

19 **Abstract**

20 **Background:** The diagnosis of rare genetic diseases is often challenging due to the
21 complexity of the genetic underpinnings of these conditions and the limited availability
22 of diagnostic tools. Machine learning (ML) algorithms have the potential to improve the
23 accuracy and speed of diagnosis by analyzing large amounts of genomic data and
24 identifying complex multiallelic patterns that may be associated with specific diseases.
25 In this systematic review, we aimed to identify the methodological trends and the ML
26 application areas in rare genetic diseases.

27 **Methods:** We performed a systematic review of the literature following the PRISMA
28 guidelines to search studies that used ML approaches to enhance the diagnosis of rare
29 genetic diseases. Studies that used DNA-based sequencing data and a variety of ML
30 algorithms were included, summarized, and analyzed using bibliometric methods,
31 visualization tools, and a feature co-occurrence analysis.

32 **Findings:** Our search identified 22 studies that met the inclusion criteria. We found that
33 exome sequencing was the most frequently used sequencing technology (59%), and rare
34 neoplastic diseases were the most prevalent disease scenario (59%). In rare neoplasms,
35 the most frequent applications of ML models were the differential diagnosis or
36 stratification of patients (38.5%) and the identification of somatic mutations (30.8%). In
37 other rare diseases, the most frequent goals were the prioritization of rare variants or
38 genes (55.5%) and the identification of biallelic or digenic inheritance (33.3%). The
39 most employed method was the random forest algorithm (54.5%). In addition, the
40 features of the datasets needed for training these algorithms were distinctive depending
41 on the goal pursued, including the mutational load in each gene for the differential
42 diagnosis of patients, or the combination of genotype features and sequence-derived
43 features (such as GC-content) for the identification of somatic mutations.

44 **Conclusions:** ML algorithms based on sequencing data are mainly used for the
45 diagnosis of rare neoplastic diseases, with random forest being the most common
46 approach. We identified key features in the datasets used for training these ML models
47 according to the objective pursued. These features can support the development of
48 future ML models in the diagnosis of rare genetic diseases.

49 **Keywords:** artificial intelligence, rare diseases, precision medicine, rare variants, DNA-
50 sequencing, genomics

51 **1. Introduction**

52 Rare diseases (RDs) continue to be a challenge to the healthcare system due to the
53 difficulty of reaching an accurate diagnosis. Although there is no uniform international
54 criteria, RDs are usually defined as those affecting fewer than 4-5 cases out of 10,000
55 individuals¹. Considering them as a whole, RDs can be regarded as a common event,
56 with 7,241 different RDs (http://www.orphadata.org/data/xml/en_product7.xml,
57 updated on June 14, 2022) with an estimated accumulated prevalence of 3.5–5.9% and
58 affecting more than 400 million people worldwide^{2,3}.

59 Most RDs appear to be caused or modified by genetic factors; up to 80% of them are
60 thought to have a genetic etiology⁴. Our current knowledge on this aspect is limited,
61 existing 3,886 RDs (53.7%) linked to, at least, a gene that cause or modify the disease
62 phenotype (http://www.orphadata.org/data/xml/en_product6.xml, updated on 14 Jun
63 22)². The improved performance and the price reduction of Next-generation sequencing
64 (NGS) technologies in recent years have made them more attractive for clinical
65 applications in RDs, increasing rapidly the number of phenotype-genotype
66 associations⁵. This has resulted in an accurate molecular diagnosis in many patients
67 suffering from monogenic RDs, which has occasionally led to personalized treatments
68 and improved disease management. Nevertheless, other patients with more complex
69 disorders receive an inconclusive genetic diagnosis, placing the diagnostic yield of
70 DNA-based NGS technologies in most studies at 40-50%^{6,7}. This is mainly caused by
71 the absence of pathogenic or likely pathogenic variants in known disease-causing genes,
72 finding instead variants of unknown significance (VUS) or variants in novel genes not
73 previously associated with the disease.

74 In this scenario of rare and complex genetic disorders where a diagnosis is not reached
75 or a prognosis is not accurate enough, more sophisticated methods should be applied to
76 analyze large-scale genomic data. The use of artificial intelligence (AI) and,
77 particularly, machine learning (ML) algorithms has raised great interest in recent years
78 due to its potential to uncover complex patterns in genomic data⁸. These ML algorithms
79 have shown the capacity to learn from and act on large, heterogeneous datasets to
80 extract new biological insights, improving the accuracy of the diagnosis of RDs⁹⁻¹².

81 Compared to previous reviews in the field of ML and RDs, such as Schaefer *et al.*⁹ or
82 Brasil *et al.*¹³, in this systematic review we used a different approach, investigating the
83 role of AI/ML algorithms in the diagnosis and prognosis of RDs using genomic data.
84 The range of options when it comes to choosing a learning algorithm or a DNA-based
85 NGS technique to address RDs is highly variable. On the one hand, ML methods are
86 usually divided into two main categories: supervised and unsupervised learning.
87 Supervised ML algorithms require labeled data to solve mainly regression and
88 classification tasks, whereas unsupervised ML algorithms address classification tasks
89 based on unlabeled data by seeking common patterns. The review from Libbrecht *et al.*
90 describes these algorithms in more detail and provides examples applied to genomic
91 data¹⁴. On the other hand, regarding NGS techniques, there are mainly two strategies: a)
92 to sequence the entirety of the DNA sequence (whole genome sequencing, WGS), or b)
93 to just sequence some regions of the DNA, such as coding regions (exome sequencing,
94 ES), or certain disease-causing genes (gene panel). Nevertheless, the raw data generated
95 in these experiments can be processed in many ways, with different workflows
96 depending on the aim of the study.

97 This systematic review presents a thorough overview of the existing evidence on the
98 application of AI/ML algorithms to the diagnosis of RDs using DNA-based sequencing

99 data. We conducted a comprehensive search of the literature and included studies that
100 used a variety of ML approaches and sequencing data sources in different research
101 settings. Our analysis focused on the evaluation of trends in the field, the ability of these
102 approaches to identify genetic variations associated with RDs, and the potential of
103 AI/ML to improve their diagnosis.

104 **2. Methods**

105 **2.1. Systematic literature search and data sources**

106 We performed a literature search using PubMed, Web of Science, and Scopus to
107 identify relevant publications on the use of AI/ML for the diagnosis and prognosis of
108 RDs using genomic data. We also used citation and hand searching to ensure that
109 potentially relevant studies were retrieved. The Preferred Reporting Items for
110 Systematic Reviews and Meta-Analyses (PRISMA) guidelines were followed to design
111 and perform this systematic review¹⁵, and its protocol was registered in PROSPERO
112 (registration number CRD42022360247).

113 A search in the selected databases using the search terms ‘rare AND ("artificial
114 intelligence" OR "machine learning" OR "deep learning") AND ((exome OR genome
115 OR panel) AND sequencing)’ and considering publications from 2012 onward resulted
116 in 296 abstracts. The citation and hand searching resulted in 10 additional records. The
117 date of the last search was September 29, 2022.

118 The list of abstracts was screened for inclusion using the following inclusion criteria: (i)
119 an application of AI/ML methods; (ii) a diagnostic or prognosis application using a
120 DNA-based NGS technique (panel, exome, or genome sequencing); and (iii) an
121 application to a RD within the orpha.net database. Non-English articles, review articles,
122 duplicate records, and studies not relevant to any RD or AI/ML were excluded. To
123 narrow our focus to clinical applications, we excluded animal studies as well as
124 publications that only reported methodological aspects of AI/ML without presenting
125 clinical data from the study population. For all articles considered relevant, the full text
126 was reviewed using the same screening procedure as in the first stage.

127 **2.2. Data extraction**

128 All the selected articles were evaluated to gather data on five main aspects: i) study
129 characteristics and study population (subjects included, RDs studied, study design, use
130 of secondary data), ii) characteristics of the applied AI/ML techniques (selected ML
131 model, programming languages used, input data, associated features, feature selection
132 methods, model evaluation), iii) information about the DNA-based NGS technology
133 used (type, sample collected, DNA sequencing kit, sequencing platform, read length,
134 mean coverage), iv) the variant discovery approach (alignment method, used
135 SNV/Indel/CNVs callers, variant annotation software, variant filtering criteria), and v)
136 authors (number of authors and institutions involved, authors' countries) and journal
137 details (name, category, journal impact factor, journal citation indicator).

138 **2.3. Data analysis**

139 The data collected from selected articles were summarized and analyzed using a variety
140 of approaches. Journal Impact Factor (JIF) scores were obtained from the Journal
141 Citation Report (JCR) database. Bibliometric networks, including data from authors and
142 abstracts, were constructed and visualized using VOSviewer¹⁶. Similarly, full-text
143 articles were analyzed using WordStat 9.0 (Provalis Research, Montreal, Quebec,
144 Canada) to extract main topics and keywords.

