

1 FHIR-DHP: A Standardized Clinical Data Harmonisation Pipeline for 2 scalable AI application deployment

3 Elena Williams¹, Manuel Kienast¹, Evelyn Medawar^{1*}, Janis Reinelt¹, Alberto Merola¹, Sophie
4 Anne Ines Klopfenstein², Anne Rike Flint², Patrick Heeren², Akira-Sebastian Poncette², Felix
5 Balzer², Julian Beimes³, Paul von Büнау³, Jonas Chromik⁴, Bert Arnrich⁴, Nico Scherf⁵,
6 Sebastian Niehaus^{1,5}

7
8 1 aicura medical GmbH, Berlin, Germany

9 2 Institute of Medical Informatics, Charité – Universitätsmedizin Berlin, Berlin, Germany

10 3 idalab GmbH, Berlin, Germany

11 4 Digital Health – Connected Healthcare, Hasso Plattner Institute, University of Potsdam,
12 Potsdam, Germany

13 5 Max Planck Institute for Human Cognitive and Brain Sciences, Leipzig, Germany

14 * corresponding author: evelyn.medawar@aicura-medical.com; Bessemerstr. 22, Berlin,
15 Germany

16

17 *Background:* Increasing digitalisation in the medical domain gives rise to large amounts of
18 healthcare data which has the potential to expand clinical knowledge and transform patient
19 care if leveraged through artificial intelligence (AI). Yet, big data and AI oftentimes cannot
20 unlock their full potential at scale, owing to non-standardised data formats, lack of technical
21 and semantic data interoperability, and limited cooperation between stakeholders in the
22 healthcare system. Despite the existence of standardised data formats for the medical
23 domain, such as Fast Healthcare Interoperability Resources (FHIR), their prevalence and
24 usability for AI remains limited.

25 *Objective:* We developed a data harmonisation pipeline (DHP) for clinical data sets relying on
26 the common FHIR data standard.

27 *Methods:* We validated the performance and usability of our FHIR-DHP with data from the
28 MIMIC IV database including > 40,000 patients admitted to an intensive care unit.

29 *Results:* We present the FHIR-DHP workflow in respect of transformation of “raw” hospital
30 records into a harmonised, AI-friendly data representation. The pipeline consists of five key
31 preprocessing steps: querying of data from hospital database, FHIR mapping, syntactic
32 validation, transfer of harmonised data into the patient-model database and export of data
33 in an AI-friendly format for further medical applications. A detailed example of FHIR-DHP
34 execution was presented for clinical diagnoses records.

35 *Conclusions:* Our approach enables scalable and needs-driven data modelling of large and
36 heterogenous clinical data sets. The FHIR-DHP is a pivotal step towards increasing
37 cooperation, interoperability and quality of patient care in the clinical routine and for
38 medical research.

39

40 **Keywords:** Data interoperability, FHIR, data standardisation pipeline, MIMIC IV

41

42 **Competing Interests:** The authors declare no competing interests.

43 **Introduction**

44

45 The increasing digitalisation of healthcare creates vast amounts of clinical data that are
46 collected and stored in an Electronic Health Record (EHR). Patient information from all
47 medical domains is captured in diverse sets of data recorded in standalone systems. With
48 the prevalent use of EHRs in healthcare organisations, there is abundant opportunity for
49 additional application of EHR data in clinical and translational research. For instance, such
50 data can be used to develop artificial intelligence (AI) algorithms which have the potential to
51 transform patient care and medical research. Resource intensive and inefficient clinical
52 workflows could be optimised by the analysis of historical data with AI applications (1,2). In
53 particular, the time-consuming and high-priced process of identifying and enrolling the right
54 patients into a clinical trial manually can be reduced significantly by automation (3,4).
55 However, the exchange of medical data remains limited due to the lack of data
56 interoperability between healthcare providers, owing to outdated IT infrastructure,
57 inconsistencies in data formats, poor data quality, inadequate data exchange solutions and
58 data silos (5,6). To achieve data interoperability, the following steps must be incorporated: i)
59 integration of isolated data silos, ii) safe exchange of data and iii) effective use of the
60 available data (7). Each of these operations includes database schema matching (8) and
61 schema mapping (9), which allow translation of the relationships between the source
62 database and the target data standard.

63 Employing a harmonised data format will facilitate the exchange of medical data, enabling
64 wide-ranging data-driven collaborations within the private and public healthcare sectors.
65 Data interoperability requires EHR data to be structured in a common format and in
66 standardised terminologies. Standardisation is often performed by adopting the Health Level
67 7 (HL7) Fast Healthcare Interoperability Resources (FHIR) model (10), which is supported by

68 numerous healthcare institutions and vendors of clinical information systems (11). FHIR is an
69 international industry standard with the benefit of integrating diverse sets of data in well-
70 defined exchangeable segments of information, which are known as FHIR resources.
71 Therefore, FHIR facilitates interoperability between healthcare organisations and allows
72 third-party developers to provide medical applications which can be easily integrated into
73 existing systems. FHIR enables the harmonization of data and thus allows standardized data
74 processing and also the rollout of AI applications across different clinics and hospitals
75 regardless of which information system they use. Therefore, FHIR forms an important
76 component for the scalable development and deployment of AI in clinics and hospitals.

77 However, to apply AI, the input data needs to be adapted to the AI algorithms. The
78 conventional AI frameworks such as Tensorflow (14) and Pytorch (15) require data to take a
79 tensor form which is a vector or matrix of n-dimensions that represents various types of data
80 (e.g., tabular, time series, image, text). FHIR facilitates the application of AI in medical
81 domain as it provides needed interoperability for a standardised access of EHR data. FHIR
82 format's multi-layered nested structure requires case-specific data pre-processing to use it
83 for AI algorithms. Depending on the AI application and the chosen data source, a custom
84 data preprocessing pipeline needs to be designed, which leads to diminished AI scalability.
85 Up to the present time, a number of studies have attempted to solve this problem. Prior
86 research addressed this problem in different forms, but focuses on individual use cases and
87 thus constrains the basic idea of FHIR to be independent of the use case.. There have been a
88 few attempts to flatten the hierarchical FHIR structure and transform it into NDJSON-based
89 data format (16) or tabular format saved in CSV files (17). Such formats are more AI-friendly
90 as they represent the data in a more accessible and standardised form for an application of
91 common AI frameworks. Nonetheless, the NDJSON-based FHIR data transformation

92 approach (16) does not provide data selection criteria and filtering capabilities. The
93 approach presented in (17) requires expert knowledge of *FHIRPath* query language.

