

Standardization of drug names in the FDA Adverse Event Reporting System: The DiAna dictionary.

Running Title: The DiAna dictionary for drug standardization

Authors: Michele Fusaroli^{1*}, Valentina Giunchi^{1*}, Vera Battini², Stefano Puligheddu¹, Charles Khouri^{3,4}, Carla Carnovale², Emanuel Raschi¹, Elisabetta Poluzzi¹

Affiliations:

¹*Unit of Pharmacology, Department of Medical and Surgical Sciences, University of Bologna, Bologna, Italy.*

²*Unit of Clinical Pharmacology, Department of Biomedical and Clinical Sciences (DIBIC), ASST Fatebenefratelli-Sacco University Hospital, Università degli Studi di Milano, Milan, Italy.*

³*Grenoble Alpes University Hospital, Pharmacovigilance Unit, Grenoble, France*

⁴*Univ. Grenoble Alpes, HP2 Laboratory, Inserm U1300, Grenoble, France*

Corresponding author: Michele Fusaroli. E-mail address: michele.fusaroli2@unibo.it

Declarations

Funding: No specific funding supported this research.

Conflicts of interest: The authors declare no conflict of interest specific for this research.

Ethics approval: not applicable. FAERS spontaneous reports are anonymous and publicly available.

Consent to participate: not applicable. FAERS spontaneous reports are anonymous and publicly available.

Consent for publication: not applicable. FAERS spontaneous reports are anonymous and publicly available.

Availability of data and material: The dictionary and the linkage to the ATC classification are available at https://osf.io/zqu89/?view_only=237d052047c142cabd5d8ea2e765efc6.

Code availability: All the processing, analyses, and visualization were obtained through R-software (version 4.2.1). The code for using the dictionary in the cleaning of the FAERS database is available at <https://github.com/fusarolimichele/DiAna>.

Author contributions: MF, VG, VB conceptualized and designed the study and developed the methodology. MF, VG implemented the automatic translation, and MF, VG, VB, SP manually checked the automatic translation. MF, VG performed the visualization. MF, VG wrote the original draft. All the authors reviewed the draft and approved the final version.

35 Acknowledgements

36 The results, discussion, and conclusions of the study are those of the authors alone and in no way
37 represent the position of the WHO Collaborating Centre for Drug Statistics Methodology on this
38 subject.

39 MF is enrolled in the PhD in General Medical and Services Sciences, Università degli studi di Bologna,
40 which supports their fellowship. V.G. is supported by EU funds (Programma Operativo Nazionale
41 Italian funds for green and innovative research based on European Structural and Investment
42 Funds). VB is enrolled in the PhD in Experimental and Clinical Pharmacological Sciences, Università
43 degli Studi di Milano, which supports her fellowship. We thank Chiara Ballarin, Margherita Bonaiuti,
44 and Laura Pierantozzi (University of Bologna) for the help in the manual revision of drug names
45 standardisation.

46 **Abstract:**

47 **Introduction:** The FDA Adverse Event Reporting System (FAERS) receives drug names in
48 various forms, including brand names, active ingredients, abbreviations, and misspellings,
49 which creates challenges in nomenclature standardization. The lack of consensus on
50 standardization strategies and of transparency hampers replicability and accuracy in
51 conducting disproportionality analysis using FAERS data.

52 **Aim:** We have developed an open-source drug-to-ingredient dictionary called the DiAna
53 dictionary (short for Disproportionality Analysis). Additionally, we have linked the DiAna
54 dictionary to the WHO Anatomic Therapeutic Chemical (ATC) classification system.

55 **Methods:** We retrieved all drug names reported to the FAERS from 2004 to December 2022.
56 Using existing dictionaries such as RxNorm and string editing techniques, we automatically
57 translated the drug names to active ingredients. Manual revision was performed to correct
58 errors and improve translation accuracy. The resulting DiAna dictionary was linked to the
59 ATC classification, proposing a primary ATC code for each ingredient.

60 **Results:** We retrieved 18,151,842 reports, with 74,143,411 drug entries. We automatically
61 translated and manually checked the first 14,832 terms, up to terms occurring at least 200
62 times (96.88% of total drug entries), to 6,282 unique active ingredients. Automatic
63 unchecked translations extend the standardization to 346,854 terms (98.94%). After linking
64 to the ATC classification, the most prominent drug classes in FAERS reports were
65 immunomodulating (37.40%) and nervous system drugs (29.19%).

66 **Conclusion:** We present the DiAna dictionary as an open-source tool and encourage experts
67 to provide input and feedbacks. Regular updates can improve research quality and promote
68 a common pharmacovigilance toolbox, ultimately advancing safety and improving study
69 interpretability.

70

71 **Key points:**

- 72 • Drug name standardization impacts signal detection accuracy.
- 73 • DiAna dictionary cleanses drugs in FAERS for improved data control.
- 74 • DiAna's transparency and flexibility improves interpretability.

75 1. Introduction

76 1.1. *The need for transparency on data preprocessing*

77 Spontaneous reporting systems (SRSs) are public and private services collecting reports
78 of suspected adverse drug reactions in order to timely detect potential issues in drug safety
79 [1,2]. However, these reports are gathered from a variety of regional, national, and
80 manufacturer databases, which have different languages, rules, and forms for data storage.
81 Furthermore, the same SRS may collect reports from different kinds of reporter (e.g.,
82 manufacturers, healthcare professionals, lawyers, consumers) and using both paper and
83 electronic forms. This heterogeneity leads to highly variable content, particularly in free text
84 fields which may contain misspellings, out-of-context information, and the use of different
85 lexicons. Furthermore, reports may be incomplete, due to the spontaneous nature of the
86 reporting, or duplicates, due to multiple reporting and system errors. For these reasons, an
87 extensive cleaning procedure is needed prior to any kind of statistical analysis.