145 Selected articles were divided into “rare neoplastic diseases” and “other rare diseases”
146 to enable comparisons. AI/ML models were categorized into three categories:
147 supervised, unsupervised, and deep learning models. Input variables that the model uses
148 to make predictions (features) were classified in 1) “clinical features”, which include
149 information about patients' clinical characteristics; 2) “phenotype-related features”,
150 including data about the association between genes and phenotypes (e.g., Human

151 Phenotype Ontology); 3) “read alignment features”, which include the properties related
152 to read mapping and sequencing quality; 4) “genotype-related features”, including
153 details of variants found in patients (e.g., variant allele frequency, count of variants in a
154 certain gene, length of indel); 5) “sequence region and structural features”, including
155 information about the region where the variant is located (e.g., gene size, GC content);
156 6) “network features”, which include details about the pathways in which a particular
157 gene is involved (e.g., number of pathway, network neighbors); 7)
158 “evolutionary/pathogenicity features”, which include pathogenicity and evolutionary
159 conservation scores of variants (e.g., CADD, PolyPhen-2); 8) “gene expression
160 features”, including data on gene expression; 9) “tissue-specific features”, including
161 features which are specific for certain types of tissues; and 10) “disease-specific
162 features”, including features which are specific for certain types of diseases. The co-
163 occurrence of these features in the datasets used for training AI/ML models was
164 examined and plotted using UpSetR¹⁷.

165 **3. Results**

166 **3.1. Included studies**

167 The literature search in databases identified 494 studies, with 296 remaining after
168 removing duplicates (**Supplementary Table 1**). Among them, 93 studies were selected
169 for full-text review, and 14 were included in the final analysis. In addition, 11 studies
170 were identified through hand and citation searching. After screening, 8 further studies
171 met the selection criteria of this systematic review. Thus, 22 studies were included in
172 the final analysis (**Supplementary Table 2**). **Figure 1** shows the PRISMA flow
173 diagram for article selection, including the reasons for excluding records.

174 **3.2. Temporal trends and bibliometrics**

175 To assess the temporal trends in the use of AI/ML methods for the diagnosis and
176 prognosis of RDs using sequencing data, meta-data from included articles was retrieved
177 (**Supplementary Table 3**). In recent years, we noticed a relative rise in the number of
178 studies that address this challenge using AI/ML (**Figure 2A**). Most of these articles
179 were published in journals belonging to the first quartile (90.9%) and within the
180 “Genetics & Hereditary” JCR category (31.8%) (**Supplementary Figure 1**). It should
181 be noted that the count for 2022 is based on studies published up to September 29,
182 2022.

183 A total of 318 authors contributed to the selected articles. The bibliometric analysis
184 showed a low level of collaboration between authors of different articles, creating 19
185 clusters where only 3 authors participated in 2 or more articles (**Supplementary Figure**
186 **2A**). The term co-occurrence analysis of abstracts found 100 relevant terms divided into
187 3 clusters that summarize the main topics of this research field. These clusters group
188 together terms mainly associated with genetics (cluster 1), cancer (cluster 2), and
189 methodology terms (cluster 3) (**Figure 2B**). The most frequently occurring terms in

190 these abstracts were “genetics” (18 occurrences), “machine learning” (15 occurrences)
191 and “whole-exome sequencing” (10 occurrences). These key terms were also among the
192 most frequently used terms in the analysis of full-text articles, where terms such as
193 “random forest” (59.1% of studies), “somatic mutations” (54.5% of studies), or “rare
194 variants” (54.5% of studies) were also in a significant proportion of studies
195 (**Supplementary Figure 2B**).

196 **3.3. Application areas for AI/ML techniques**

197 The most common disease scenario was rare neoplastic diseases (59%). The remaining
198 studies investigated different kinds of RDs, such as developmental, neurological, or
199 circulatory diseases (**Figure 3A**). Exome sequencing was the most used NGS method in
200 both rare neoplastic diseases (61.5%) and other RDs (55.5%) (**Figure 3B**). Of note,
201 63.6% (14/22) of the studies employed sequencing data stored in external databases,
202 primarily The Cancer Genome Atlas (TCGA), but also the Myocardial Genetics
203 Consortium (MIGEN), or the Undiagnosed Diseases Network (UDN). These studies
204 showed larger sample sizes than those using their own cohorts (**Supplementary Figure**
205 **3**), but they also showed higher intra-method variability, as seen by the mixed sample
206 processing methods they employed (**Supplementary Figure 4**). **Supplementary Table**
207 **4** summarizes the NGS-related and sequencing data processing methods in detail.

208 Supervised machine learning methods were chosen in 86.3% of the studies, with
209 Random Forest (RF) being the most employed algorithm within this group (54.5%)
210 (**Figure 3C**). One study discarded the genetic features after the feature selection
211 process, and three studies did not describe the selected features in detail, one of which
212 was due to a commercial interest (**Supplementary Table 5**).

213 **3.4. AI/ML in the study of rare genetic diseases**

214 The objectives of AI/ML approaches in the different studies were investigated. It was
215 found that the primary goal of using AI/ML in rare neoplastic diseases was the
216 differential diagnosis of patients (5/13), followed by the identification of somatic
217 mutations when a matched normal tissue was not available (4/13). In contrast, the major
218 goals in other RDs were to prioritize variants and candidate genes (5/9) and to identify
219 biallelic or digenic inheritance (3/9). To date, the use of AI/ML for the differential
220 diagnosis of patients with non-neoplastic diseases is uncommon (1/9) (**Figure 4A**).

221 Looking at the types of instances (labels) and features (attributes) of datasets used for
222 training these AI/ML models, we found that they were distinctive and different
223 depending on the goal pursued (**Table 1** and **Figure 4B**). For the differential diagnosis
224 of patients, most datasets included only features related to the genotype of patients.
225 These features primarily contained mutational load data for each gene or genomic
226 window using collapsing methods. Models trained to predict the prognosis of RDs
227 included clinical features (e.g., sex, age, exposure to certain substances) in addition to
228 genotype features. The four AI/ML models aimed at finding possible pathogenic
229 combinations of genes (digenic) or variants (biallelic) shared the usage of features
230 related to biological networks or pathways (e.g., the associated pathway of each gene in
231 KEGG or Reactome, network neighbors). Datasets focused on training models for
232 variant or gene prioritization were distinguished by using features linked to predictors
233 of variant pathogenicity at protein level and conservation across the genome of different
234 species. Finally, for the identification of somatic mutations without a matched normal
235 sample, the AI/ML models combined genotype features (e.g., variant allele frequency)
236 with characteristics of the genome region where the variant is located (e.g., GC-content)
237 or sequencing and mapping quality scores (e.g., coverage). **Supplementary Table 6**
238 contains further information regarding the types of features mentioned above.

239 **3.5. Data and code access for reproducibility**

240 When it comes to studies that define ML models, reproducibility is a key factor. Of the
241 selected articles, 16 studies (72.7%) provided access to the data used during the
242 analysis; 3 studies did so only upon data request; and 3 did not explicitly declare in the
243 text that data were available, one of which was due to commercial confidentiality. In
244 terms of the code of AI/ML models, 16 studies (72.7%) had made it publicly available.
245 With respect to the variant discovery approaches, all studies specified the software used
246 for sequence alignment; 21 studies (95.5%) included information about the variant
247 calling step; 17 studies (77.3%) did not mention the use of copy number variations
248 (CNVs) during the analysis, and 3 studies did not state how the variants were annotated.
249 **Supplementary Table 7** summarizes data availability and reproducibility information.

250 **4. Discussion**

251 AI/ML involve the use of algorithms to process and gain insights from data with the aim
252 of making predictions or decisions that can be applied to a wide range of fields,
253 including healthcare and genetics. In this systematic review, we have evaluated the
254 latest developments in AI/ML when it comes to rare genetic conditions and examined
255 the ways in which the use of DNA sequencing data can improve their diagnosis. In
256 addition, we have identified some challenges and opportunities for future research in
257 this area.

258 **4.1. Exome sequencing and rare neoplastic diseases as main topics**

259 Although to a lesser extent than in other types of diagnostic methods, such as medical
260 imaging, AI/ML are increasingly being used in the field of RDs^{9,18,19}. This trend was
261 also found when focusing only on those studies that use DNA sequencing data to
262 improve the diagnostic process. Through the bibliometric study carried out in this
263 review, and the subsequent manual analyses, we found that exome sequencing was the
264 most prevalent sequencing approach in the field, and that rare neoplastic diseases were
265 the most prevalent clinical scenario. Exome sequencing continues to be a good starting
266 point for the genetic diagnosis of RDs, as it provides a cost-effective and efficient way
267 to identify disease-causing variants²⁰. However, depending on the specific rare disease
268 context, genome sequencing may be necessary to provide a complete diagnosis,
269 including the analysis of non-coding variations, CNVs, or chromosomal
270 rearrangements^{21,22}.

271 Rare neoplastic diseases generally have a worse diagnosis and higher funding
272 opportunities than other RDs, making them the type of rare disease in which AI/ML are
273 used the most^{19,23}. This is also due to the existence of public databases such as TCGA,

274 which allow researchers to access a large amount of genomic data and use AI/ML
275 techniques to identify patterns and make predictions²⁴. When we analyzed the data on
276 which these AI/ML models were trained, we saw that many of them (63.6%) were based
277 on sequencing data from external sources, such as TCGA. These studies showed larger
278 sample sizes, but also a greater diversity in sequencing technology characteristics such
279 as read depth, different length of reads or different sequencing kits and platforms.
280 Mixing sequencing data from different technologies, qualities, and batches can have
281 several biases that can influence the variant calling results, affecting in turn the results
282 of downstream analyses, and making difficult to draw accurate conclusions from the
283 data²⁵. The precision when taking clinical decisions must be maximized, so these
284 studies must have control over these factors²⁶. Different studies have shown how to
285 approach this process^{25,27}.