94 In this paper, we address the challenge of data interoperability in the healthcare sector by
95 proposing a FHIR Data Harmonisation Pipeline (DHP) that provides EHR data in an AI-friendly
96 format. The newly developed FHIR-DHP represents a data workflow solution that includes
97 the aforementioned operations such as data exchange, mapping, and export. Data privacy is
98 a delicate topic in healthcare and is of great ethical concern (18). Given the degree of
99 automation, such pipeline should allow preprocessing of unseen data in an isolated hospital
100 environment, which makes the harmonisation privacy-preserving. In this setting, direct
101 access to the sensitive data would not be required to run the standardisation pipeline. FHIR-
102 based data preprocessing pipelines have already been implemented in different contexts: as
103 electronic data capture (12), as a natural language processing tool (13) and as a
104 standardisation protocol based on the Resource Description Framework (RDF) (6). Despite
105 their immense benefit of processing EHR data, existing approaches are limited to specific use
106 cases or require considerable data preparation to perform standardisation. Moreover, their
107 final output is not easily accessible by common data preprocessing tools and thus hinders
108 the application of AI.

109

110 **Methods**

111

112 *FHIR-DHP Development*

113 In our work, we propose a generic solution to harmonise hospital EHR data. The FHIR-DHP
114 was designed based on the Extract-Transform-Load (ETL) framework (19) in which the data is
115 pulled out (i.e. queried) from diverse sources, processed into the desired format and loaded
116 into a data warehouse, namely the "patient-model DB". As the hospital database (DB)

117 contains highly sensitive patient data, it is located behind the hospital's security
118 infrastructure and is completely isolated from outside access. Therefore, an edge-
119 computation solution was designed, bringing the FHIR-DHP into the hospital's own
120 infrastructure. The edge-computation solution represents a set of frameworks which
121 perform data querying, preprocessing, storage and export. In this setting, direct access to
122 the sensitive data is not required to run the standardisation pipeline. The queries to the data
123 are defined beforehand based on the database documentation.

124 To bring the data into a harmonised form we used Fast Healthcare Interoperability
125 Resources (FHIR) data model which is applied by mapping the relationships between the
126 source database and the desired data standard. The FHIR standard is straightforward to
127 implement because it provides a choice of JavaScript Object Notation (JSON), Extensible
128 Markup Language (XML), or Resource Description Format (RDF) for data representation. The
129 mapping pipeline was developed in Python programming language to translate queried
130 hospital data into matching FHIR concepts and save the resulting resources in JSON format.
131 The conversion to FHIR was designed to only support a core standard of the FHIR format to
132 allow generic data preprocessing.

133 Syntactic validation of FHIR resources is necessary in the remote data standardisation
134 scenario to prevent errors. For instance, conversion of data types can sometimes lead to
135 wrong values, especially with date features. Automatic syntactic validation allows logging of
136 occurred errors and improvement of standardisation pipeline when working with unseen
137 data. After the mapped data is validated, FHIR resources should be sent to the database for
138 storage to allow fast and easy retrieval of preprocessed data for AI applications.

139 In the final stage of data export, we designed the output that provides the benefits of the
140 original FHIR format with a high level of clinical detail, yet which is also easily accessible for
141 computational tools. Moreover, we wanted to restructure the data representation in a way
142 which supports effortless data selection and filtering capabilities and which would not require
143 knowledge of *FHIRPath* query language. Consequently, such output format would enable
144 smooth conversion of data into a “tensor” format required by conventional AI frameworks.

145 *FHIR-DHP Validation*

146 To demonstrate and evaluate how the FHIR-DHP works, we used the openly available
147 Medical Information Mart for Intensive Care IV (MIMIC IV) database (20). MIMIC IV includes
148 patient data from over 40,000 individuals admitted to intensive care units at a tertiary
149 academic medical center in Boston, MA. We selected a wide range of tables from MIMIC IV
150 which cover most of the events occurring during the hospital stay as well as core patient
151 details, information about admissions and hospital transfers (further referred as core tables).
152 The event tables include laboratory results, diagnoses, prescriptions and other details as
153 shown in **Table 1**. MIMIC IV includes so-called reference tables containing matching
154 dictionaries with medical terms which are used in the hospital records.

155

156 **Table 1.** The table lists selected core and event MIMIC IV tables as well as the reference dictionary tables
157 which were merged together with core and event tables for FHIR mapping.

| Selected core and event MIMIC IV tables | Selected MIMIC IV reference tables |
|--|---------------------------------------|
| Patient | |
| Admissions | |
| Transfers | |
| Chartevents | d_items |
| Labevents | d_labitems |
| Procedureevents | d_items |
| Prescriptions | |
| Inputevents | d_items |
| Microbiologyevents | |
| Outputevents | d_items |
| Procedures_icd | d_icd_procedures |
| Diagnoses_icd | d_icd_diagnoses |

158

159 The selected tables were mapped to FHIR standard. Automatic semantic validation is
160 unfeasible, so two of the authors manually validated the mapping semantics independently
161 of each other. There are many tools which perform automatic syntactic validation, such as
162 the Python-based package `fhir.resources` used herein (21). To evaluate the exporting of data
163 from the patient-model DB, we retrieved diagnoses records.

164

165 **Results**

166 *FHIR-DHP Architecture*

167

168 The approach presented herein represents a scalable protocol for harmonising hospital EHR
169 datasets based on five stages from data query to data export in a standardised format.

170

171 *1. Querying data from the hospital database*

172 To connect the FHIR-DHP pipeline to the hospital DB, a communication server is employed.
173 This server runs all necessary queries to retrieve the patient data. The query execution can
174 be run at regular intervals as well as in batches of patients, so as not to overload the data
175 pipeline. Furthermore, the queries pre-structure the data according to their semantic
176 relations before proceeding to data mapping.

177 *2. Mapping data to FHIR*

178 FHIR allows describing data formats and elements which are recorded as "resources" and an
179 application of a programming interface (API) for exchanging EHRs. To perform the mappings,
180 semantics of features from the source database and FHIR concepts are explored as well as
181 relationships between the data tables. Consequently, the mappings between the database

182 tables and FHIR resources are defined. Features where a matching FHIR concept is not found
183 are excluded. The resulting FHIR resources are then saved in JSON format.

