reusable. The modelled information would have multiple uses, which may not be known to the
336 creator of information.

337 *“So, if they use the wrong tool to model something, you don’t have the appropriate dataset to it [...]*
338 *when you go into your authoring tool, do I use a slab tool or do I go and use a generic solid*
339 *modelling tools then try and attach a dataset to it[...]if they’re not, therefore we have to go in there*
340 *and say, well I can’t just say there’s a slab now, that’s just a piece of geometry” (Digital engineer,*
341 *C3I3)*

342 For example, a slab could be modelled as a generic solid model with a dataset attached to it or as a slab component.
343 From the human point of view, the information contained in both models is the same. However, during automated
344 quantity take-off, the slab modelled as a generic solid would not be considered as the computer cannot classify it
345 as a slab.

346 *“For tunneling purposes, when you try to extract 2D drawings from 3D BIM models, those drawings*
347 *are not as correct and as detailed as they used to be. They have glitches, they have errors”*
348 *(Technical manager, C3I5)*

349 This issue was observed predominantly in the water project. Limited observation of this issue in the student
350 apartment and the metro project can be attributed to the absence of multiple firms in modelling the data.

351 Organizational divisions in large projects lead to lower machine readability and data quality because of
352 differences in modelling approaches and software tools used. Software tools offer different approaches to model
353 the same information at the same level of detail. Moreover, modelling approaches used by the firms are guided
354 by the norms and practices followed in the firm. These norms may be different for the firms who use the data.
355 Whilst examining the cases of the student apartment and metro project, where the information modelling is
356 performed by a single firm, the issues such as interoperability and improper modelling of information were limited.
357 However, in the water project, where the information modelling spanned over different firms, there was evidence

358 of issues of interoperability of data and improper information modelling. Therefore, the information created could
359 be used as a digital submission but limits further use. Even when the information is present in the model in the
360 correct format, the level of detail of the modelled information is less than desired, making the information not fit
361 for further use. This problem reduces the quality of the data concerning accuracy and completeness. The fact that
362 the data exists but not in the way it was supposed to be a case of syntactic inaccuracy.

363 **4.2 Information sharing**

364 This section presents the codification challenges in construction, such as unstructured information
365 sharing, drawing and file-based sharing, document control issues, and lack of process change. The implications
366 of these challenges on the data quality dimensions are presented in this section.

367 **4.2.1 Unstructured information sharing**

368 Information shared over modes such as meetings, reports, e-mails, etc., contains relevant data for decision
369 making. This information is embedded in documents, is shared in a human-readable format (documents,
370 PowerPoint presentations, drawings, etc.), and is in formats which are both human and machine readable (e.g.
371 spreadsheet, BIM, etc.). The main limitation of the information shared through these unstructured channels is its
372 accessibility, which is limited to people involved in the meeting or e-mail conversation.

373 *“A guy made a design on a spreadsheet for quantities. Some people knew about it; he logged it as*
374 *well. And I didn’t know that at all. So, at the very end of the day, on the eleventh hour when I have*
375 *finished everything, by the way, we have this spread sheet, and it’s exactly what I wanted to do”*
376 *(Technical manager, C3I5)*

377 In this case, the information required for a task was already available. However, it was not accessible for the
378 person who needed it, who was not part of the group within which the information was shared. This led to the
379 recreation of the information and loss of productive time. The unstructured information-sharing issue was
380 observed in all cases. The office visits in the student apartment (C1S1), meetings with project participants in the
381 metro project (C2I1, C2I2, C2I3), and interviews with participants in the water project (Table 3) revealed the
382 problem of unstructured information in the projects.

383 <<Insert Table 3>>

384 Even when there are structured workflows for information sharing, project participants find it easier to use
385 the unstructured channels of communication. They often find structured information flow through the CDE slow
386 and complicated. Despite being more traceable and accountable compared to unstructured information-exchange
387 practices, the complexity of the new structured methods for information sharing and the poor understanding of

388 workflows across the teams lead to the use of a combination of structured and unstructured channels for
389 information sharing. This issue, however, has an implication on data quality, lowering it with respect to the
390 accessibility dimension since data is not available in a common repository. Instead, it is distributed in different
391 silos and e-mail databases, and the access is limited. In addition, it introduces inconsistency as the same
392 information is distributed amongst different databases which are not connected or synchronized. Tracing the
393 source of the data and its history is also difficult when using unstructured channels for information sharing, thereby
394 reducing the quality of data associated with the provenance dimension. These issues limit the data science as the
395 datasets for drawing inferences are siloed and disconnected.

396 **4.2.2 Drawings and file-based sharing**

397 Even though the projects follow BIM level 2, the engineers interviewed are more comfortable performing
398 submissions and approvals using drawings rather than model-based information sharing. This is mainly because
399 they find it intuitive to use drawings.

400 *“I use, I’m not very good at, but I use all the navigator tools that we’ve got here. But I prefer to use*
401 *AutoCAD because I find it a lot easier” (Project engineer, C3I1)*

402 Despite having a BIM coordination tool, Bentley Navigator, the project engineers resort to using the CAD tool
403 because they find it easier. In the contractor’s offices of the projects studied, the engineers had drawings on their
404 desks and the CAD software opened on the monitors. If they find errors in the drawings, they modify and edit
405 them first on paper and then on the computer.

406 *“Because I have to open up every drawing individually and print them all, or even if you do it the*
407 *other way around, it’s very slow anyway. And then, once I’ve printed them, reviewed them all, you*
408 *could do it on there [computer] but it’s not the best way because we haven’t got the technology. I*
409 *haven’t got a big screen” (Project engineer, C3I1)*

410 The engineers use laptops with small screens. Some of the hot desks have an additional screen; however, these
411 were also small (less than 23 inches). This is highly inconvenient when reviewing large drawings as they pan and
412 zoom to detect mistakes. This issue was also observed in the student apartment and metro projects.

413 During the visit to the contractor’s office (C1S1) in the student apartment, the first author observed multiple
414 discussions between engineers using drawings as a common representation medium. On the site, a tablet-based
415 application was used to open the drawing. Similarly, in the metro project, the workflow for design, review, and
416 approval (C2D2) presents how drawings are reviewed and processed using the CDE. Such evidence points towards
417 the drawing and file-based information sharing in the construction phase.

418 The file-based sharing impacts the data quality dimensions associated with temporality and consistency.
419 Most of the file-based manipulations happen within the computer and are uploaded only when complete. Hence,
420 there is a mismatch between the rate of the volatility of data and currency of data. The volatility is high as the data
421 is manipulated on the users' desks (for example, they are manipulating the information in a printed drawing).
422 However, the currency of the data is low as it is uploaded as a batch. That means the data is updated at a lesser
423 speed than it is varied. This has an implication on the data quality associated with the timeliness dimension as the
424 data in the CDE is not the latest version. File-based sharing impacts consistency too. The files act as individual
425 entities and have information from related files in them. Unlike a model, this information is not connected.
426 Therefore, when the source is updated, the information in the file may not be updated, which introduces
427 inconsistency problems.

428 **4.2.3 Document control bottlenecks**

429 Document control plays a vital role in the flow of structured information in the projects, and document
430 control professionals are tasked with managing the access, version control, and availability of documents. Before
431 being uploaded to the CDE, such documents must be approved by the relevant authority (depending on the
432 document). The quantity of documents uploaded to CDEs in the projects studied is enormous, and in each project,
433 processing information becomes a significant task, with bottlenecks in the process leading to workarounds and
434 data quality issues. As the authorization of documents is limited to specific individuals, they tend to get asked for
435 a huge volume of authorizations, and this slows down the information flow.

436 *“So, there’s certain people who are responsible for issuing information or authorize certain*
437 *communications, and if they’re not available then things can stop. Or they may have a high volume*
438 *of these authorizations to do that it takes them a lot of time to get through” (Project planner, C3I2)*

439 This process of authorization means that multiple versions of designs can be circulating in different parts of the
440 project. For example, whilst one design is in use by the construction team, in the meantime, the designers could
441 have progressed the design, and the latest design is not uploaded as it is in the queue to get approved.

442 *“Design teams and checkers and approvers could be progressing designs but then it’d be held up*
443 *when someone say high up needed to actually approve the whole design” (Digital engineer, C3I3)*

444 *“It’s just making sure that I’m getting the latest information and no-one’s updating it in the*
445 *background and then the right versions are going onto [CDE1] [...] it’s no longer the most current*
446 *version anymore by the time I’m reviewing it” (Project engineer, C3I1)*

447 “[...]when this revision has been updated from the designer to revision 10 and I sit here on my desk
448 checking the revision 1 and the designer has the revision ten, then that revision 10 is internal [...] because he keeps on updating but he hasn't he hasn't put it on the [CDE].” (Technical manager,
450 C315)

451 This shows that the information available in the CDE is not the latest version, thereby reducing data quality
452 associated with the timeliness dimension. The CDE has provisions for labelling status of a document as work-in-
453 progress. However, even with that, it is hard to ascertain whether the information at hand is the latest as the work-
454 in-progress documents can only be accessed by internal teams. This accessibility becomes a dimension of data
455 quality which such approval processes make challenging. Additionally, even when the information is internally
456 approved, it does not get stored on the CDE. Another approval is necessary to share the information with other
457 stakeholders, thus limiting other stakeholders and algorithms from accessing this information for further decision
458 making. In the metro project, an innovative approach to address this issue for drawings was implemented, placing
459 a quick response (QR) code (a matrix barcode) in the document and a mobile app to scan the QR code and inform
460 the user whether the drawing is the latest version or not. However, the document status must be continuously
461 updated to make this useful and is limited to the issued drawings and not the work-in-progress drawings.