286 **4.2. AI/ML algorithms and feature selection in genetic studies**

287 Most of the methods utilized in the selected studies fall into the category of supervised
288 learning (86.7%), with RF being the most common algorithm among them (73.7%). RF
289 algorithm offers a combination of properties that makes it one of the most widely used
290 and suitable algorithms for the study of genetic variants^{28,29}. RF combines multiple
291 decision trees (forest) that can handle high-dimensional data, capturing interactions and
292 complex relationships between features by creating random subsets of both, data and
293 features, at each tree. In addition, RF also allows to compute feature importances, which
294 can be used to identify the most relevant features for the prediction task, providing
295 interpretable models³⁰. All this makes RF well suited for complex genetic problems and
296 explains its popularity among the genetic studies.

297 The structure of the dataset is a fundamental and key aspect of any AI/ML model, as it
298 is the data that the model uses to learn and make predictions. The processes of feature
299 selection and feature engineering can have a substantial effect on the performance of the
300 model; hence, it is essential that the final features possess relevance to the problem at
301 hand³¹. In this systematic review, we have identified the features used by each of the
302 selected studies and found that these features were specific to each of the objectives
303 pursued. This insight can be valuable in understanding the current state of research in
304 the field, and it can serve as a starting point for creating new datasets in future studies.

305 The results suggested that collapsing or burden methods seem to be crucial for setting
306 up the features of datasets used to train models for the stratification or differential
307 diagnosis of patients. These methods divide the genome into portions (bins or genes)
308 and summarize the information contained in these segments into a burden value, which
309 can be calculated in different ways^{32,33}. This approach has shown its usefulness in
310 finding candidate genes in different complex RDs with both genome and exome
311 sequencing data³⁴⁻³⁶. Thus, applied to AI/ML tasks, this process helps to decrease the
312 dimensionality of datasets based on genetic variants by grouping them into one value
313 per gene or bin, which helps to reduce the curse of dimensionality and improve
314 interpretability³⁷.

315 On the other hand, models focused on predicting patient prognosis integrate clinical
316 and genomic data to obtain a more complete picture of the patient and assess the risk of
317 disease progression. Previous studies, particularly in cancer, have shown how this
318 integration of data provides a more comprehensive and accurate assessment of patient
319 outcome^{38,39}. Alternatively, models aimed at predicting possible pathogenic
320 combinations of genes use features that summarize the association of these genes with
321 the biological pathways in which they participate. The use of these features is supported

322 by the fact that digenic diseases are usually caused by variants in genes that are
323 functionally related and have a common pathway^{40,41}.

324 **4.3. Future challenges**

325 From the results of this review, we identified some challenges that need to be addressed
326 in future studies. When we analyzed the type of genomic data used to train the AI/ML
327 models reviewed, we realized that most of them (77.3%) were based exclusively on
328 single nucleotide variants or short indels, not including the analysis of CNVs. CNVs are
329 a significant source of genetic diversity in humans that has remained understudied due
330 to the difficulty of detection. However, today there are different algorithms for CNV
331 detection that simplify the task considerably, as well as guidelines that help us to
332 interpret them^{42,43}. This allows the possibility of evaluating its effect on the
333 pathogenesis and outcome of RD. On the other hand, when we examine the goals
334 pursued in the analysis of neoplastic RDs, we can see that the differential diagnosis or
335 stratification of patients stands out above the other objectives. This is totally different in
336 other RDs, where, in fact, this objective is the least pursued of the 3 objectives
337 identified, and, therefore, a field where the contribution of genetic variation to the
338 phenotype is not well understood. The use of AI/ML algorithms on rare disease
339 sequencing data can support the identification of novel genetic interactions, uncovering
340 patterns and relationships that may not be immediately apparent and providing a better
341 understanding of the regulatory mechanisms mediated by these variants in the
342 phenotype. The use of unsupervised methods would be a possible first approach to
343 achieve the objective of identifying clusters of patients according to their genetic
344 background⁴⁴.

345 **4.4. Limitations**

346 Our review is limited by the design of the systematic search and the exclusion of purely
347 methodological articles, focusing only on those studies with clinical applications.
348 Because of the limited number of studies available on the topic, and although it has
349 been studied, articles have not been discarded because of the quality of the journal in
350 which they were published (i.e., JIF), and this may have influenced, in some way, the
351 results of this review. In addition, to reduce variability in study methodology and
352 facilitate the analysis, we have only focused on those studies using DNA-based
353 sequencing, not including other NGS methodologies such as RNA-seq, which are
354 widely used in conjunction with AI/ML methodologies⁴⁵⁻⁴⁷.

355 **5. Conclusions**

356 We have conducted a systematic review of ML algorithms to the diagnosis of RDs
357 using DNA-based sequencing data, providing an overview of the current state of the
358 field and the potential of these methods to improve diagnostic accuracy. Exome
359 sequencing is the most widely used sequencing technology and rare neoplastic diseases
360 are the most common disease scenario. On the other hand, the goals of AI/ML
361 algorithms in RDs using sequencing data are broad, ranging from patient stratification
362 to the identification of possible pathogenic combinations of variants. However, we
363 found common patterns in these goals when configuring the datasets with which these
364 models are trained, identifying key features for each of the objectives. Finally, we
365 identified possible future challenges, such as the use of CNV to train the AI/ML
366 models, or the application of AI/ML for the stratification of patients with non-neoplastic
367 RDs. Thus, this systematic review can be used as a reference for further studies,
368 supporting the development of future ML models in the diagnosis of rare genetic
369 diseases

370 **6. Fundings**

371 JALE has received funds from Instituto de Salud Carlos III (Grant# PI20-1126),
372 CIBERER (Grant# PIT21_GCV21), Andalusian University, Research and Innovation
373 Department (PY20-00303, EPIMEN), Andalusian Health Department (Grant# PI027-
374 2020), Asociación Síndrome de Meniere España (ASMES) and Meniere's Society, UK.
375 PRNV is supported by PY20-00303 Grant (EPIMEN). AMPP is a PhD student in the
376 Biomedicine Program at Universidad de Granada and his salary was supported by
377 Andalusian University, Research and Innovation Department (Grant#
378 PREDOC2021/00343).

379 **7. Declaration of Competing Interest**

380 The authors declare that they have no competing interests. The research was conducted
381 independent of any commercial or financial relationships that could be construed as a
382 potential conflict of interest.

383 8. References

- 384 1. Richter T, Nestler-Parr S, Babela R, et al. Rare Disease Terminology and Definitions—A
385 Systematic Global Review: Report of the ISPOR Rare Disease Special Interest Group. *Value*
386 *Health*. 2015;18(6):906-914. doi:10.1016/j.jval.2015.05.008
- 387 2. Orphadata: Free access data from Orphanet. © INSERM 1999. Available on
388 <http://www.orphadata.org>. Data version (XML data version).
- 389 3. Nguengang Wakap S, Lambert DM, Olry A, et al. Estimating cumulative point prevalence of
390 rare diseases: analysis of the Orphanet database. *Eur J Hum Genet*. 2020;28(2):165-173.
391 doi:10.1038/s41431-019-0508-0
- 392 4. 100,000 Genomes Project Pilot Investigators, Smedley D, Smith KR, et al. 100,000
393 Genomes Pilot on Rare-Disease Diagnosis in Health Care - Preliminary Report. *N Engl J*
394 *Med*. 2021;385(20):1868-1880. doi:10.1056/NEJMoa2035790
- 395 5. Wise AL, Manolio TA, Mensah GA, et al. Genomic Medicine for Undiagnosed Diseases.
396 *Lancet Lond Engl*. 2019;394(10197):533-540. doi:10.1016/S0140-6736(19)31274-7
- 397 6. Vinkšl M, Writzl K, Maver A, Peterlin B. Improving diagnostics of rare genetic diseases
398 with NGS approaches. *J Community Genet*. 2021;12(2):247-256. doi:10.1007/s12687-020-
399 00500-5
- 400 7. Dai P, Honda A, Ewans L, et al. Recommendations for next generation sequencing data
401 reanalysis of unsolved cases with suspected Mendelian disorders: A systematic review and
402 meta-analysis. *Genet Med*. Published online May 14, 2022. doi:10.1016/j.gim.2022.04.021
- 403 8. Routhier E, Mozziconacci J. Genomics enters the deep learning era. *PeerJ*.
404 2022;10:e13613. doi:10.7717/peerj.13613
- 405 9. Schaefer J, Lehne M, Schepers J, Prasser F, Thun S. The use of machine learning in rare
406 diseases: a scoping review. *Orphanet J Rare Dis*. 2020;15(1):145. doi:10.1186/s13023-020-
407 01424-6
- 408 10. Setty ST, Scott-Boyer MP, Cuppens T, Droit A. New Developments and Possibilities in
409 Reanalysis and Reinterpretation of Whole Exome Sequencing Datasets for Unsolved Rare
410 Diseases Using Machine Learning Approaches. *Int J Mol Sci*. 2022;23(12):6792.
411 doi:10.3390/ijms23126792
- 412 11. Cohen ASA, Farrow EG, Abdelmoity AT, et al. Genomic answers for children: Dynamic
413 analyses of >1000 pediatric rare disease genomes. *Genet Med Off J Am Coll Med Genet*.
414 2022;24(6):1336-1348. doi:10.1016/j.gim.2022.02.007
- 415 12. Okazaki A, Ott J. Machine learning approaches to explore digenic inheritance. *Trends*
416 *Genet TIG*. Published online May 14, 2022;S0168-9525(22)00105-6.
417 doi:10.1016/j.tig.2022.04.009
- 418 13. Brasil S, Pascoal C, Francisco R, dos Reis Ferreira V, A. Videira P, Valadão G. Artificial
419 Intelligence (AI) in Rare Diseases: Is the Future Brighter? *Genes*. 2019;10(12):978.
420 doi:10.3390/genes10120978