184

185 *3. Syntactic validation of FHIR mappings*

186 During validation, mapped data is ensured to have the correct data types as well as the
187 syntactic format where the hierarchy is maintained and entries follow FHIR standard
188 specifications. All mappings are validated first during the development stage to identify
189 structural errors and data type inconsistencies. A validation algorithm is incorporated into
190 the pipeline to confirm the correctness of transformed data in the remote data
191 standardisation scenario.

192

193 *4. Transferring FHIR resources to patient-model DB*

194 The database of choice for the patient-model is Postgres (22) which is an open-source
195 relational database management system (RDBMS) featuring SQL compliance and storage of
196 JSON documents. PostgreSQL allows handling both small and large workloads. The database
197 for FHIR resources is used to harmonise the locally available data only once to allow further
198 application of various medical AI-based solutions. The data is stored according to FHIR
199 resource type where each resource is saved in a separate JSON structure.

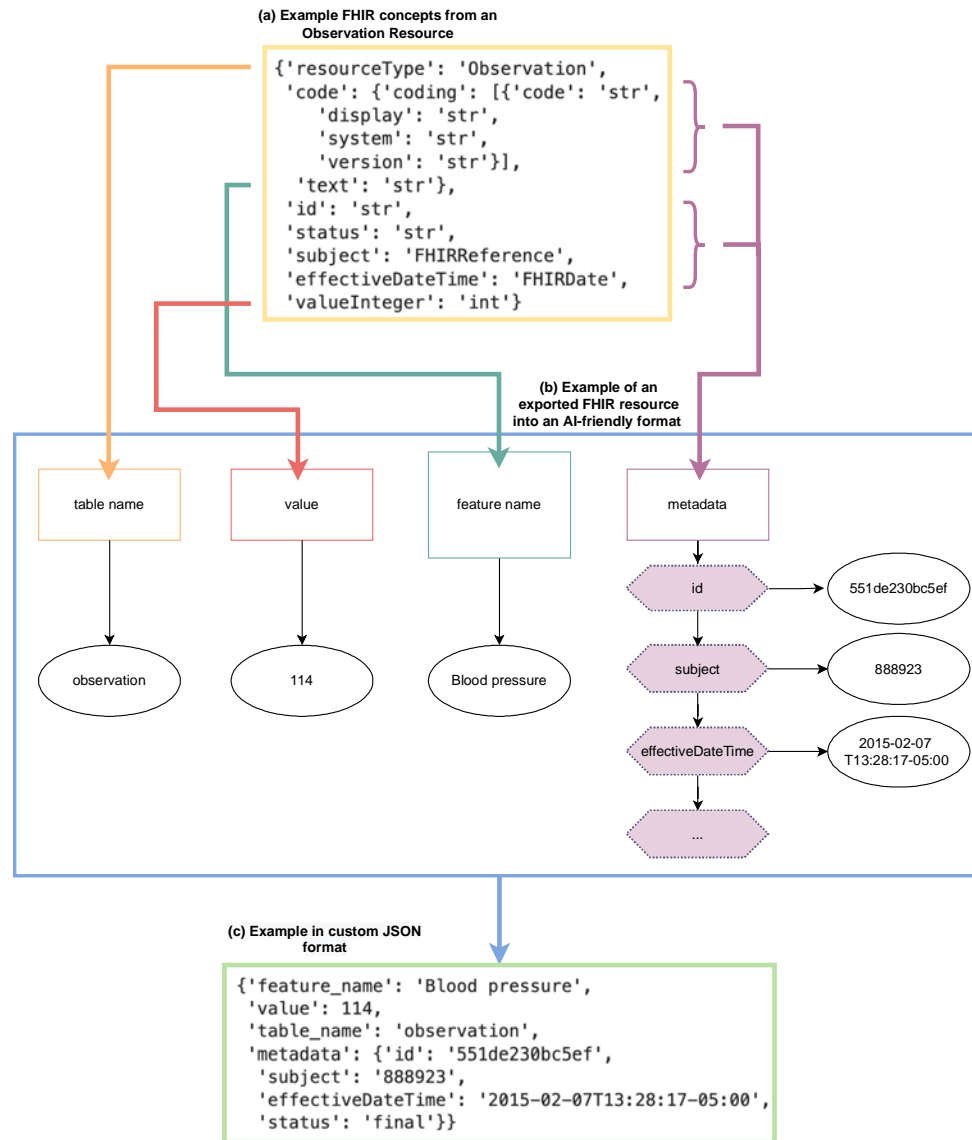
200 *5. Exporting data into custom JSON format*

201 To export the data from the patient-model DB, the selection is performed by outlining the
202 tables and features of interest in a configuration file which is then used to determine which
203 data is queried from the patient-model DB. Following that, the data is exported into the
204 custom JSON file adhering to specific formatting rules in respect of its key-value structure. To
205 create a custom JSON structure, *FHIRPath* queries were written to retrieve all elements from
206 FHIR resources. Such transformation flattens the hierarchical structure of FHIR resources and

207 makes the data more accessible for common data preprocessing tools. The final flattened
208 output does not require expert knowledge of *FHIRPath* query language and supports
209 effortless data selection and filtering. The resulting file allows uncomplicated conversion of
210 data into a “tensor” format required by conventional AI frameworks and fast data selection
211 based on four keys: feature_name, table_name, value and metadata.

212 In **Figure 1**, we demonstrate how the FHIR-DHP recodes nested FHIR syntax to more
213 accessible features in an AI-friendly format. Example FHIR concepts from an Observation
214 resource are given in **Figure 1a** where the code’s entity “text” defines the record or
215 measurement label. The entity “text” is often duplicated in the item “display”. However,
216 depending on the coding system this “display” item can change, whereas “text” always stays
217 the same and is therefore used as a feature name. The information from the FHIR resource is
218 grouped into four concept-keys such as feature name (ex. “Blood pressure”), value (ex.
219 “114”), table name (ex. “observation”) and metadata (**Figure 1b**). For a given FHIR resource
220 type, the metadata may include concepts such as dates, references, coding system details,
221 resource ID amongst other things. As an output, feature names together with a
222 corresponding value and available metadata are provided in a custom JSON structure (**Figure**
223 **1c**). The defined format allows uncomplicated data selection and aggregation based on
224 resource type (ex. “table_name”), feature name and value. Additional information in a
225 standardised format can be easily accessed from the metadata key and allows further data
226 manipulation.

227



228
229
230
231

Figure 1. Conceptual overview for an exemplary FHIR structure and hospital record which are transformed from FHIR standard to an AI-friendly format.