88 The United States (US) Food and Drug Administration (FDA) Adverse Event Reporting
89 System (FAERS) is one of the most used SRS in signal detection, because of its free access
90 and its large catchment area mirroring the entire world (even if with a large representativity
91 for the US). There are two ways the FDA provides free access to the FAERS data, with
92 notable differences.

93 First, already pre-processed data through an online public dashboard [3], developed for
94 transparency reasons and the promotion of higher quality reporting. However, this tool is
95 obtained from spontaneous reports through a cleaning procedure that is, allegedly, both
96 undisclosed and partial, for example because it does not attempt any duplicates detection
97 and it does not provide all the information reported in individual case safety reports (ICSRs).
98 Therefore, it is not suitable for research purpose, namely complex analyses beyond simple
99 exploratory analyses.

100 Second, raw quarterly data (both in ASCII and XML format) [4] that allow to knowingly
101 perform and document the entire pre-processing procedure. This cleaning and
102 normalization procedure requires the researchers a conspicuous effort and faces them with
103 the need to make multiple operative choices: for example how to deal with duplicates, how
104 to deal with dates that up to 2012 were completed automatically when partial, how to deal
105 with unclear entries (e.g., in 2019 “RN” was often recorded as a reporter type; while this

106 entry may refer to registered nurses, it is not documented in the “readme” file of the
107 FAERS).

108 These choices are seldom driven by objectivity alone and must be considered both in
109 the design and the interpretation, because the same database cleaned with different
110 procedures may give different results [5]. The subsequent lack of replicability heavily
111 impacts on the credibility of SRSs, already diminished due to their inherent bias, which
112 hinders any interpretation of disproportionality analysis beyond the generation of
113 hypotheses [6].

114 To sum up, SRSs’ data are extremely raw and heterogeneous, necessitating a
115 meticulous cleaning process guided by clinical and pharmacological reasoning before
116 conducting any analysis [7]. Throughout each stage of this process, it is crucial to uphold
117 collaboration among multiple professionals and maintain an understanding of the data
118 collection features and the relevant underlying theory for the phenomenon under
119 investigation [8]. Additionally, it is of utmost importance, for researchers, not only to be in
120 control and knowledgeable of the pre-processing procedure, but also to be transparent
121 about their operational choices, allowing for external assessment and interpretation.

122

123 1.2 Drug nomenclature issues

124 Among the information requiring standardization prior to case retrieval and analysis,
125 a particular effort should be focused on the standardization of drug names, which in the
126 FAERS are collected as verbatim (free) text. A drug may be recorded in the FAERS using
127 either the brand name or the active ingredients, with the international nonproprietary term
128 (INN, defined by the WHO) or the United States adopted name (USAN, defined by the USAN
129 council), with the full name or abbreviations. Misspellings might easily occur, and the drug
130 name is sometimes followed by dose, route, and formulation details.

131 An objective standardization of these entries is unattainable: for example, the same
132 brand name may refer to different compositions in different countries, or two brand names
133 may be just one letter apart thus being extremely susceptible to misspellings. The
134 inconsistencies that derive from the multiple operative options and researcher’s personal
135 choices can affect case-retrieval and impair replicability among studies. Nonetheless, no
136 consensus on the best operative procedures has been achieved, and already-published
137 analyses using SRSs’ data are rarely transparent on the cleaning choices adopted, lacking

138 documentation on whether and how the FAERS was prepared for statistical analyses. The
139 FAERS system itself, which does have a formal dictionary through which it cleans the data
140 for the public dashboard, does not make it publicly available. Among the 17 studies
141 conducted using the FAERS quarterly data and published in February 2023 (accessed on
142 PubMed on March 15th, 2023), ten did not state any drug standardization process, six used
143 an automatic translation via dictionaries, one performed a manual translation but did not
144 make it publicly available (Table S1). This lack of transparency also affects some free ready-
145 to-use pharmacovigilance databases which provide already pre-processed FAERS data
146 [9,10]. Moreover, these databases standardize drugs only according to US dictionaries of
147 drug names (e.g., RX-norm, orange book), thus failing to identify foreign drug names,
148 misspellings, and other free text issues that were described above. The WHODrug Global is a
149 dictionary that compiles extensive drug information from across the globe [11]. However, it
150 does not consider the challenges posed by free text issues, and it is not available as an open-
151 source resource. Other tools attempt to edit drug names using an automatic translation of
152 potential misspellings, which may result in mistranslations due to the existence of similar
153 drug names referring to formulations with different active ingredients [12]. For this reason,
154 already in 2015 Wong et al. produced a manually revised translation of the LAERS drug
155 names (the FAERS system up to 2012), with a transparent explanation of the operative
156 choices [13]. Still, their translation was not publicly available and did not come into use.