- 421 14. Libbrecht MW, Noble WS. Machine learning applications in genetics and genomics. *Nat*
422 *Rev Genet.* 2015;16(6):321-332. doi:10.1038/nrg3920
- 423 15. Page MJ, McKenzie JE, Bossuyt PM, et al. The PRISMA 2020 statement: An updated
424 guideline for reporting systematic reviews. *PLOS Med.* 2021;18(3):e1003583.
425 doi:10.1371/journal.pmed.1003583
- 426 16. van Eck NJ, Waltman L. Software survey: VOSviewer, a computer program for bibliometric
427 mapping. *Scientometrics.* 2010;84(2):523-538. doi:10.1007/s11192-009-0146-3
- 428 17. Lex A, Gehlenborg N, Strobel H, Vuillemot R, Pfister H. UpSet: Visualization of Intersecting
429 Sets. *IEEE Trans Vis Comput Graph.* 2014;20(12):1983-1992.
430 doi:10.1109/TVCG.2014.2346248
- 431 18. Oren O, Gersh BJ, Bhatt DL. Artificial intelligence in medical imaging: switching from
432 radiographic pathological data to clinically meaningful endpoints. *Lancet Digit Health.*
433 2020;2(9):e486-e488. doi:10.1016/S2589-7500(20)30160-6
- 434 19. Lee J, Liu C, Kim J, et al. Deep learning for rare disease: A scoping review. *J Biomed Inform.*
435 2022;135:104227. doi:10.1016/j.jbi.2022.104227
- 436 20. Klau J, Abou Jamra R, Radtke M, et al. Exome first approach to reduce diagnostic costs and
437 time – retrospective analysis of 111 individuals with rare neurodevelopmental disorders.
438 *Eur J Hum Genet.* 2022;30(1):117-125. doi:10.1038/s41431-021-00981-z
- 439 21. Marwaha S, Knowles JW, Ashley EA. A guide for the diagnosis of rare and undiagnosed
440 disease: beyond the exome. *Genome Med.* 2022;14(1):23. doi:10.1186/s13073-022-01026-
441 w
- 442 22. Souche E, Beltran S, Brosens E, et al. Recommendations for whole genome sequencing in
443 diagnostics for rare diseases. *Eur J Hum Genet.* 2022;30(9):1017-1021.
444 doi:10.1038/s41431-022-01113-x
- 445 23. Dlamini Z, Francies FZ, Hull R, Marima R. Artificial intelligence (AI) and big data in cancer
446 and precision oncology. *Comput Struct Biotechnol J.* 2020;18:2300-2311.
447 doi:10.1016/j.csbj.2020.08.019
- 448 24. Weinstein JN, Collisson EA, Mills GB, et al. The Cancer Genome Atlas Pan-Cancer Analysis
449 Project. *Nat Genet.* 2013;45(10):1113-1120. doi:10.1038/ng.2764
- 450 25. De-Kayne R, Frei D, Greenway R, Mendes SL, Retel C, Feulner PGD. Sequencing platform
451 shifts provide opportunities but pose challenges for combining genomic data sets. *Mol*
452 *Ecol Resour.* 2021;21(3):653-660. doi:10.1111/1755-0998.13309
- 453 26. Goldfeder RL, Priest JR, Zook JM, et al. Medical implications of technical accuracy in
454 genome sequencing. *Genome Med.* 2016;8(1):24. doi:10.1186/s13073-016-0269-0
- 455 27. Ellrott K, Bailey MH, Saksena G, et al. Scalable Open Science Approach for Mutation Calling
456 of Tumor Exomes Using Multiple Genomic Pipelines. *Cell Syst.* 2018;6(3):271-281.e7.
457 doi:10.1016/j.cels.2018.03.002
- 458 28. Goldstein BA, Polley EC, Briggs FBS. Random Forests for Genetic Association Studies. *Stat*
459 *Appl Genet Mol Biol.* 2011;10(1):32. doi:10.2202/1544-6115.1691

- 460 29. Chen X, Ishwaran H. Random Forests for Genomic Data Analysis. *Genomics*.
461 2012;99(6):323-329. doi:10.1016/j.ygeno.2012.04.003
- 462 30. Breiman L. Random Forests. *Mach Learn*. 2001;45(1):5-32. doi:10.1023/A:1010933404324
- 463 31. Pudjihartono N, Fadason T, Kempa-Liehr AW, O'Sullivan JM. A Review of Feature Selection
464 Methods for Machine Learning-Based Disease Risk Prediction. *Front Bioinforma*. 2022;2.
465 Accessed January 21, 2023.
466 <https://www.frontiersin.org/articles/10.3389/fbinf.2022.927312>
- 467 32. Dering C, König IR, Ramsey LB, Relling MV, Yang W, Ziegler A. A comprehensive evaluation of
468 collapsing methods using simulated and real data: excellent annotation of functionality
469 and large sample sizes required. *Front Genet*. 2014;5. Accessed January 21, 2023.
470 <https://www.frontiersin.org/articles/10.3389/fgene.2014.00323>
- 471 33. Nicolae DL. Association Tests for Rare Variants. *Annu Rev Genomics Hum Genet*.
472 2016;17(1):117-130. doi:10.1146/annurev-genom-083115-022609
- 473 34. Roman-Naranjo P, Gallego-Martinez A, Soto-Varela A, et al. Burden of Rare Variants in the
474 OTOG Gene in Familial Meniere's Disease. *Ear Hear*. 2020;41(6):1598-1605.
475 doi:10.1097/AUD.0000000000000878
- 476 35. Dillioott AA, Abdelhady A, Sunderland KM, et al. Contribution of rare variant associations to
477 neurodegenerative disease presentation. *NPJ Genomic Med*. 2021;6:80.
478 doi:10.1038/s41525-021-00243-3
- 479 36. Lin J, Li C, Cui Y, et al. Rare variants in IMPDH2 cause autosomal dominant dystonia in
480 Chinese population. *J Neurol*. Published online January 17, 2023. doi:10.1007/s00415-023-
481 11564-x
- 482 37. Altman N, Krzywinski M. The curse(s) of dimensionality. *Nat Methods*. 2018;15(6):399-400.
483 doi:10.1038/s41592-018-0019-x
- 484 38. Lobato-Delgado B, Priego-Torres B, Sanchez-Morillo D. Combining Molecular, Imaging, and
485 Clinical Data Analysis for Predicting Cancer Prognosis. *Cancers*. 2022;14(13):3215.
486 doi:10.3390/cancers14133215
- 487 39. Gonzalez-Bosquet J, Gabrilovich S, McDonald ME, et al. Integration of Genomic and Clinical
488 Retrospective Data to Predict Endometrioid Endometrial Cancer Recurrence. *Int J Mol Sci*.
489 2022;23(24):16014. doi:10.3390/ijms232416014
- 490 40. Gazzo A, Raimondi D, Daneels D, et al. Understanding mutational effects in digenic
491 diseases. *Nucleic Acids Res*. 2017;45(15):e140. doi:10.1093/nar/gkx557
- 492 41. Schäffer AA. Digenic inheritance in medical genetics. *J Med Genet*. 2013;50(10):641-652.
493 doi:10.1136/jmedgenet-2013-101713
- 494 42. Riggs ER, Andersen EF, Cherry AM, et al. Technical standards for the interpretation and
495 reporting of constitutional copy-number variants: a joint consensus recommendation of
496 the American College of Medical Genetics and Genomics (ACMG) and the Clinical Genome
497 Resource (ClinGen). *Genet Med*. 2020;22(2):245-257. doi:10.1038/s41436-019-0686-8