232 FHIR-DHP Validation

233 The MIMIC IV data was queried accordingly to the defined FHIR mappings. The core and
234 event MIMIC IV tables were merged with reference tables to contain complete description
235 of the hospital records. As a result, the data was grouped and restructured into the
236 information blocks required in FHIR standard. Manual independent validation of the
237 mapping semantics resulted in slight discrepancies which were subsequently resolved to

238 adhere closely to the FHIR standard. The automatic syntactic validation allowed prompt
239 verification of standardisation operations.

240 **Table 2** shows to which FHIR resources the MIMIC IV tables were mapped. The largest
241 proportion of tables (4 out of 12 tables) were mapped to the *Observation* FHIR resource type
242 which included lab, microbiology, output and charted events collected throughout the
243 patient stay. The information on admissions and transfers was translated into the *Encounter*
244 FHIR resource (2 out of 12 tables). Procedure events and ICD codes (2 out of 12 tables) were
245 stored in the *Procedure* FHIR resource. Given that the prescriptions table contains
246 medication requests (1 out of 12 tables) and inpuvents table holds records of medication
247 administration (1 out of 12 tables), these tables were mapped to corresponding FHIR
248 resource types. Finally, the *Condition* FHIR resource was used to map the table with patients'
249 diagnoses details (1 out of 12 tables).

250 Table 2. Overview of mappings performed on the selected MIMIC DB tables to FHIR resource types.

| MIMIC IV DB | FHIR Resource Type |
|--------------------|--------------------------|
| Patients | Patient |
| Admissions | Encounter |
| Transfers | Encounter |
| Chartevents | Observation |
| Labevents | Observation |
| Procedureevents | Procedure |
| Prescriptions | MedicationRequest |
| Inputevents | MedicationAdministration |
| Microbiologyevents | Observation |
| Outputevents | Observation |
| Procedure_icd | Procedure |
| Diagnoses_icd | Condition |

251

252 In **Table 3**, we demonstrate how the mapping of the MIMIC IV “diagnoses_icd” table to
253 *Condition* FHIR resource was conducted. Multiple columns of the “diagnoses_icd” table such
254 as “icd_code”, “icd_version” and “long_title” were mapped to FHIR “condition.code”
255 concept, which has a nested structure and provides keys to store the exact ICD code, version

256 of the coding system and the code title. The full diagnosis title was mapped both to the
257 “display” and “text” entities.

258 Table 3. Mapping of “diagnoses_icd” table to Condition FHIR resource.

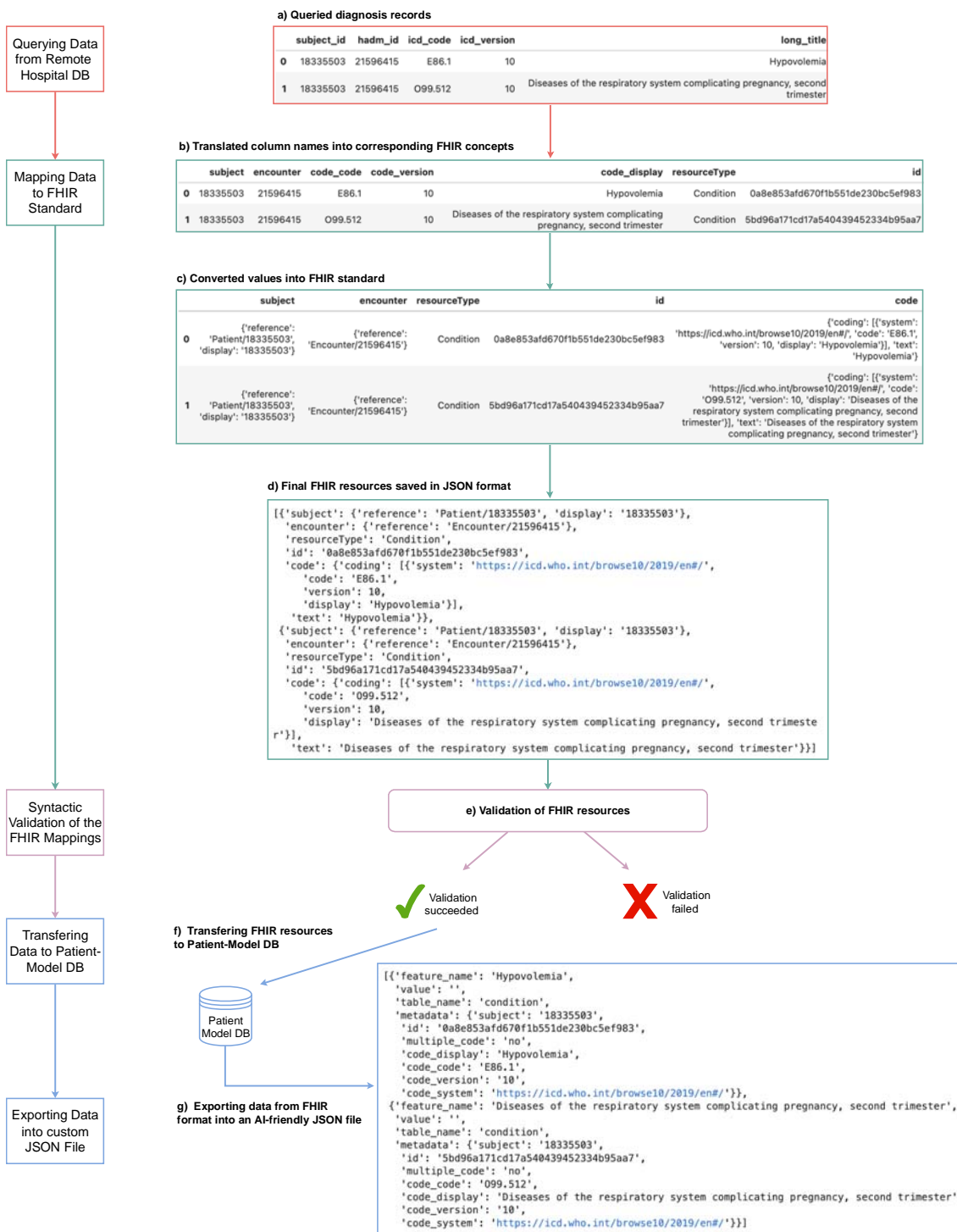
| MIMIC format | FHIR resource format |
|---------------------------------|-----------------------------|
| mimic.diagnoses_icd.subject_id | fhir.condition.subject |
| mimic.diagnoses_icd.hadm_id | fhir.condition.encounter |
| mimic.diagnoses_icd.icd_code | fhir.condition.code_code |
| mimic.diagnoses_icd.icd_version | fhir.condition.code_version |
| mimic.diagnoses_icd.long_title | fhir.condition.code_display |
| mimic.diagnoses_icd.long_title | fhir.condition.code_text |