462 Although it seems to be straightforward from the outset, many users find document control frustrating
463 because submission ends up being a long process even when all the attributes are correct. For example, engineers
464 submit the packages to the document controller in their firm, who sends it to the document controller in the other
465 firm (receiving end), which is then sent to the team lead and, finally, to the user who would get the useful
466 information out of that package. This is a long process with checks and iterative cycles involved in each stage.
467 The document controller makes sure that the files in the CDE have the relevant attributes before they are published
468 in the CDE. If there are missing attributes, the submission is rejected. In the case of a Request For Information
469 (RFI), if the document controller does not understand the request, it gets rejected. At times, it takes more than two
470 to three weeks for the document to reach the recipient whilst following the document control workflow. This is
471 essentially slowing down the whole information flow by implementing a system which was supposed to speed up
472 the process.

473 *“It took us three weeks to actually get the package with all the right documents and revisions in
474 there, and that's a long time; again [...] I gave it to my document control and my document control
475 sent it to the other company's document control, the document control there sends it to whoever the
476 lead is, and the lead then sends it on to whoever's doing the work – and that may take a week or*

477 *two. And that's completely wasted time, and no one in the middle of that process has done any*
478 *work[...] And actually, by the time it gets to the people who are reviewing the actual technical data,*
479 *it may be three or four weeks later [...] In fact, I did it yesterday, I sent a load of RFIs through to*
480 *[designer] informally, five minutes after I'd sent it through my document control.” (Principal*
481 *engineer, C318)*

482 To bypass this obstacle, workers send the information through an unstructured channel in addition to the structured
483 workflows. This is because of poor understanding of document control workflows amongst the project participants
484 regarding the requirement of these structured workflows and the CDE. The slow processes and the need for
485 completing the task before deadlines force employees not to follow document control workflows.

486 *“There is generally a poor understanding of document control requirements, certification*
487 *requirements [...] we're finding that general good practice that people should have brought with*
488 *them from other projects is being conveniently put to one side for the purposes of expediting the*
489 *work that people are being asked to do” (Information manager, C314)*

490 In addition, the workflows in the CDE are complex and not intuitive, making it difficult for users to follow the
491 protocols for document control.

492 *“Just because of the way they need to store it in certain places and stuff like that, and it can't be*
493 *done... The way that [CDE1] is set up here, I believe it is not easy to use” (Project engineer, C311)*

494 *“It seems very complex, you open them up, there's lots of things going on [...] I just want to know*
495 *where I can get my latest drawing” (Digital engineer, C313)*

496 This has resulted in employees bypassing the workflow, which leads to system conflicts and further delays in the
497 processes and, at times, in the information being stored on the CDE.

498 *“I think someone within the doc management system had obviously circumnavigated it somehow, to*
499 *get the drawings out. And then when we were trying to get the said revisions for our set out, the*
500 *system wouldn't allow it because directory hadn't been properly created.” (Technical manager,*
501 *C315)*

502 Not following the document control workflows leads to information loss in the CDE. This is a major setback to
503 codification as the data is stored in an unstructured way which is difficult to access. This problem was found in
504 the metro project as well. Conversations with the project manager (C211), chief site engineer (C212), and BIM
505 consultants (C214, C254, C216) revealed that the project team weren't exposed to structured information flow
506 used with digital technologies in the past. This made the implementation of CDE-based structured workflows

507 difficult despite the training given to the participants, resulting in a combination of structured and unstructured
508 workflows in the project.

509 The complexity of workflows and document control measures has implications for data quality. There are
510 shared norms, values, and expectations for the users regarding the tools, such as speed, easy communication,
511 transparency, and so on, which were developed based on their previous experiences of collaboration. When the
512 new tool does not meet the expected qualities, it reduces their productivity, and users move back to the older ways
513 of information sharing to expedite the task. When users find it difficult to utilize the CDE for structured
514 information flow, they bypass the workflows to get the work expedited. This leads to the loss of metadata,
515 document trails, and information dependencies as these unstructured communication channels offer limited or no
516 codification. In addition, the document control workflow itself makes the process slow. Document control
517 bottlenecks have multiple implications for data quality. Firstly, the value for the timeliness dimension is lowered
518 as the data which is published might not be the latest. Hence, the inferences are based on old data, which leads to
519 false interpretations. Secondly, this lowers the semantic accuracy of the data as its attributes might no longer be
520 true. This also introduces the problem of consistency. Depending on which database employees look at, they see
521 different values. For instance, one CDE to which the information was packaged would have the latest value, whilst
522 the one which must go through another document controller would have a different value. This tampers with the
523 idea of a ‘single source of truth’. In addition, when users circumnavigate the workflow, there are more data quality
524 issues related to unstructured information sharing such as accessibility and provenance.

525 ***4.2.4 Lack of process change***

526 Even though structured information flow is digitized through the introduction of a CDE, the process
527 enabling information flow remains unchanged. For example, for a piece of information to be approved, it must be
528 printed, associated with a cover sheet, and signed.

529 *“We’re actually going to export that out of the CDE, we’re going to print it out, we’re going to*
530 *staple it together, we’re going to put our own cover sheet on the front of it, with the exact same*
531 *details on the back and we’re going to go off and go and get three signatures, scan it back in, put it*
532 *back into [CDE1] and submit it.”(Project engineer, C311)*

533 Printing and scanning the document results in loss of metadata. A scanned document in a PDF format has little
534 machine-readable information in it. Inferring the contents from a scanned document is also resource intensive
535 when compared to its original form. In the process of printing and scanning, the content becomes digitally
536 accessible but not machine readable, thereby limiting the application of data science.

537 *“What I'm finding now is it's not actually speeding everything up, it's sort of making everything a*
538 *lot slower; which I find very frustrating” (Principal engineer, C318)*

539 The lack of change in the processes reduces the value in the adoption of a CDE as it increases the time to do these
540 tasks rather than a total reduction in time.

541 In the water project, three CDEs were used for the project, which created issues such as double handling,
542 data inconsistencies, and so on. The presence of multiple CDEs in a project is another example of the lack of
543 process change. Multiple CDEs resembles the paper-based workflows such as document flows between the
544 designer and contractor, another set of document flows within the contractor's office, and another set of document
545 flows to the clients for approvals.

546 *‘Our client prefers [CDE1], and we have the designer who stores things in [CDE2], so we have*
547 *both of those tools, and we have to balance between those two. That can be very confusing when we*
548 *have two platforms’ (technical manager, C315).*

549 For the information submissions to the clients, CDE1 was used; for the information from the design consultants
550 to the contractor, CDE2 was used; and for internal file handling and sharing with the contractor, CDE3 was used.
551 CDE1 and CDE3 came from the same vendor. However, CDE2 was from a different vendor.

552 *“Everything had to be taken out of one data environment and pushed into another. One of the issues*
553 *with that is the consistency or the compliance or knowing the latest versions of information” (Digital*
554 *engineer, C316)*

555 The existence of multiple CDEs within a project introduces the problem of data inconsistencies. Documents must
556 be taken out of one CDE and placed in another. When the volume of information is huge, with each file having
557 multiple versions, it is difficult to maintain consistency of documents across multiple CDEs. This means that
558 information in a CDE might not be accurate and up to date, leading to incorrect interpretations when data analytics
559 is performed on it.

560 *“As the contractor, then we have to deliver it to a completely separate, disconnected CDE [...]*
561 *we're double-handling” (Digital engineer, C313)*

562 When the CDEs are disconnected, the document trail is lost when a document is moved from one CDE to another,
563 leading to the loss of traceability.

564 Lack of process change has multiple implications on data quality as well. Printing and scanning remove
565 metadata and data relationships from the files. A scanned version of the file would also have very little machine-
566 readable information embedded in it and would require resource-intensive methods to extract insights from it.

567 This reduces the data quality dimensions such as accessibility (as the metadata and data relationships are removed),
568 completeness (information is not complete), and provenance (document trail is lost in the process). In addition,
569 having multiple CDEs introduces the data quality issue associated with provenance as files residing in multiple
570 CDEs are disconnected, and the relationships of that particular file with another file are lost in the transition
571 process. There are further issues with synchronization of the information when it is distributed in multiple CDEs.
572 This introduces the data quality issues associated with consistency, which limits the quality of findings made by
573 inference algorithms.