157

158 *1.3 Aim of the study*

159 In this work, we follow Wong et al. efforts and extend their work to consider
160 previously unattended nomenclature issues. We propose an open-source drug name-to-
161 ingredient dictionary for standardizing the FAERS updated to December 2022 together with
162 a transparent report of the data cleaning protocol to identify and resolve drug
163 nomenclature issues. This pharmacovigilance tool, that we define as DiAna
164 (Disproportionality Analyses) dictionary, with its linkage to the Anatomic Therapeutic
165 Chemical (ATC) classification, aims to support pharmacovigilance researchers towards a
166 greater control on the FAERS data, a higher replicability and accuracy of disproportionality
167 analyses, and a more appropriate interpretation of their results.

168

169 2. Methods

170

171 *2.1 The FAERS database*

172 We downloaded FAERS Quarterly Data (trimestral) Extract Files[4] in ASCII format from
173 04Q1 to 22Q4. These files are composed of five tables linked through a primary key
174 (“primaryid”) identifying a specific version of a report: DEMO (demographic and
175 administrative information), DRUG (information on reported medications), REAC (adverse
176 events), OUTC (outcomes), RPSR (report sources).

177 DRUG is also linked through “primaryid” and a secondary key (“drug_seq”), identifying a
178 specific medication within a report, to other two tables: THER (start dates and end dates for
179 the reported medications), INDI (indications for using the reported medications).

180

181 *2.2 Automatic set up of the dictionary*

182 We combined all DRUG quarters into one database. We focused on three columns used
183 to identify the medicinal product:

- 184 • Drugname, recording the name of the medicinal product.
- 185 • Prod_ai, recording the product's active ingredients, when available.
- 186 • Val_vbm, recording whether the source of drugname was a validated trade name
187 (value = 1) or a verbatim name (value = 2).

188 Since our aim was the translation to active ingredients, we did not consider the
189 column “val_vbm”. We instead retrieved all the unique terms from the other two columns
190 (i.e., Prod_ai and Drugname), lowered upper cases, and removed multiple spaces, leading
191 and trailing spaces and punctuation, and spaces between parentheses and included text.
192 We merged the unique terms with RX-norm⁷ and WHO-ATC substances⁸ to create a
193 dictionary with automatic translations to active ingredients. The merging process was also
194 repeated after several rounds of text editing, during which we removed leading or trailing
195 spaces and specific terms or symbols such as chirality indicators (e.g., "+", "-", "d", "s") and
196 text between brackets or caret symbols.

197

198 *2.3 Manual revision*

199 We then manually revised all the automatic translations starting with the most
200 frequently reported ones up to the terms recorded in a minimum of 200 reports (and
201 beyond, ongoing). We started translating drug names to active ingredients included in the
202 2023 update of the ATC classification of the WHO, integrating it whenever we met new
203 active ingredients. Hand search of foreign drug names was performed using online
204 databases (e.g., DrugBank.com [14] and Drugs.com [15]), manufacturer websites, and
205 websites storing information from foreign package labels (e.g., Kusuri-no-Shiori –drug
206 information sheets– from the Japanese regulatory agency, accessed at <https://www.radar.or.jp/siori/english/>).

208

209 *2.4 Nomenclature issues*

210 Multiple issues were identified in the process of translating drug names, including brand
211 names and abbreviations) to active ingredients (e.g., “Zantac” was translated into
212 “raniditine”):

- 213 • A drug may include multiple ingredients. We translated the drug to all its ingredients
214 and ordered them alphabetically, separating them by a semicolon. For example,
215 “Entresto” was translated into “sacubitril;valsartan”.
- 216 • The spelling of an active ingredient can be different between the United States
217 Adopted Name (USAN) and the International Nonproprietary Name (INN): for
218 example, acetaminophen (USAN) = paracetamol (INN); amphetamine (USAN) =
219 amfetamine (INN); dimethicone/simethicone (USAN) = simeticone (INN); cysteamine
220 (USAN) = mercaptamine (INN). We gave preference to the INN.
- 221 • The active ingredient may be recorded in languages different from English (e.g.,
222 acide folique). We translated everything to the English INN.
- 223 • Typing mistakes can occur (e.g., “zopiclone” instead of “zopiclone”; “Diavan®”
224 instead of “Diovan®”). We manually fixed the mistakes taking into account the INN.
- 225 • The same drug name may contain different ingredients in different countries (e.g.,
226 Gaster® contains famotidine in the US, omeprazole in Japan, cromoglicic acid in Italy;
227 Previscan® contains fluindione in the US, pentoxifylline in Italy and Spain; Furix®
228 contains furosemide in the US, cefuroxime in India). In these cases, we translated the

229 brand name to the active ingredient contained in the US packaging, assuming that
230 US is more represented in the FAERS rather than other countries. When the brand
231 name of interest was not sold in the US, we checked the most reported country in
232 the FAERS for that specific brand name.

- 233 • The drug name may be missing or underspecified:
 - 234 ○ When there was no medication, we translated the drug name field to "no
235 medication".
 - 236 ○ When the medication was unspecified, we translated the drug name field as
237 "unspecified".
 - 238 ○ Unspecified drug-class terms were translated to the most specific term
239 possible (e.g., "water pills" as "diuretics, unspecified", "antihypertensives" as
240 "antihypertensives, unspecified").
- 241 • We specified when drug (or placebo) consumption occurred in a clinical trial (such as
242 when blinded was specified, or when the investigational name was used –e.g., cc-
243 223 for onatasertib) translating the drug name as “active ingredient, trial”, “placebo,
244 trial”, or “unspecified, trial”.
- 245 • Additionally, we decided to standardize terms other than drugs to broader
246 categories, since specific details are seldom provided: minerals (e.g., calcium),
247 vitamins (e.g., vitamin b5, independent of the route of administration; vitamin b12,
248 independent of the form—e.g., cyanocobalamin, mecobalamin—), devices (e.g.,
249 intrauterine contraceptive device), vaccines (e.g., COVID-19 vaccine), and
250 phytotherapies (e.g., plantago spp).
- 251 • If a drug name was reported followed by a non-coherent active ingredient in square
252 brackets, we assumed that an error was made during the compilation by the
253 pharmacovigilance expert. In this case, we translated the entry based on the drug
254 name alone, considering incorrect the active ingredient listed in square brackets.