259

260 **Figure 2** shows an example of how queried diagnoses records are harmonised to an AI-
261 friendly format. The standardisation follows the described FHIR-DHP stages. At first, the raw
262 data from tables “diagnoses_icd” and “d_icd_diagnoses” is queried (**Figure 2a**) and merged
263 accordingly to the defined FHIR mappings. Then the features are renamed as defined in
264 **Table 3** for FHIR Condition resource and required entities such as “resourceType” and “id”
265 are created (**Figure 2b**). Finally, the values are placed into a nested FHIR structure (**Figure**
266 **2c**), and subsequently the data is transformed into JSON format (**Figure 2d**), which can be
267 automatically validated (**Figure 2e**) and saved in the patient-model DB. When the resource is
268 not approved in terms of its syntactic quality, e.g. data type, nested structure or cardinality,
269 an error is raised which prevents further saving of this resource in the patient-model DB
270 (**Figure 2e**). Otherwise, the resource is transferred into a storage (**Figure 2f**) and the
271 requested data is exported in a custom AI-friendly JSON format (**Figure 2g**).

272 We provide an example of further two-step transformation of harmonised example
273 diagnoses data to a “tensor” format in Supplementary Material, chapter A.

274



275
276
277
278
279
280
281
282
283
284

Figure 2. Flow chart showing an example diagnoses data being processed through the five stages in FHIR-DHP. The first stage includes querying of the diagnoses records (a), at the second stage the data is mapped to FHIR standard (b-c), and the third stage carries out the syntactic resource validation. If the FHIR resource is successfully validated, it is being transferred into the patient-model DB (f) and then exported in a custom AI-friendly JSON format (g).

285 **Discussion**

286

287 Harmonisation of EHR data is a crucial step towards increasing cooperation, interoperability
288 and quality of patient care in the clinical routine and medical research. To drive
289 harmonisation of medical data forward, we developed the FHIR-DHP and evaluated it on key
290 MIMIC IV tables. A detailed example of data standardisation was presented for clinical
291 diagnoses records from the MIMIC IV database. The FHIR-DHP allows querying of health data
292 in an isolated environment employing an edge-computation solution and a communication
293 server which retrieve patient data and pre-structure it for further mapping to the FHIR
294 standard. A validation step ensures syntactic compliance and initiates transfer of formatted
295 data to the patient-model DB. The data export provides FHIR resources in a custom JSON file
296 format.

297 Owing to the FHIR format's multi-layered nested structure, its accessibility for AI algorithms
298 is low as it requires transformation into a format compatible with common data
299 preprocessing tools. Up to the present time, a number of studies have attempted to solve
300 this problem. However, the final output of these studies has not supported data selection
301 criteria and filtering capabilities (16) and requires expert knowledge of *FHIRPath* query
302 language (17). Here, we introduce a custom JSON format which represents a higher level of
303 abstraction to support easier data selection based on four keys: `feature_name`, `table_name`,
304 `value` and `metadata`. Moreover, the newly developed JSON structure fits the expected data
305 format of common data preprocessing frameworks, which are designed to work efficiently
306 with tabular data. As a result, the output presented facilitates generic and fast deployment of
307 AI and patient cohort identification algorithms.

308 In comparison to (12,13), the details of FHIR-DHP execution inside the hospital environment
309 to protect data privacy are discussed. This step, though crucial, is often omitted and left out
310 of the published standardisation protocols. The edge-computation solution sets up the FHIR-
311 DHP in a privacy-preserving way where preprocessing of the patient-related data is
312 performed inside the hospital and is completely isolated from outside access. So-called
313 federated learning (FL) framework (23) can be integrated into FHIR-DHP workflow to run
314 algorithms locally using the data on the on-premise component in the respective hospitals
315 and to merge model parameters centrally in the cloud without any patient data leaving the
316 hospital. The FL framework requires data to be in a consistent format across various hospital
317 systems. The developed pipeline achieves such a format and enables scaling of AI
318 applications. Furthermore, given the degree of automation, the setup of the pipeline
319 facilitates preprocessing of unseen data in an isolated hospital environment, which makes
320 the harmonisation privacy-preserving.

321 To the date of publication of this paper, there are only two studies attempting to perform
322 mapping of MIMIC IV database (24,25). In (24), the mapping was performed on fewer tables
323 than our approach (8 versus 12 tables). The FHIR mappings from (25) have been recently
324 released and were not yet widely validated. Similarly to (12,13,24), FHIR-DHP includes
325 verification of the performed FHIR mapping which is essential to ensure validity of data
326 transformation. An automated syntactic verification of translated to FHIR data is crucial to
327 adhere to FHIR version updates. Moreover, in comparison to (12,13,24), FHIR-DHP
328 represents a generic approach to standardise EHR data and can be applied to various
329 hospital database systems.

330 The FHIR-DHP allows integration into the hospital data management system which facilitates
331 the development and application of advanced AI and patient cohort identification algorithms

332 without compromising on data privacy protection laws. With the introduction of the FHIR-
333 DHP into the hospital environment, a number of patient stay parameters can be potentially
334 optimised using AI-based algorithms. For example, the length of stay as well as mortality
335 could be reduced (26) and patients suitable for trial treatment could be automatically and
336 efficiently identified (27). In consequence the financial impact on medical providers in
337 respect of personnel time and resources would decrease considerably. The FHIR-DHP aims to
338 bring healthcare closer to digital transformation and thus towards Healthcare 4.0 (28) by
339 making EHR data usable “from bedside-to-bench”. By inverting the idea of translational
340 research, in contrast to “from bench-to-bedside”, accessing the full potential of medical big
341 data with AI will further inform and advance basic research.