574 **4.3 Construction process information**

575 This section presents the codification challenges related to construction process information. The data
576 analysis shows that product information is relatively well structured as BIM models, analysis models, and CAD
577 drawings. However, construction process information is relatively less structured and detailed when compared to
578 product information. Process information is codified as Gantt chart models in scheduling software and then linked
579 to the BIM model. The detailed process information is not structured into a model. Instead, it is shared as method
580 statements in less structured PowerPoint presentations and PDF documents. Sharing information in these
581 unstructured formats has implications for data quality, which are presented in this section.

582 **4.3.1 Loss of constraint information**

583 In construction, the constraints for any activity execution are discussed during team meetings as the
584 constraints span between different teams. For example, logistics constraints span amongst prefabrication, logistics,
585 and site teams. These discussions lead to the removal of constraints by rearranging the start and end times for the
586 activities. These are then translated to Gantt charts as an output.

587 *“So, from an engineering perspective, we have to interpret engineering information, whether it be*
588 *drawing or written constraints, written narratives and interpret those into a Gantt chart. So, we*
589 *physically need that information to know what we’re building and what the constraints in building it are.”*
590 *(Project planner, C3I2).*

591 The above statement from the project planner provides evidence on the processes used to convert the constraints
592 into a Gantt chart for communication. However, during this process, many of the constraints themselves, and thus
593 context information for rearranging the activities, are lost. This is because Gantt charts can hold only precedence
594 constraints. Other constraints, such as disjunctive (where activities cannot overlap) and logical constraints, are not
595 embedded into the Gantt model. Instead, they are retained only as tacit information by project participants
596 involved in team meetings. This is a case of incomplete information within the dataset as this information is only

597 accessible to the meeting participants. For example, one of the meetings in the water project had an issue with
598 piling, where the pile-driving equipment did not have access to the site for a specific date as there was another
599 activity going on which limited the width of the site access road. At the meeting, this was raised:

600 *“Access chamber works will conflict with access road for pile work, piling work package has to be moved*
601 *back 2 weeks.” (Progress review meeting, C3M3)*

602 Here, there is a dependency between access chamber works, and the piling work package as the access chamber
603 work would reduce the road width. Therefore, the piling activity was delayed to a later date. The constraint was
604 removed. However, the knowledge that there was a constraint is not recorded, and thus the presence of that
605 constraint is not codified. This means such constraints are not machine readable as the access to this information
606 was limited to the participants of a particular meeting. If an automatic scheduler is used to reschedule these
607 activities, this rescheduling activity would not have access to this information, resulting in an unrealistic schedule.

608 Similar issues were observed in all cases. This issue introduces data quality problems associated with
609 accessibility (information is limited to people who attended the meeting) and data completeness (the model does
610 not include any constraint information; hence, it is incomplete).

611 **4.3.2 Low level of detail**

612 The precedence information codified into Gantt charts is linked with BIM to create 4D BIM simulations.
613 However, the lack of detail in work packaging and associated information (such as constraints and resources)
614 results in misinterpretation from the 4D models. The metro project follows a 5D BIM workflow (digital project
615 management PPT—C2D2), where the schedule is linked to a BIM model, Bill of Quantities (BOQ), and an
616 Enterprise Resource Planning (ERP) system to compare the cost based on the quantities versus the cost stated in
617 the work orders from the subcontractors. The progress information is also linked to this model to ensure that the
618 work is done before sanctioning the bills for the work orders.

619 *“Work package for three spans were linked to a work order. Model showed the deck for a span was*
620 *completed before the pier supporting it was completed because the work package for the first span*
621 *was reported as completed.” (Field notes- BIM Consultant 2, C215)*

622 <<Insert Figure 2>>

623 The deck of the metro can be completed only when the pier supporting it is completed, as shown in Figure 2.
624 However, the system recorded the deck assembly to be completed when the pier was not completed. This was
625 because work packages for the deck and pier were different as they were done by different subcontractors, and
626 the level of detail of the work package is low. A subcontractor who dealt with deck assembly had a part of the

627 work package completed, but the lack of detail in work packaging triggered the computer to record the whole
628 work package as completed, resulting in the error. This is a clear case of lack of detail in the model leading to
629 wrong inferencing.

630 Similar issues were observed in the student apartment (Case 1) by examining the schedule data in the
631 construction programme (C1D1) and the water project (Case 3) by examining the construction programme update
632 (C3D4) and observing the review meeting (C3M3). These issues are caused by low data quality due to incomplete
633 information related to the completeness dimension.

634

635 **5. Discussion**

636 This section discusses the software usage, information sharing, and construction process information
637 codification challenges which limit the uptake of data science in construction, drawing on the evidence from the
638 empirical study. The discussion relates the findings to the literature on BIM use in practice (e.g. Dossick and Neff,
639 2010; Harty and Whyte, 2010) and other strands of research on data quality, machine readability, and BIM
640 adoption and implementation to articulate how these new analyses contribute by extending understanding of
641 codification challenges. Furthermore, building on and extending Batini and Scannapieco (2016), it shows how
642 these codification challenges are then mapped to their data quality dimensions, such as accuracy, completeness,
643 timeliness, consistency, accessibility, and data provenance.

644 **5.1 Software usage**

645 The findings on software usage show that, despite significant digitization of work processes, data remains
646 fragmented into different domains and formats because of the multiple software tools in use across the
647 organizations involved in construction. Work in the construction information technology community is pioneering
648 new data management solutions to improve interoperability (Hu et al., 2016; Pauwels et al., 2010; Pazlar and Turk,
649 2008; Redmond et al., 2012), and it is disappointing to find that construction projects still suffer from poor quality
650 data as a result of problems of interoperability caused by the existence of multiple domain-specific tools and
651 modelling practices. In their work, Dossick and Neff (2010) have previously shown how the organizational and
652 cultural divisions between the designers and builders, contractors, and subcontractors stifle collaborative work.
653 This paper shows these issues are not resolved. In the projects studied, organizational and cultural divisions
654 between the firms involved in the late design and construction stages of projects cause software usage problems
655 (interoperability, information loss during format conversion, multiple modelling techniques). Whilst there may be
656 shared norms and tools within a firm for modelling information, these norms differ across the firms which

657 modelled project information. Multiple modelling techniques (as described in 4.1.3) and data created using
658 different software (as described in 4.1.1) result in datasets which are not interoperable and require format
659 conversions, resulting in loss of information and low machine readability. The water project (Case 3) had multiple
660 firms working on the data over different phases of the project, with interoperability problems more prevalent in
661 this case in comparison with the student apartment (Case 1), in which a single leading firm was involved with the
662 creation and use of the model. Although the metro project (Case 2) had different firms over the different phases
663 of the project, digital data creation was handled through a single owner support organization, limiting the impact
664 of this problem. Software usage problems are found to lead to challenges of codification for data science; hence,
665 this work extends prior insights by Dossick and Neff (2010) to show how organizational and cultural divisions
666 between designers and builders not only stifle collaborative work and joint problem-solving but also result in
667 fragmented datasets in construction, leading to data silos and data loss and, thus, resulting in poor data quality
668 which is more difficult to use in data science.

669 As it is relatively unusual and potentially undesirable to have one firm or owner with overall control of
670 the model, to enable more distributed working, developers of new tools or digitally enabled processes should
671 consider the implication of organizational separation in the sector in addition to the technical requirements. In
672 their work, Dossick and Neff (2010) have described the influence of strong leadership to hold people together on
673 a project to improve collaboration despite professional segregation. Similarly, a set of common practices and a
674 larger vision for the data creation and management should be laid out in the project to ensure the data meets the
675 necessary quality to enable its use without loss of information in between. To achieve better-quality data in
676 projects, practitioners must focus beyond the individual scope of their multiple firms towards the common goals
677 of the project.

678 **5.2 Information sharing**

679 The analyses suggest that the construction sector has not yet made the transition from document-based
680 to model-based ways of organizing digital data. The use of drawings and file-based sharing, unstructured
681 information sharing, printing and scanning of documents, multiple CDEs, and so forth in information sharing has
682 a significant impact on the machine readability of data. Paper-based practices are institutionalized in the sector,
683 and while they are being replaced by digital ways of working, this change is slow, with users of construction
684 information still conditioned to work with drawings and PDFs and unstructured information sharing. Even in
685 projects which are championing newer BIM-based workflows using CDEs, this work finds it is difficult to replace
686 these practices, as evidenced by the problems associated with information sharing (section 4.2). The complexity

687 and long processing times involved in these workflows force users to shift back to existing practices and
688 workarounds to expedite their work. The findings from this paper also support the previous characterization of
689 users in construction combining new structured methods of information sharing along with the prior practice of
690 unstructured information sharing when they were hindered by bottlenecks in processes, such as document control.
691 Thus, we can characterize the project participants use a range of new and existing practices together as ‘hybrid
692 practices’ (Harty and Whyte, 2010), and their circumnavigation of workflows results in unstructured information
693 sharing (as shown in a previous discussion in Whyte et al. (2016)). However, this paper goes beyond such studies
694 to characterize the implications for data quality and to highlight, building on Hartmann (2008), the potential to
695 develop newer workflows and digitally enabled processes which address the challenges faced by practitioners.