- 498 43. Gordeeva V, Sharova E, Arapidi G. Progress in Methods for Copy Number Variation
499 Profiling. *Int J Mol Sci*. 2022;23(4):2143. doi:10.3390/ijms23042143
- 500 44. Basile AO, Ritchie MD. Informatics and Machine Learning to Define the Phenotype. *Expert*
501 *Rev Mol Diagn*. 2018;18(3):219-226. doi:10.1080/14737159.2018.1439380
- 502 45. Wang L, Xi Y, Sung S, Qiao H. RNA-seq assistant: machine learning based methods to
503 identify more transcriptional regulated genes. *BMC Genomics*. 2018;19(1):546.
504 doi:10.1186/s12864-018-4932-2
- 505 46. Gunavathi C, Sivasubramanian K, Keerthika P, Paramasivam C. A review on convolutional
506 neural network based deep learning methods in gene expression data for disease
507 diagnosis. *Mater Today Proc*. 2021;45:2282-2285. doi:10.1016/j.matpr.2020.10.263
- 508 47. Figgett WA, Monaghan K, Ng M, et al. Machine learning applied to whole-blood RNA-
509 sequencing data uncovers distinct subsets of patients with systemic lupus erythematosus.
510 *Clin Transl Immunol*. 2019;8(12):e01093. doi:10.1002/cti2.1093
- 511 48. Parida L, Haferlach C, Rhrissorrakrai K, et al. Dark-matter matters: Discriminating subtle
512 blood cancers using the darkest DNA. *PLoS Comput Biol*. 2019;15(8):e1007332.
513 doi:10.1371/journal.pcbi.1007332
- 514 49. Parvande S, Donehower LA, Panagiotis K, et al. EPIMUTESTR: a nearest neighbor machine
515 learning approach to predict cancer driver genes from the evolutionary action of coding
516 variants. *Nucleic Acids Res*. 2022;50(12):e70. doi:10.1093/nar/gkac215
- 517 50. Peneder P, Stütz AM, Surdez D, et al. Multimodal analysis of cell-free DNA whole-genome
518 sequencing for pediatric cancers with low mutational burden. *Nat Commun*.
519 2021;12(1):3230. doi:10.1038/s41467-021-23445-w
- 520 51. Li Y, Luo Y. Performance-weighted-voting model: An ensemble machine learning method
521 for cancer type classification using whole-exome sequencing mutation. *Quant Biol Beijing*
522 *China*. 2020;8(4):347-358. doi:10.1007/s40484-020-0226-1
- 523 52. Aguiar-Pulido V, Wolujewicz P, Martinez-Fundichely A, et al. Systems biology analysis of
524 human genomes points to key pathways conferring spina bifida risk. *Proc Natl Acad Sci U S*
525 *A*. 2021;118(51):e2106844118. doi:10.1073/pnas.2106844118
- 526 53. Chaix MA, Parmar N, Kinnear C, et al. Machine Learning Identifies Clinical and Genetic
527 Factors Associated With Anthracycline Cardiotoxicity in Pediatric Cancer Survivors. *JACC*
528 *CardioOncology*. 2020;2(5):690-706. doi:10.1016/j.jacc.2020.11.004
- 529 54. Zauderer MG, Martin A, Egger J, et al. The use of a next-generation sequencing-derived
530 machine-learning risk-prediction model (OncoCast-MPM) for malignant pleural
531 mesothelioma: a retrospective study. *Lancet Digit Health*. 2021;3(9):e565-e576.
532 doi:10.1016/S2589-7500(21)00104-7
- 533 55. Carter H, Douville C, Stenson PD, Cooper DN, Karchin R. Identifying Mendelian disease
534 genes with the variant effect scoring tool. *BMC Genomics*. 2013;14 Suppl 3:S3.
535 doi:10.1186/1471-2164-14-S3-S3

- 536 56. Vitsios D, Petrovski S. Mantis-ml: Disease-Agnostic Gene Prioritization from High-
537 Throughput Genomic Screens by Stochastic Semi-supervised Learning. *Am J Hum Genet.*
538 2020;106(5):659-678. doi:10.1016/j.ajhg.2020.03.012
- 539 57. Majithia AR, Tsuda B, Agostini M, et al. Prospective functional classification of all possible
540 missense variants in PPARG. *Nat Genet.* 2016;48(12):1570-1575. doi:10.1038/ng.3700
- 541 58. Carss KJ, Baranowska AA, Armisen J, et al. Spontaneous Coronary Artery Dissection:
542 Insights on Rare Genetic Variation From Genome Sequencing. *Circ Genomic Precis Med.*
543 2020;13(6):e003030. doi:10.1161/CIRCGEN.120.003030
- 544 59. Davis NA, Lareau CA, White BC, et al. Encore: Genetic Association Interaction Network
545 centrality pipeline and application to SLE exome data. *Genet Epidemiol.* 2013;37(6):614-
546 621. doi:10.1002/gepi.21739
- 547 60. Mukherjee S, Cogan JD, Newman JH, et al. Identifying digenic disease genes via machine
548 learning in the Undiagnosed Diseases Network. *Am J Hum Genet.* 2021;108(10):1946-1963.
549 doi:10.1016/j.ajhg.2021.08.010
- 550 61. Laan M, Kasak L, Timinskas K, et al. NR5A1 c.991-1G > C splice-site variant causes familial
551 46,XY partial gonadal dysgenesis with incomplete penetrance. *Clin Endocrinol (Oxf).*
552 2021;94(4):656-666. doi:10.1111/cen.14381
- 553 62. Ainscough BJ, Barnell EK, Ronning P, et al. A deep learning approach to automate
554 refinement of somatic variant calling from cancer sequencing data. *Nat Genet.*
555 2018;50(12):1735-1743. doi:10.1038/s41588-018-0257-y
- 556 63. Wood DE, White JR, Georgiadis A, et al. A machine learning approach for somatic mutation
557 discovery. *Sci Transl Med.* 2018;10(457):eaar7939. doi:10.1126/scitranslmed.aar7939
- 558 64. Kalatskaya I, Trinh QM, Spears M, McPherson JD, Bartlett JMS, Stein L. ISOWN: accurate
559 somatic mutation identification in the absence of normal tissue controls. *Genome Med.*
560 2017;9(1):59. doi:10.1186/s13073-017-0446-9
- 561

562 **Figures**

563 **Figure 1:** PRISMA flow diagram for the identification, screening and selection of
564 genetic studies using AI/ML for the diagnosis of rare diseases.

565 **Figure 2:** Visualization of temporal trends and bibliometrics. Panel **A)** shows the
566 selected studies distributed per year and divided into deciles (**D**) according to journal
567 impact factors (**JIF**). Panel **B)** displays a keyword co-occurrence network using abstracts
568 of selected studies.

569 **Figure 3:** Methods and areas of application. Panel **A)** displays the distribution of rare
570 disease identified in selected studies. Panel **B)** shows the next-generation sequencing
571 (NGS) methods used in studies targeting rare neoplastic diseases and other rare
572 diseases. Panel **C)** summarizes the types of machine learning algorithms applied in
573 selected studies.

574 Footer: KNN: K-Nearest Neighbors; LDA: Linear Discriminant Analysis; FNN:
575 Feedforward neural network; CNN: Convolutional Neural Networks.

576 **Figure 4:** Objectives and settings of AI/ML models. Panel **A)** displays the goals of
577 AI/ML models in rare neoplastic diseases (blue) and other rare diseases (orange). Panel
578 **B)** contains an upset plot showing the different combinations of features in the training
579 datasets of AI/ML models depending on the objective pursued.

580 **Table**

581 **Table 1:** Distinctive features identified in the datasets used for training the ML models of included studies, based on the specific goal being

Objective AI/ML algorithm	Type of instances	Distinctive dataset feature/s	Example of feature	Use cases (ref)
Stratification/Differential diagnosis	Patients	Genotype features	Burden value	[48–52]
Prognosis of patients	Patients	Genotype + Clinical features	Burden value + age	[53,54]
Variant/Gene prioritization	Genes/variants	Pathogenicity features	CADD score	[55–58]
Identification of digenic/biallelic combinations	Pair of genes/variants	Network features	Number of pathways shared	[59–61]
Identification of somatic mutations	Variants	Genotype + Sequence features	VAF + GC-content	[62–64]

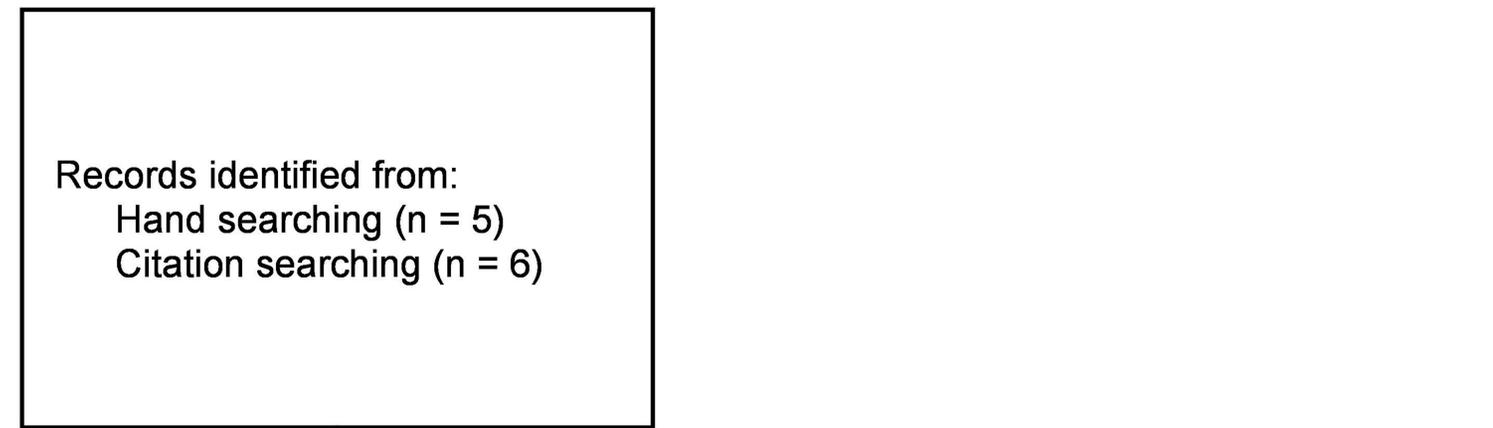
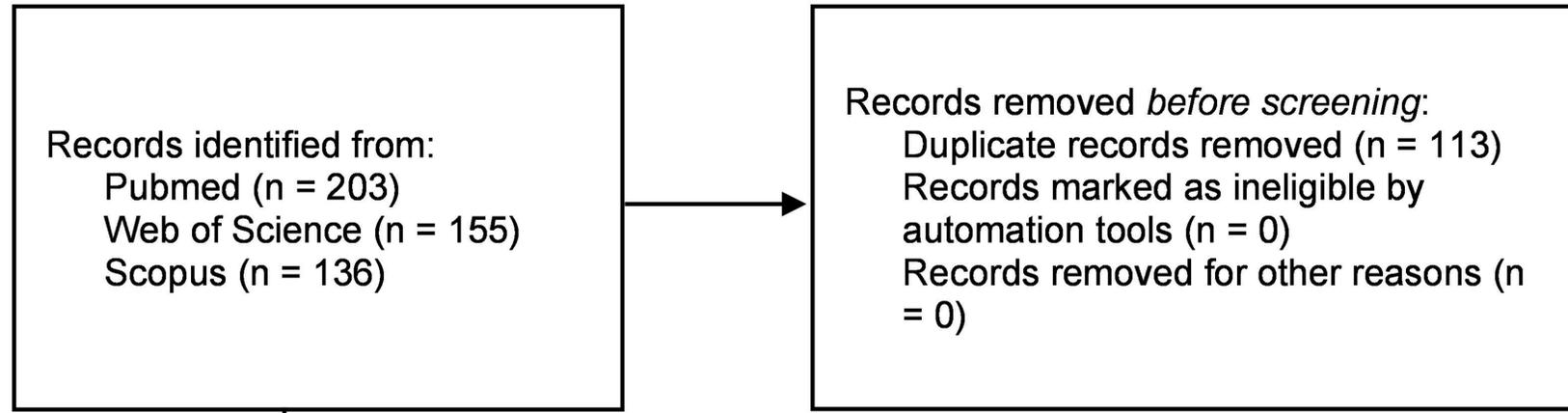
582 pursued.