255

256 The whole list of standardized names is available in the open-source repository
257 [https://osf.io/zqu89/?view_only=237d052047c142cabd5d8ea2e765efc6].

258

259 *2.5 Linkage to the ATC classification*

260 Furthermore, we performed data-linkage between the DiAna dictionary and the
261 hierarchical ATC classification, which was downloaded from the WHO Collaborating Centre
262 for Drug Statistics and Methodology website [16] using the R package “rvest”. Since this
263 classification is mainly a tool for drug utilization research, the same active ingredient may be
264 given more than one ATC code if it is available in multiple strengths or routes of
265 administration with clearly different therapeutic uses [17]. We linked the final list of
266 individual active ingredients from our translation with the ATC classification, manually
267 integrating for different choices in the nomenclature (e.g., “vitamin b9”–folic acid– to
268 B03BB01), for classes of drugs (e.g., “antihypertensives, unspecified” to C02), and for drugs
269 recorded in the ATC only in combination (glecaprevir and pibrentasvir both to J05AP57). We
270 have here linked each active ingredient to all its ATC codes (most of the time it is not
271 possible to discriminate between the different ATC codes based only on the drug name) but,
272 since sometimes it is important to count each ingredient only once, we also proposed a
273 unique primary ATC code for each ingredient. To this end, we prioritized the first level in the
274 following order (“H”, “J”, “P”, “L”, “M”, “N”, “C”, “G”, “R”, “B”, “D”, “A”, “S”, “V”). We furtherly
275 moved vitamin c from G01AD03 to A11GA01, sex hormones having both a genitourinal and
276 an immunomodulating code to the genitourinal, and sodium and calcium chloride to the
277 alimentary ATC code instead of the blood-related. This data linkage, while not useful to
278 identify specific formulations, may be used to define the drugs of interest or for visualization
279 purposes in the implementation of disproportionality analyses.

280

281 **3. Results**

282 We downloaded the FAERS quarterly data up to 22Q4 and retrieved 18,151,842
283 ICSRs, for a total of 74,143,411 drug entries (92.81% allegedly recorded using a validated
284 trade name) and 955,778 unique drugname and prod_ai terms (see **Figure 1**). After the
285 initial formatting, we reduced them to 793,274 unique entries. We automatically translated
286 346,854 terms (98.94% of total drug entries) and manually checked the first 14,832 terms
287 (96.88%) up to 174 occurrences (<0.00015%, ongoing).

288 A total of 6,282 unique ingredients were included in the DiAna dictionary, of which
289 3,209 were linked to the ATC classification. The most common primary ATC classes in the
290 FAERS, after translation with the DiAna dictionary, were immunomodulating (reported in

291 37.40% FAERS reports), nervous system (29.19%), alimentary tract (25.18%), and
292 cardiovascular agents (20.17%) (see **Figure 2**). Most frequently reported medicinal products
293 were paracetamol (5.45%), acetylsalicylic acid (4.62%), adalimumab (3.81%), etanercept
294 (3.35%), levothyroxine (3.17%), ranitidine (3.13%) (see **Table 1**). When compared with the
295 untranslated formatted FAERS and with the FAERS translated according to RxNorm, the
296 translation based on DiAna dictionary showed clear advantages in case retrieval (98.94% of
297 total drug entries against 76.32% by RxNorm). Among the most reported medicinal
298 products, DiAna allowed to retrieve more cases than RxNorm, from a ratio of 1.01 for
299 etanercept (638,427 vs 632,130), to a ratio of 8.55 for ranitidine (69,883 vs 597,604). Due to
300 differences in nomenclature some ratios were not calculated (e.g., paracetamol is translated
301 to acetaminophen and acetylsalicylic acid to aspirin by RxNorm). For some drugs the added
302 value of DiAna translation for case retrieval was extremely high: for example, rimegepant
303 (ratio = 277.91; 6,392 vs 23), adapalene (122.60; 174,711 vs 1,425), drospirenone (108.49;
304 86,356 vs 796), and umeclidinium (105.66; 45,751 vs 433; not shown in the table).

305 The translation also took account information about placebo and experiments, thus
306 identifying 50,967 reports as generated within trials (0.28%).