696 **5.3 Construction process information**

697 Regarding construction process information, this paper shows that the process codification is limited to
698 the master planning or phase planning level and lacks the level of detail and linkage required for the application
699 of data science tools. Whilst product information is relatively well codified in BIM, process information is less
700 well detailed, and there is a lack of constraint information. These challenges are identified by researchers
701 implementing 4D BIM. For example, Han and Golparvar-Fard (2017) have stated that the process modelling
702 methods fail to document field issues to be made available for further analysis; for instance, the 4D BIM’s “Model
703 Breakdown Structure typically does not match operational details or require creating complicated namespaces
704 which, without visual representations, are difficult to communicate” (p. 1733). Giretti et al. (2012) have further
705 reported the lack of correlation between the resources employed hourly and work progress. This led to the
706 decomposition of tasks into sub-tasks to determine causal relationships between the involved variables so the
707 whole progress could be determined. This study suggests that to overcome such reported issues, the methods for
708 codifying construction process information must be more detailed. The institutionalized practice of planning being
709 limited to master planning and phase planning, without the focus on granular planning such as look-ahead planning
710 and weekly planning, is causing this codification challenge, with the lack of semantic relationships embedded in
711 the model limiting the application of automatic schedulers. These issues suggest both a change in the modelling
712 of process information in construction, with the need to develop tools which support modelling of complex
713 constraint information, and also a change in the practice to codify the process information in greater detail so that
714 data science could be employed to augment decision making in construction.

715 **5.4 Machine readability of construction datasets**

716 This section discusses the machine readability of the construction datasets. Common construction
717 datasets are classified based on the set of ‘rules’ for creating structured data described by Berners-Lee (2006) in
718 Table 4.

719 <<Insert Table 4>>

720 Most of the construction information observed from the cases satisfies the requirement for a one-star
721 category. The observed projects use a CDE for storing and managing project data, resulting in indexing the data
722 and storing it on online servers, resulting in one-star data. CDE makes the data easier to retrieve for the computers
723 to make inferences on them. However, the complexity of the new structured methods for information sharing
724 using CDE, and the poor understanding of workflows across the teams, leads to the use of a combination of
725 structured and unstructured channels for information sharing, as discussed in section 4.2. This aspect reduces the
726 machine readability of the information distributed over unstructured channels as the information is not indexed
727 nor available on a common server. The same issue occurs when the users circumnavigate the workflows to get the
728 work expedited. Similarly, codification challenges associated with construction process information also lower
729 the machine readability as the information is not recorded (lack of detail and loss of construction information)
730 and, hence, not indexed or stored in online servers. These issues make the information inaccessible for inferencing.

731 With regard to the structure of the construction information, construction datasets in the form of BIM,
732 project management information in project management software (Primavera P6, Asta powerproject, etc.), outputs
733 from Microsoft tools such as Excel, and so forth are structured, satisfying the requirements for two-star data.
734 However, construction data are also unstructured in the forms of PDF documents, drawings, and other file-based
735 formats, as described in sections 4.2.2 and 4.2.4. The lack of structure in the datasets makes inferencing from
736 them difficult, leading to the need for complex algorithms. Where the construction data is structured, the data
737 structure is in proprietary formats which require the APIs to access the semantic relationships within the data. The
738 observed projects do not use open formats or open standards for publishing the data. Proprietary tools for the
739 authorship of construction data are far more advanced and easier to use than the open-source tools. Hence,
740 construction projects resort to using the tested and robust proprietary tools, resulting in issues associated with
741 interoperability and loss of information presented in section 4.1. Thus, the construction information rarely
742 achieves three-star classification, as mentioned by Berners-Lee (2006). In conclusion, the maximum level of
743 machine readability of the construction datasets in the observed projects is two-star, with most of the information
744 with a one-star rating.

745 Low machine readability of the data has implications on data quality. When the construction information
746 is not stored on servers due to information sharing issues or lack of detail in the models, the accessibility of that
747 data is affected, thus reducing the accessibility dimension of the data quality. When information is stored in CDE
748 (satisfying conditions for one-star data) as PDF documents, the structure of the data is not maintained, resulting
749 in data quality issues associated with syntactic accuracy and consistency. Lack of a data structure removes the
750 context from the information, thus resulting in the need for complex algorithms to infer the contexts and infer
751 from data. This issue also introduces problems associated with consistency as the data value for a field might be
752 different in different files, and the lack of context limits the computer programs to detect it. This problem is further
753 worsened as the datasets are not linked since the links are lost when a file is moved from one CDE to another.
754 Storing the information in proprietary formats also reduces the accessibility dimension for data quality.

755 **5.5 Implications for data quality**

756 The codification challenges discussed earlier have many significant implications for data quality. To
757 unpack these in this section, they are mapped onto the different quality dimensions.

758 <<Insert Table 5>>

759 *Accuracy:* The organizational and cultural divisions between different teams results in problems
760 associated with multiple modelling techniques, leading to data quality issues concerning the syntactic accuracy of
761 the data. This was evident from the dataset when different people had different perceptions of the model, as in the
762 case of the slab example in 4.1.3. Similarly, the hybrid practices associated with information sharing lead to
763 lowering the semantic accuracy of the data as the data with which inferences are made are not accurate due to
764 inefficiencies in information sharing. When there are syntactic errors in the data, this leads to incorrect insights
765 (for example, if a slab is modelled as a geometric object with attributes attached to it and when a software tool is
766 used to compute quantity take-off from the model for all the slabs). The output would be zero as quantity required
767 as the program fails to identify the geometric object as a slab. If this software is integrated with a costing tool used
768 for cash flow analytics, this error gets propagated into that tool. These errors can be removed to an extent using
769 semantic enrichment programs. However, even the accuracy of inferences of semantic enrichment programs is
770 dependent on the quality of input datasets (Sacks et al., 2017).

771 *Completeness:* Concerning the completeness of the data, organizational and cultural divisions between
772 the teams resulting in problems such as interoperability, format conversions, multiple modelling techniques and
773 the implication of hybrid practices such as printing and scanning the documents play a role in reducing the data
774 quality. For example, the software usage issues caused by the organizational divisions leads to format conversion

775 resulting data loss leading incomplete dataset. Lack of process change also leads to similar problems such as loss
776 of metadata when documents are printed and scanned. The institutionalized practices of process modelling with
777 low levels of detail and the practice of not codifying constraints in the model result in incomplete data. Inferring
778 insights from incomplete datasets reduces the quality of the output. For example, if the constraints are not codified
779 in a schedule, an automatic scheduler would create an unrealistic schedule. This leads to further problems down
780 the line. In the case of a product model, an incomplete dataset used for a structural capacity prediction would give
781 incorrect results.

782 *Timeliness:* This dimension of data quality is mostly affected by the information-sharing practices in
783 construction using hybrid practices by using old and new practices simultaneously. File-based sharing (resulting
784 in slowing down the data updates compared to model-based sharing), document control bottlenecks delaying the
785 submission of data into the system, and the use of multiple CDEs requiring data transfer from one to another
786 (resulting in a delay in fetching the data) all constitute outdated data-skewing analytics. When outdated data is
787 used for analytics, the resulting inference is not time appropriate, and decisions taken with those inferences lead
788 to problems. For example, if the construction plan is made using resource availability data, during the actual
789 construction date, the assigned resource might not be available. When this is managed manually, the planner
790 makes sure this does not happen. However, when performed automatically, it is necessary that the datasets are
791 updated close to real time.

792 *Consistency:* Hybrid practices in information sharing, resulting in document control bottlenecks, multiple
793 CDEs, and circumnavigating workflows, induce inconsistency in the data. The data in one CDE might be different
794 from that in another CDE. Similarly, the document control bottlenecks in publishing the data result in different
795 teams and their databases having different versions of the data. This creates a problem during the data analysis
796 phase. For example, for the same content, different values might exist, of which only one is true. This induces
797 problems of semantic accuracy if the computer takes the incorrect value for analysis.

798 *Accessibility:* This dimension of data quality is affected by organizational and cultural divisions between
799 teams, hybrid practices, and institutionalized practices in process modelling. Organizational and cultural divisions
800 between teams, resulting in the use of multiple software, multiple modelling techniques, and so forth, affect the
801 accessibility of data as employees must convert the data to access it, and there are instances which show
802 conversions result in loss of data. Hence, the data is only accessible in full content to the team, which created it.
803 Similarly, the hybrid practices, particularly printing and subsequent scanning of the documents, remove the
804 metadata as well as semantic information within those data, resulting in a less-efficient analysis. In addition, the

805 lack of constraint codification restricts the information regarding constraints, making it accessible to the few
806 people who attended the meeting. It is not embedded in the model and, hence, not accessible to any analytic
807 algorithms. This reduces the capability of such algorithms to infer accurate information and determine causalities.