342 There are several limitations that we would like to emphasise. FHIR-DHP only works with a
343 core standard of the FHIR format. Those core FHIR resource types have a bounded set of
344 concepts which presents a constraint to mapping accuracy. Although the standard resources
345 can be expanded using profiling technique or FHIR extensions, the use of those would make
346 the FHIR-DHP less generic. Hence, we implemented the mapping using only the standard
347 FHIR resources and omitted some of the MIMIC IV data features which did not have a
348 matching concept in FHIR. Additionally, the FHIR mapping step is subject to the extent of the
349 detail of the database documentation used to infer semantic and syntactic properties of the
350 data. A solution for an automatic concept recognition can potentially solve this problem. The
351 existing approach in (6) is limited to a small number of FHIR resources and requires an
352 extensive data preparation. Further experiments in this direction could alleviate the concept
353 matching problem and the requirement for a detailed database description. Moreover, the
354 validation and robustness of FHIR-DHP needs to be tested on other EHR datasets to evaluate

355 its generic setup. In addition, to validate the FHIR-DHP compatibility with machine learning
356 pipelines, further experiments are needed.

357 The proposed FHIR-DHP pipeline highlights the therein featured essential data
358 standardisation stages and holds the potential to becoming an interoperable harmonisation
359 system with an AI-friendly data format. FHIR-DHP enables interoperability and cooperation
360 between clinical institutions, rapid patient cohort identification for clinical trials and unlocks
361 the potential of big medical data.

362 **Conclusions**

363

364 We provide a comprehensive approach to transforming unstandardised EHR data into a
365 harmonised multi-layered nested FHIR format and then to a more readable, more efficient
366 AI-friendly JSON structure. We developed a five-stage data harmonisation pipeline, which
367 includes validation checks. The AI-friendly format of patient data allows generic and fast
368 integration of both AI and patient cohort identification algorithms. Harmonised and
369 standardised health care data is of great value to advancing efficiency in big data processing,
370 cooperation and multi-center data exchange in the clinical sector, in order to boost medical
371 research, patient care and clinical trial cohort identification. The next steps would include
372 validating our approach in a hospital environment and applying privacy-preserving FL
373 framework to make use of advanced AI deployment.

374

375 **Acknowledgements**

376 Not applicable.

377

378 **List of abbreviations**

379

| Abbreviation | Full name |
|--------------|-----------------------------|
| DHP | Data Harmonisation Pipeline |

| | |
|-------|---|
| EHR | Electronic Health Record |
| FHIR | Fast Healthcare Interoperability Resources |
| JSON | JavaScript Object Notation |
| MIMIC | Medical Information Mart for Intensive Care |
| RDF | Resource Description Format |
| XML | Extensible Markup Language |

380

381

382

383 **Declarations**

384

385 *Availability of data and materials*

386 MIMIC IV database which was used in this study is openly available to credentialed users
387 who sign “Data Use Agreement” at PhysioNet website (20). The code is not publicly available
388 due to privacy but a demo is available from the corresponding author on request.