307

308 4. Discussion

309 4.1 The DiAna Dictionary

310 The sensitivity of case retrieval and the relevant disproportionality analysis results may
311 vary depending on the drug cleaning procedures used in SRSs. Disproportionality analysis is
312 mostly performed on public dashboards or other analytical tools with no access to
313 underlying data, ready-to-use databases with partial or non-transparent translation, or
314 individually cured undisclosed databases. While these tools provide easy access to
315 disproportionality analysis, they also pose a risk of inappropriate analyses and
316 interpretation due to users' unawareness on the nature of data. Common drug translation
317 procedures involve automatic linkage to existing dictionaries (offering only partial
318 translation) and automatic algorithms dealing with misspellings (potentially introducing
319 errors). To address these concerns, a dictionary for drug name-to-ingredient translation was
320 developed through an automatic procedure that was manually checked and extended. This
321 dictionary, called DiAna dictionary, required a time-consuming effort and is made available
322 opensource for everyone to use it and propose changes. The use of the DiAna dictionary will

323 allow the pharmacovigilance community to agree on the best possible translation. A greater
324 control on data cleaning will result in improved replicability and accuracy of signals, and
325 more conscious and appropriate interpretation of results, with relevant benefit for the
326 scientific community.

327

328 *4.2 Better retrieval for higher sensitivity*

329 We were able to translate 98.94% of total drug entries to 6,282 unique active
330 ingredients using the DiAna dictionary, compared to 76.32% using only RxNorm. When
331 considering unique drug entries, we translated 346,854 terms over 793,274 (43.72%). We
332 manually checked the first 14,832 terms (up to 174 occurrences), which were responsible
333 for the translation of 96.88% of total drug entries. We believe that this is a good starting
334 point to share our work with the pharmacovigilance community and enable more
335 participative use and development of the DiAna dictionary. In contrast to the previous work
336 by Wong et al.[13], made on the FAERS up to 2012, we made our dictionary (up to 2022)
337 open source. We chose to design the translation so that a new column is produced with only
338 active ingredients, while keeping the original verbatim text in a separate column for more
339 in-depth analyses. We have also decided not to translate to salts as this is rarely taken into
340 account in disproportionality analysis and can lead to confusion about whether the same
341 ingredient with unspecified salt should be considered among cases or non-cases. Instead,
342 we have included the linkage to the ATC classification, and provided translation also to
343 higher ATC classes such as “antihypertensives, unspecified”, as this information can be
344 important for adjusting the analysis and assessing individual cases.

345 The DiAna dictionary translates a higher proportion of the database, enabling a higher
346 sensitivity in case retrieval, and a higher number of identified cases. This results in better
347 specificity in the definition of non-cases and higher accuracy in signal detection. This means
348 earlier and clearer signals, as in some specific products the number of reports retrieved
349 significantly increased. For example, for rimegepant the DiAna dictionary identifies 278
350 times more reports than RxNorm alone.

351 In addition to identifying active ingredients, the drug name information enabled us to
352 identify reports derived from clinical trials (0.28% of total reports), as they recorded
353 placebo, blinding, or drug codes. This information can help researchers exclude evidence

354 already taken into account in other steps of drug safety characterization from the
355 disproportionality analysis.

356 Finally, the linkage between the DiAna dictionary and the ATC classification can help in
357 the retrieval of drug classes and in visualization. The information on the distribution of drug
358 classes in the database is particularly useful for the design of future disproportionality
359 analyses, as it provides insight into the representativeness of the population chosen as
360 comparison. Over one-third of the database consists of reports documenting the utilisation
361 of immunomodulating drugs. Recent observations, specifically in the context of the
362 extensive rollout of COVID-19 vaccines, have reignited the attention to the possibility that
363 this uneven distribution of drugs in the SRS may lead to masking/cloaking bias, hiding
364 disproportionality signals [18].

365 The DiAna dictionary and its linkage to the ATC classification are freely available online
366 for everyone to use
367 (https://osf.io/zqu89/view_only=237d052047c142cabd5d8ea2e765efc6), and can be
368 corrected and expanded by experts in the field. Changes can be proposed in the GitHub
369 repository (<https://github.com/fusarolimichele/DiAna>) under the issue DiAna dictionary and
370 will be periodically validated and integrated into the existing dictionary. This collaborative
371 effort will improve the quality and reproducibility of pharmacovigilance research. The
372 dictionary can be downloaded in Excel and csv formats and can be imported into any data
373 management software, such as R, to automatically translate drug names to active
374 ingredients before conducting analyses. Users can also easily modify the translation of
375 specific terms for their analyses, which is not possible with ready-to-use FAERS databases.

376

377 *4.3 Limitations, Strengths, and Further Goals*

378 The DiAna dictionary is not designed as a static dictionary: it will require ongoing efforts to
379 keep up with new drugs and terms. We are recursively extending our translation to reach
380 and maintain a full checked translation of any entry with over 100 occurrences. Users of the
381 DiAna dictionary should be aware of this limitation (which is even more impairing in other
382 disproportionality tools), especially with less frequent terms that may not be included in the
383 dictionary. It's recommended that before any research on a specific drug, inherent terms
384 are checked in the dictionary and any new translations are shared to integrate into the
385 DiAna dictionary for everyone to benefit. The translation will plausibly never be complete,

386 since some terms are not easily translated (e.g., “chinese food”) and many choices are partly
387 subjective. However, these choices can be defined in agreement with the entire
388 pharmacovigilance community.

389 The translation of ambiguous terms was also noted as a challenge, especially with over-
390 the-counter cold, cough, and flu agents (multiple ingredients changing over the years).
391 When we were not certain, we used the higher-level term (e.g., “cough preparations,
392 unspecified”). The lack of expertise in supplements and phytotherapies may have resulted in
393 the dictionary being excessively generic (for example referring to *plantago* spp instead of
394 individual species, and to covid 19 vaccines instead of specific types), and it could benefit
395 from expert refinement for higher specificity and coverage of entries provided to other
396 spontaneous report databases (CAERS and VAERS are more appropriate to investigate the
397 safety profile of these medicinal products).