583 CADD: Combined Annotation Dependent Depletion; VAF: Variant allelic frequency

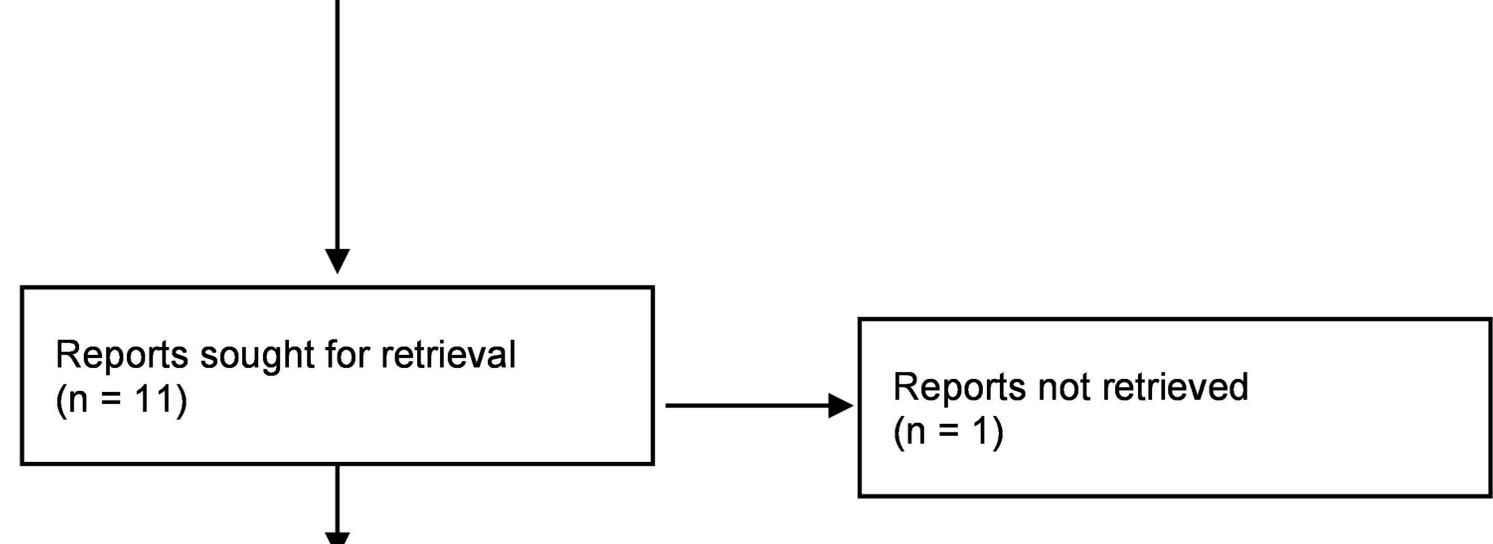
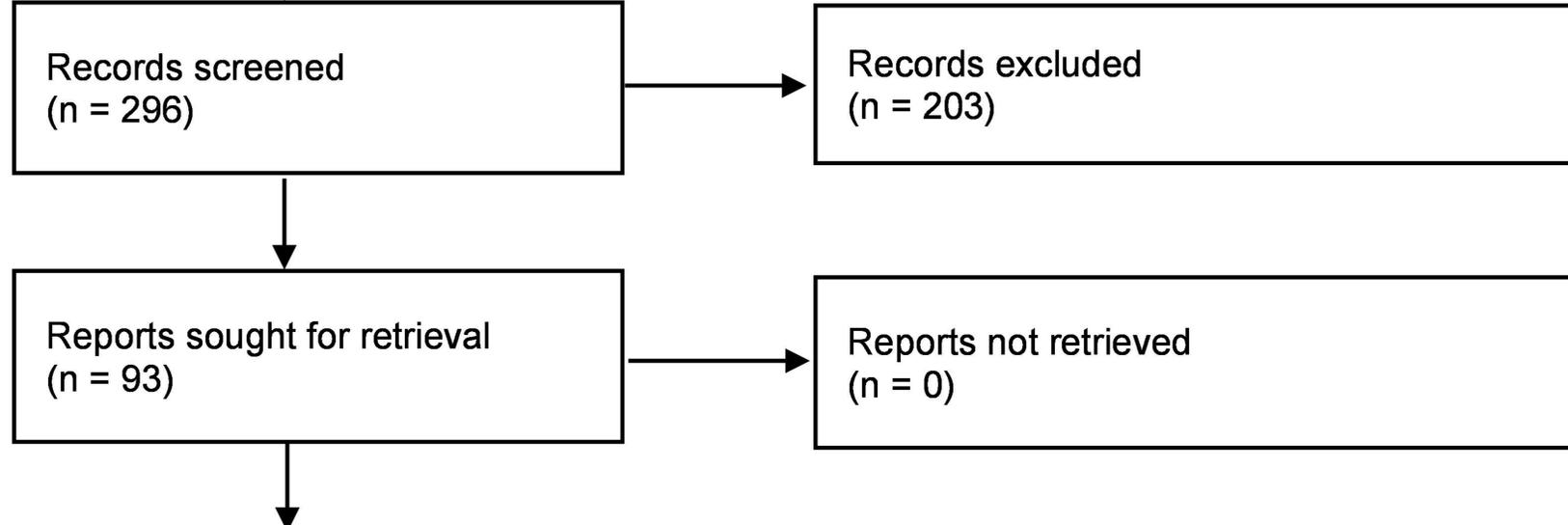
Identification of studies via databases and registers