808 *Provenance:* This dimension of the data quality is affected by both organizational and cultural divisions
809 in the industry as well as hybrid practices. Fragmentation in the industry leads to many conversions and much
810 manipulation of data to suit purposes, in the process removing the traceability and origins of the data. Similarly,
811 multiple CDEs, printing and scanning of documents and unstructured information sharing also lead to loss of
812 metadata and ends up removing information regarding origins of the data. This limits algorithms from making
813 inferences based on data origins and their incremental manipulation.

814 This section described the challenges of codification as causes of data quality issues, which have implications
815 for different data quality dimensions and the output of data analytics. While big data techniques suggest future
816 opportunities to use data science with an increasing variety of data of different levels of veracity, data cleaning is
817 a resource-intensive process, and significant training models are required. To improve the uptake of data science
818 in construction, high-quality data is necessary, and achieving this requires overcoming the codification challenges
819 identified here.

820

821 **6. Conclusions and future work**

822 Codification challenges for data science in construction are found to be related to: 1) software usage—
823 interoperability, information loss during conversion, and multiple modelling techniques, 2) information sharing—
824 unstructured information sharing, drawing and file-based sharing document control, and lack of process change,
825 and 3) construction process information—loss of constraints and low levels of detail. The implication of these
826 challenges was discussed by mapping them to data quality dimensions such as accuracy, completeness, timeliness,
827 consistency, and provenance. Through the identification of the codification challenges in the late design and
828 construction phase of the projects and their mapping to the data quality dimensions, this paper extends the
829 knowledge on data quality issues in construction. It shows how data quality arises from organizational as well as
830 technical practices. The persistence of organizational and cultural divisions, paper-and document-based (as well
831 as digital- and model-based) ways of working, and institutionalized practices of construction process modelling
832 are challenges for the uptake of data science as they lead to the partial codification of information in machine-
833 readable formats. Whilst the fragmented nature of the sector is well understood, this work shows how codification
834 challenges arise because of the different digital workflows and working practices across projects, and how these

835 lead to fragmented data as a result of the use of multiple software packages, poor information sharing and only
836 partially captured construction process information.

837 While all of the projects studied were BIM enabled, within these projects, information sharing issues and
838 software usage issues emerge as a result of their document-based, and sometimes paper-based, practices. Project-
839 wide standards and policies are created to streamline information sharing through structured workflows, and these
840 workflows are aimed at improving collaboration, but the process of document control raises issues. In addition,
841 this study identifies inefficiencies (such as long processing times, complexity, etc.) in these workflows, which
842 pushes the users to bypass the structured methods in the current workflows. This forces the users to revert to older
843 methods or use a combination of old and new methods, resulting in the generation of unstructured information.
844 These cause data quality issues relating to inaccessibility of data, timeliness of data, and data consistency.

845 Limited codification in construction process information is a result of the institutionalized practice of
846 scheduling focusing on codifying just precedence relationships into a Gantt chart. The advancements in 4D and
847 5D BIM have attempted to address this problem to an extent by integrating the resource and cost information to a
848 single platform. However, the data herein suggest that the level of detail of the scheduling is still at a macro-
849 planning level. In addition, the constraint relationships between resources and processes, resources and the site
850 conditions, and site conditions and the processes discussed in weekly meetings are not codified into a model. This
851 limits the data science due to data quality issues associated with completeness and accessibility.

852 How might codification challenges be overcome to enable greater uptake of data science in construction?
853 The evidence from this paper points towards a need for change in the policy and practice to ensure machine
854 readability and better-quality construction information. The current policies on information modelling are aimed
855 at improving collaboration amongst the project participants. Building on the existing policies, newer policies
856 should ensure that the machine readability and quality of the data are maintained during information sharing. Care
857 should be taken to ensure that newer policies suit the existing work practices in construction so that the workflows
858 recommended by such policies are not bypassed by the users. Newer policies also ensure the transition from file-
859 based information sharing to model-based information sharing, where the model data is shared based on the
860 principles stated by Berners-Lee (2006) to ensure machine readability. Standards under development, such as
861 ISO/DIS 21597—information container for data drop, exchange specification parts 1 and 2—are moving in this
862 direction, where linked data technologies are used to define the relationships between documents and datasets,
863 thereby automating the document control processes (International Organization for Standardization, 2019). These
864 standards would have a positive impact on data quality pertaining to dimensions, timeliness, accessibility, and

865 consistency. Even though policy interventions would have a positive effect on data quality, the construction sector
866 is known for institutionalized practices and resistance to change. Thus, policies must influence practices to enable
867 greater uptake of data science.

868 Thus, in addition to the policies, it is, therefore, necessary to adopt newer methods in practice to make
869 construction information machine readable. This study examined a student apartment building project constructed
870 by a leading contractor in the United Kingdom using state-of-the-art offsite manufacturing techniques (Case 1), a
871 metro project pioneering digital transformation in India (Case 2), and an innovative water megaproject in the
872 United Kingdom (Case 3). Despite these projects having advanced workflows and innovative approaches to ensure
873 structured information, lack of process change after digitization of workflows was evident in the datasets, which
874 included multiple CDEs, printing and scanning of documents, need for wet signatures, and such. These practices
875 continue to occur due to the lack of trust in digital workflows; thus, there must be a change in this mindset, and
876 trust in digital data must be developed in the construction practice. There is also a need for a shift to model-based
877 information sharing from the file-based information sharing between teams. To address this issue, related to a lack
878 of detail in process information, there is a need for change in the look-ahead scheduling practice to incorporate
879 the codification of different constraint relationships into information. This might involve a significant change in
880 the planning of a progress review process by integrating BIM-based scheduling tools within the planning and
881 progress meetings and modelling different constraints within the software environment to capture it from the
882 discussions.

883 The current study has identified the codification challenges by examining information sharing in the late
884 design and construction phase and mapped it to the data quality dimensions. The findings from this study inform
885 the researchers, who are developing frameworks and methods to codify construction information, of the
886 organizational issues to be considered in their work. This paper also provides the data quality implications of
887 issues associated with BIM implementation, motivating the researchers focusing on implementation studies to
888 widen their scope from collaboration to include data quality and machine readability. Future research can build
889 on this study to develop recommendations for ensuring the machine readability of construction information
890 generated throughout different stages of the lifecycle of the project. To achieve this, future research should focus
891 on different aspects. Firstly, researchers should elaborate on the organizational issues identified in the current
892 study across different phases of the project and amongst different stakeholders involved in the project. Secondly,
893 accounting for the complexity of information-use trends in construction, there is a need for research in
894 fundamental data science to pave the way for the integration of the disconnected information in the construction

895 sector. This includes developing newer information modelling approaches which adapt to the current work
896 processes as well as support codification of information, such as algorithms for crawling through the disconnected
897 information to draw insights and learning algorithms to detect discrepancies in data (such as with accuracy,
898 completeness, timeliness, consistency, and provenance) and predict their consequences. Finally, researchers and
899 managers of construction projects should work together towards developing workflows and information sharing
900 practices which ensure the machine readability of construction information whilst considering the issues of
901 fragmentation, hybrid practices and institutionalized practices. This paper has provided a foundation for such
902 future research by extending the knowledge on data quality issues in construction through the identification of
903 codification challenges, considering the wider practice of model and document-based information sharing.

904 **7. Data availability statement**

905 To comply with the research ethics process set out by Imperial College London
906 (<https://www.imperial.ac.uk/research-ethics-committee/application-process/>), an agreement was made with the
907 participants of this study. Sharing of the data, that support the findings of this study, is restricted by the agreement
908 made with the research participants. As a result, the data cannot be shared without the permission of the
909 participants.

910 **8. Acknowledgements**

911 The authors are grateful to the research participants in the three case studies. The PhD research of the
912 first author is co-funded by Bentley Systems UK and through a Skempton Scholarship from the Department of
913 Civil and Environmental Engineering, Imperial College, London. During the development of this paper, this
914 author was supported by the PhD enrichment scholarship from the Alan Turing Institute, the United Kingdom's
915 National Institute for Data Science and AI. The second author acknowledges the support of Laing O'Rourke and
916 the Royal Academy of Engineering for cosponsoring her Professorship.