398 Since lack of completeness is a known problem in spontaneous reports, and other
399 information is not always available, we implemented sharp-cut operative choices to retrieve
400 active ingredients based only on the drug name. The use of additional columns such as
401 country, year of occurrence, dose, indication, and route of administration, could help
402 discriminate between mistranslations when the same drug name may be translated to
403 multiple active ingredients. Moreover, information from the drug name column could be
404 used to impute information into other columns. For example, “nizoral a-d” is translated to
405 ketoconazole and refers specifically to an anti-dandruff shampoo (i.e., the indication,
406 formulation, and route of administration could be imputed if missing), while “hypersal”
407 refers to a sodium chloride nebulizer solution, and “jinarc” refers to a formulation of
408 tolvaptan specifically indicated for autosomal dominant polycystic disease. By incorporating
409 a drug name-to-product translation feature, for example referring to the WHO Drug Global
410 or to the Identification of Medicinal Products (IDMP), we could streamline the process of
411 imputation of structured fields using free text, thereby enhancing the value of the DiAna
412 dictionary.

413 Linking INN names to ATC codes was a complex task due to the existence of combination
414 products (e.g., glecaprevir and pibrentasvir), medicinal products with ingredients that do
415 not have an ATC code yet, and experimental substances which are missing even the INN.
416 The linkage will be annually updated according to changes in the ATC classification to
417 preserve its utility.

418

419 5. Conclusion

420 We offer the DiAna dictionary as open-source tool for the pharmacovigilance community to
421 standardize drug names in the FAERS database. Its public accessibility, transparency, and
422 flexibility provide a foundation for ongoing improvement and refinement through input
423 from experts in the field. With periodic updates, this open-source project can drive a
424 common effort towards a more transparent and cleaner shared FAERS database, leading to
425 more replicable and higher quality research in pharmacovigilance. Ultimately, by sharing
426 and mutually enriching our knowledge, we can develop a common pharmacovigilance
427 toolbox that advances safety and improves the accuracy, replicability and reliability of
428 pharmacovigilance studies.

429 **BIBLIOGRAPHY**

- 430 1. Poluzzi E, Raschi E, Piccinni C, Ponti FD. Data Mining Techniques in Pharmacovigilance:
431 Analysis of the Publicly Accessible FDA Adverse Event Reporting System (AERS). Data Mining
432 Applications in Engineering and Medicine. IntechOpen; 2012;
- 433 2. Raschi E, Moretti U, Salvo F, Pariente A, Antonazzo IC, Ponti FD, et al. Evolving Roles of
434 Spontaneous Reporting Systems to Assess and Monitor Drug Safety. Pharmacovigilance
435 [Internet]. 2018 [cited 2019 Feb 3]; Available from: [https://www.intechopen.com/online-](https://www.intechopen.com/online-first/evolving-roles-of-spontaneous-reporting-systems-to-assess-and-monitor-drug-safety)
436 [first/evolving-roles-of-spontaneous-reporting-systems-to-assess-and-monitor-drug-safety](https://www.intechopen.com/online-first/evolving-roles-of-spontaneous-reporting-systems-to-assess-and-monitor-drug-safety)
- 437 3. FDA. FDA Adverse Event Reporting System (FAERS) Public Dashboard | FDA [Internet].
438 [cited 2022 Dec 14]. Available from: [https://www.fda.gov/drugs/questions-and-answers-](https://www.fda.gov/drugs/questions-and-answers-fdas-adverse-event-reporting-system-faers/fda-adverse-event-reporting-system-faers-public-dashboard)
439 [fdas-adverse-event-reporting-system-faers/fda-adverse-event-reporting-system-faers-](https://www.fda.gov/drugs/questions-and-answers-fdas-adverse-event-reporting-system-faers/fda-adverse-event-reporting-system-faers-public-dashboard)
440 [public-dashboard](https://www.fda.gov/drugs/questions-and-answers-fdas-adverse-event-reporting-system-faers/fda-adverse-event-reporting-system-faers-public-dashboard)
- 441 4. Center for Drug Evaluation and Research. FDA Adverse Event Reporting System - Latest
442 Quarterly Data Files [Internet]. FDA. 2019 [cited 2019 Jul 28]. Available from:
443 [http://www.fda.gov/drugs/fda-adverse-event-reporting-system-faers/fda-adverse-event-](http://www.fda.gov/drugs/fda-adverse-event-reporting-system-faers/fda-adverse-event-reporting-system-faers-latest-quarterly-data-files)
444 [reporting-system-faers-latest-quarterly-data-files](http://www.fda.gov/drugs/fda-adverse-event-reporting-system-faers/fda-adverse-event-reporting-system-faers-latest-quarterly-data-files)
- 445 5. Hauben M, Reich L, Gerrits CM, Younus M. Illusions of objectivity and a recommendation
446 for reporting data mining results. *Eur J Clin Pharmacol*. 2007;63:517–21.
- 447 6. Mouffak A, Lepelley M, Revol B, Bernardeau C, Salvo F, Pariente A, et al. High prevalence
448 of spin was found in pharmacovigilance studies using disproportionality analyses to detect
449 safety signals: a meta-epidemiological study. *Journal of Clinical Epidemiology*. 2021;138:73–
450 9.
- 451 7. Rocca E, Grundmark B. Monitoring the Safety of Medicines and Vaccines in Times of
452 Pandemic: Practical, Conceptual, and Ethical Challenges in Pharmacovigilance [Special
453 Issue]. *Argumenta*. 2021;7:127–46.
- 454 8. Leonelli S. The challenges of big data biology. *eLife*. 8:e47381.
- 455 9. Banda JM, Evans L, Vanguri RS, Tatonetti NP, Ryan PB, Shah NH. A curated and
456 standardized adverse drug event resource to accelerate drug safety research. *Sci Data*
457 [Internet]. 2016 [cited 2020 Dec 17];3. Available from:
458 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4872271/>
- 459 10. Khaleel MA, Khan AH, Ghadzi SMS, Adnan AS, Abdallah QM. A Standardized Dataset of a
460 Spontaneous Adverse Event Reporting System. *Healthcare. Multidisciplinary Digital*
461 *Publishing Institute*; 2022;10:420.
- 462 11. Lagerlund O, Strese S, Fladvad M, Lindquist M. WHODrug: A Global, Validated and
463 Updated Dictionary for Medicinal Information. *Ther Innov Regul Sci*. 2020;54:1116–22.
- 464 12. Stanford T. The fuzzyfaers package [Internet]. 2022 [cited 2022 Dec 24]. Available from:
465 <https://github.com/tystan/fuzzyfaers>