389

390 *Conflict of interest*

391 The authors declare that they have no competing interests.

392

393 *Funding*

394 This work was partially funded by the German Federal Ministry of Education and Research
395 under Grant 16SV8559.

396

397 *Authors' contributions*

398 Study conception: EW, SN, MK, JR, AM

399 Data analysis: EW, MK

400 Figures: EW, SN, EM

401 Methods: EW, MK, AM, SN

402 Writing: EW, EM, JR, SAIK

403 Revising: BA, JB, PVB, JC, ARF, ASP, NS

404

405

406

407 References

408

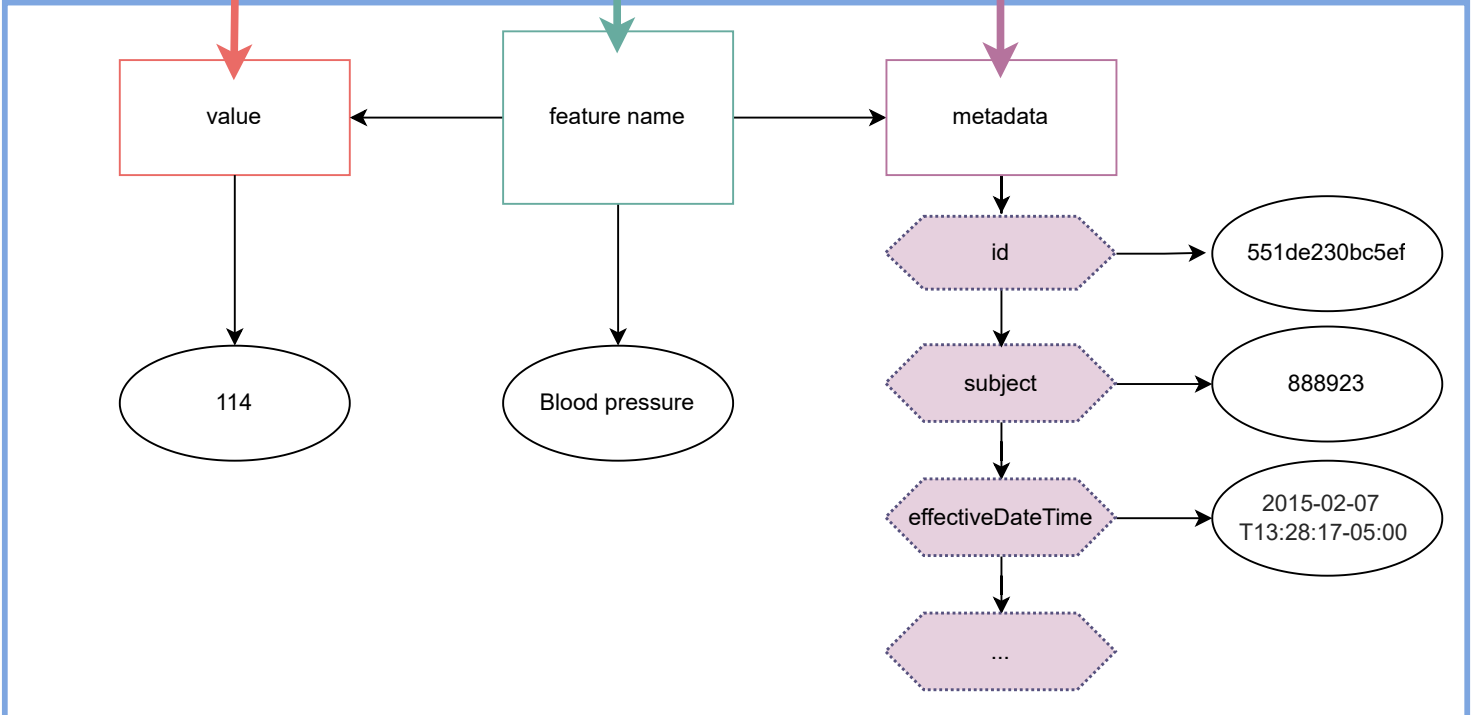
- 409 1. Au-Yeung WTM, Sahani AK, Isselbacher EM, Armoundas AA. Reduction of false
410 alarms in the intensive care unit using an optimized machine learning based
411 approach. NPJ Digit Med [Internet]. 2019;2:86. Available from:
412 <https://europepmc.org/articles/PMC6728371>
- 413 2. Desautels T, Calvert J, Hoffman J, Jay M, Kerem Y, Shieh L, et al. Prediction of
414 Sepsis in the Intensive Care Unit With Minimal Electronic Health Record Data: A
415 Machine Learning Approach. JMIR Medical Informatics. 2016 Jun;4.
- 416 3. Maier C, Kapsner L, Mate S, Prokosch HU, Kraus S. Patient Cohort Identification
417 on Time Series Data Using the OMOP Common Data Model. Applied Clinical
418 Informatics. 2021 Jun;12.
- 419 4. Ni Y, Wright J, Perentesis J, Lingren T, Deleger L, Kaiser M, et al. Increasing the
420 efficiency of trial-patient matching: automated clinical trial eligibility Pre-
421 screening for pediatric oncology patients. BMC Medical Informatics and
422 Decision Making [Internet]. 2015;15(1):28. Available from:
423 <https://doi.org/10.1186/s12911-015-0149-3>
- 424 5. de Mello BH, Rigo SJ, da Costa CA, da Rosa Righi R, Donida B, Bez MR, et al.
425 Semantic interoperability in health records standards: a systematic literature
426 review. Health and Technology [Internet]. 2022;12(2):255–72. Available from:
427 <https://doi.org/10.1007/s12553-022-00639-w>
- 428 6. Kiourtis A, Mavrogiorgou A, Menychtas A, Maglogiannis I, Kyriazis D.
429 Structurally Mapping Healthcare Data to HL7 FHIR through Ontology Alignment.
430 Journal of Medical Systems. 2019 Jun;43.
- 431 7. Pagano P, Candela L, Castelli D. Data Interoperability. Data Sci J.
432 2013;12:GRD119–25.
- 433 8. Rahm E, Bernstein P. A Survey of Approaches to Automatic Schema Matching.
434 VLDB J. 2001 Jun;10:334–50.
- 435 9. Kolaitis P. Schema mappings, data exchange, and metadata management. In:
436 Proceedings of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of
437 Database Systems. 2005. p. 61–75.
- 438 10. HL7 FHIR. <https://www.hl7.org/fhir/>.
- 439 11. Vorisek C, Lehne M, Klopfenstein S, Bartschke A, Haese T, Thun S. Fast
440 Healthcare Interoperability Resources (FHIR) for Interoperability in Health
441 Research: A Systematic Review (Preprint). JMIR Medical Informatics. 2021 Jul;
- 442 12. Zong N, Wen A, Stone DJ, Sharma DK, Wang C, Yu Y, et al. Developing an FHIR-
443 Based Computational Pipeline for Automatic Population of Case Report Forms
444 for Colorectal Cancer Clinical Trials Using Electronic Health Records. JCO Clinical
445 Cancer Informatics. 2020;4.
- 446 13. Hong N, Wen A, Shen F, Sohn S, Wang C, Liu H, et al. Developing a scalable
447 FHIR-based clinical data normalization pipeline for standardizing and
448 integrating unstructured and structured electronic health record data. JAMIA
449 Open. 2019 Jun;2.
- 450 14. Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen,
451 Craig Citro, et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous
452 Systems [Internet]. 2015. Available from: <https://www.tensorflow.org/>
- 453 15. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An
454 Imperative Style, High-Performance Deep Learning Library. In: Advances in

- 455 Neural Information Processing Systems 32 [Internet]. Curran Associates, Inc.;
456 2019. p. 8024–35. Available from: [http://papers.neurips.cc/paper/9015-](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
457 [pytorch-an-imperative-style-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
- 458 16. Liu D, Sahu R, Ignatov V, Gottlieb D, Mandl K. High Performance Computing on
459 Flat FHIR Files Created with the New SMART/HL7 Bulk Data Access Standard.
460 AMIA Annu Symp Proc. 2020 Mar 4;2019:592–6.
- 461 17. Oehm J, Storck M, Fechner M, Brix T, Yildirim K, Dugas M. FhirExtinguisher: A
462 FHIR Resource Flattening Tool Using FHIRPath. In: Studies in health technology
463 and informatics. 2021.
- 464 18. Mittelstadt B, Floridi L. The Ethics of Big Data: Current and Foreseeable Issues in
465 Biomedical Contexts. *Sci Eng Ethics*. 2015 Jul;
- 466 19. Denney M, Long D, Armistead M, Anderson J, Conway B. Validating the Extract,
467 Transform, Load Process Used to Populate a Large Clinical Research Database:
468 *International Journal of Medical Informatics*. 2016 Jun;94.
- 469 20. Johnson A, Bulgarelli L, Pollard T, Horng S, Celi LA, Roger M. MIMIC-IV (version
470 2.0). *PhysioNet*. 2022.
- 471 21. PostgreSQL, PostgreSQL Global Development Group.
472 <https://www.postgresql.org>. Accessed 15 June 2022.
- 473 22. Islam N. FHIR® Resources. <https://github.com/nazrulworld/fhir.resources>.
474 Accessed 20 May 2022.
- 475 23. Konecný J, McMahan HB, Yu FX, Richtárik P, Suresh AT, Bacon D. Federated
476 Learning: Strategies for Improving Communication Efficiency. *ArXiv*.
477 2016;abs/1610.05492.
- 478 24. Ulrich H, Behrend P, Wiedekopf J, Drenkhahn C, Kock-Schoppenhauer AK,
479 Ingenerf J. Hands on the Medical Informatics Initiative Core Data Set — Lessons
480 Learned from Converting the MIMIC-IV. In: Studies in health technology and
481 informatics. 2021.
- 482 25. Bennett A, Wiedekopf J, Ulrich H, Johnson A. MIMIC-IV Clinical Database Demo
483 on FHIR (version 2.0). . *PhysioNet*. 2022.
- 484 26. Shimabukuro DW, Barton CW, Feldman MD, Mataraso SJ, Das R. Effect of a
485 machine learning-based severe sepsis prediction algorithm on patient survival
486 and hospital length of stay: a randomised clinical trial. *BMJ Open Respiratory*
487 *Research* [Internet]. 2017;4(1). Available from:
488 <https://bmjopenrespres.bmj.com/content/4/1/e000234>
- 489 27. Sarmiento RF, Deroncourt F. Improving Patient Cohort Identification Using
490 Natural Language Processing. In: Data MITC, editor. *Secondary Analysis of*
491 *Electronic Health Records* [Internet]. Cham: Springer International Publishing;
492 2016. p. 405–17. Available from: [https://doi.org/10.1007/978-3-319-43742-](https://doi.org/10.1007/978-3-319-43742-2_28)
493 [2_28](https://doi.org/10.1007/978-3-319-43742-2_28)
- 494 28. Li J, Carayon P. Health Care 4.0: A vision for smart and connected health care.
495 *IISE Transactions on Healthcare Systems Engineering* [Internet].
496 2021;11(3):171–80. Available from:
497 <https://doi.org/10.1080/24725579.2021.1884627>
- 498 29. McKinney W. Data Structures for Statistical Computing in Python. In 2010. p.
499 56–61.
- 500