References

- 917
- 918 Akintola, A., Venkatachalam, S., and Root, D. (2017). "New BIM Roles' Legitimacy and Changing Power
919 Dynamics on BIM-Enabled Projects." *Journal of Construction Engineering and Management*, 143(9),
920 04017066.
- 921 Batini, C., and Scannapieco, M. 2016. *Data and information quality*. Cham, Switzerland: Springer International
922 Publishing.
- 923 Beetz, J., Berlo, v. L. L., Laat, d. R., and Helm, v. d. P. (2010). "BIMSERVER.ORG – An open source IFC Model
924 Server." In *Proc., CIB W78 2010: 27th International Conference*, Cairo, Egypt.
- 925 Berners-Lee, T. 2006. "Linked Data - Design Issues." Accessed: 19/03/2019.
926 <https://www.w3.org/DesignIssues/LinkedData.html>.
- 927 Bilal, M., Oyedele, L. O., Qadir, J., Munir, K., Ajayi, S. O., Akinade, O. O., Owolabi, H. A., Alaka, H. A., and
928 Pasha, M. (2016). "Big Data in the construction industry: A review of present status, opportunities, and future
929 trends." *Advanced Engineering Informatics*, 30(3), 500-521.
- 930 Bolton, A., Butler, L., Dabson, I., Enzer, M., Evans, M., Fenemore, T., Harradence, F., Keaney, E., Kemp, A.,
931 Luck, A., Pawsey, N., Saville, S., Schooling, J., Sharp, M., Smith, T., Tennison, J., Whyte, J., Wilson, A., and
932 Makri, C. (2018). "The Gemini Principles: Guiding values for the national digital twin and information
933 management framework." Centre for Digital Built Britain and Digital Framework Task Group, Cambridge,
934 UK.
- 935 British Standards Institution (2018). "ISO19650 - 1 & 2 : Organization and digitization of information about
936 buildings and civil engineering works, including building information modelling (BIM). Information
937 management using building information modelling." Retrieved from:
938 <https://bsol.bsigroup.com/Bibliographic/BibliographicInfoData/00000000030333754>.
- 939 Cai, L., and Zhu, Y. (2015). "The Challenges of Data Quality and Data Quality Assessment in the Big Data Era."
940 *Data Science Journal*, 14(2), 1-10.
- 941 Cao, D., Li, H., and Wang, G. (2014). "Impacts of Isomorphic Pressures on BIM Adoption in Construction
942 Projects." *Journal of Construction Engineering and Management*, 140(12), 04014056
- 943 Carrillo, P., Harding, J., and Choudhary, A. (2011). "Knowledge discovery from post-project reviews."
944 *Construction Management and Economics*, 29(7), 713-723.

945 Chang, C.-Y., Pan, W., and Howard, R. (2017). "Impact of Building Information Modeling Implementation on
946 the Acceptance of Integrated Delivery Systems: Structural Equation Modeling Analysis." *Journal of*
947 *Construction Engineering and Management*, 143(8), 04017044.

948 Chegu Badrinath, A., and Hsieh, S.-h. (2019). "Empirical Approach to Identify Operational Critical Success
949 Factors for BIM Projects." *Journal of Construction Engineering and Management*, 145(3), 04018140.

950 Choi, B., Lee, H.-S., Park, M., Cho, Y. K., and Kim, H. (2014). "Framework for Work-Space Planning Using
951 Four-Dimensional BIM in Construction Projects." *Journal of Construction Engineering and Management*,
952 140(9), 04014041.

953 Delone, W., and McLean, E. (2014). "The DeLone and McLean Model of Information Systems Success: A Ten-
954 Year Update." *Journal of Management Information Systems*, 19(4), 9-30.

955 Donato, V., Lo Turco, M., and Bocconcinio, M. M. (2017). "BIM-QA/QC in the architectural design process."
956 *Architectural Engineering and Design Management*, 14(3), 239-254.

957 Dossick, C. S., and Neff, G. (2010). "Organizational Divisions in BIM-Enabled Commercial Construction."
958 *Journal of Construction Engineering and Management*, 136(4), 459-467.

959 Eastman, C., Teicholz, P., Sacks, R., and Liston, K. 2008. *BIM Handbook: A Guide to Building Information*
960 *Modeling for Owners, Managers, Designers, Engineers, and Contractors*. Hoboken: John Wiley & Sons, Inc.

961 Fan, W. (2015). "Data Quality: From theory to practice." *ACM SIGMOD Record*, 44(3), 7-18.

962 Farias, T. M. d., Roxin, A., and Nicolle, C. (2018). "A rule-based methodology to extract building model views."
963 *Automation in Construction*, 92, 214-229.

964 Giretti, A., Carbonari, A., Novembri, G., and Robuffo, F. (2012). "Estimation of job-site work progress through
965 on-site monitoring." In *Proc., 9th International Symposium of Automation and Robotics in Construction*,
966 *ISARC 2012*, Eindhoven, The Netherlands.

967 Goedert, J. D., and Meadati, P. (2008). "Integrating Construction Process Documentation into Building
968 Information Modeling." *Journal of Construction Engineering and Management*, 134(7), 509-516.

969 Gu, N., and London, K. (2010). "Understanding and facilitating BIM adoption in the AEC industry." *Automation*
970 *in Construction*, 19(8), 988-999.

971 Han, K. K., and Golparvar-Fard, M. (2017). "Potential of big visual data and building information modeling for
972 construction performance analytics: An exploratory study." *Automation in Construction*, 73, 184-198.

973 Hartmann, T. (2008). "A grassroots model of decision support system implications by construction project teams."
974 Ph.D., Stanford University, California, USA.

975 Harty, C., and Whyte, J. (2010). "Emerging Hybrid Practices in Construction Design Work: Role of Mixed
976 Media." *Journal of Construction Engineering and Management*, 136(4), 468-476.

977 Hendler, J., and Pardo, T. A. 2012. "A Primer on Machine Readability for Online Documents and
978 Data." Accessed: 10th October 2019. [https://www.data.gov/developers/blog/primer-machine-readability-
979 online-documents-and-data](https://www.data.gov/developers/blog/primer-machine-readability-online-documents-and-data).

980 Hu, Z.-Z., Zhang, X.-Y., Wang, H.-W., and Kassem, M. (2016). "Improving interoperability between architectural
981 and structural design models: An industry foundation classes-based approach with web-based tools."
982 *Automation in Construction*, 66, 29-42.

983 Hwang, B.-g., Zhao, X., and Yang, K. W. (2019). "Effect of BIM on Rework in Construction Projects in
984 Singapore: Status Quo, Magnitude, Impact, and Strategies." *Journal of Construction Engineering and
985 Management*, 145(2), 04018125.

986 International Organization for Standardization (2019). "ISO/DIS 21597 Information container for data drop --
987 Exchange specification", Retrieved from <https://www.iso.org/standard/74389.html>

988 Janssen, M., van der Voort, H., and Wahyudi, A. (2017). "Factors influencing big data decision-making quality."
989 *Journal of Business Research*, 70(1), 338-345.

990 Jaradat, S., Whyte, J., and Luck, R. (2013). "Professionalism in digitally mediated project work." *Building
991 Research & Information*, 41(1), 51-59.

992 Jayawardene, V., Sadiq, S., and Indulska, M. (2015). "An Analysis of Data Quality Dimensions" *ITEE Technical
993 Report*, The University of Queensland.

994 Jordani, D. A. (2010). "BIM and FM: The Portal to Lifecycle Facility Management." *Journal of Building
995 Information Modeling*, 13-16.

996 Kassem, M., and Succar, B. (2017). "Macro BIM adoption: Comparative market analysis." *Automation in
997 Construction*, 81, 286-299.

998 Kim, K., Kim, H., Kim, W., Kim, C., Kim, J., and Yu, J. (2018). "Integration of ifc objects and facility
999 management work information using Semantic Web." *Automation in Construction*, 87, 173-187.

1000 Krijnen, T., and Beetz, J. (2017). "An IFC schema extension and binary serialization format to efficiently integrate
1001 point cloud data into building models." *Advanced Engineering Informatics*, 33, 473-490.

1002 Kwon, O. S., Park, C. S., and Lim, C. R. (2014). "A defect management system for reinforced concrete work
1003 utilizing BIM, image-matching and augmented reality." *Automation in Construction*, 46, 74-81.

1004 Lee, Y.-C., Eastman, C. M., and Solihin, W. (2018). "Logic for ensuring the data exchange integrity of building
1005 information models." *Automation in Construction*, 85, 249-262.

1006 Loshin, D. 2010. *The Practitioner's Guide to Data Quality Improvement*. Burlington, MA: Morgan Kaufmann.

1007 Mirarchi, C., and Pavan, A. (2019). "Building information models are dirty." In *Proc., 2019 European Conference*
1008 *of Computing in Construction (2019 EC³)*, Chania, Greece.

1009 Naumann, F., and Rolker, C. (2000). "Assessment Methods for Information Quality Criteria." In *Proc., Fifth*
1010 *Conference on Information Quality (IQ 2000)*, Cambridge, US.

1011 Oti, A. H., Tah, J. H. M., and Abanda, F. H. (2018). "Integration of Lessons Learned Knowledge in Building
1012 Information Modeling." *Journal of Construction Engineering and Management*, 144(9), 04018081.