- 466 13. Wong CK, Ho SS, Saini B, Hibbs DE, Fois RA. Standardisation of the FAERS database: a
467 systematic approach to manually recoding drug name variants. *Pharmacoepidemiology and*
468 *Drug Safety*. 2015;24:731–7.
- 469 14. Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, et al. DrugBank 5.0: a major
470 update to the DrugBank database for 2018. *Nucleic Acids Research*. 2018;46:D1074–82.
- 471 15. Drugs.com | Prescription Drug Information, Interactions & Side Effects [Internet].
472 Drugs.com. [cited 2022 Dec 24]. Available from: <https://www.drugs.com/>
- 473 16. WHOCC - ATC/DDD Index [Internet]. [cited 2023 May 9]. Available from:
474 https://www.whocc.no/atc_ddd_index/
- 475 17. WHO Collaborating Centre for Drug Statistics Methodology. Guidelines for ATC
476 classification and DDD assignment 2023 [Internet]. [cited 2023 May 9]. Available from:
477 https://www.whocc.no/atc_ddd_index_and_guidelines/guidelines/
- 478 18. Harpaz R, DuMouchel W, Van Manen R, Nip A, Bright S, Szarfman A, et al. Signaling
479 COVID-19 Vaccine Adverse Events. *Drug Saf*. 2022;45:765–80.
- 480

481 **Figure 1 – Translation pipeline.** Flowchart showing the procedure to translate drug names to active
482 ingredients. Some examples of the processing of entries are provided in the background. The color
483 remains constant across the steps, and within each step the dimension is proportional to the
484 number of occurrences.

485

486 **Figure 2 – Distribution of medicinal products in the FAERS.** Drugs most frequently reported in the FAERS, after
487 translation, according to ATC class. Each step is a first level, starting from the most reported one. Within each
488 level, a tree map shows how ATC levels 2 and 4 are reported in FAERS reports. The 3 most reported active
489 ingredients of each 1st level are also shown.

490

491

492 **Table 1 – Performance of DiAna translation.** Drugs most frequently reported in the FAERS, after DiAna
 493 translation, relative to simple formatting and to the merging with RxNorm. The number of occurrences in the
 494 three translations are reported together with the ratio of occurrences between DiAna and the others. In some
 495 cases, differences in the nomenclature resulted in empty cells.

Active substance	Formatting (n. occurrences)	RxNorm (n. occurrences)	DiAna (n. occurrences)	<i>DiAna/RxNorm</i>	<i>DiAna/Formatting</i>
paracetamol	112,939	(-)	1,040,051	(-)	9.21
acetylsalicylic acid	61,652	(-)	942,051	(-)	15.28
adalimumab	21,403	699,331	727,730	1.04	34.00
etanercept	19,332	632,130	638,427	1.01	33.02
levothyroxine	55,115	288,916	604,688	2.09	10.97
ranitidine	69,874	69,883	597,604	8.55	8.55
methotrexate	224,467	230,866	553,011	2.40	2.46
prednisone	177,872	181,092	552,531	3.05	3.11
omeprazole	132,069	273,775	539,760	1.97	4.09
insulin	73,619	(-)	539,088	(-)	7.32
metformin	279,461	320,567	534,234	1.67	1.91
atorvastatin	210,528	421,702	529,453	1.26	2.51
calcium	150,507	(-)	513,772	(-)	3.41
amlodipine	250,232	347,454	508,501	1.46	2.03
furosemide	99,789	269,463	472,999	1.76	4.74
oxycodone	98,512	287,818	462,109	1.61	4.69
salbutamol	40,326	(-)	435,816	(-)	10.81
pantoprazole	182,322	288,211	421,710	1.46	2.31
metoprolol	76,981	104,124	420,331	4.04	5.46
magnesium	68,548	(-)	414,166	(-)	6.04
fluticasone	12,549	91,204	397,614	4.36	31.68
lenalidomide	27,201	370,133	380,151	1.03	13.98
hydrochlorothiazide	72,069	(-)	380,068	(-)	5.27
gabapentin	77,949	170,801	372,316	2.18	4.78
dexamethasone	91,353	131,730	365,042	2.77	4.00
lisinopril	125,376	150,113	361,929	2.41	2.89
vitamin b9	100	(-)	356,564	(-)	3,565.64
simvastatin	101,474	163,417	341,174	2.09	3.36
vitamin d3	142,967	(-)	327,770	(-)	2.29

hydrocodone	55,669	56,140	325,467	5.80	5.85
-------------	--------	--------	---------	------	------