Identification of studies via other methods

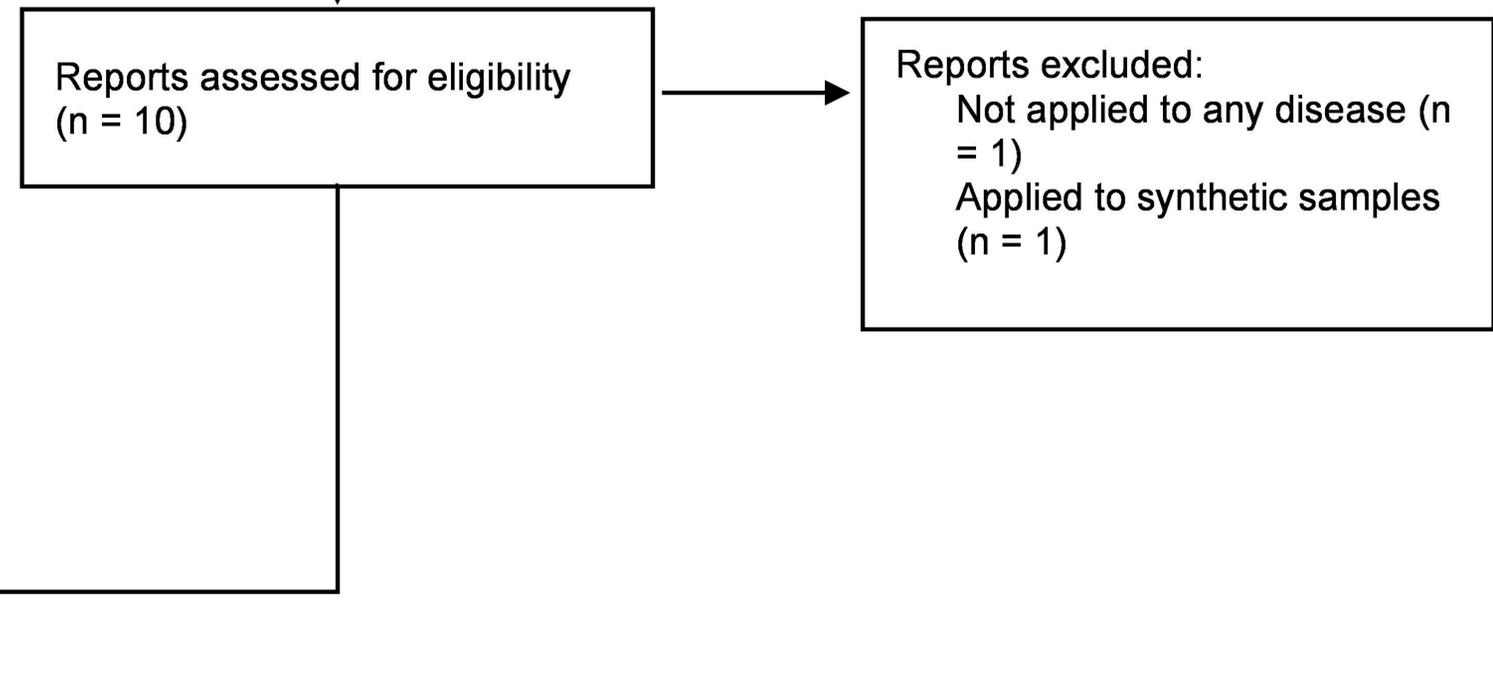
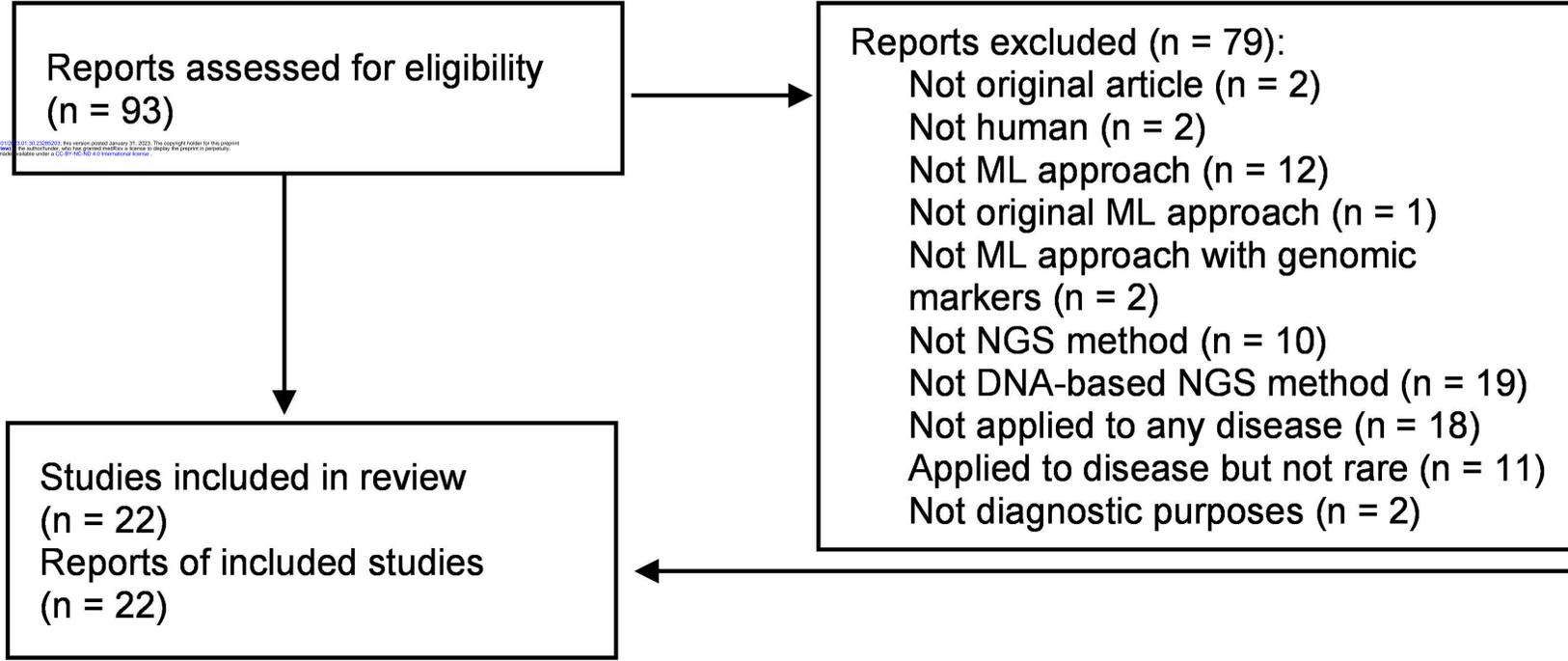
Identification

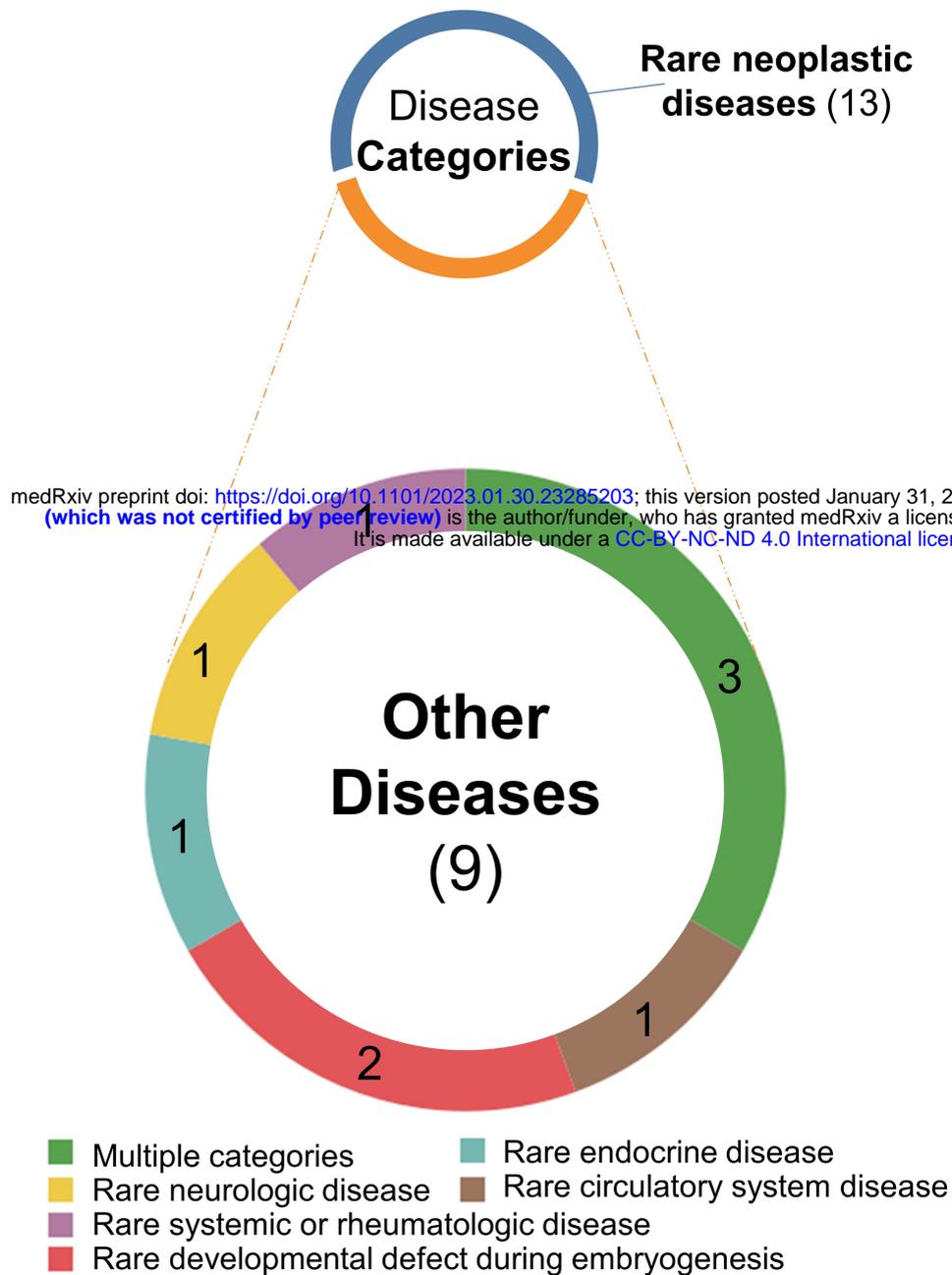
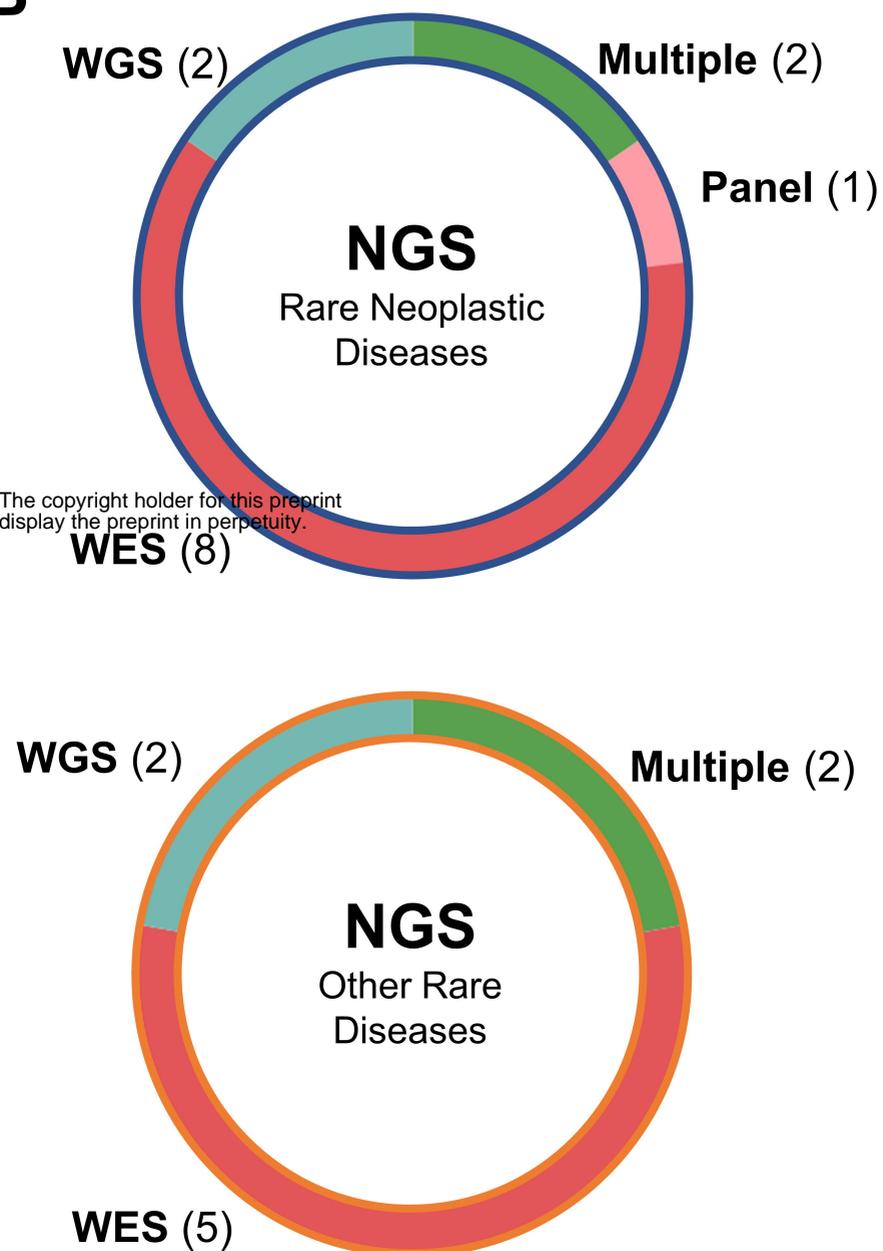
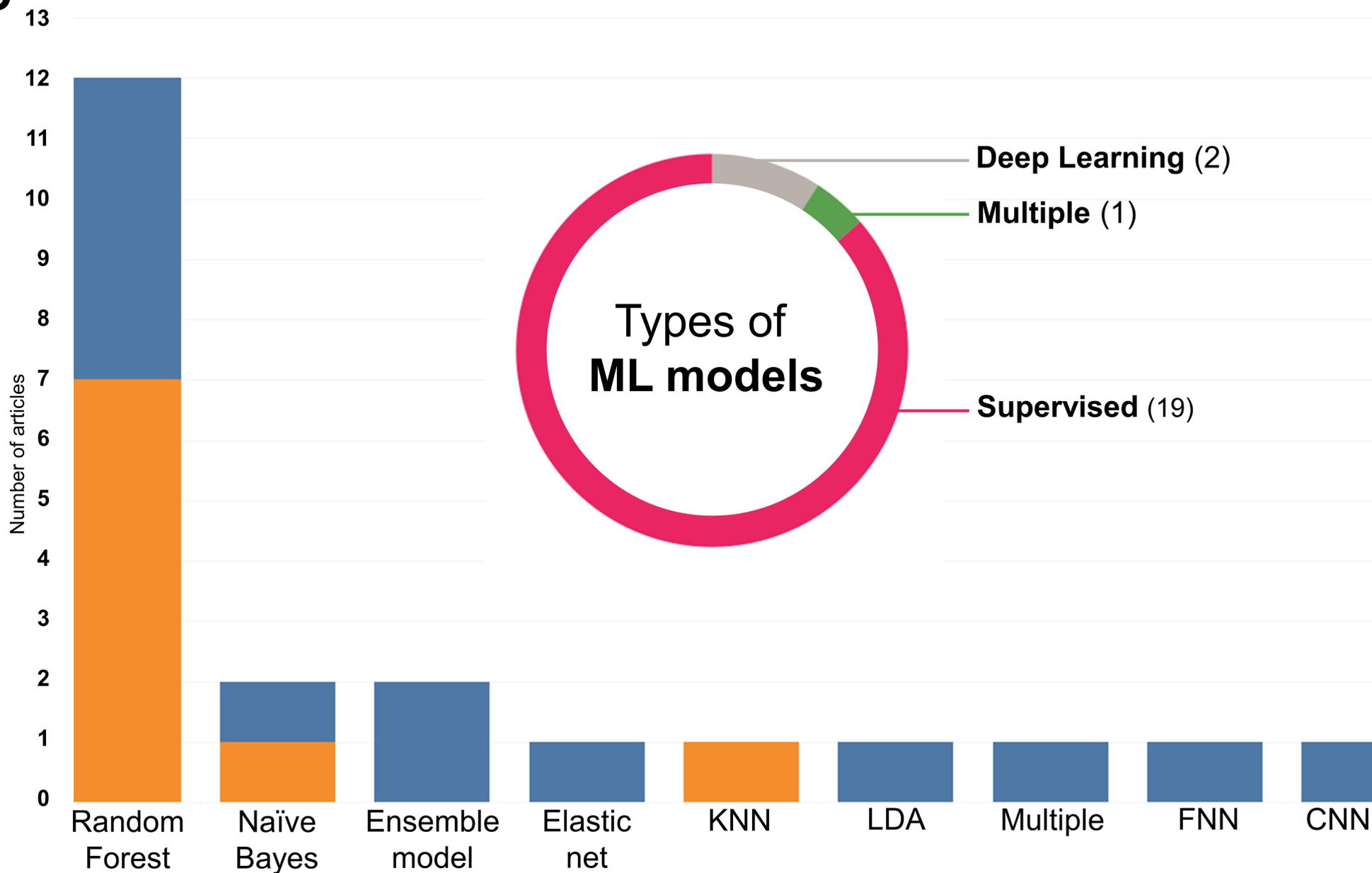