1013 Pauwels, P., De Meyer, R., Van Campenhout, J., Meyer, R. D., and Campenhout, J. V. (2010). "Interoperability
1014 for the Design and Construction Industry through Semantic Web Technology." *Lecture Notes in Computer*
1015 *Science*, 6725(1), 143-158.

1016 Pauwels, P., Törmä, S., Beetz, J., Weise, M., and Liebich, T. (2015). "Linked Data in Architecture and
1017 Construction." *Automation in Construction*, 57, 175-177.

1018 Pazlar, T., and Turk, Ž. (2008). "Interoperability in practice: Geometric data exchange using the IFC standard."
1019 *Journal of Information Technology in Construction*, 13(1), 362-380.

1020 Pedro, A., Lee, D. Y., Hussain, R., and Park, C. S. "Linked Data System for Sharing Construction Safety
1021 Information." In *Proc., 34th International Symposium on Automation and Robotics in Construction (ISARC 2017)*,
1022 Taipei, Taiwan.

1023 Preidel, C., Borrmann, A., Oberender, C., and Tretheway, M. (2016). "Seamless Integration of Common Data
1024 Environment Access into BIM Authoring Applications: the BIM Integration Framework." In *Proc., 11th*
1025 *European Conference on Product and Process Modelling (ECPPM 2016)*, Limassol, Cyprus, 119.

1026 Quintero, D., Genovese, W., Kim, K., Li, M., Martins, F., Nainwal, A., Smolej, D., Tabinowski, M., and Tiwary,
1027 A. 2015. *IBM software defined environment*. IBM Redbooks.

1028 Redmond, A., Hore, A., Alshawi, M., and West, R. (2012). "Exploring how information exchanges can be
1029 enhanced through Cloud BIM." *Automation in Construction*, 24, 175-183.

1030 Sacks, R., Eastman, C., Lee, G., and Teicholz, P. 2018. *BIM handbook: a guide to building information modeling*
1031 *for owners, designers, engineers, contractors, and facility managers*. New Jersey, USA: John Wiley & Sons.

1032 Sacks, R., Ma, L., Yosef, R., Borrmann, A., Daum, S., and Kattel, U. (2017). "Semantic Enrichment for Building
1033 Information Modeling: Procedure for Compiling Inference Rules and Operators for Complex Geometry."
1034 *Journal of Computing in Civil Engineering*, 31(6), 04017062.

1035 Sebastian, R. (2011). "Changing roles of the clients, architects and contractors through BIM." *Engineering,*
1036 *Construction and Architectural Management*, 18(2), 176-187.

1037 Singh, V., Gu, N., and Wang, X. (2011). "A theoretical framework of a BIM-based multi-disciplinary
1038 collaboration platform." *Automation in Construction*, 20(2), 134-144.

1039 Solihin, W., Eastman, C., and Lee, Y.-C. (2015). "Toward robust and quantifiable automated IFC quality
1040 validation." *Advanced Engineering Informatics*, 29(3), 739-756.

1041 Taylor, J. E. (2007). "Antecedents of Successful Three-Dimensional Computer-Aided Design Implementation in
1042 Design and Construction Networks." *Journal of Construction Engineering and Management*, 133(12), 993-
1043 1002.

1044 Wang, L. (2017). "Heterogeneous Data and Big Data Analytics." *Automatic Control and Information Sciences*, 3,
1045 8-15.

1046 Wang, R. Y., and Strong, D. M. (1996). "Beyond Accuracy: What Data Quality Means to Data Consumers."
1047 *Journal of Management Information Systems*, 12(4), 5-33.

1048 Whyte, J. (2011). "Managing digital coordination of design: emerging hybrid practices in an institutionalized
1049 project setting." *Engineering Project Organization Journal*, 1(3), 159-168.

1050 Whyte, J., Stasis, A., and Lindkvist, C. (2016). "Managing change in the delivery of complex projects:
1051 Configuration management, asset information and 'big data'." *International Journal of Project Management*,
1052 34(2), 339-351.

1053 Won, J., Lee, G., Dossick, C., and Messner, J. (2013). "Where to Focus for Successful Adoption of Building
1054 Information Modeling within Organization." *Journal of Construction Engineering and Management*, 139(11),
1055 04013014.

1056 Younas, M. (2019). "Research challenges of big data." *Service Oriented Computing and Applications*, 13(2), 105-
1057 107.

1058 Zadeh, P. A., Wang, G., Cavka, H. B., Staub-French, S., and Pottinger, R. (2017). "Information Quality
1059 Assessment for Facility Management." *Advanced Engineering Informatics*, 33, 181-205.

1060 Zhang, C., and Beetz, J. (2015). "Model Checking on the Semantic Web: IFC Validation Using Modularized and
1061 Distributed Constraints." In *Proc., 32nd CIB W78 Conference 2015, 27th-29th October 2015*, Eindhoven, The
1062 Netherlands, 819-827.

1063 Zhang, J., Liu, Q., Hu, Z., Lin, J., and Yu, F. (2017). "A multi-server information-sharing environment for cross-
1064 party collaboration on a private cloud." *Automation in Construction*, 81, 180-195.

1065 Zhang, S., Pan, F., Wang, C., Sun, Y., and Wang, H. (2017). "BIM-Based Collaboration Platform for the
1066 Management of EPC Projects in Hydropower Engineering." *Journal of Construction Engineering and*
1067 *Management*, 143(12), 04017087.

1068

1069 **Fig. 1.** Visual representation showing the effect of the low level of detail in process information
1070 **Table 1.** Level of machine readability the data based on linked data principles set out by Berners-Lee
1071 (2006)
1072 **Table 2.** Codification challenges across cases
1073 **Table 3.** Information sharing between different teams
1074 **Table 4:** Construction data sets categorized based on levels of machine readability
1075 **Table 5.** Mapping between codification challenges and data quality dimensions
1076

1077

1078 **Table 1:** Level of machine readability the data based on linked data principles set out by Berners-Lee
1079 (2006)

Quality of data	Principles for publishing a machine-readable data set
1-Star	Data is available on the web
2-Star	1-star data structured in a proprietary format
3-Star	1-star data structured in a non-proprietary format
4-Star	3-star data that is published using open standards
5-Star	4-star data with links to other 4 star datasets

1080

1081

1082 **Table 2:** Codification challenges across the cases

Codification challenges	Observations	The student apartment	The metro project	The water project
Software usage	Interoperability	X	X	X
	Information loss during conversions	X		X
	Modelling technique			X
Information sharing	Unstructured information sharing	X	X	X
	Drawing and file-based sharing	X	X	X
	Document control bottleneck		X	X
	Lack of process change	X	X	X
Process information	Loss of constraints	X	X	X
	Low level of detail	X	X	X

1083

1084 **Table 3:** Information sharing between different teams

Media	Example evidence from the dataset
Common Data Environments (CDE1, CDE2 from different vendors)	<i>Formal submissions, drawing receipts, design and temporary works:</i> “if it’s formal document submission, we do it through [CDE1] [...] When I receive drawings from [design consultant], I get them through [CDE2]. And quite a lot of the designers use [CDE2] and[...] Well, we try get all the design functions, including temporary works to use [CDE2].” (Project engineer, C3I1)
Reports	<i>Spreadsheets and documents:</i> “There’s lots of reporting on the project[...] And then that gets out into various outputs, so that could be just a schedule in Excel. Lots of Excel outputs as well, huge amount of Excel outputs. And, if it’s a commercial discussion there may need to be some narrative around it, so using Microsoft Word to develop a narrative.” (Project Planner lead, C3I2)
Meetings	<p><i>Design review meetings:</i> “You could do a design meeting, review something and then say, write comment on that[...]to understand what information, they’re going to require at a particular stage. So that may consist of meetings; that might consist of face-to-face conversations, emails, etc.” (Digital engineering lead, C3I3)</p> <p><i>Buildability meetings</i> “attendance to buildability meetings and trying to get out of them what sort of temporary works may be needed to build something” (Project engineer, C3I1)</p> <p><i>Client meetings:</i> “I will be going off-site to attend meetings with the client” (Information manager, C3I4)</p> <p><i>Design for Manufacture and Assembly (DfMA) input:</i> “I’m trying to attend meetings and troubleshoot and try to help and provide technical input into the design and assisting the designer and support team” (Technical manager/DfMA coordinator, C3I5)</p>
Email	<p><i>Highly used:</i> “So, obviously we do use emails a lot.” (Project Planner lead, C3I2)</p> <p><i>A normal type of communication:</i> “Then emails, meetings, usually types of communication.” (Project engineer, C3I1)</p> <p><i>Provides remote precision:</i> “If I’m communicating over longer distances or if I think to myself, I’d better make a precise request, then it’ll be emails. We don’t use a communicator-type facility in [CDE1]”. (Information manager, C3I4)</p> <p>“Stakeholders, yes, it’s certainly meetings and emails. Most of our stakeholders don’t want to use [CDE1], because of the admin that comes with that” (Principal engineer, C3I8)</p>
Remote conversations	<p><i>Online meetings:</i> “Generally, like a Skype, conference calls, linked meetings” (Senior digital engineer, C3I6)</p> <p><i>Online communications and records of design logs, discussion points, online forums:</i> “Yes, so all the design data is held within a common data environment, which was [CDE2]. And so, I managed that area and access to that area. Then all communications were stored on SharePoint on Microsoft online, so everyone had access to registers or design logs or discussion points, almost used as an online forum where anyone could ask questions” (Senior digital engineer, C3I6)</p> <p><i>Telephone calls with design consultants:</i> “So, I’ll start with between us and design consultants: there’s meetings, emails, and phone calls. I prefer meetings and phone calls” (Principal engineer, C3I8)</p>