496

Retrieval

ACETAMINOPHEN HYDROCODONE
 BITARTRATE
ZANTAC VITAMINS
RANITIDINE
 CALCIUM CHLORIDE DEXTROSE MAGNESIUM
 CHLORIDE SODIUM CHLORIDE SODIUM LACTATE
 COUMADIN
TOFACITINIB CITRATE

DRUG rows
 N = 74,143,411

Unique drug names
 N = 942,534

Unique prod_ai
 N = 13,244

FISH OIL
 CYTAA 7 (MAG) (MAG)
 UNSPECIFIED INGREDIENT
 ALPHA-1-PROTEINASE INHIBITOR HUMAN
 BAKERSCH THERAPY
 CYANOCOBALAMIN
 PROXYTOX (CYCLOPHOSPHAMIDE) 300MG/VAL
ADALIMUMAB
ASPIRIN
PROACTIV MD ADAPALENE ACNE TREATMENT
ACETAMINOPHEN

Formatting

acetaminophen hydrocodone bitartrate
 vitamins bitartrate cosmetic
zantac ranitidine
 calcium chloride dextrose magnesium chloride
 sodium chloride sodium lactate
 coumadin
tofacitinib citrate

955,778
 unique entries

1. lower capitalisation
2. remove excess spaces
3. remove starting and trailing punctuation

fish oil
 alpha-1-proteinase inhibitor human
 cyanocobalamin
adalimumab
aspirin
 proactiv md adapalene acne treatment
acetaminophen

Standardization

hydrocodone; paracetamol
 ranitidine tofacitinib
 calcium chloride glucose magnesium sodium chloride sodium lactate
ranitidine warfarin
 vitamins, unspecified

793,274
 unique entries

1. automatic translation through linkage with existing dictionaries
2. recursive text editing
3. manual revision and integration



adalimumab
 adapalene vitamin b9
 unsaturated fatty acids
 paracetamol
acetylsalicylic acid
 drosiprenone ethinylestradiol vitamin b12
 cyclophosphamide

ATC linkage

R05DA03; N02BE01
 A02BA03
A02BA02 A11
 A12AA07; B05C001; A12CC; A12CA01; A12CA
A02BA02 A02BA03
 M03AX01
B01AA03 L04AA29

346,854 translated (43.72%)
 (98.94%)
 14,832 checked (1.84%)
 (96.88%)

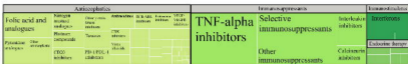
6,282
 unique ingredients

3,209
 ingredients with ATC
 (3,744 ATC codes)

C10AX06 A06AA01; N03AX01
L04AB04 B03BB01
 N03CF01
D10AD03 N02BE01
 B02AB02
N02BA01 B03BA
 G03AC10; G03CA01
L01AA01

Immunomodulating (I)
(37.40%)

adalimumab (3.81%)
etanercept (3.35%)
methotrexate (2.90%)

**Nervous (N)**

(29.19%)
paracetamol (5.45%)
acetilsalicylic acid (4.94%)
oxycodone (2.42%)

**Alimentary (A)**

(25.18%)
ranitidine (3.13%)
insulin (2.83%)
omeprazole (2.83%)

**Cardiovascular (C)**

(20.17%)
atorvastatin (2.78%)
amlodipine (2.67%)
furosemide (2.48%)

**Blood (B)**

(12.72%)
vitamin b9 (1.87%)
warfarin (1.40%)
iron (1.29%)

Hormones (H)

(12.17%)
levothyroxine (3.17%)
prednisone (2.90%)
dexamethasone (1.91%)

Respiratory (R)

(11.36%)
salbutamol (2.28%)
fluticasone (2.08%)
hydrocodone (1.71%)

Antifungals (J)

(10.87%)
tenofovir (0.84%)
emtricitabine (0.65%)
trimethoprim (0.96%)

Musculo-Skeletal (M)

(10.52%)
ibuprofen (1.63%)
naproxen (1.08%)
denosumab (1.01%)

Genitourinary (G)

(8.12%)
levonorgestrel (1.15%)
etinylestradiol (1.02%)
tamsulosin (0.75%)

Dermatologicals (D)

(4.52%)
adapalene (0.92%)
dupilumab (0.87%)
fumaric acid (0.67%)

Various (V)

(2.08%)
folic acid (0.35%)
oxygen (0.22%)
deferasirox (0.19%)

Herbals & Probiotics

(1.38%)
senna spp (0.43%)
lactobacillus spp (0.12%)
curcuma spp (0.09%)

Antiparasitic (P)

(1.34%)
hydroxychloroquine (0.92%)
simeticone (0.13%)
quinine (0.07%)

Sensory (S)

(1.27%)
latanoprost (0.28%)
bimatoprost (0.18%)
ranibizumab (0.16%)

Unclassified (-)

(8.95%)
untranslated (4.66%)
contraceptive IUD (0.99%)
skin care (0.86%)