(a) Example FHIR concepts from an Observation Resource

```
{'code': {'coding': [{'code': 'str',  
  'display': 'str',  
  'system': 'str',  
  'version': 'str'}],  
  'text': 'str'},  
  'id': 'str',  
  'status': 'str',  
  'subject': 'FHIRReference',  
  'effectiveDateTime': 'FHIRDate',  
  'valueInteger': 'int'}
```

(b) Example of an exported FHIR resource into an AI-friendly format



(c) Example in custom AI-friendly JSON format

```
{'feature_name': 'Blood pressure',  
  'value': 114,  
  'table_name': 'observation',  
  'metadata': {'id': '551de230bc5ef',  
    'subject': '888923',  
    'effectiveDateTime': '2015-02-07T13:28:17-05:00',  
    'status': 'final'}}
```

Querying Data from Remote Hospital DB

a) Queried diagnosis records

| | subject_id | hadm_id | icd_code | icd_version | long_title |
|---|------------|----------|----------|-------------|---|
| 0 | 18335503 | 21596415 | E86.1 | 10 | Hypovolemia |
| 1 | 18335503 | 21596415 | O99.512 | 10 | Diseases of the respiratory system complicating pregnancy, second trimester |

Mapping Data to FHIR Standard

b) Translated column names into corresponding FHIR concepts

| | subject | encounter | code_code | code_version | code_display | resourceType | id |
|---|----------|-----------|-----------|--------------|---|--------------|----------------------------------|
| 0 | 18335503 | 21596415 | E86.1 | 10 | Hypovolemia | Condition | 0a8e853afd670f1b551de230bc5ef983 |
| 1 | 18335503 | 21596415 | O99.512 | 10 | Diseases of the respiratory system complicating pregnancy, second trimester | Condition | 5bd96a171cd17a540439452334b95aa7 |

c) Converted values into FHIR standard

| | subject | encounter | resourceType | id | code |
|---|--|-------------------------------------|--------------|----------------------------------|---|
| 0 | {'reference': 'Patient/18335503', 'display': '18335503'} | {'reference': 'Encounter/21596415'} | Condition | 0a8e853afd670f1b551de230bc5ef983 | {'coding': [{'system': 'https://icd.who.int/browse10/2019/en#/', 'code': 'E86.1', 'version': 10, 'display': 'Hypovolemia'}], 'text': 'Hypovolemia'} |
| 1 | {'reference': 'Patient/18335503', 'display': '18335503'} | {'reference': 'Encounter/21596415'} | Condition | 5bd96a171cd17a540439452334b95aa7 | {'coding': [{'system': 'https://icd.who.int/browse10/2019/en#/', 'code': 'O99.512', 'version': 10, 'display': 'Diseases of the respiratory system complicating pregnancy, second trimester'}], 'text': 'Diseases of the respiratory system complicating pregnancy, second trimester'} |

d) Final FHIR resources saved in JSON format

```
[{'subject': {'reference': 'Patient/18335503', 'display': '18335503'},
'encounter': {'reference': 'Encounter/21596415'},
'resourceType': 'Condition',
'id': '0a8e853afd670f1b551de230bc5ef983',
'code': {'coding': [{'system': 'https://icd.who.int/browse10/2019/en#/',
'code': 'E86.1',
'version': 10,
'display': 'Hypovolemia'}],
'text': 'Hypovolemia'}},
{'subject': {'reference': 'Patient/18335503', 'display': '18335503'},
'encounter': {'reference': 'Encounter/21596415'},
'resourceType': 'Condition',
'id': '5bd96a171cd17a540439452334b95aa7',
'code': {'coding': [{'system': 'https://icd.who.int/browse10/2019/en#/',
'code': 'O99.512',
'version': 10,
'display': 'Diseases of the respiratory system complicating pregnancy, second trimester'}],
'text': 'Diseases of the respiratory system complicating pregnancy, second trimester'}}
```

Syntactic Validation of the FHIR Mappings

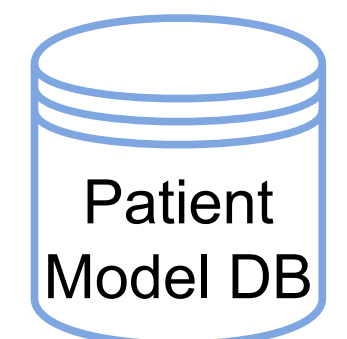
e) Validation of FHIR resources

✓ Validation succeeded

✗ Validation failed

Transferring Data to Patient-Model DB

f) Transferring FHIR resources to Patient-Model DB



Exporting Data into custom JSON File

g) Exporting data from FHIR format into an AI-friendly JSON file

```
[{'feature_name': 'Hypovolemia',
'value': '',
'table_name': 'condition',
'metadata': {'subject': '18335503',
'id': '0a8e853afd670f1b551de230bc5ef983',
'multiple_code': 'no',
'code_display': 'Hypovolemia',
'code_code': 'E86.1',
'code_version': '10',
'code_system': 'https://icd.who.int/browse10/2019/en#/'}},
{'feature_name': 'Diseases of the respiratory system complicating pregnancy, second trimester',
'value': '',
'table_name': 'condition',
'metadata': {'subject': '18335503',
'id': '5bd96a171cd17a540439452334b95aa7',
'multiple_code': 'no',
'code_code': 'O99.512',
'code_display': 'Diseases of the respiratory system complicating pregnancy, second trimester',
'code_version': '10',
'code_system': 'https://icd.who.int/browse10/2019/en#/'}}]
```

medRxiv preprint doi: <https://doi.org/10.1101/2022.11.07.22281564>; this version posted November 13, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.