1085

1086

1087 **Table 4:** Construction data sets categorized based on levels of machine readability

Quality of data	Principles for publishing a machine-readable data set	Construction data sets
1-Star	Data is available on the web	Files and Models uploaded in the common data environment
2-Star	1-star data structured in a proprietary format	BIM files in proprietary formats (Revit files, Microstation files, etc.), project management information (Asta power project, Primavera P6, Microsoft project), design rationale and associated information (in Microsoft Excel), etc.
3-Star	1-star data structured in a non-proprietary format	BIM files in IFC format, CSV data etc.
4-Star	3-star data that is published using open standards	BIM files published using open standards such as ifcOWL, BOT ontology etc.
5-Star	4-star data with links to other 4 star datasets	BIM files published using open standards linked to other such files (BIM files, GIS data etc).

1088

1089

Data quality dimension	Codes and data
Accuracy	<p><i>Multiple modelling techniques:</i> “use the wrong tool to model something [...] I can’t just say there’s a slab now, that’s just a piece of geometry” (Digital engineer, C3I3); “when you try to extract 2D drawings from 3D BIM models, those drawings are not as correct and as detailed as they used to be” (Technical manager, C3I5)</p> <p><i>Document control bottlenecks:</i> “it’s no longer the most current version anymore by the time I’m reviewing it” (Project engineer, C3I1) “he keeps on updating but he hasn’t he hasn’t put it on the [CDE1].” (Technical manager, C3I5) “I just want to know where I can get my latest drawing” (Digital engineer, C3I3) “I think someone within the doc management system had obviously circumnavigated it somehow, to get the drawings out. And then when we were trying to get the said revisions for our set out, the system wouldn’t allow it because directory hadn’t been properly created.” (Technical manager, C3I5)</p>
Completeness	<p><i>Interoperability:</i> “transferring things [...], you lose data” (Technical manager, C3I5)</p> <p><i>Information loss during conversion:</i> “when you upload a PDF.” (Information manager, C3I4)</p> <p><i>Multiple modelling techniques:</i> “use the wrong tool to model something [...] I can’t just say there’s a slab now, that’s just a piece of geometry” (Digital engineer, C3I3); “when you try to extract 2D drawings from 3D BIM models, those drawings are not as correct and as detailed as they used to be” (Technical manager, C3I5)</p> <p><i>Lack of process change:</i> “We’re going to print it out, we’re going to staple it together [...] get three signatures, scan it back in, put it back into [CDE1] and submit it.”(Project engineer, C3I1) “it’s not actually speeding everything up, it’s sort of making everything a lot slower; which I find very frustrating” (Principal engineer, C3I8) “can be very confusing when we have two platforms” (Technical manager, C3I5) “Everything had to be taken out of one data environment and pushed into another. One of the issues with that is the consistency or the compliance or knowing the latest versions of information” (Digital engineer, C3I6) “As the contractor, then we have to deliver it to a completely separate, disconnected CDE [...] we’re double-handling” (Digital engineer, C3I3)</p> <p><i>Loss of constraint information:</i> “we physically need that information to know what we’re building and what the constraints in building it are.” (Project planner, C3I2). “Access chamber works will conflict with access road for pile work, piling work package has to be moved back 2 weeks.” (Progress review meeting, C3M3); <i>Low level of detail:</i> “Work package for three spans were linked to a work order. Model showed the deck for a span was completed before the pier supporting it was completed because the work package for the first span was reported as completed.” (Field notes- BIM Consultant 2, C2I5)</p>
Timeliness	<p><i>Unstructured information sharing:</i> “when I have finished everything - by the way, we have this spread sheet” (Technical manager, C3I5) <i>Document control bottlenecks:</i> “it’s no longer the most current version anymore by the time I’m reviewing it” (Project engineer, C3I1) “he keeps on updating but he hasn’t he hasn’t put it on the [CDE1].” (Technical manager, C3I5) “I just want to know where I can get my latest drawing” (Digital engineer, C3I3) “I think someone within the doc management system had obviously circumnavigated it somehow, to get the drawings out. And then when we were trying to get the said revisions for our set out, the system wouldn’t allow it because directory hadn’t been properly created.” (Technical manager, C3I5)</p> <p><i>Lack of process change:</i> “We’re going to print it out, we’re going to staple it together [...] get three signatures, scan it back in, put it back into [CDE1] and submit it.”(Project engineer, C3I1) “it’s not actually speeding everything up, it’s sort of making everything a lot slower; which I find very frustrating” (Principal engineer, C3I8) “can be very confusing when we have two</p>

	platforms” (Technical manager, C3I5) “Everything had to be taken out of one data environment and pushed into another. One of the issues with that is the consistency or the compliance or knowing the latest versions of information” (Digital engineer, C3I6) “As the contractor, then we have to deliver it to a completely separate, disconnected CDE [...] we’re double-handling” (Digital engineer, C3I3)
Consistency	<i>Document control bottlenecks:</i> “it’s no longer the most current version anymore by the time I’m reviewing it” (Project engineer, C3I1) “he keeps on updating but he hasn’t he hasn’t put it on the [CDE1].” (Technical manager, C3I5) “I just want to know where I can get my latest drawing” (Digital engineer, C3I3) “I think someone within the doc management system had obviously circumnavigated it somehow, to get the drawings out. And then when we were trying to get the said revisions for our set out, the system wouldn’t allow it because directory hadn’t been properly created.” (Technical manager, C3I5)
Accessibility	<i>Interoperability:</i> “transferring things [...], you lose data” (Technical manager, C3I5); <i>Unstructured information sharing:</i> “when I have finished everything - by the way, we have this spread sheet” (Technical manager, C3I5); <i>Drawings and file-based sharing:</i> “I use all the navigator tools that we’ve got here. But I prefer to use AutoCAD because I find it a lot easier” (Project engineer, C3I1) “It’s not the best way because we haven’t got the technology. I haven’t got a big screen” (Project engineer, C3I1) <i>Lack of process change:</i> “We’re going to print it out, we’re going to staple it together [...] get three signatures, scan it back in, put it back into [CDE1] and submit it.”(Project engineer, C3I1) “it’s not actually speeding everything up, it’s sort of making everything a lot slower; which I find very frustrating” (Principal engineer, C3I8) “can be very confusing when we have two platforms” (Technical manager, C3I5) “Everything had to be taken out of one data environment and pushed into another. One of the issues with that is the consistency or the compliance or knowing the latest versions of information” (Digital engineer, C3I6) “As the contractor, then we have to deliver it to a completely separate, disconnected CDE [...] we’re double-handling” (Digital engineer, C3I3) <i>Loss of constraint information:</i> “we physically need that information to know what we’re building and what the constraints in building it are.” (Project planner, C3I2). “Access chamber works will conflict with access road for pile work, piling work package has to be moved back 2 weeks.” (Progress review meeting, C3M3); <i>Low level of detail:</i> “Work package for three spans were linked to a work order. Model showed the deck for a span was completed before the pier supporting it was completed because the work package for the first span was reported as completed.” (Field notes- BIM Consultant 2, C2I5)
Data Provenance	<i>Information loss during conversion:</i> “when you upload a PDF.” (Information manager, C3I4); <i>Drawings and file-based sharing:</i> “I use all the navigator tools that we’ve got here. But I prefer to use AutoCAD because I find it a lot easier” (Project engineer, C3I1) “It’s not the best way because we haven’t got the technology. I haven’t got a big screen” (Project engineer, C3I1) <i>Lack of process change:</i> “We’re going to print it out, we’re going to staple it together [...] get three signatures, scan it back in, put it back into [CDE1] and submit it.”(Project engineer, C3I1) “it’s not actually speeding everything up, it’s sort of making everything a lot slower; which I find very frustrating” (Principal engineer, C3I8) “can be very confusing when we have two platforms” (Technical manager, C3I5) “Everything had to be taken out of one data environment and pushed into another. One of the issues with that is the consistency or the compliance or knowing the latest versions of information” (Digital engineer, C3I6) “As the contractor, then we have to deliver it to a completely separate, disconnected CDE [...] we’re double-handling” (Digital engineer, C3I3)

1093

1094