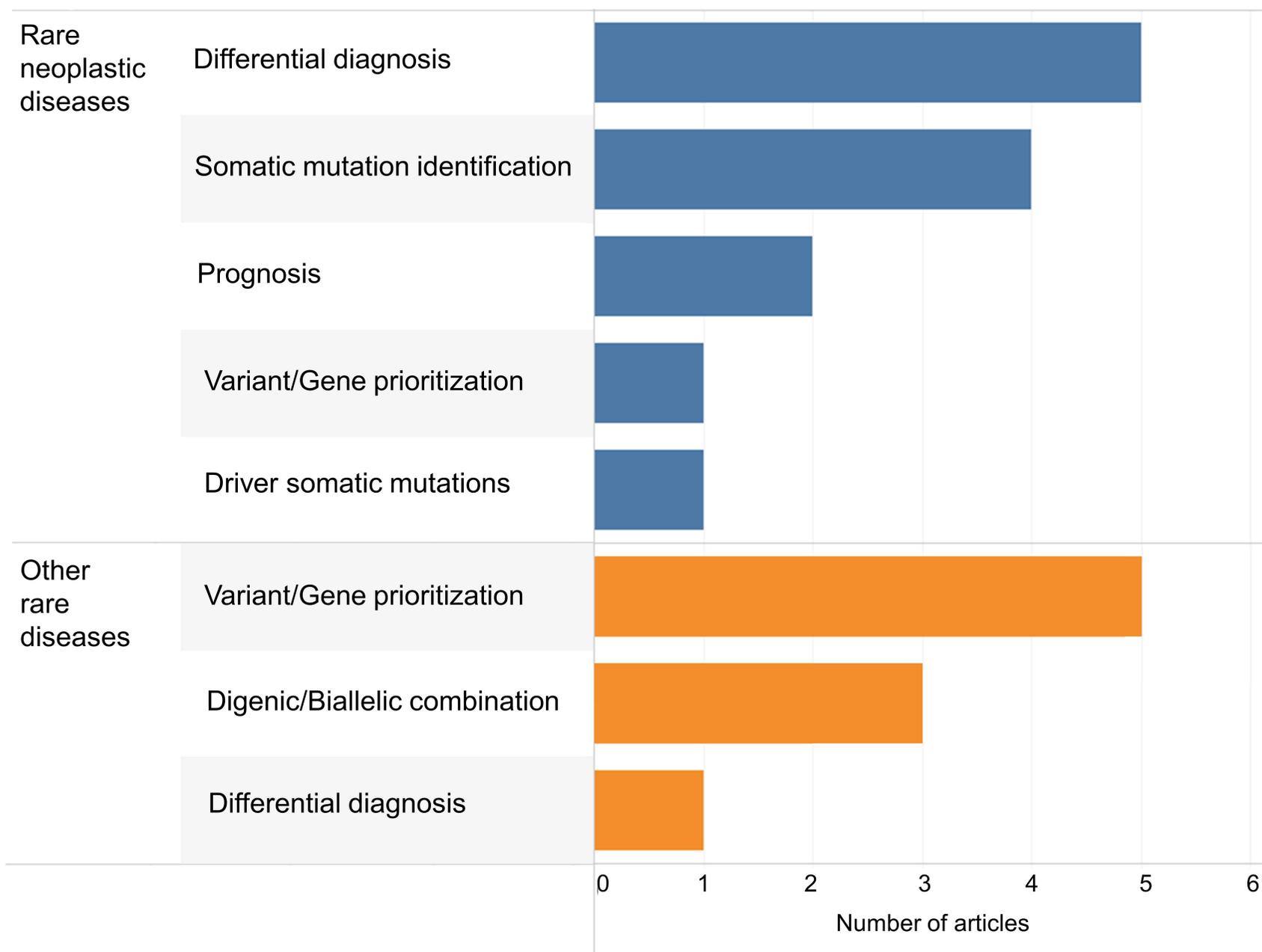
Screening



Included



A**B****C**

A**B**